

# Dynamical Friction Strikes Back


$$\frac{dv_M}{dt} = -16\pi^2 \ln \Lambda G^2 m(M+m) \int_0^{v_M} f(v_m) v_m^2 dv_m v_M$$

**Proceedings of the 34th Symposium on Celestial Mechanics**  
**March 11-13, 2002 at Hakone-Onsen, Kanagawa, Japan**  
**Eiichiro Kokubo, Takashi Ito & Hideyoshi Arakida (eds.)**



天体力学N体力学研究会 (2002年3月11日 - 3月13日 文部科学省共済組合箱根宿泊所 静雲荘にて)



# 箱根天体力学N体力学研究会集録

平成 14 年 3 月 11 日–13 日

神奈川県足柄郡箱根町強羅 静雲荘

## Dynamical Friction Strikes Back

## いまさら力学的摩擦？

Proceedings of the 34th Symposium on Celestial Mechanics

March 11–13, 2002 at Hakone-Onsen, Kanagawa, Japan

Editors: E. Kokubo, T. Ito, and H. Arakida

## Preface/序文

2001 年度の天体力学  $N$  体力学研究会 (通称箱根  $N$  体) は、2002 年 3 月 11 日から 13 日にかけてまだ温泉の温かさがうれしい早春の箱根温泉静雲荘にて滞在型研究会の形で開催されました。口頭発表が 14 件、ポスター発表が 17 件あり、参加者は大学生からシニア研究者までの 47 名を数えました。

箱根  $N$  体のメインテーマは、Dynamical Friction Strikes Back – いまさら力学的摩擦? – ということで、恒星系力学の基礎の 1 つである力学的摩擦でした。力学的摩擦は銀河団から惑星系までさまざまなスケールの天体現象の中で重要な役割を果たしています。例えば、銀河中心における巨大ブラックホール形成、球状星団での恒星の質量分離、惑星系形成時の惑星の移動や軌道の円軌道化などをあげることができます。今回は力学的摩擦について 2 つの招待講演を企画しました。まず恒星系、特に球状の系での力学的摩擦についての基礎的な理論について東京大学の牧野淳一郎氏にレビューしていただきました。そして惑星系、つまり円盤系に適応した場合の力学的摩擦について東京工業大学の田中秀和氏にレビューをしてもらいました。興味深くわかりやすいレビューをしていただいた両氏にはこの場を借りて感謝したいと思います。力学的摩擦の奥の深さを感じてもらえませんでしたでしょうか。「いまこそ力学的摩擦!」ということを感じていただけたなら世話人としてうれしい限りです。

研究会運営にあたっては静雲荘の職員の方々に大変お世話になりました。厚くお礼申し上げます。また、国立天文台の木下宙氏、谷川清隆氏、福島登志夫氏には集録の出版費用を提供していただきました。厚くお礼申し上げます。

平成 14 年初秋 世話人代表 小久保英一郎

## Editors/世話人

- |                            |  |
|----------------------------|--|
| 小久保英一郎 (Kokubo, Eiichiro)  | 国立天文台理論天文学研究系<br>kokubo@th.nao.ac.jp     |
| 伊藤孝士 (Ito, Takashi)        | 国立天文台天文学データ解析計算センター<br>tito@cc.nao.ac.jp |
| 荒木田英禎 (Arakida, Hideyoshi) | 総合研究大学院大学<br>h.arakida@nao.ac.jp         |



# Table of Contents/目次

## “Dynamical Friction Strikes Back” Feature Articles

Dynamical friction: Overview

*Shigeru Ida* ..... 1

Dynamical friction in stellar systems and their simulations — The reality and the fiction

*Junichiro Makino* ..... 2

Dynamical friction in gravitating disk systems and radial migration

*Hidekazu Tanaka* ..... 25

## Stellar Dynamics

Life expectancy of the large Magellanic cloud

*Yoko Funato, Yoshikazu Hashimoto, and Junichiro Makino* ..... 30

New limits on the mass of the Milky Way

*Tsuyoshi Sakamoto and Masashi Chiba* ..... 38

Gravothermal catastrophe and Tsallis’ generalized entropy of self-gravitating systems: Thermodynamic properties of stellar polytrope

*Atsushi Taruya and Masa-aki Sakagami* ..... 78

Stationary state in N-body system with power law interaction

*Osamu Iguchi* ..... 104

The effect of tidal force and mass loss in star clusters sinking process

*Tatsushi Matsubayashi and Toshikazu Ebisuzaki* ..... 112

Dynamical friction between lopsided disks and dark halos

*Makoto Ideta* ..... 134

## Formation of Planetary Systems

Formation of terrestrial planets in a dissipating gas disk

*Junko Kominami and Shigeru Ida* ..... 152

The evidence of a stellar encounter on the distribution Edgeworth-Kuiper belt object

*Hiroshi Kobayashi, Shigeru Ida, and Hidekazu Tanaka* ..... 162

Formation of low-mass multiple satellites

*Takaaki Takeda* ..... 189

Orbital stability of a protoplanet system in the nebular gas: Dependence on the masses of protoplanets

*Kazunori Iwasaki, Hiroshi Emori, Hidekazu Tanaka, and Kiyoshi Nakazawa* 193

## Extrasolar Planetary Systems

Dynamical stability of planetary system of GJ876

*Hiroshi Kinoshita and Hiroshi Nakai* ..... 199

## Solar System Dynamics

A comprehension for the solution due to Schwarzschild

*Takeshi Inoue* ..... 226

Subsystems in a stable planetary system I. A classification

*Kiyotaka Tanikawa and Takashi Ito* ..... 236

Evolution of obliquity of a terrestrial planet due to gravitational perturbation by a giant planet

*Keiko Atobe, Takashi Ito, and Shigeru Ida* ..... 255

Motion around triangular Lagrange points perturbed by other bodies

*Hideyoshi Arakida and Toshio Fukushima* ..... 265

Minor planet's orbits in or near mean motion resonances with Jupiter

*Hiroshi Nakai and Hiroshi Kinoshita* ..... 289



Orbital theory of a highly eccentric satellite disturbed by a massive inner satellite <i>Yoshimitsu Masaki and Hiroshi Kinoshita</i> .....	303
Size and spatial distributions of sub-km main-belt asteroids <i>Fumi Yoshida and Tsuko Nakamura</i> .....	330

## Theory of Dynamical Systems

Dynamical systems which produce the Lévy flights <i>Tomoshige Miyaguchi and Yoji Aizawa</i> .....	353
Dynamical ordering of non-Birkhoff orbits and topological entropy in the standard mapping <i>Yoshihiro Yamaguchi and Kiyotaka Tanikawa</i> .....	359
Non-Birkhoff periodic orbits in a circle mapping <i>Yoshihiro Yamaguchi and Kiyotaka Tanikawa</i> .....	395
Relaxation in Hamiltonian systems with long-range interactions <i>Yoshiyuki Y. Yamaguchi</i> .....	407

## Numerical Techniques

A modified Hermite integrator for planetary dynamics <i>Eiichiro Kokubo and Junichiro Makino</i> .....	415
Secular numerical error in $H = T(p) + V(q)$ symplectic integrator: simple analysis for error reduction <i>Takashi Ito and Kiyotaka Tanikawa</i> .....	437
Construction of heterogeneous computer system with GRAPE-5 and VPP5000 by using IMPI <i>Mitsuru Hayashi, Takashi Ito, Eiichiro Kokubo, Hiroshi Koyama, Koji Tomisaka, Keiichi Wada, Nobuhiro Uchida, Noboru Asai, Eiji Uemura, and Kazutaka Sugimoto</i> .....	528
Non-linear harmonic analysis of the time ephemeris of the Earth <i>Wataru Harada and Toshio Fukushima</i> .....	533

Problems of symmetric multistep methods for Keplerian motion  
*Tadato Yamamoto and Toshio Fukushima* ..... 541

Other Topics

Response of lifespan of organisms to secularly changing environment using a new dynamical model  
*Toshihiro Handa, Kiyotaka Tanikawa, and Takashi Ito* ..... 550

Symposium program ..... 565

Author index and participant list ..... 568



# いまさら(いまこそ)力学的摩擦 -Introduction-

井田茂 [東工大・地惑]

- 力学的摩擦は銀河ダイナミクス、惑星形成における Key process のひとつ
  - 銀河ダイナミクス (牧野/東大)
  - 銀河中心部巨大ブラックホール形成 (松林/東工大、理研)
  - 惑星落下問題 (田中/東工大)
  - 地球型惑星形成 (小南/東工大)



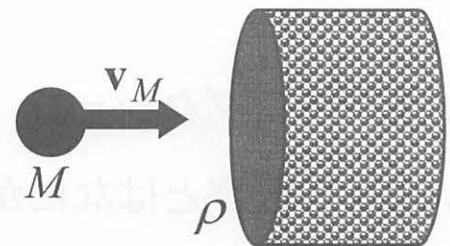
- 「力学的摩擦」はいろいろな意味で使われる
  - Exotic particle の速度の減衰
  - 相対的に質量の大きな天体の \*\*速度\*\* の減衰
  - Energy equipartition
  - Mass segregation, \*\*中心\*\* への落下
  - ガス成分との重力相互作用

## 力学的摩擦はちゃんとわかっているのか？

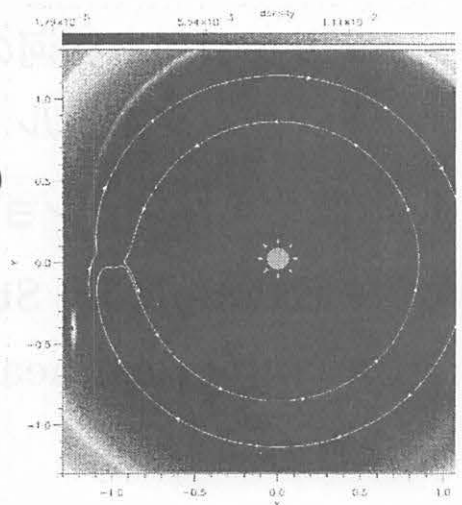
Chandrasekhar's formula

$$M \frac{dv_M}{dt} \approx -\rho v_M \cdot \pi \left( \frac{GM}{v_M^2} \right)^2 \ln \Lambda \cdot v_M \quad \text{for } v_M > v_m$$

$$M \frac{dv_M}{dt} \approx -\rho v_M \cdot \pi \left( \frac{GM}{v_m^2} \right)^2 \ln \Lambda \cdot v_m \quad \text{for } v_M < v_m$$



- $v_m$  はどこまで減衰? (equipartition?)
- 直線で飛んでいないとき? 媒質に系統的運動があるとき? (例 ケプラー粒子)
  - 軌道離心率、軌道長半径の減衰?
- $\rho$  に比例: 媒質の粒子の大きさによらない
  - 相手がガスでもOK?  $v_m = c_s$ ?
  - 粒子描像 vs 流体描像



# 恒星系及び恒星系N体シミュレーションに おける力学的摩擦 — 現実と虚構

牧野淳一郎

## 概要

1. イントロダクション
2. 力学的摩擦とはなにか？
3. 「現実の」系における力学的摩擦
  - サテライト銀河の進化
  - ブラックホール
4. 数値シミュレーションにおける力学的摩擦
  - Example by Steimetz and White
  - Numerical heating of thin disk
5. まとめ



# イントロダクション

この講演ではとりあえずダイナミカルフリクションとはどんなもので、何故そういうものを考えないといけないかという話を主に恒星系を例にとってする。

- 力学的摩擦とはなにか
- どう計算されるか
- その計算のしかたは「正しい」か

## 力学的摩擦とはなにか

基本的には、自己重力多体系のなかで、平均の運動エネルギーよりも高いエネルギーをもった粒子が受ける抵抗。

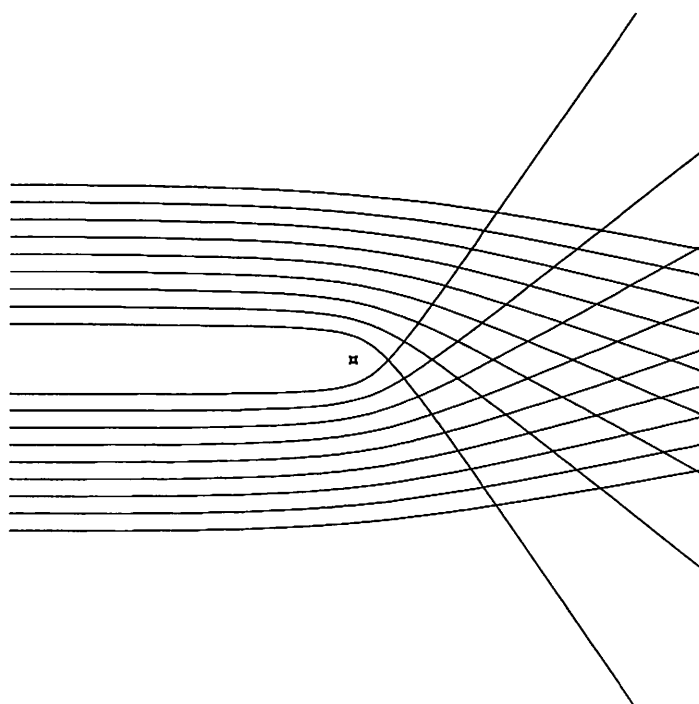
抵抗を受ける「理由」：「熱平衡状態に近づくため」といってもなんだか分からないので、もうちょっと簡単な場合を考える。

# もっとも簡単なモデル

今、温度0（だと、本当はジーンズ不安定が起きるわけだがこれはとりあえず考えない）の、無限に一様な物質分布の中を、適当な大きさを持った球対称なポテンシャルの摂動（質点によるものでもOK）が動いているとする。

まわりの粒子は質点からの力を受けて速度を得る。

## 回りの粒子の軌道



## 抵抗になる理由

もともとの止まっていた物質分布に固定された座標系で考える:

散乱されたものは、左向きと中心向きの速度をもらうことになり、ネットに加速されている。

エネルギーをもらっている。

## まっとうな導出

分布している質点の質量を  $m$ 、数密度を  $n$  とする。テスト粒子が一つの粒子から距離（インパクトパラメータ）  $p$  を相対速度  $V = v_t - v_f$  で通った時に曲がる角度:

$$\tan \theta = \frac{2p}{(p/p_0)^2 - 1}$$
$$p_0 = \frac{G(m_t + m_f)}{V^2}$$

で与えられる。

## 回りがとまっている場合 (1)

回りはまだ速度 0 とすると:

$$\Delta v_{vert} = \frac{m_f}{m_t + m_f} V \sin \theta = 2V \frac{m_f}{m_t + m_f} \frac{p/p_0}{1 + (p/p_0)^2}$$

$$\Delta v_{para} = \frac{m_f}{m_t + m_f} V (1 - \cos \theta) = -2V \frac{m_f}{m_t + m_f} \frac{1}{1 + (p/p_0)^2}$$

## 回りがとまっている場合 (2)

単位時間当たりの衝突回数  $2\pi p n_f V dp$  を掛けて積分すると:

$$\langle \Delta v_{vert}^2 \rangle = \frac{2n_f \Gamma}{V}$$

$$\langle \Delta v_{para} \rangle = - \left( 1 + \frac{m_t}{m_f} \right) \frac{n_f \Gamma}{V^2}$$

$$\langle \Delta v_{para}^2 \rangle = \frac{n_f \Gamma}{V \ln \Lambda}$$

ここで  $\Gamma$  は

$$\Gamma = 4\pi G^2 m_f^2 \ln \Lambda \quad (1)$$

である。

## $\ln \Lambda$ って?

インパクトパラメータ  $p$  での積分:  $p \rightarrow \infty$  で形式的に発散する。

$$\Delta v_{para} = -2V \frac{m_f}{m_t + m_f} \frac{1}{1 + (p/p_0)^2}$$

に  $2\pi n p dp$  をかけて積分するから。  $p \rightarrow 0$  は (質点粒子でも) 発散しない。

$p$  の上限

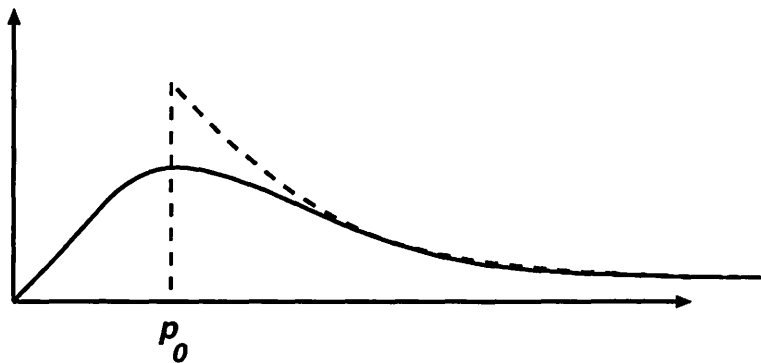
- 系の「大きさ」くらい
- 密度のスケールハイト

もちろん  $n$  の  $p$  への依存性を書き下せればそれを使うべき?

## 積分の下限 (1)

質点粒子:  $\Delta v_{para}$  をまともに積分すればよい。

実際には  $\Delta v_{para} \sim 1/p$  として  $p = p_0$  で積分を打ち切るとあんまり変わらない。



## 積分の下限(2)

粒子が広がっている場合 (衛星銀河とか):

$\Delta v_{para}$  を軌道をといてちゃんと計算すればいい。

粒子サイズ  $\gg p_0$  の時は適当な近似ができる。

この場合にも、実際には  $\Delta v_{para} \sim 1/p$  として  $p$  が典型的な半径 (half-mass radius とか) で積分を打ち切るので普通は大丈夫。

## 力学的摩擦の性質

$$\langle \Delta v_{para} \rangle = - \left( 1 + \frac{m_t}{m_f} \right) \frac{4\pi G^2 m_f^2 \ln \Lambda n_f}{V^2}$$

- $m_t$  に比例 ( $m_t \gg m_f$  で)
- $m_f$  によらない ( $\rho = m_f n_f$  一定なら)
- $V^2$  に反比例

割合変な力。重い粒子の  $\Delta v$  がその粒子の質量に比例。



## 回りが動いているとき

導出は面倒なので結果だけ。速度分布が等方的なら

$$\begin{aligned} F_n(v) &= \int_0^v \left(\frac{v_f}{v}\right)^n f(v_f) dv_f \\ E_n(v) &= \int_v^\infty \left(\frac{v_f}{v}\right)^n f(v_f) dv_f \end{aligned} \quad (2)$$

というものを考えると、

$$\langle \Delta v_{para} \rangle = -4\pi\Gamma \left(1 + \frac{m}{m_f}\right) F_2(v)$$

## 回りが熱平衡なら

$$f_0(v) = \frac{n_f}{(2\pi\sigma^2)^{3/2}} \exp\left(\frac{-v^2/2}{\sigma^2}\right)$$

$$\langle \Delta v_{para} \rangle = -4\frac{n_f\Gamma}{\sigma^2} \left(1 + \frac{m}{m_f}\right) G(x)$$

ここで erf は誤差関数であり、

$$G(x) = \frac{\operatorname{erf}(x) - x\operatorname{erf}'(x)}{2x^2}$$

また  $x = v_t/(\sqrt{2}\sigma)$ 。

## 回りが止まってるときとの違い

結局、

$$\text{erf}(x) - x\text{erf}'(x)$$

の分だけ弱くなる。

$v_t \rightarrow 0$  の極限:  $v_t$  に比例

$v_t \rightarrow \infty$  の極限:  $v_t^2$  に反比例

中間: 回りがとまっているときよりファクターで弱い(あんまり変わらない)

## 実際の恒星系での力学的摩擦

様々なところに現れる。

基本的には熱平衡でなければ必ず等分配からずれた粒子がある。

それは必ず力学的摩擦を受ける/およぼす。

自己重力系は原理的に熱平衡でありえない。

## 具体例

- 衛星銀河の進化 (船渡講演)
- 銀河中心の複数ブラックホール系
- 若い高密度星団

## 衛星銀河の進化

衛星銀河:親銀河のダークハロー (あれば) の力学的摩擦を受ける。

Binney & Tremaine に書いてある式

$$t_{fric} \simeq 6 \times 10^8 \left( \frac{r}{1\text{kpc}} \right)^2 \left( \frac{v_c}{100\text{kms}^{-1}} \right) \left( \frac{5 \times 10^6 M_\odot}{m} \right) \text{yr}$$

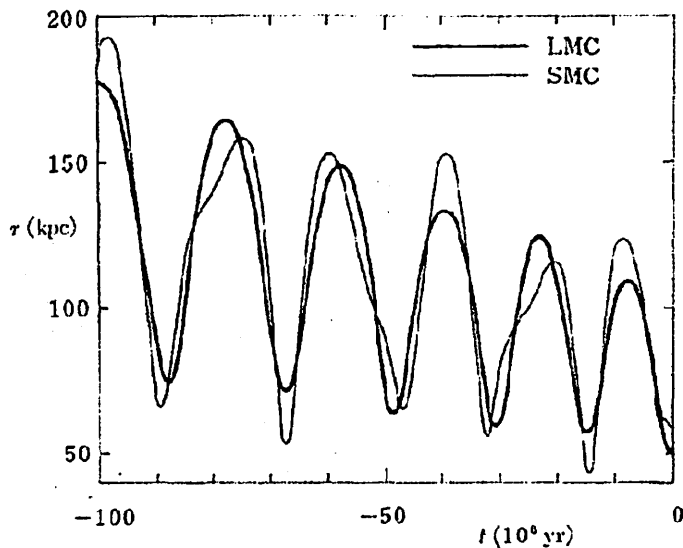
これは銀河中心の星団用なのでちょっと短い、、、  $\log \Lambda$  は 5 くらいをいれたはず。

まあ、質量が親銀河の数パーセント以上あればタイムスケールは結構短い。

# LMC-SMC

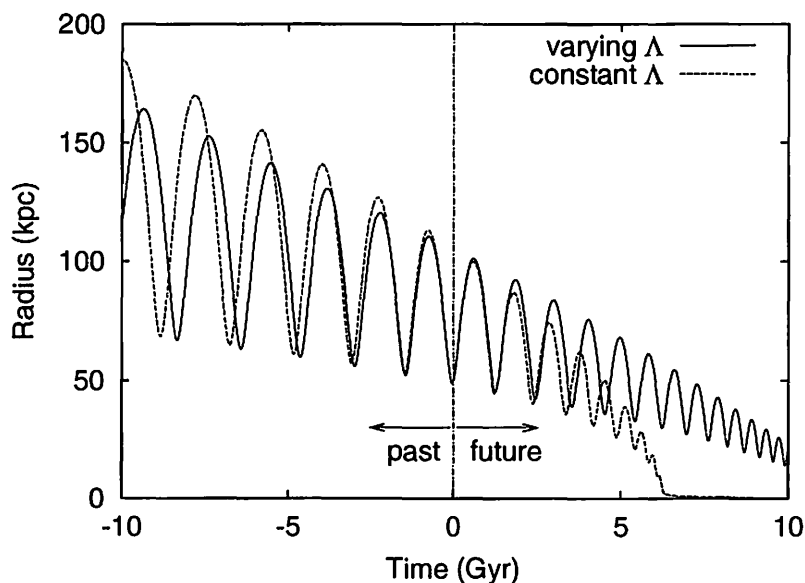
古典的計算: 村井・藤本

論文読んでもどうやって計算したのかよくわからない、、、



これは合ってるか?というのは実は問題

## $\ln \Lambda$ のとりかた

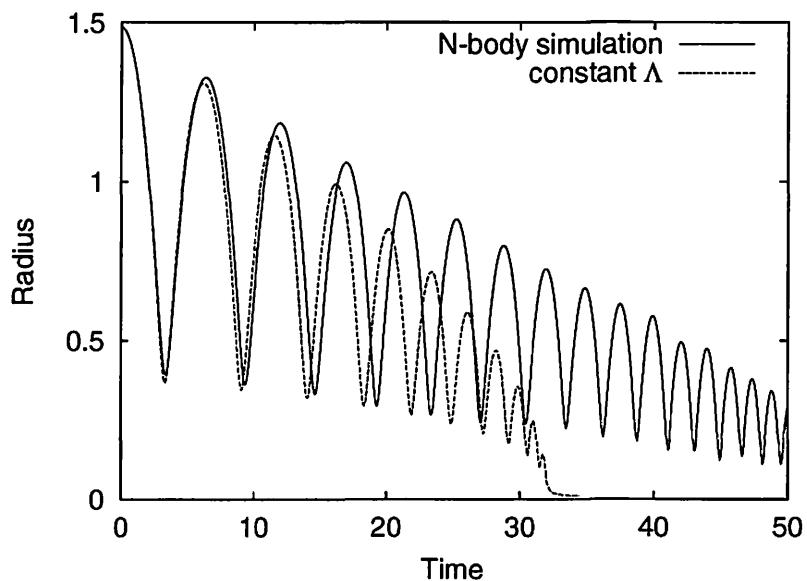


(Hashimoto, et al. in preparation)

破線:  $b_{max} = R_{halo}$  実線:  $b_{max} = r_s$  (衛星の位置)

破線が村井・藤本によく合う。

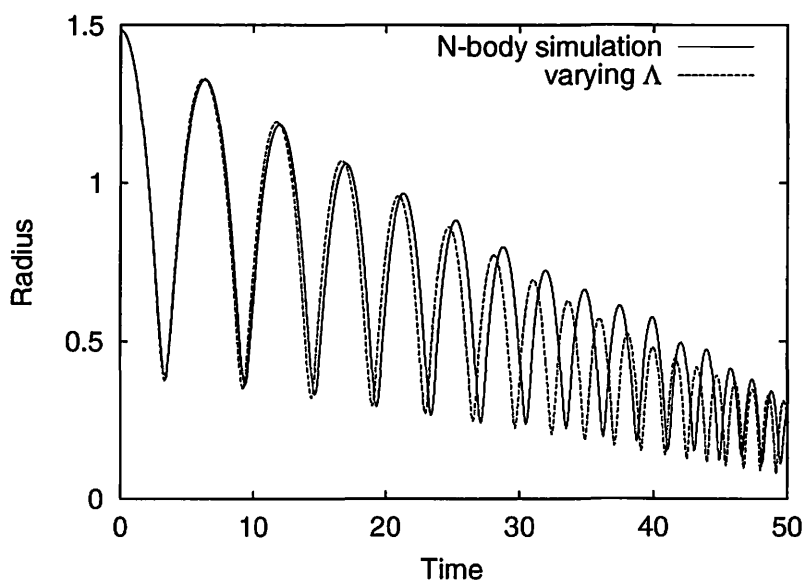
## $N$ 体との比較 (1)



$$b_{max} = R_{\text{halo}}$$

$N$  体より進化速い 軌道丸くなる

## $N$ 体との比較 (2)



$$b_{max} = r_s$$

$N$  体とよく合う 軌道丸くならない

## $N$ 体との比較まとめ)

$\ln \Lambda$  の絶対値だけの問題ではないことに注意

軌道が丸くなるかならないか

= 近点で DF がどれくらい有効か

$b_{max} = r_s$  とすると近点で  $\ln \Lambda$  が非常に小さくなる

= circularization が抑えられる

## $\ln \Lambda$ の問題

DF の公式を使って衛星銀河の進化を調べた論文は無数にあるが、、、

これまで発見できたなかでは Tremaine (1976) だけが

$b_{max} = r_s$  採用。

一番悪い論文:  $\Lambda = M_H/M_s$  これは衛星銀河を質点と仮定しているのと同じ。

こういう困った論文もいっぱいある。

# チャンドラセカル公式の限界

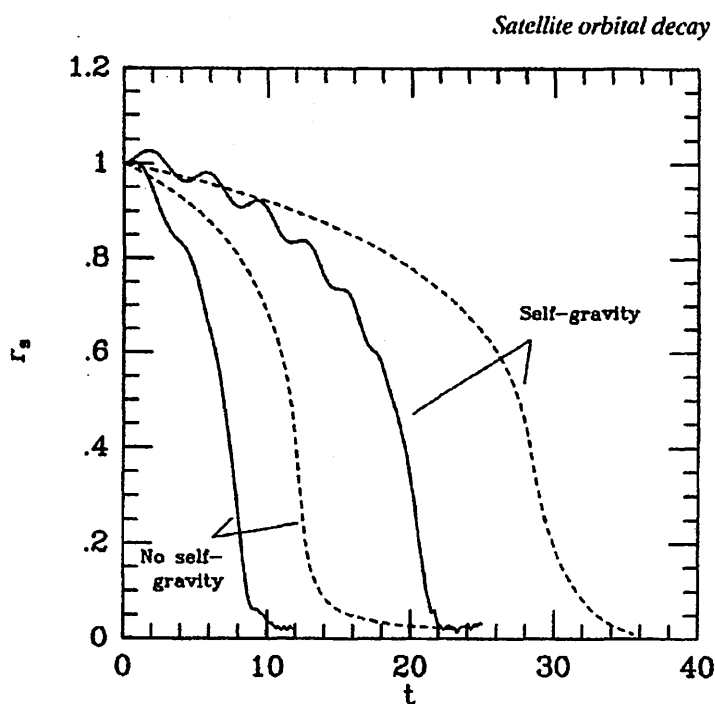
- テスト粒子、フィールド粒子共に直線運動をすると仮定
- フィールド粒子同士の自己重力は無視

テスト粒子の近くではそんなに悪くない

親銀河の大きさ (テスト粒子の銀河中心からの距離) くらい離れたところではまったくなりたない仮定。  
(楽観的な期待: 現実的な効果は DF を小さくする)

## 「もっと精密な」 計算法

分布関数のグローバルな応答を計算 (Weinberg 1989)





## 「もっと精密な」 計算法 (2)

精密にしたらよくあうというものでもない？

この不一致の理由は結局よくわかってない(?)

## デモンストレーション

サテライト銀河の進化

親銀河: King ( $W_0 = 9$ ) モデル。 Heggie Unit

衛星銀河: King ( $W_0 = 9$ ) モデル。 質量  $1/8, 1/16, 1/32$

初期位置 (5,0,0)

初期速度 (0,0.3,0)

粒子数  $32768+65536$ , GRAPE-6 direct method

# 複数ブラックホール系

(単一ブラックホールだと中心に落ちて終わりなのでつまらない)

2 個ではなにが起きるか？

- 中心にブラックホールを持つ銀河同士の合体
- 中間質量ブラックホール ???

# 複数ブラックホール系

中心に沈むのはまあそうなるとして、、、

それからどうなるか？

= ブラックホール連星系の進化

フィールド粒子と高速回転するブラックホール連星の相互作用

物理 (統計力学) としては力学的摩擦に似ている。

実際の連星系の進化: いろいろ分からないことがある。

# $N$ 体計算

(Makino 1997)

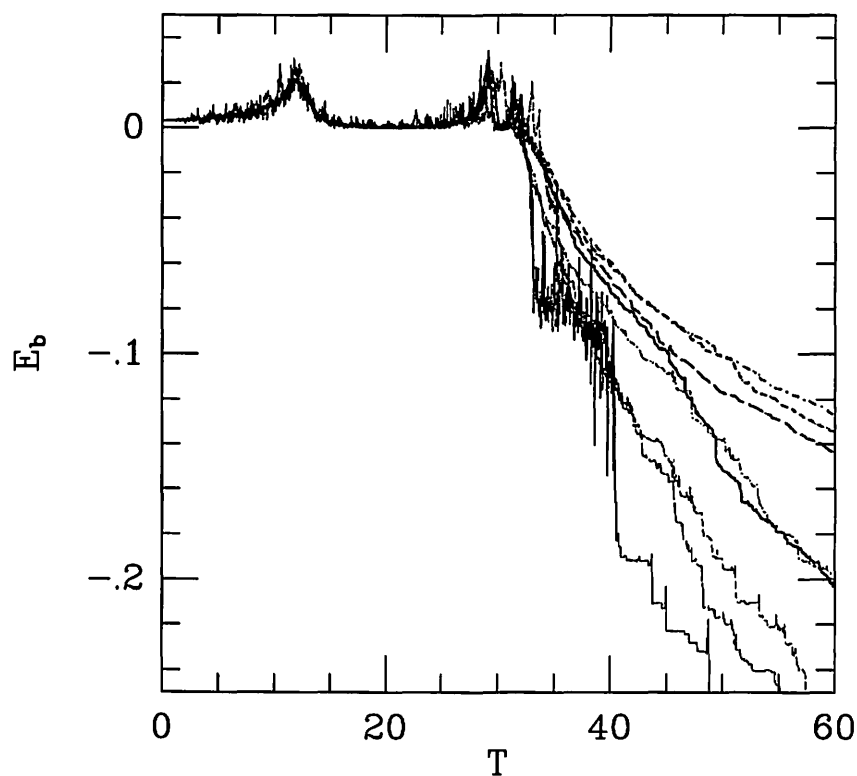
理屈では ブラックホール連星の進化タイムスケールは緩和時間に比例するはず

相互作用できる星を弾き飛ばす。2体緩和でまた拡散してくる。

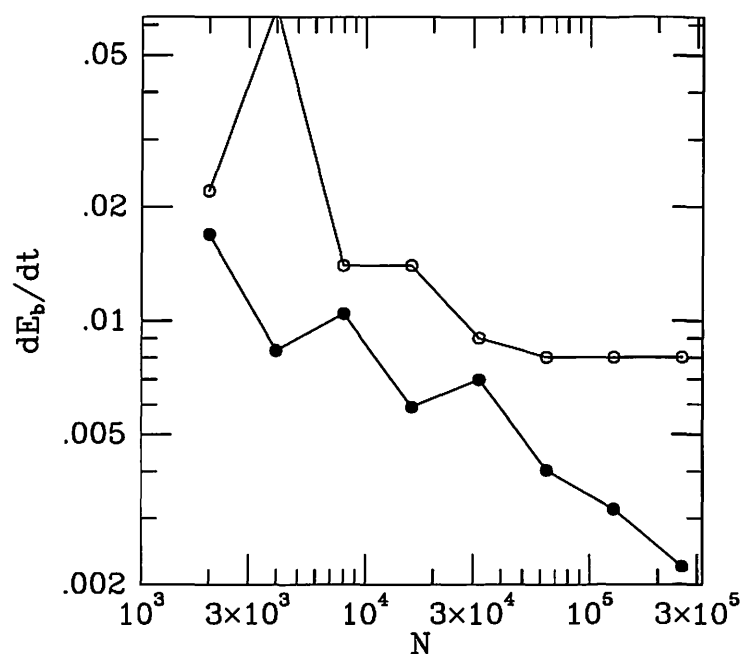
粒子数をかなり広く変えて計算してみた (2k — 256k)

結果はなんだかよく分からない。

## ブラックホール連星: 計算結果



# 成長率



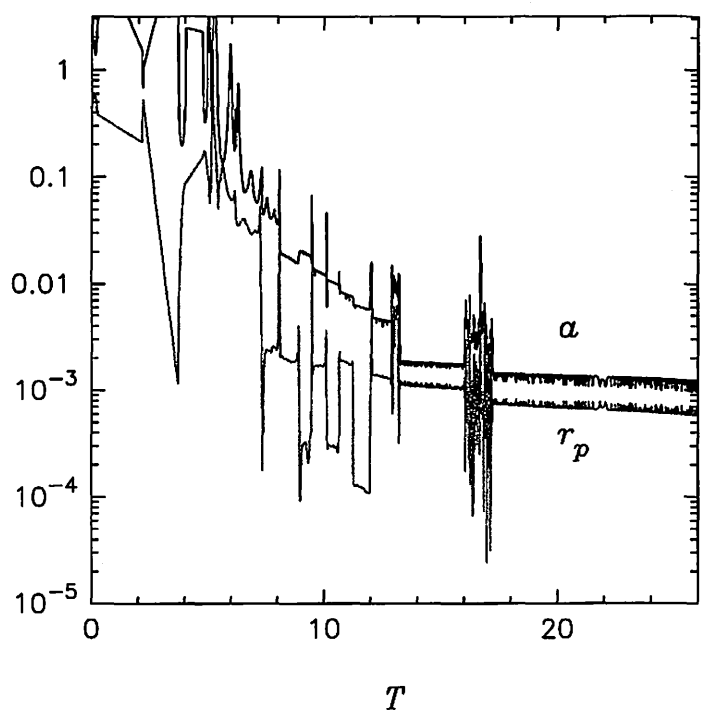
白丸: 最初のほう  
黒丸: 後のほう

Timescale:

初期:  
収束?  
後期:  
 $\propto N^{1/3}???$

# BH 3個

(SC2001 proceedings paper)



軌道長半径と離心率の進化  
近点での距離は極めて小さくなれる — 割合簡単に重力波で合体?

## 「現実」の系での力学的摩擦

- 割合いろんなところにててくる
- (正しく使えば) チャンドラセカールの公式はかなり正確
- 重要なのは  $\ln \Lambda$  の推定

## $N$ 体計算における数値的な力学的摩擦

これがいろんなところにててきては本当はいけない。

が、、、

というわけで、ちょっとそういう方面の話。

# どういう場合に問題になるか

主に「現実」は無衝突系の場合。

- シミュレーションは粒子数がずっと少ない
- 粒子が等質量とは限らない

注意していないと2体緩和や力学的摩擦の影響が極めて大きいことがある

## 例:円盤銀河

質量の 90%以上がダークハロー

ディスクは 10% 以下

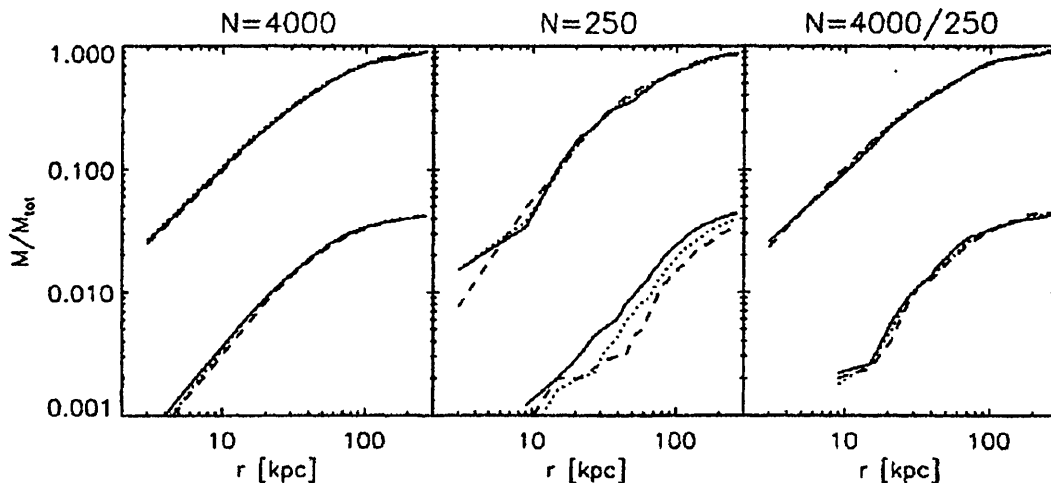
バルジは1% 程度

ディスク、バルジに分解能をもたせるにはハローに使う粒子数を節約したいような気がする。

これをやるとなかなか楽しいことが起きる。

# Steinmetz and White によるデモンストレーション

Steinmetz and White 1996



ハロー、ガス共に球状。 40 および 80 億年後  
ハローに 4000 粒子もあればまあまあ、、、

## Thin disk の場合

一昨日急いで作った例なのであんまりよくないですが、、、  
ダークハロー: Plummer model, heggie unit  
厚さ、温度 0 の (ほぼ) テスト粒子ディスク。8192 体  
半径 0.4, 一様分布、(ほぼ) つりあった回転を与える。  
ハロー粒子 1024-131072 体



## わかること

10万体くらい使ってもディスクはすぐに膨らむ。  
数千体では全く論外。

## ディスク粒子が $z$ 方向の速度を得る時間 スケール

緩和時間:  $z$  方向の速度が回転速度とコンパラ。既にディスクではない

緩和時間の  $1/10$  :  $z$  方向の速度が回転速度の  $1/3$ 。

緩和時間の  $1/100$  :  $z$  方向の速度が回転速度の  $1/10$ 。

緩和時間の  $1/1000$  :  $z$  方向の速度が回転速度の  $1/30$ 。

ディスク数回転の間厚さを 3% くらいに保つためにはハローに 10 万粒子以上必要。

# まとめ

## 数値的な力学的摩擦

- 主に問題になるのは高温 (ランダム) な部分系によっては低温 (規則運動) 部分系が加熱されること。
- 薄いディスクを作ろうと思うとハロー (のディスクがあるあたり) だけに数十万粒子

## 全体のまとめ

- 現実の系での力学的摩擦はどこにでもある。
- 古典的なチャンドラセカル公式は正しく使えば十分使える。
- 数値的なものはもっとどこにでもある。

# Dynamical Friction in Gravitating Disk Systems and Radial Migration

## 円盤重力系における力学的摩擦と動径方向移動

Hidekazu TANAKA

*Dept. of Earth and Planetary Sciences, Tokyo Institute of Technology*

田中秀和

東京工業大学 大学院理工学研究科 地球惑星科学専攻

### Abstract

Studies on dynamical friction in gravitating particle disks are reviewed. Dynamical friction plays an important role in “the velocity relaxation process” and “radial migration”. In disk systems, the velocity relaxation can be considered as a local relaxation process and its characteristic time is given by “Chandrasekhar’s relaxation time”. On the other hand, radial migration in disk systems is considered as global evolution of disks and the evolution time is much longer than Chandrasekhar’s relaxation time. The gravitational interaction between a particle and a gaseous disk is also described. the interaction with gaseous disks is very similar to that with particle disks. The relaxation time in gaseous disks is given by the same formula as particle disks if the sound velocity is taken as the relative velocity.

### 0. はじめに

粒子系の重力緩和過程や動径方向移動は、銀河、惑星形成、惑星リングなどを考える上で重要な素過程であり、今日まで多くの研究がなされてきた。これら重力緩和過程と動径方向移動は、広い意味での力学的摩擦によるものと説明されている。

一方、(原始) 惑星系や銀河では、構成粒子は円盤状に分布しているが、この円盤重力系の重力緩和過程は、重力多体系に見られるいわゆる長距離相互作用による困難を含まず、比較的簡単に理論化すること可能であった。今日では、円盤系での局所的な重力緩和過程については、理論的に良く理解されていると言えるであろう。

本稿では、「円盤重力系における力学的摩擦」に関連した従来の研究を、まとめた上で紹介していこう。取り上げる内容は、以下のようものである。

1. 力学的摩擦とは何か？
2. 円盤系の特徴
3. 円盤系での粒子の速度進化
4. 円盤系での粒子の動径方向移動

ガス円盤と惑星の間の重力相互作用は、惑星を落下させるなどの効果を持ち、粒子同士の重力相互作用と共に惑星形成過程において重要な働きをする。一見、ガス円盤との重力相互作用は、粒子同士の重力相互作用とは全

く別のものと思われるが、これら2つの重力相互作用は多くの共通点を持っている。本稿ではこれらの共通点についても議論していく。

## 1. 力学的摩擦とは何か？

力学的摩擦といっても、研究者によって色々な意味で用いられている。その点をまず整理しておこう。

1つの粒子が流体の中をある相対速度で運動する場合、その粒子は流体から抵抗力を受ける。チャンドラセカールは、流体中ではなく粒子集団の中を粒子が運動する際にも、重力相互作用によって同様な抵抗力が発生することを示し、この効果を「力学的摩擦 (dynamical friction)」と命名した (Chandrasekhar 1943ab)。これが通常の見方である。

力学的摩擦による抵抗力で粒子が減速する時間は、「チャンドラセカールの緩和時間」で与えられる。この「緩和時間」は、重力多体系のより一般的な緩和過程を記述している。例えば、重力多体系の速度分布は、マクスウェル分布に緩和していくが、これに要する時間も「緩和時間」で与えられる。また、異なる粒子質量を持つ2つの粒子系が混ざりあっている場合、「緩和時間」程度で2つの粒子系は、エネルギー等分配の状態に近づいていく。惑星科学の分野では、力学的摩擦をより広い意味で捉えて、「エネルギー等分配に近づけるもの」という意味で用いる場合が多い (e.g., Hornung et al. 1985; Ida 1990)。

ここで、「チャンドラセカールの緩和時間、 $T_{\text{relax}}$ 」の表式について説明しておこう。質量  $M$  をもつ粒子が、質量  $m$  の粒子の集団に対して相対速度  $v$  で運動している場合、チャンドラセカールの緩和時間は次式で与えられる：

$$T_{\text{relax}} = \frac{M}{m} \frac{1}{n \sigma v}. \quad (1)$$

ここで、 $n$  は、集団の粒子数密度である。また、

重力散乱の断面積  $\sigma$  は

$$\sigma = \pi (GM/v^2)^2 \ln \left( \frac{Lv^2}{GM} \right) \quad (2)$$

で与えられる。上式で、 $G$  は重力定数、 $L$  は系の特徴的な長さである。円盤系においては、特徴的な長さ  $L$  は円盤の厚さである。(1) 式では、粒子集団の熱速度 (速度分散の平方根)、 $v_m$  が相対速度  $v$  に比べて小さいことが仮定されている。逆に、 $v_m$  の方が大きい場合には、上の表式で  $v$  の代わりに  $v_m$  を用いられれば良い。この様に表された緩和時間によって、多くの重力緩和過程を理解することができるのである。

銀河や (原始) 惑星系のような粒子円盤系において、大きな粒子は他の粒子との重力相互作用により、(多くの場合) 円盤中心に向かって落下していく。この粒子の動径方向移動を重い粒子が重力ポテンシャルのより低い所に向かうために起こると考えれば、これを重力緩和過程の1つと解釈できるであろう。多くの研究者は、この「円盤系での動径方向移動」も力学的摩擦による効果と解釈している (e.g., Donner & Sundelius 1993; Wahde et al. 1996)。しかし、円盤系においてはすべての粒子達が同じ回転速度で運動しているので、平均的に見るとそれらは相対速度を持たない。そのため、回転方向には通常の意味での力学的摩擦による抵抗力は働かないのである。円盤系における動径方向移動を考える場合には、このことを気をつける必要があるであろう。

以上のように、力学的摩擦には、「抵抗力」、「速度分布を緩和させるもの」、「動径方向移動させるもの」という3つの意味がある。以下では、力学的摩擦が引き起こすとされている「速度緩和」と「動径方向移動」の2つについて着目し、両者の比較を通じてこれらを理解していこう。

## 2. 円盤系の特徴

円盤重力系には、銀河、惑星系円盤、惑星リングなどがあり、これらは(1)「薄い」、(2)「粒子はほぼ円運動をしている」という特徴を持っている。又、惑星系や惑星リングでは、粒子は中心天体の周りをほぼケプラー運動をしている。上の(1)、(2)の特徴は、軌道の離心率  $e$  や傾斜角  $i$  が小さいことを意味している。

相互作用する2粒子間の相対速度は、 $e$  と  $i$  で決まりこれらとともに大きくなる。(この粒子相対速度はランダム速度とも呼ばれる。) よって、速度緩和は、 $e$  や  $i$  の分布の緩和であるといえる。一方、粒子の動径方向移動は、回転速度のエネルギー変化と関係している。

円盤系では、 $e, i$  が小さいため、相対速度(又はランダム速度)のエネルギーは回転速度のエネルギーに比べ圧倒的に小さい。このため、エネルギーの小さいランダム速度の緩和に比べ、「動径方向の移動」は長時間で進行する。「速度緩和」と「動径方向の移動」が異なる時間スケールで進行することは、円盤系の進化における重要な特徴である。円盤系では、局所的な速度緩和が円盤の各部分で達成された後に、粒子の動径方向の大幅な移動という系全体の進化がゆっくりと進んでいくのである。

## 3. 円盤系での粒子の速度進化

局所的緩和である速度分布の緩和を、惑星系円盤の場合を例にとり具体的に見ていこう。

ランダム速度のもととなる  $e$  や  $i$  の進化は、次の3つの素過程で決まる。1つ目は、力学的摩擦であり、各粒子をエネルギー等分配(又は熱平衡)の状態へ向かわせる。2つ目は、粘性加熱である。これは、差動回転円盤内で粒子が相互作用することにより、回転運動のエネルギーをランダム運動(熱運動)のエネルギーへと変換する効果であり、 $e, i$  を増加させる。3つ目は、粒子に働くガス抵抗であ

り、これはランダム運動を抑える働きをする。 $e, i$  の大きさは、粘性加熱とガス抵抗の釣り合いによって決まる。

これら3つの素過程が働く時間は、どれもチャンドラセカールの緩和時間で与えられ、局所的な緩和はこの緩和時間で進行していく。惑星系円盤において、チャンドラセカールの緩和時間は惑星形成時間に比べ短いため、惑星形成時において「速度の平衡」は常に実現されていると考えて良い。

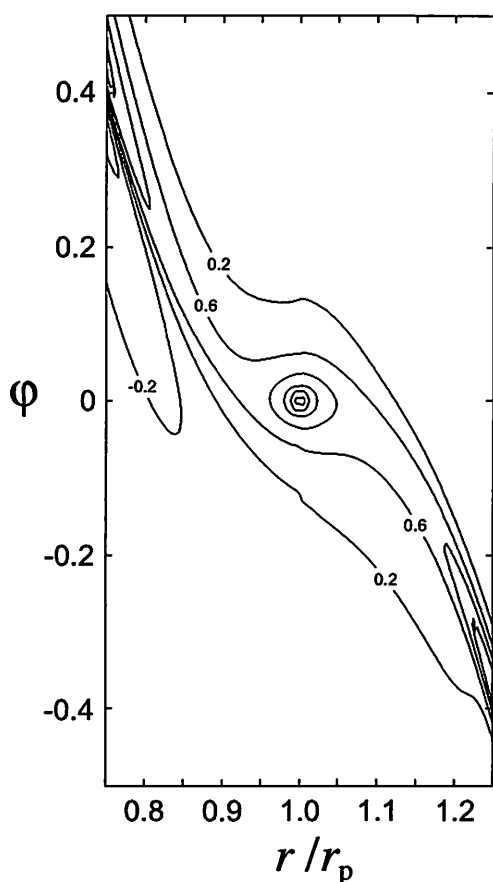


図1: 惑星がガス円盤にたてる密度波の面密度等高線。惑星は図の中心に位置する。惑星軌道の内側と外側に一本ずつ腕状の波が作られている。惑星は内側の腕からは正のトルク、外側の腕からは負のトルクを受け、正味としては負のトルクを受ける(Tanaka et al. 2002より)。惑星がガス円盤から密度波の励起を通して受ける抵抗力は、粒子円盤から力学的摩擦により受ける抵抗力と非常に良く似ている。

粒子に働くガス抵抗についてももう少し詳しく述べておこう。惑星系円盤の場合、月程度より小さい天体は、流体力学的な通常のガス抵抗を受ける。一方、月より大きい天体では、ガス円盤との重力による相互作用が卓越する。この場合のガス抵抗力は、惑星がガス円盤に密度波を励起することにより発生する(図1参照)。惑星のランダム運動のエネルギーが、密度波のエネルギーに変換されるため、「抵抗力」が働くのである(e.g., Goldreich & Tremaine 1980, Artymowicz 1993, Tanaka et al. 2000)。

興味深いことに、密度波励起による抵抗力というメカニズムの違いにもかかわらず、この抵抗力は相対速度を音速とした場合の力学的摩擦の公式で大体説明することができる。このような粒子円盤とガス円盤の間の類似性は、速度緩和の場合だけでなく、後で述べる「動径方向移動」においても成り立っている。

## 4. 円盤系での粒子の動径方向移動

「動径方向の移動」は、粒子が円盤からトルクを受けることにより起こる。もし円盤から正のトルクを受ければ、粒子は角運動量を得て外側に移動することになる。この円盤と粒子のトルクのやりとりは、ガス円盤の場合に対しては Goldreich & Tremaine (1979) に始まり、詳しく調べられている。そこで先ず、ガス円盤との相互作用を説明しよう。粒子円盤の場合に対しては、ガス円盤の場合の類推により、移動速度の見積りをすることができる。

### 4.1. ガス円盤の場合

一般に、粒子は、その軌道に対し円盤の内側部分からと外側部分からは反対のトルクを受ける。このため、粒子が受ける正味のトルクは、内側部分と外側部分との非対称性から生まれる。ガス円盤との相互作用の場合、粒子は図1にみられる密度波からトルクをうけ

る。粒子軌道に対し内側の腕は粒子回転を加速させ正のトルクを、外側の腕は負のトルクをおよぼす。これらの打ち消し合いの結果、正味のトルクが決まる。

正味のトルクを生む非対称には、(1) 円盤の曲率の効果と、(2) 面密度、圧力、温度の動径方向の変化とがある。このような非対称性により天体に働く正味のトルクは、Ward (1986) によって見積もられていた。彼の結果によると、(1) の曲率の効果は粒子に負のトルクをおよぼす。また (2) の効果については、通常ガス円盤で面密度、圧力、温度は、負の動径方向の勾配を持っており、その場合やはり負のトルクを粒子に与える。従って、通常ガス円盤に対して粒子は常に中心へ落下していく。

おおよその移動時間は、「チャンドラセカールの緩和時間」 $T_{\text{relax}}$  を用いて表すことができる。密度波の特徴的な長さは、ガス円盤の厚さ  $h$  で与えられる。今、トルクの打ち消し合いを考えず、外側の円盤からのトルクだけで粒子が移動するとすると、ガス円盤の厚さ  $h$  程度移動するのにかかる時間は、 $T_{\text{relax}}$  で与えられる。実際には、正味のトルクは打ち消し合いの結果ファクター  $h/r$  ( $r$  は粒子の軌道半径) 程度小さくなるので、粒子が軌道半径程度を移動する時間、 $T_{\text{migration}}$  は、

$$T_{\text{migration}} \sim T_{\text{relax}} \times (r/h)^2 \quad (3)$$

で与えられることになる。

地球が原始惑星系円盤内にあった場合、この移動時間は10万年程度になる。この移動時間は惑星形成時間に比べて短く、惑星形成理論において重大な問題となっている。

### 4.2. 粒子円盤の場合

粒子円盤においても、同様な方法で移動時間を見積もることができる。実際に、粒子円盤の場合も、粒子は円盤の内側部分からと外側部分からは反対のトルクを受け、正味のト

ルクは、内側部分と外側部分との非対称性から生まれる。粒子円盤の場合、相互作用する範囲の特徴的長さは、 $e \times r$ で与えられる。これより、粒子円盤での移動時間の見積りにおいては、ガス円盤での厚さ  $h$  の代わりに  $er$  を用いれば良い。即ち、粒子円盤での移動時間に対して次式を得る：

$$T_{\text{migration}} \sim T_{\text{relax}} \times (1/e)^2. \quad (4)$$

前に述べたように、円盤系では離心率  $e$  は 1 に比べて小さいため、移動時間は  $T_{\text{relax}}$  に比べて長くなる。このように得られた移動時間に含まれるファクター  $1/e^2$  は、緩和に必要なエネルギー変化量からも理解することができる。粒子軌道が移動する際には、その軌道エネルギーが大きく変化する。 $e$  や  $i$  に関係したランダム速度のエネルギーに比べて、粒子の軌道エネルギーは、ファクター  $1/e^2$  だけ大きいので、これを変化させるためにはその分長い時間を要する。そのため移動時間には、ファクター  $1/e^2$  が含まれるのである。

以上では簡単な議論によりおおよその移動時間を見積もったが、移動方向の決定は簡単ではない。移動方向は粒子円盤の様子に依存し変わり得るからである。例えば、粒子円盤の面密度が半径とともに増加する場合は、粒子が外側に移動することもあり得る。方向についての詳細なことは、粒子円盤での「曲率の効果」が調べられていないため、まだ明らかになっていないというのが現状である。円盤系においては、「重い粒子が中心に落ち込み、軽い粒子が外側に移動する」のが当然というところまでには至っていないのである。

## 5. まとめ

本稿では、粒子円盤とガス円盤に対して、重力緩和過程としての「粒子速度進化」と「動径方向の移動」を論じた。

速度進化である  $e, i$  の進化は、円盤での局

所的緩和と考えることができ、その緩和過程は「チャンドラセカールの緩和時間」で進行する。一方、「粒子の動径方向の移動」は、系全体での進化と考えることができる。この系全体の進化は、それに伴うエネルギー変化の大きさゆえに、「チャンドラセカールの緩和時間」 $T_{\text{relax}}$  に比べて、ファクター  $1/e^2$  だけ長くなる。

又、ガス円盤と粒子の重力相互作用は、粒子円盤との重力相互作用と非常に似通っており、力学的摩擦やチャンドラセカールの緩和時間でおおよそ説明することができる。

粒子円盤での「粒子の動径方向の移動」の研究は、本稿での大雑把な見積りより正確なものはほとんどない。また、移動の方向を決めるには、曲率の効果などの詳細な研究が必要がある。これらの研究は、今後なされていくべきものであろう。

## 参考文献

- Artymowicz, P. 1993, ApJ, 419, 166
- Chandrasekhar, S. 1943a, ApJ, 97, 255
- Chandrasekhar, S. 1943b, Principle of Stellar Dynamics (New York: Dover)
- Donner, K.J. & Sundelius, B. 1993, MNRAS, 265, 88
- Goldreich, P. & Tremaine, S. 1979, ApJ, 233, 857
- Goldreich, P. & Tremaine, S. 1980, ApJ, 241, 425
- Hornung, P., Pellat, R. & Barge, P. 1985, Icarus, 64, 295
- Ida, S. 1990, Icarus, 88, 129
- Tanaka, H., Takeuchi, T. & Ward, W.R. 2000, Proceedings of 31st LPSC
- Tanaka, H., Takeuchi, T. & Ward, W.R. 2002, ApJ, 565, 1257
- Wahde, M., Donner, K.J., Sundelius, B. 1996, MNRAS, 281, 1165
- Ward, W.R. 1986, Icarus, 67, 164
- Ward, W.R. 1988, Icarus, 73, 330



# Life Expectancy of the Large Magellanic Cloud

Yoko Funato<sup>1</sup>, Yoshikazu Hashimoto<sup>2</sup>, Junichiro Makino<sup>2</sup>

(1) *General Systems Studies, Graduate Division of International and Interdisciplinary Studies,  
University of Tokyo, 3-8-1, Komaba, Meguroku, Tokyo, 153, Japan*

funato@chianti.c.u-tokyo.ac.jp

(2) *Department of Astronomy, Graduate School of Science,  
the University of Tokyo*

hashimoto@astron.s.u-tokyo.ac.jp, makino@astron.s.u-tokyo.ac.jp

## ABSTRACT

We investigated the orbital evolution of satellite galaxies using numerical simulations. It has been long believed that the orbit suffers circularization due to the dynamical friction from the galactic halo during orbital decay. This circularization was confirmed by numerous simulations where dynamical friction is added as external force. However, some of the recent  $N$ -body simulations demonstrated that circularization is much slower than expected from approximate calculations. In this study we will show that

- (1) *The discrepancy really exists, in other words, it is not any of error caused during numerical simulations.*
- (2) *The dominant reason for the discrepancy is the assumption that Coulomb logarithm  $\log \Lambda$  is constant, which has been used in practically all recent calculations.*

Since the size of the satellite is relatively large, accurate determination of the outer cutoff radius is crucial to obtain good estimate for the dynamical friction. An excellent agreement between  $N$ -body simulations and approximate calculations was observed when the outer cutoff radius is taken to be the distance of the satellite to the center of the galaxy. When satellite is at the periastron, the distance to the center is smaller and therefore  $\log \Lambda$  becomes smaller. As a result, the dynamical friction becomes less effective.

- (3) *Applying our result to orbital evolution of the Large Magellanic cloud, the expected lifetime of the LMC is twice as long as that would be predicted with previous calculations.*

Previous study predicts that the LMC will merge into the Milky Way after 7 G years, while we found that the merging will take place after 14 G years from now. Our result suggests that generally satellites formed around a galaxy have longer lifetime than previous estimates.

*Subject headings:* celestial mechanics, stellar dynamics — Galaxy: kinematics and dynamics — galaxies: Magellanic Clouds — Local Group — methods: numerical

## 1. Introduction

Recent observations have revealed that there are many satellite galaxies around the Milky Way. In the hierarchical clustering scenario, it is expected many of such dwarf satellites are formed. In fact, one of the most serious problems with the present hierarchical clustering scenario is that it

predict too many satellite galaxies, about a factor of 10 more than the number observed in the Local group (*e.g.*, Moore *et al.*, 1999). A number of explanations, including exotic theories which relies on hot or self-interacting dark matter, have been proposed.

In this paper, we go back to the basic problem: how long are the satellites lives? In other words,

how do the orbits of satellites evolve through interaction with the gravitational field of its parent galaxy? The dominant driving force of the evolution is the dynamical friction. For satellites like the LMC-SMC pair and the Sagittarius dwarf, there are many detailed studies of their orbital evolution, in which the dynamical friction is included as the external force operating on the center-of-mass motion of the satellite. Well known works include Murai and Fujimoto (MCs) and Ibata and Lewis (Sagittarius). In both of these studies, and in all other studies where the dynamical friction formula is used, significant circularization of the orbit of the satellite is observed. This circularization is the natural result of the fact that the dynamical friction is proportional to the local density of the background stars, and therefore the strongest at the periastron.

However, recent  $N$ -body simulations of the orbital evolution of satellites resulted in rather counter-intuitive result. Van den Bosch et al (1999, hereafter BLLS) performed the  $N$ -body simulation of the satellite, where the parent galaxy is modeled directly as self-consistent  $N$ -body system. The satellite is modeled as one massive particle with spline potential softening used in PKDGRAV (Dikaiakos & Stadel, 1996). They investigated the evolution of the orbit for wide variety of model parameters such as the mass of the satellite and initial orbital eccentricity. They observed practically no circularization in any of their simulations.

Jiang and Binney (2000, hereafter JB) performed fully self-consistent simulation of the satellite, where both the parent galaxy and the satellite are modeled as self-consistent  $N$ -body systems. They compared their result with the result of approximate model in which the usual dynamical friction formula is used. Though they argued that the agreement is good, from their figure 3 it is clear that approximate models suffer strong circularization and evolve faster than their  $N$ -body counterpart.

Neither of above two papers discussed the reason of this rather serious discrepancy between the result of  $N$ -body simulations and previous analytic prediction. The purpose of this paper is to understand its cause. In section 2, we describe our model experiment designed to reproduce the discrepancy observed by BLLS and JB. In section 3

we show our result. Our result is consistent with both of the previous works.  $N$ -body simulation showed only marginal circularization but approximate calculation using dynamical friction formula showed strong circularization. In section 4, we investigate the reason. There are several possible candidates for the reason. We consider a few of them, and found that a simple modification of the conventional form of the dynamical friction formula results in a quite remarkable improvement of the agreement between  $N$ -body and approximate calculations. In section 5 we apply our formalism to the LMC. In this cases, orbital evolution becomes significantly slower than prediction by previous calculations using conventional formula. The lifetime of LMC was 7 Gyr with conventional formula, but is 14 Gyr with our formalism. We also discuss the implication of our result to the so-called “dwarf problem”.

## 2. Numerical Simulation

We carried out a set of numerical simulations to see whether the results obtained by BLLS and JB are really true or not. In this section, we describe the models we used.

### 2.1. $N$ -body simulation

We performed  $N$ -body simulations of the evolution of a satellite orbiting in a massive dark halo of a galaxy.

The massive halo is composed by  $N$  equal mass particles, while the satellite dwarf is modeled by a single particle with a certain softening length. The softening is used to mimic the finite size of the satellite.

We adopted a King model of the concentration ratio  $\Psi_0 = 9$  as a model of the galactic halo. The system of units is the Heggie unit (Heggie and Mathieu (1986)) where the gravitational constant  $G$  is 1, the mass and the binding energy are 1 and 0.25, respectively.

JB used a composite disk+halo model in which the halo is expressed by particles and the disk is assumed to be rigid. BLLS used a single spherical halo. In both works, the halo density profile has the form

$$\rho = \rho_0 \frac{r_c^2 \exp[-(r/r_0)^k]}{r_c^2 + r^2}, \quad (1)$$

where  $r_c$  and  $r_t$  are the core radius and the outer scale radius of the halo and  $\rho_0$  is the central density of the halo. BLLS adopted  $k = 2$  while JB adopted  $k = 1$ .

We did not follow the models in their works. The standard dynamical friction formula is derived for the case of field stars with the Maxwell distribution. However, the distribution function associated with eq. (1) is rather different from the Maxwell distribution. This may cause difference in the effect of the dynamical friction. Also, the distribution function would relax to the Maxwellian through two-body relaxation, causing a small change in both the distribution function and the density profile.

In addition, the range of radius for which the density slope is approximately  $-2$  is rather narrow with this model, since the slope is noticeably shallower than  $-2$  for  $r \leq 10r_c$ .

The distribution function of the King model is a simple lowered Maxwellian. Therefore the agreement with the true Maxwellian is very good within the half-mass radius. Also, since the distribution function is practically as close as the true Maxwellian as we can make, thermal relaxation is minimized, though it still present (see e.g., Quinlan 1996?). Also, the King model with  $W_0 = 9$  has fairly wide range of radius in which the slope of the density is approximately  $-2$ . So it is a fairly good model for a spherical halo with flat rotation.

The satellite galaxy is modeled by a single particle with mass  $M_s$  and softening length  $\epsilon_s$ . The force on the satellite from a particle in the halo is calculated as follows

$$\mathbf{F} = -\frac{GmM_s(\mathbf{r}_{\text{sat}} - \mathbf{r}_{\text{halo}})}{(|\mathbf{r}_{\text{sat}} - \mathbf{r}_{\text{halo}}|^2 + \epsilon_s^2 + \epsilon_{\text{halo}}^2)^{\frac{3}{2}}}. \quad (2)$$

Here  $\epsilon_{\text{halo}}$  is the softening length for the particle in the galactic halo. The value of the gravitational constant  $G$  is 1 in the standard units.

In our simulations, equations of motions of all particles in a dark halo and the satellite, *i.e.*,  $N+1$  particles, are integrated self-consistently. In other words, the dynamical friction effect from halo particles to the satellite is included naturally.

We used GRAPE6 to calculate the acceleration. We adopted simple  $O(N^2)$  direct summation, to avoid any possible numerical artifact caused by the approximations made in force cal-

culation. BLLS used the treecode and JB used a composite grid-based code. We do not think the numerical method caused the difference, but we want to be absolutely sure that our  $N$ -body simulation is as accurate as possible. The number of particles  $N$  used in the simulations shown in this paper is 32768. We varied  $N$  from 8192 to 32768, and found any noticeable difference in the orbit of the satellite. We integrated the orbits of the satellites and halo particles using the standard leapfrog scheme.

## 2.2. Semi-analytic Integration

We performed semi-analytic calculations to follow the evolution of satellite orbits.

In these calculations, the model of the satellite is the same as in the  $N$ -body simulations, *i.e.*, a single particle with mass  $M_s$  and the softening length  $\epsilon_s$ .

Instead of being represented by  $N$  particles, the potential of the galactic halo is evaluated by using the gravitational potential of King 9 model with the same mass and scales as those adopted in  $N$ -body simulations.

In this integration, the force to the satellite due to the dynamical friction from the halo is evaluated by using an analytical formula, too.

For the dynamical friction formula, we follow JB (and also Murai and Fujimoto) to use the standard ‘‘Chandrasekhar’s dynamical friction formula’’. It is expressed as

$$\frac{d\mathbf{v}}{dt} = -16\pi^2 G^2 m(M_s + m) \ln \Lambda \frac{\int_0^{v_{\text{max}}} f(v) d\mathbf{v}}{|\mathbf{v}|^3} \quad (\S)$$

where  $M_H$  and  $M_{\text{sat}}$  are the masses of host galaxy and its satellite galaxy (Chandrasekhar 1943; Binney and Tremaine 1987). Here  $\ln \Lambda$  is the Coulomb logarithm

$$\ln \Lambda = \ln(R_{\text{halo}}/\epsilon_s V_s^2), \quad (4)$$

where  $R_{\text{halo}}$  is the scale length of the galactic halo. This formula has been adopted by many semi-analytic studies of the orbital evolution of satellite galaxies (*e.g.*, Murai and Fujimoto 1980; Helmi and White, 1999; Johnston *et al.*, 1995). It is also used in cosmological studies of galaxy formation in order to estimate the merging time scale of satellite galaxies (*e.g.*, Kauffmann, *et al.*, 1994).

### 3. Result

Figure 1 shows the orbital evolution of a model satellite galaxy. The ordinate and abscissa are the distance of the satellite from the center of the galaxy and time in the  $N$ -body units. The solid and dashed curves correspond to result of  $N$ -body simulation and that of semi-analytic model with standard dynamical friction formula (3).

In Figure 1 two curves are in good agreement only for a first few dynamical times. After a few orbits, two curves deviate from each other. Figure 1 shows that the orbital decay calculated with formula (3) is faster than that obtained by  $N$ -body simulation. If one measure the orbital eccentricity, it is clear that  $N$ -body result shows only a small change in the eccentricity, while semi-analytic result shows significant circularization.

Thus, even though we used completely different initial models and numerical method, we confirmed previous results by BLLS and JB that  $N$ -body simulation shows little circularization while semi-analytical calculation with standard dynamical friction formula shows strong circularization. In the next section, we discuss the possible causes of this discrepancy.

### 4. Possible causes of discrepancy

Since we have obtained quite different results with  $N$ -body and semi-analytic models, *at least* one of them must be wrong. Since  $N$ -body calculation can suffer many numerical problems due to limited resolution and particle noise, one might think  $N$ -body result is probably wrong. However, additional tests with different number of particles and different sizes of timestep showed very good agreement (Hashimoto *et al.*, 2002). Therefore it seems our  $N$ -body result is sound. In addition, as we stated in the previous section, our  $N$ -body result is in good agreement with BLLS and JB. Though it is not impossible, it is certainly unlikely that all of these three works are wrong.

So let us now consider the possibility that the standard dynamical friction formula is wrong.

The standard dynamical friction formula is obtained under the assumption that the massive object moves straight in a uniform and isotropic distribution of field particles. Field particles are also assumed to be moving straight, and any interac-

tion between field particles is ignored. Clearly, the satellite does not move straight, but circle around the center of the parent galaxy. The distribution of field stars within the parent galaxy is far from uniform, and field stars also circle around in the parent galaxy. Thus, it is not really surprising that the naive use of the dynamical friction formula gives rather bad result.

One obvious way to improve the accuracy of the dynamical friction formula is to calculate the linear response of the global distribution function of the parent galaxy to the presence and the orbit of the satellite (Weinberg, 1995). This approach would certainly give accurate and reliable result which agrees well with  $N$ -body result (Hernquist and Weinberg 1989). However, since the global response depends on the distribution function itself, the result cannot be expressed in a compact and form. So here we consider the possibility to improve the standard formula.

As we noted above, there are at least two problems with the standard formula. First, it assumes that both the satellite and field stars move straight. Second, it assumes that the density of the field star is the same everywhere.

The first assumption is clearly wrong, but its effect is difficult to estimate. Let us consider the effect of the second assumption, which is much easier to evaluate. In previous works, the outer cutoff radius of the Coulomb logarithm is taken to be the scale length of the halo, while the representative density of the field stars is taken to be the local density around the satellite. This would clearly cause an overestimate of the Coulomb integral, for the case of the singular isothermal sphere (or the King model we used), since the stellar density drops off as fast as  $1/r^2$ . This means the logarithmic divergence of the Coulomb integral does not actually occur if we takes into account the effect of the density gradient.

To correctly take into account the effect of the density gradient is a tricky problem, since for encounters with impact parameter comparable or larger than  $R_s$ , the distance to the center of the galaxy, we cannot really use the straight line approximation. On the other hand, just to ignore any encounter with impact parameters  $R_s$  might not be too bad assumption, since density drops off rapidly and realistic effect is unlikely to enhance the effect of the encounter (except for the small

fraction of the orbits in resonance with the orbit of satellite).

Thus, it might be more sensible to use  $R_s$  as the outer cutoff radius for the Coulomb logarithm, that to use the traditional  $R_{halo}$ . In fact, this use of  $R_s$  is first proposed by a pioneering work by Tremaine (1976) on the effect of the dynamical friction to the orbit of LMC-SMC pair.

To use  $\epsilon_s$  as the inner cutoff is okay as an order-of-magnitude estimate, but can be improved by actually integrating the effect of all encounters with small impact parameters for Plummer potential, following the treatment by White (1976). For Plummer model, the integration can be performed analytically and the result is that inner cutoff radius is  $r_{in} = 1.4\epsilon_s$ .

Figure 2 is the same as Figure 1 but for the above discussed choice of the Coulomb logarithm

$$\ln \Lambda = \ln \left( \frac{R_s}{1.4\epsilon_s} \right). \quad (5)$$

When the  $R_s$  becomes smaller than  $1.4\epsilon_s$ , we simply put the dynamical friction term to be zero, since it is clearly unphysical to apply dynamical “acceleration”.

The agreement between the  $N$ -body result and semi-analytic treatment is quite remarkable.

Figure 3 shows evolution of eccentricities. In Figure 3, solid, thin dashed and thick dashed curves corresponds to the result of  $N$ -body simulation, that of semianalytic formula with constant  $\Lambda$  and that with varying  $\Lambda$ , respectively. The results of  $N$ -body simulation and that obtained using varying  $\Lambda$  formula demonstrate good agreement, while the result of calculation using a constant  $\Lambda$  does not.

In Figure 4, evolution of eccentricities are plotted against that of apogalactic distance  $R_{max}$ . In Figure 4, curves are the same as those in Figure 3. The results of  $N$ -body simulation and that obtained using varying  $\Lambda$  formula demonstrate good agreement, while the result of calculation using a constant  $\Lambda$  does not. Figure 4 shows that quick circularization appearing in result of semi-analytic integration using constant  $\Lambda$  is not a matter of time-scale of orbital evolution. Instead, there is a qualitative difference in the understanding and use of dynamical friction.

Figures 2 and 3 show that the discrepancy

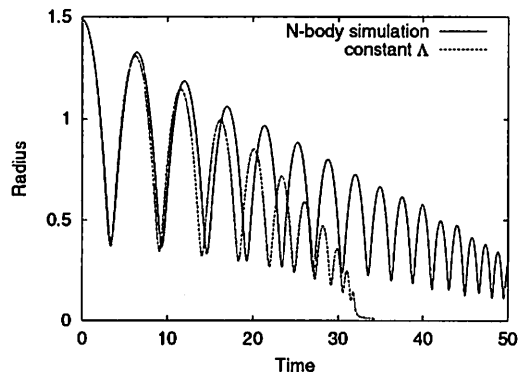


Fig. 1.— Time evolution of radius of satellite position from the galaxy center. Solid: result of  $N$ -body simulation. Dashed: semi-analytical integration using constant  $\Lambda$ .

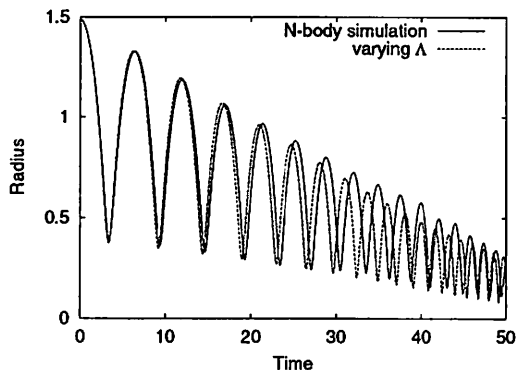


Fig. 2.— Same as Figure 1, but for the variable  $\Lambda$ .

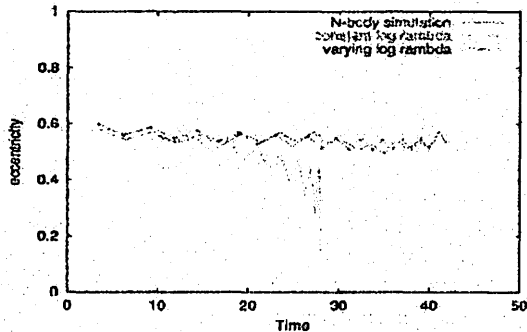


Fig. 3.— Time evolution of eccentricities. Result of semi-analytic integration using varying  $\Lambda$  is in good agreement with that of  $N$ -body simulation, while that using constant  $\Lambda$  is not.

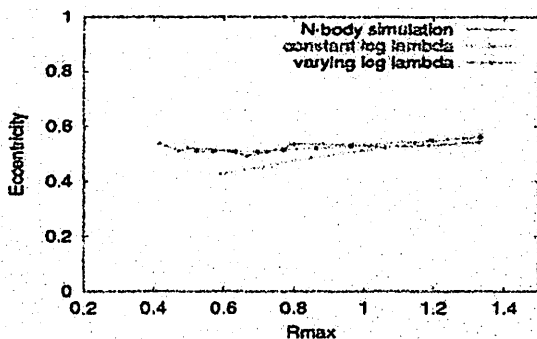


Fig. 4.— Evolution of eccentricities are plotted against that of apogalactic distance. Result of semi-analytic integration using varying  $\Lambda$  is in good agreement with that of  $N$ -body simulation, while that using constant  $\Lambda$  is not.

shown in Figure 1 is caused by an inadequate estimate of  $\Lambda$ . Figure 4 shows that the difference between Figure 1 and Figure 2 and the reason for the discrepancy in Figure 1 are never a matter of time-scale of orbital evolution. Other possible reasons, such as the effect of the global response of the distribution function, might still be important, but they are clearly not the prime reason of the discrepancy between  $N$ -body and semi-analytic works which we discussed in the introduction and section 3.

The improved agreement with the  $N$ -body result is explained as follows. With  $b_{max} = R_{cut}$ , the semi-analytical treatment causes strong circularization and faster orbital evolution. This implies that the semi-analytical treatment overestimated the dynamical friction around the periastron. Around the apoastron, the error might exist, but relatively small compared to that at the perigalacticon. The use of variable  $b_{max}$  reduces the value of  $\ln \Lambda$  both at perigalacticon and apogalacticon, but by a much larger factor at the perigalacticon simply because  $R_p$  is smaller. Thus, effectively we reduced the dynamical friction around the periastron, which resulted in the improvement in the agreement with the  $N$ -body result.

In hindsight, it looks too obvious that the traditional use of the dynamical friction formula was inappropriate. Theoretically, it is clearly not justifiable to assume that the stellar density is the same up to the outer cutoff radius of the halo. From comparison between the  $N$ -body result and those of semi-analytic treatment, it also is clear that previous semi-analytic treatment overestimates deceleration due to the dynamical friction around the perigalacticon.

To summarize our result, the orbital decay of satellites is slower than ever estimated, the eccentricity of orbit of revolution of a satellite around the host galaxy is almost constant. The reason why previous estimates are wrong is that previous studies overestimated the effect of dynamical friction at the perigalacticon.

## 5. Summary and Applications

We performed  $N$ -body simulations of satellite orbits. We found that the circularization of the orbit due to the dynamical friction is much slower than commonly believed. This discrepancy was

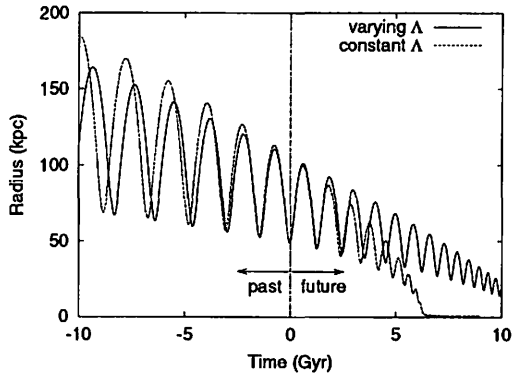


Fig. 5.— Radial Evolution of LMC. From -10 G years to 10 G years.

also reported by BLLS, and we can see the same tendency from the numerical result reported by JB.

Previous studies of satellite orbits used the outer cutoff radius of the dark halo as  $b_{max}$ . We found that the effective  $b_{max}$  should be of the order of  $R_s$ , the distance of the satellite from the center of the galaxy, which varies as the satellite orbits around the galaxy. Our formula results in a greatly improved agreement with the  $N$ -body result.

### 5.1. Application to the Milky Way

We investigated the orbital evolution of the Large Magellanic Cloud and the Sagittarius; two of the most famous satellites of the MW.

#### 5.1.1. the Large Magellanic Cloud

The Large Magellanic Cloud is the most famous satellite of Milky Way. Its orbit has been investigated from both observation and numerical simulations (*e.g.*, Toomre, 1970; Tremaine, 1976; Lin and Lynden-Bell, 1977; Murai and Fujimoto, 1980). The importance of the effect of dynamical friction from the galactic halo on the orbit evolution LMC is first emphasized by Tremaine (1976).

By using numerical simulation, Murai and Fujimoto (1980) (hereafter MF80) determined the orbital elements and the present phase of the LMC. They performed a number of backward numerical integrations of the orbits of the LMC and SMC from various initial conditions, and integrated or-

bits of test particles in the LMC and SMC for each condition. Comparing the result of distribution of test particles and the observed Magellanic stream, they chose the initial condition which gives the best fit.

In their numerical integration, they assumed a halo expressed by a singular isothermal sphere, which is a simple flat-rotation halo. In their paper, it is not clear what assumption is used for  $\ln \Lambda$ , since there is no discussion on how they determined  $\ln \Lambda$  though it appeared in their equation (13). However, the fact that the orbit of the LMC obtained in their calculation shows significant circularization strongly suggest that they treated  $\ln \Lambda$  as constant.

In order to see the effect of changing  $\ln \Lambda$ , we integrated the orbit of LMC both forward and backward in time, using both the constant  $\Lambda$  and variable  $\Lambda$  ( $b_{max} = R_s$ ). In this study, we express LMC as a single Plummer-softened particle with mass  $2 \times 10^{10} M_\odot$  and softening length 5 kpc. The rotation velocity of the halo is 250 km/s, same as what is used by MF. We simulated the orbit of the LMC only, since our purpose here is to demonstrate the effect of  $\Lambda$  and not the accurate determination of the orbits of the Clouds.

The solid curve in Figure 5 corresponds to the orbit obtained when the dynamical friction is calculated using equation (5). The dashed curve in Figure 5 correspond to the orbit obtained using the formula (3). The backward part of this dashed curve is in very good agreement with the result of MF, indicating that what MF used is indeed a constant  $\Lambda$ .

Figure 5 shows that real evolution of the orbit of LMC (with variable  $\Lambda$ ) is significantly smaller than what is obtained by MF. 10 Gyrs ago, the "true" apocenter was only 160 kpc, while the solution by MF was 180 kpc.

A more remarkable difference is in the future of the LMC. With the constant  $\Lambda$ . The LMC will fall to the galactic center in only 7 Gyrs with constant  $\Lambda$ , while our result suggests that it will take more than 14 G years for the LMC to fall to the galactic center.

### 5.2. Eccentricity Distribution of Satellites

Our study shows that the time evolution of the eccentricity of satellites is rather small. Thus, we

may assume that the distribution of eccentricities of satellite galaxies at present directly reflects that at the formation epoch of the Galaxy. Therefore the distribution of eccentricities of satellites galaxies can be an important clue to the formation of the Galaxy.

### 5.3. Number Evolution of Faint Galaxies

Cosmological studies on galaxy formation are based on this estimate and discuss the number evolution of galaxies.

The lifetime of the satellite is estimated using the dynamical friction timescale with  $\ln \Lambda$  taken to be  $M_H/M_s$  (Lacey and Cole 1993; Kauffmann et al. 1994). This would cause a quite serious overestimate in the dynamical friction timescale, since the factor one should use is the ratio between the size of the halo and the size of the satellite. If we assume  $M \propto \sigma^4$ , we have  $R \propto M^{1/2}$ . Thus, there is at least a factor of two difference in the value of  $\ln \Lambda$ . Since there are too many other uncertainties in the semi-analytic modelling of the galaxy number evolution, how serious this difference is is not clear. However, it certainly affects the estimate of presently observed satellites rather strongly. A more detailed study on this aspect is clearly necessary.

## REFERENCES

- Binney, J., and Tremaine S. Galactic Dynamics, 1987, Princeton University Press
- Chandrasekhar, S. 1943, ApJ, 97, 255.
- Dikaiakos, M., & Stadel, J. 1996, in Proc. Industrial Computing Soc. Conf. (Research Triangle Park, NC: Instrument Soc. Am.)
- Hashimoto, Y., and Makino, J., 2002, in proceedings of IAU symp. 208.
- Hashimoto, Y., Funato, Y., and Makino, J., 2002, in preparation.
- Heggie, D. C. and Mathieu, R. D., 1986, "Standardised Units and Time Scales", in The Use of Supercomputers in Stellar Dynamics, p233.
- Helmi, A., White, S. D. M., de Zeeuw, P. T. & Zhao, H. ,1999, Nature 402, 53-55.
- Hernquist, L. and Weinberg, M. D., 1989, MNRAS, 238, 407.
- Ibata, R. A., & Lewis, G. F. 1998, ApJ, 500, 575.
- Jiang, I.-G., & Binney, J. 2000, MNRAS, 314, 468.
- Johnston, K. V., Spergel, D. N. & Hernquist, L. 1995, ApJ, 451, 598.
- Kauffmann, G., Guiderdoni, B., & White, S. D. M. 1994, MNRAS, 267, 981.
- Lacey, C., & Cole, S. 1993, MNRAS, 262, 627.
- Moore, B., Ghigna, S., Governato, F., Lake, G., Quinn, T., Stadel, J. & Tozzi, P. 1999, ApJ, 524, L19.
- Murai, T. and Fujimoto, M. 1980, PASJ, 32, 581.
- Navarro, J. F., Frenk, C. S., & White, S. D. M. 1995, MNRAS, 275, 720.
- Van den Bosch, F. C., Lewis, G. F., Lake, G., and Stadel, J. 1999, ApJ, 515, 50.
- Weinberg, M. D. 1995, ApJ, 455, L31.
- White, S. D. M., 1976, MNRAS, 174, 467.

---

This 2-column preprint was prepared with the AAS L<sup>A</sup>T<sub>E</sub>X macros v5.0.



# New Limits on the Mass of the Milky Way

Tsuyoshi Sakamoto

Department of Astronomical Science, The Graduate University for Advanced Studies,  
Mitaka, Tokyo 181-8588, Japan  
email: sakamoto@pluto.mtk.nao.ac.jp

Masashi Chiba

National Astronomical Observatory, Mitaka, Tokyo 181-8588, Japan  
email: chibams@gala.mtk.nao.ac.jp

Timothy C. Beers

Department of Physics and Astronomy, Michigan State University, East Lansing, MI 48824  
email: beers@pa.msu.edu

## ABSTRACT

We set new limits on the mass of the Milky Way, making use of the latest kinematic information for Galactic satellites and halo objects. Our sample consists of 11 satellite galaxies, 137 globular clusters, and 413 field horizontal-branch stars at large distances from the sun. Roughly half of the objects in this sample have measured proper motions, permitting the use of their full space motions in our analysis. Two alternative methods of mass estimation are explored in this paper. First, the constraint that rest-frame velocities of the sample objects be lower than their escape velocities at their estimated distances, provided by prescribed Galactic potentials, provides a lower limit on the total mass of the Galaxy of  $1.3 \sim 1.4 \times 10^{12} M_{\odot}$ . We demonstrate that this mass estimate is basically determined by the motions of seven high-velocity objects (Leo I, Pal 3, Draco, and four horizontal-branch stars), *not* by a single object alone (such as Leo I), as has often been the case in past analyses. We also find that a gravitational potential that gives rise to a declining rotation curve is insufficient to bind many of our sample objects to the Galaxy. Second, for a family of phase-space distributions in a potential with a flat rotation curve, a Bayesian likelihood approach is used to reproduce the observed distribution of current positions and motions of the sample. This method enables a search for the most likely total mass of the Galaxy, without suffering a large influence in the final result due to the presence or absence of Leo I, provided that both radial velocities and proper motions are used. Although the best mass estimate depends somewhat on the model assumptions, such as the unknown prior

probabilities for the model parameters, the resultant systematic change in the mass estimate is confined to a relatively narrow range of a few times  $10^{11} M_{\odot}$ . The most likely total mass derived from this method is  $2.5^{+0.5}_{-1.0} \times 10^{12} M_{\odot}$  (including Leo I), and  $1.8^{+0.4}_{-0.7} \times 10^{12} M_{\odot}$  (excluding Leo I). The mass estimate within the distance to the Large Magellanic Cloud ( $\sim 50$  kpc) is essentially independent of the model parameters, yielding  $5.5^{+0.0}_{-0.2} \times 10^{11} M_{\odot}$  (including Leo I) and  $5.4^{+0.1}_{-0.4} \times 10^{11} M_{\odot}$  (excluding Leo I). Implications for the origin of halo microlensing events (e.g., the possibility of brown dwarfs as the origin of the microlensing events toward the LMC may be excluded by our lower mass limit) and prospects for more accurate estimates of the total mass are also discussed.

*Subject headings:* Galaxy: halo — Galaxy: fundamental parameters — Galaxy: kinematics and dynamics — stars: horizontal-branch

## 1. INTRODUCTION

Over the past decades, various lines of evidence have revealed that the mass density in the Milky Way is largely dominated by unseen dark matter, from the solar neighborhood to the outer reaches of the halo (e.g., Fich & Tremaine 1991). Moreover, the presence of a dark component similar to that found in our own Galaxy appears to be a generic feature in external galaxies, as inferred from, e.g., flat rotation curves in their outer parts, the presence of (a gravitationally bound) hot plasma in early-type galaxies, and the observed gravitational lensing of background sources (e.g., Binney & Tremaine 1987). A determination of the extent over which such dark-matter-dominated mass distributions apply for most galaxies, including our own, is of great importance for understanding the role of dark matter in galaxy formation and dynamical evolution. In particular, the mass estimate of the Galaxy is closely relevant to understanding the origin of the microlensing events toward the Large Magellanic Cloud (LMC) (e.g., Alcock et al. 2000; Alcock et al. 2001).

While mass estimates of external galaxies can (in principle) be obtained in a relatively straightforward fashion using various dynamical probes, the total mass of the Galaxy remains rather uncertain, primarily due to the lack of accurate observational information for its outer regions, where the dark matter dominates. The precise shape of the outer rotation curve, as deduced from H II regions and/or H I gas clouds (e.g., Honma & Sofue 1997), is still uncertain because its determination requires knowledge of accurate distances to these tracers (Fich & Tremaine 1991). Also, interstellar gas can be traced only up to

$\sim 20$  kpc from the Galactic Center, and hence provides no information concerning the large amount of dark matter beyond this distance.

The most suitable tracers for determination of the mass distribution in the outer halo of the Galaxy are the distant luminous objects, such as satellite galaxies, globular clusters, and halo stars on orbits that explore its farthest reaches (e.g., Miyamoto, Satoh, & Ohashi 1980; Little & Tremaine 1987; Zaritsky et al. 1989; Kochanek 1996; Wilkinson & Evans 1999, hereafter WE99). However, the limited amount of data presently available on the full space motions of these tracers, and the small size of the available samples, have stymied their use for an accurate determination of the Galaxy's mass. In particular, most previous mass estimates (except for WE99, see below) depend quite sensitively on whether or not a distant satellite, Leo I, is bound to the Galaxy. Leo I has one of the largest radial velocities of the known satellites, despite its being the second most distant from the Galaxy (Mateo 1998; Held et al. 2001). As a consequence, estimates of the total mass of the Galaxy are much more uncertain (by as much as an order of magnitude) than, for instance, the value of the circular speed in the solar neighborhood (Kerr & Lynden-Bell 1986; Fich & Tremaine 1991; Miyamoto & Zhu 1998; Méndez et al. 1999).

Recently, by making use of both the observed radial velocities and proper motions of six distant objects, WE99 demonstrated that the use of full space motions can provide a reliable mass estimate of the Galaxy without being largely affected by the presence or absence of Leo I. They also argued that the primary uncertainties in their mass estimate arose from the small size of the data and the measurement errors in the full space motions, especially the proper motions. This work motivated us to investigate a much larger data set, with more accurate kinematic information, to set tighter limits on the mass of the Galaxy. Specifically, as we show below, there are two objects among the WE99 sample (Draco and Pal 3) that have relatively large velocity errors, yet still play crucial roles in a determination of the Galaxy's mass, so the addition of more (and better data) is important.

Over the past few years, the number of distant satellite galaxies and globular clusters with available proper motions has gradually increased (e.g., Mateo 1998; Dinescu, Gerard, & van Altena 1999; Dinescu et al. 2000; Dinescu et al. 2001). In addition, another tracer population that is suitable for exploring mass estimates of the Galaxy has become available from the extensive compilation of A-type metal-poor stars by Wilhelm et al. (1999b), which provided radial velocity measurements, as well as estimates of the physical parameters of these stars (e.g.,  $[\text{Fe}/\text{H}]$ ,  $T_{\text{eff}}$ ,  $\log g$ ). Among the Wilhelm et al. sample, the luminous field horizontal-branch (FHB) stars are the most useful mass tracers, both because of their intrinsic brightness, and the fact that accurate distance determinations can be inferred from their absolute magnitudes on the horizontal branch (e.g., Carretta et al. 2000). Moreover,

there exist proper-motion measurements for many of these stars, provided by both ground- and space-based proper-motion catalogs (Klemola, Hanson, & Jones 1994; Röser 1996; Platais et al. 1998; Hog et al. 2000), from which full space motions may be derived.

In this paper we re-visit the mass determination of the Galaxy, based on a sample of 11 satellite galaxies, 137 globular clusters, and 413 FHB stars, out of which 5 satellite galaxies, 41 globular clusters, and 211 FHB stars have measured proper motions. We adopt two different methods for obtaining this mass estimate: (1) A method based on the requirement that the rest-frame velocities of observed sample objects be less than their escape velocities at their present distance from the Galactic center (e.g., Miyamoto et al. 1980; Carney, Laird, & Latham 1988), and (2) A method, based on a Bayesian likelihood analysis, that seeks to reproduce both the current positions and velocities of the sample objects (e.g., Little & Tremaine 1987; Kochanek 1996; WE99). Because our present sample of tracers is, by far, the largest and most accurate one available, it is possible to place more reliable limits on the total mass of the Galaxy. In § 2 we describe our sample objects and the assembly of their kinematic data. In § 3 and § 4 the results on the mass estimates of the Galaxy are presented. § 3 is devoted to the method of mass estimation based on escape velocities; in § 4 we adopt a Bayesian likelihood method to obtain the most likely total mass of the Galaxy. In § 5 we discuss implications for the origin of the halo microlensing events toward the LMC and the mass estimate of the Local Group, and consider the prospects for more obtaining more accurate estimates of the total mass of the Galaxy in the near future.

## 2. DATA

We consider a sample of objects that serve as tracers of the Galactic mass distribution consisting of 11 satellite galaxies, 137 globular clusters, and 413 FHB stars. In the case of the satellite galaxies, all of the basic information for our kinematic analysis, i.e., positions, heliocentric distances, and heliocentric radial velocities, are taken from the compilation of Mateo (1998). For the globular clusters, we adopt the information provided by Harris (1996), including their positions and heliocentric radial velocities, their metal abundances,  $[\text{Fe}/\text{H}]$ , and the apparent magnitude of the clusters' horizontal branch (HB). The catalog of Wilhelm et al. (1999b) is our source of similar information for the FHB stars. We derive an internally consistent set of distance estimates for the globular clusters and the FHB stars from the recently derived relationship between the absolute magnitude of the HB,  $M_V(\text{HB})$ , and  $[\text{Fe}/\text{H}]$ , by Carretta et al. (2000),

$$M_V(\text{HB}) = (0.18 \pm 0.09)([\text{Fe}/\text{H}] + 1.5) + (0.63 \pm 0.07) . \quad (1)$$

Clearly, we have assumed that there is no large offset between the absolute magnitudes of FHB stars and their counterpart HB stars in the globular clusters (a view also supported by the recent work of Carretta, Gratton, & Clemintini 2000). Figure 1 shows the spatial distribution of the globular clusters, satellite galaxies, and FHB stars on the plane perpendicular to the Galactic disk, where  $X$  axis connects the Galactic center ( $X=0$ ) and the sun ( $X=8.0$  kpc). The filled and open symbols denote the objects with and without proper-motion measurements, respectively. Satellite galaxies are the most distant tracers, with Galactocentric distances  $r$  greater than 50 kpc. The globular clusters extend out to almost  $r = 40$  kpc, while the present sample of FHB stars are confined to locations within 10 kpc of the sun. Thus, our sample objects are widely, though not uniformly, distributed throughout the volume of the Galaxy.

Among these sample objects there exist proper-motion measurements for 5 of the satellite galaxies, 41 of the globular clusters, and 211 of the FHB stars. The proper motion data for LMC, Sculptor, and Ursa Minor are taken from WE99, whereas those for Sagittarius and Draco are taken from Irwin et al. (1996) and Scholz & Irwin (1994), respectively. The proper motions for most of the globular clusters have been compiled by Dinescu et al. (1999). We adopt the data from this source, except for two globular clusters with recently revised proper-motion measurements (NGC 6254: Chen et al. 2000; NGC 4147: Wang et al. 2000), and for three additional globular clusters compiled recently (Pal 13: Siegel et al. 2000; Pal 12: Dinescu et al. 2000; NGC 7006: Dinescu et al. 2001). Proper motions for 211 of the FHB stars in the Wilhelm et al. (1999b) sample are available from one or more existing proper-motion catalogs. These include the STARNET Catalog (Röser 1996), the Yale-San Juan Southern Proper Motion Catalog (SPM 2.0: Platais et al. 1998), the Lick Northern Proper Motion Catalog (NPM1: Klemola, Hanson, & Jones 1994), and the TYCHO-2 Catalog (Høg et al. 2000). Many of these FHB stars have been independently measured in two or more catalogs, so that by combining all measurements one can reduce the statistical errors, as well as minimize any small remaining systematic errors in the individual catalogs, as was argued in Martin & Morrison (1998) and Beers et al. (2000).

We estimate average proper motions,  $\langle \mu \rangle$ , and their errors,  $\langle \sigma_\mu \rangle$ , weighted by the inverse variances, as

$$\langle \mu \rangle = \left( \sum_{i=1}^n \mu_i / \sigma_{\mu_i}^2 \right) / \left( \sum_{i=1}^n 1 / \sigma_{\mu_i}^2 \right), \quad (2)$$

$$\langle \sigma_\mu \rangle = \left( \sum_{i=1}^n 1 / \sigma_{\mu_i}^2 \right)^{-1/2}, \quad (3)$$

where  $n$  denotes the number of measurements. Table 1 lists these compilations, as well as the estimated distances to the FHB stars, where  $r$  and RV denote the Galactocentric distances and heliocentric radial velocities, respectively. Typical errors in the reported proper-motion measurements range from  $1 \sim 5 \text{ mas yr}^{-1}$  for individual field stars, whereas those for satellite galaxies and globular clusters are about  $0.3 \text{ mas yr}^{-1}$  and  $1 \text{ mas yr}^{-1}$ , respectively.

We assume a circular speed of  $V_{LSR} = 220 \text{ km s}^{-1}$  at the location of the sun (i.e.  $R_{\odot} = 8.0 \text{ kpc}$  along the disk plane) and a solar motion of  $(U, V, W) = (-9, 12, 7) \text{ km s}^{-1}$  (Mihalas & Binney 1981), where  $U$  is directed outward from the Galactic Center,  $V$  is positive in the direction of Galactic rotation, and  $W$  is positive toward the North Galactic Pole. We then calculate the space motions and their errors, fully taking into account the reported measurement errors in the radial velocities of the individual satellite galaxies (typically a few  $\text{km s}^{-1}$ ), adopting a typical radial-velocity error for other objects ( $10 \text{ km s}^{-1}$ ), the measurement errors assigned to the proper motions of each object (when available, adopting a mean error for the source catalog when not), and distance errors for the satellite galaxies (10 % relative to the measured ones), or as obtained from eq. 1 for the globular clusters and FHB stars.

It is worth noting that the reported proper motions of the FHB stars in our sample may yet contain unknown systematics with respect to their absolute motions in a proper reference frame; this caution applies to the globular clusters and satellite galaxies as well. It is an important goal to make efforts to reduce the systematic, as well as random, errors in the proper motions upon which studies of Galactic structure and kinematic studies are based, using much higher precision astrometric observations than have been obtained to date.

### 3. MASS DETERMINATION BASED ON ESCAPE VELOCITIES

#### 3.1. Methods and Mass Model

If we model the Galaxy as an isolated, stationary mass distribution, and assume that all of our tracer objects are gravitationally bound to it, then the rest-frame velocities of all objects,  $V_{RF}$ , must be less than their escape velocities,  $V_{esc} = \sqrt{2\psi}$ , where  $\psi$  denotes the gravitational potential of the Galaxy. A number of previous researchers have adopted this method for obtaining mass estimates of the Galaxy (e.g., Fricke 1949; Miyamoto et al. 1980; Carney et al. 1988; Leonard & Tremaine 1990; Dauphole & Colin 1995). We first follow this procedure using the sample described in §2.

Here we adopt two different mass models, in order to investigate the difference in estimates of the Galaxy's mass obtained by the use of different potentials. Our models, hereafter referred to as Model A and B, are the same as those adopted in WE99 and Johnston, Spergel, & Hernquist (1995) (and also used by Dinescu et al. 1999), respectively.

Model A has spherical symmetry, and results in a flat rotation curve in the inner regions of the Galaxy. The gravitational potential and mass density are given as

$$\psi(r) = \frac{GM}{a} \log \left( \frac{\sqrt{r^2 + a^2} + a}{r} \right), \quad \rho(r) = \frac{M}{4\pi} \frac{a^2}{r^2 (r^2 + a^2)^{3/2}}, \quad (4)$$

where  $a$  is the scale length of the mass distribution, and  $M$  is the total mass of the system. The central density of this model is cusped (like  $r^{-2}$ ) and falls off as  $r^{-5}$  for  $r \gg a$ . As  $M$  is derived by integrating  $\rho(r)$  from  $r = 0$  to  $\infty$ , this model contains one free parameter,  $a$ .

Model B consists of realistic axisymmetric potentials with three components (the bulge, disk, and dark halo) that reproduce the shape of the Galactic rotation curve (Johnston et al. 1995). The bulge and disk components are represented by Hernquist (1990) and Miyamoto & Nagai (1975) potentials, respectively. All of the parameters included in these potentials are taken from Dinescu et al. (1999) (see their Table 4). In order to obtain a finite total mass, we assume the following modified logarithmic potential (corresponding to an isothermal-like density distribution) for the dark halo component:

$$\psi_{halo}(r) = \begin{cases} v_0^2 \log[1 + (r/d)^2] - \psi_0, & \text{at } r < r_{cut} \\ -2v_0^2 \frac{r_{cut}}{r} \frac{c}{1+c}, & \text{at } r \geq r_{cut}, \end{cases} \quad (5)$$

$$\rho(r) = \frac{2v_0^2}{4\pi G d^2} \frac{3 + r/d}{(1 + r/d)^3}, \quad (6)$$

where  $\psi_0$  is defined as

$$\psi_0 = v_0^2 [\log(1 + c) + 2c/(1 + c)], \quad c = (r_{cut}/d)^2, \quad (7)$$

and we adopt  $v_0 = 128 \text{ km s}^{-1}$  and  $d = 12 \text{ kpc}$  (Dinescu et al. 1999). This model contains one free parameter, namely the cutoff radius of the dark halo,  $r_{cut}$ . Figure 2 shows the rotation curves for  $0 \leq R \leq 20 \text{ kpc}$ , provided by Model A with  $a = 200 \text{ kpc}$  (thick solid line) and Model B with  $r_{cut} = 170 \text{ kpc}$  (thin solid line), where both curves shown at  $R \leq 20 \text{ kpc}$  remain unchanged as long as  $a, r_{cut} \gg 20 \text{ kpc}$ . The circular speed at  $R = R_\odot$  is  $220 \text{ km s}^{-1}$  for both mass models. Also shown is the declining rotation curve with increasing radius, as obtained from Model A with  $a = 20 \text{ kpc}$  (dashed line).

### 3.2. Results

Figures 3a and 3b show the relationship between the derived escape velocities,  $V_{esc}$ , and the rest-frame velocities,  $V_{RF}$ , when we adopt Model A with  $a = 195$  kpc and Model B with  $r_{cut} = 295$  kpc, respectively. For the objects without available proper motions (open symbols), we adopt the radial velocities alone as measures of  $V_{RF}$ , hence their estimated space velocities are only lower limits. The solid line denotes the boundary between the objects that are bound (below the line) and unbound (above the line) to the Galaxy, respectively. By selecting the smallest scale length,  $a$ , that places the sample objects inside the bound region it is possible to set lower limits on the total mass of the Galaxy.

It is worth noting that this mass determination (the enclosed mass) is basically provided by the high-velocity objects located near the boundary line at each respective radius (or corresponding  $\psi$ ). For determination of the total mass, these include Leo I (for which only radial velocity information is available), Draco, Pal 3, and four FHB stars (shown inside the rectangular region). Table 2 summarizes the basic observational data for these particular objects, where columns (6) and (7) list the Galactocentric distances and heliocentric radial velocities, respectively, and the other columns are obvious. Inspection of Figure 3 highlights the following important properties of the mass determination: (1) If the proper motions of all objects are unavailable, then the mass estimate sensitively depends on the presence or absence of Leo I, as has been noted in previous studies. (2) Compared to case (1), if the available proper motions of the satellite galaxies and globular clusters are taken into account, the constraint provided by Draco and Pal 3 is basically the same as that provided by Leo I. This may explain the result of WE99, which showed that the mass determination is made insensitive to Leo I if the proper motion data of satellite galaxies and globular clusters are taken into account. However, as Figure 3 indicates, the velocity errors for Draco and Pal 3 are quite large, so these objects place only weak constraints on the mass estimate. (3) If we consider the proper motions of FHB stars, then the four FHB stars having high velocities (one of which exhibits a rather small velocity error) provide the basically the same constraint on the Galaxy's mass as Leo I, Draco, and Pal 3. These properties suggest that the inclusion of FHB stars with available proper motions is crucial, and provides constraints on the mass limit of the Galaxy that depend on *neither* the inclusion or absence of Leo I *nor* on the large velocity errors for Draco and Pal 3.

We compute the boundary line provided by Pal 3, Draco, and the four FHB stars inside the rectangular region in Figure 3, based on a weighted least-squares fitting procedure (weights being inversely proportional to the velocity errors). This exercise yields  $a = 195^{+160}_{-85}$  kpc for Model A, and  $r_{cut} = 295^{+335}_{-145}$  kpc for Model B. Using these values, we arrive at the most likely lower limits to the total mass  $M$  of the Galaxy as  $2.2^{+1.8}_{-1.0} \times 10^{12} M_{\odot}$  for Model



A and  $2.2^{+2.6}_{-1.1} \times 10^{12} M_{\odot}$  for Model B, respectively. Thus, the difference in the derived mass limits is not significant, as long as the rotation curve at outer radii is approximately constant at the adopted value of  $220 \text{ km s}^{-1}$ . It also suggests that the flattened nature of the Model B potential, due to the presence of the disk component, does not affect the results significantly – the high-velocity tracers are located at large Galactocentric distances and/or their orbits largely deviate from the disk plane.

We note that this method of mass determination, based on escape velocities, inevitably depends on the selection of a few *apparently* high-velocity objects from a much larger sample of tracers. However, we also point out that the lower mass limit obtained here is also influenced by the inclusion of additional FHB stars with  $V_{RF} \sim 500 \text{ km s}^{-1}$ , or the consideration of Draco alone, which possesses the the highest  $V_{RF}$  relative to  $V_{esc}$ . An anonymous referee echoed a concern of ours, that mass estimates obtained from tracers that exhibit extreme properties, such as high *inferred* space motions, may simply be reflecting the tail of an error distribution in the observables, e.g., the proper motions, possibly amplified (particularly in the case of the satellite galaxies) by systematic errors in distance estimates. Our principal goal, at present, is not to obtain the exact value of the lower mass limit, but to highlight the significance of considering FHB stars, which set basically the same mass limit as can be obtained from Leo I, Draco, and Pal3. The great advantage of the FHB stars is that *their* number can be expanded quickly in future studies, while the number of satellite galaxies will forever remain small.

In addition to the above experiments, we also considered a mass model that yields a declining rotation curve at outer radii, as was proposed by Honma & Sofue (1997) from their H I observations. We adopt Model A with  $a = 20 \text{ kpc}$ , which gives rise to  $V_{LSR} = 211 \text{ km s}^{-1}$  at  $R = R_{\odot}$ . The corresponding rotation curve, being reminiscent of the result in Honma & Sofue (1997), is shown as the dashed line in Figure 2. Figure 4 shows the  $V_{RF}$  vs.  $V_{esc}$  relationship that follows from adoption of this model. As is evident, the total mass obtained from a model that leads to a declining rotation curve is quite insufficient to bind many of our sample objects to the Galaxy.

## 4. MASS DETERMINATION BASED ON A BAYESIAN LIKELIHOOD METHOD

### 4.1. Method

As a second method for mass estimation of the Galaxy, we examine an alternative that takes into account *all of the positional and kinematic information* of the sample objects,

in contrast to the use of the high-velocity tracers alone, as in the previous section. In this approach, a phase-space distribution function of tracers,  $F$ , is prescribed for a specifically chosen  $\psi$ , and the model parameters included in  $F$  and  $\psi$  are derived so as to reproduce the presently observed positions and velocities of the tracers in the (statistically) most significant manner. The optimal deduced parameters relevant to  $\psi$  then allow us to estimate the total mass of the Galaxy. This method was originally proposed by Little & Tremaine (1987), and further developed by Kochanek (1996) and WE99.

Based on the results presented in the previous section, we take Model A with spherical symmetry as the mass distribution of the Galaxy, which is sufficient for the following analysis. For the sake of simplicity, and also for ease of comparison with the previous studies by Kochanek (1996) and WE99, the phase-space distribution function is taken to have the same anisotropic form as that adopted in these studies. That is, it depends on the binding energy per unit mass,  $\varepsilon$  ( $\equiv \psi - v^2/2$ ), and the angular momentum per unit mass,  $l$ , in the following way,

$$F(\varepsilon, l) = l^{-2\beta} f(\varepsilon) , \quad (8)$$

where

$$\begin{aligned} f(\varepsilon) &= \frac{2^{\beta-3/2}}{\pi^{3/2}\Gamma[m-1/2+\beta]\Gamma[1-\beta]} \frac{d}{d\varepsilon} \\ &\times \int_0^\varepsilon d\psi \frac{d^m r^{2\beta} \rho_s}{d\psi^m} (\varepsilon - \psi)^{\beta-3/2+m} , \end{aligned} \quad (9)$$

where  $\rho_s$  is the tracer density distribution,  $\Gamma$  is the gamma function, and  $m$  is an integer whose value is chosen such that the integral in eq. (9) converges (e.g., Dejonghe 1986; Kochanek 1996). In the spherical model, this form of the distribution function yields equal velocity dispersions in the orthogonal angular directions,  $\langle v_\theta^2 \rangle = \langle v_\phi^2 \rangle$ , and a constant anisotropy  $\beta = 1 - \langle v_\theta^2 \rangle / \langle v_r^2 \rangle$  everywhere in the Galaxy. Our choice of  $m = 2$  in eq. (9) (to be in accord with the WE99 work) limits the allowed range for the velocity anisotropy to  $-1.5 \leq \beta \leq 1$  when proper motion data are considered, while the use of radial velocities alone sets no limit for tangential anisotropy  $[-\infty, 1]$ .

For  $\rho_s$ , we consider WE99's two models: (a) Shadow tracers following the mass density distribution obtained from Model A (eq. 4), and (b) a power-law distribution as a function of  $r$ . The shadow-tracer model is given as

$$\rho_s(r) \propto \frac{a_s^2}{r^2(r^2 + a_s^2)^{3/2}} , \quad (10)$$

where  $a_s$  is the scale length. The power-law model with index  $\gamma$  is given as

$$\rho_s(r) \propto \frac{1}{r^\gamma} . \quad (11)$$

Here, since shadow tracers may be truncated at the distance below the scale length of the mass distribution, the scale length of the tracers,  $a_s$ , is generally different from the scale length of the Galaxy's mass,  $a$ .

Using the 27 objects (satellite galaxies and globular clusters) at  $r > 20$  kpc, WE99 derived  $a_s = 100$  kpc and  $\gamma = 3.4$  as the best fitting parameters for their spatial distribution. We re-examine  $a_s$  and  $\gamma$  using our sample of all satellite galaxies and globular clusters. Note that the FHB stars are excluded in this determination of  $a_s$  and  $\gamma$ , as they have not (yet) been completely surveyed over the Galactic volume. We obtain  $a_s = 10$  kpc and  $\gamma = 3.3$  as the best fitting values, based on a simple K-S test of the observed vs. predicted distribution functions (see Figure 5). If we exclude the globular clusters at  $r \leq 10$  kpc, for which the spherical symmetry assumption may be questionable due to the presence of the disk globular clusters, we obtain  $a_s = 50$  kpc and  $\gamma = 3.4$ . Thus,  $a_s$  depends sensitively on the selected range of radius (or in other words the selection of the sample), whereas  $\gamma$  basically remains unchanged. Therefore, we focus our attention on the results using the power-law representation for the tracer population, but the shadow-tracer population is also examined for the purpose of comparison with WE99. To see the dependence of the mass estimate on these parameters, we obtain estimates for two values of  $\gamma$  (3.4 and 4.0) and  $a_s$  (100 kpc and the scale length of the mass distribution,  $a$ ), respectively. We note that the FHB stars are also expected to follow a power-law form with  $\gamma \simeq 3.4$ , as inferred from other halo field stars (e.g., Preston, Shectman, & Beers 1991; Chiba & Beers 2001).

We calculate the likelihood of a particular set of model parameters (the scale length of the mass distribution,  $a$ , and the anisotropy parameter,  $\beta$ ) given the positions,  $r_i$ , and radial velocities,  $v_{r,i}$ , or space velocities,  $v_i$ , using Bayes' theorem. The probability that the model parameters take the values  $a$  and  $\beta$ , given the data  $(r_i, v_{(r)i})$  and prior information  $I$ , is

$$P(a, \beta | r_i, v_{(r)i}, I) = \frac{1}{N} P(a) P(\beta) \prod_{i=1}^N P(r_i, v_{(r)i} | a, \beta) , \quad (12)$$

where  $N$  is the normalization factor (Kochanek 1996; WE99). The probabilities  $P(a)$  and  $P(\beta)$  denote the prior probability distributions in  $a$  and  $\beta$ , respectively. Here,  $P(r_i, v_{(r)i} | a, \beta)$  corresponds to the probability of finding an object at position  $r_i$  moving with radial velocity  $v_{r,i}$  or space velocity  $v_i$  for a particular set of model parameters  $a$  and

$\beta$ . The complete expressions for  $P(r_i, v_{(r)i}|a, \beta)$  are shown in Table 1 of WE99. To calculate this probability for the objects with full space velocities, we take into account their large errors relative to radial velocities alone (due to the observed proper-motion errors), by multiplying by an error convolution function of the form

$$P(r_i, v_i|a, \beta) = \int \int dv_\alpha dv_\delta E_1(v_\alpha) E_1(v_\delta) P(r_i, v_{i,obs}(v_\alpha, v_\delta)|a, \beta) , \quad (13)$$

where  $(v_\alpha, v_\delta)$  are the tangential velocities along the right ascension and declination coordinates, respectively, and  $E_1$  is the Lorentzian error convolution function, defined as

$$E_1(v) = \frac{1}{\sqrt{2\pi}\sigma_1} \frac{2\sigma_1^2}{2\sigma_1^2 + (v - v_{obs})^2} , \quad (14)$$

where  $\sigma_1$  is defined as  $\sigma_1 = 0.477\sigma$  for the calibrated error estimate,  $\sigma$  (see WE99).

The prior probability in the velocity anisotropy,  $\beta$ , is taken to be of the form  $P(\beta) \propto 1/(3 - 2\beta)^n$ , where  $n = 0$  and  $2$  correspond to a uniform prior and uniform energy prior, respectively (Kochanek 1996; WE99). Larger values of  $n$  give a larger weight towards radial anisotropy. For the prior probability in  $a$ ,  $P(a)$ , we adopt  $1/a$  and  $1/a^2$  (WE99).

Using the routine AMOEBA in Numerical Recipes (Press et al. 1992), we search for a set of model parameters,  $a$  and  $\beta$ , that maximize the probability  $P(a, \beta|r_i, v_{(r)i}, I)$ . The total mass of the Galaxy,  $M$ , is then derived from the parameter  $a$ .

## 4.2. Results

Initially, we apply the Bayesian likelihood method, making use of only the radial velocities of the objects, setting aside for the moment the available proper-motion information. Specifically, we focus on the difference in the mass estimate arising from the presence or absence of Leo I. Figure 6 shows the likelihood contours in the mass-anisotropy  $(M - \beta)$  plane for the case of a power-law tracer population with  $\gamma = 3.4$ , where  $\beta$  is limited to the range of  $-1.5 \leq \beta \leq 1$ . The solid and dashed lines denote the presence and absence of Leo I, respectively. As is evident, the mass estimate sensitively depends on whether or not Leo I is bound to the Galaxy, as has been noted in previous studies. Inclusion of Leo I yields a likely total mass that is *an order of magnitude greater* than the case without Leo I. Over the range of  $\beta$  we consider, the most likely value of  $M$  with Leo I is  $21.0 \times 10^{11} M_\odot$ , corresponding to a scale length  $a = 185$  kpc, whereas excluding Leo I yields  $M = 9.6 \times 10^{11} M_\odot$  and  $a = 85$  kpc. We note that the role of Leo I in the Galaxy's

mass estimate is also understandable from the escape-velocity argument – if only the sample radial velocities are taken into account, Leo I alone determines the best-fit boundary line  $V_{RF} = V_{esc}$  in the  $V_{RF}$  vs.  $V_{esc}$  diagram (Figure 3).

As is seen in Figure 6, the high-probability region is biased toward the line  $\beta = -1.5$ . This bias arises from the specific form of the phase-space distribution function  $F(\varepsilon, l)$  given in equation (8), where the probability  $P(a, \beta | r_i, v_{(r)i}, I)$  is high at large  $F$ . We plot  $F$  in Figure 7 for a set of  $r$  and  $\beta$  (solid and dotted lines for  $\beta = -1$  and  $1$ , respectively). It follows that  $F$  at high  $\varepsilon$  is larger for smaller  $\beta$ , whereas  $F$  at low  $\varepsilon$  is larger for larger  $\beta$ . The range of  $\varepsilon$  corresponding to these two different cases depends on  $r$ , as can be deduced from the comparison between panel (a) and (b) in Figure 7. Since our sample objects are mainly distributed in the region of higher  $\varepsilon$  (solid histograms for the sample with radial velocities), the probability is highest at smallest  $\beta$ .

Following the above experiments, we drop the lower bound of  $-1.5$  for  $\beta$ , and search for the maximum probability at smaller  $\beta$ . No maximum is found up to  $\beta = -20$ , although the large discrepancy in  $M$  between the cases with and without Leo I remains. When we confine ourselves to the sample at  $r > 10$  kpc, there exists a maximum probability at  $\beta = -2.75$  (with Leo I), with a corresponding mass  $32.0 \times 10^{11} M_{\odot}$ . For the sample at  $r > 20$  kpc, we obtain  $11.4 \times 10^{11} M_{\odot}$  at  $\beta = 0.8$ . This clearly suggests that the best-fitting  $\beta$ , obtained from the analysis when only radial velocities are considered, is rather sensitive to the range of  $r$  for the sample selection. This in turn affects the number distribution  $N(\varepsilon)$ , which is relevant to the likely range of  $F$  (Figure 7).

With these unavoidable limitations of the present sample in mind, Table 3 summarizes the likelihood results for the limited range of  $-1.5 \leq \beta \leq 1$ , obtained for power-law and shadow tracers using a variety of different priors on  $a$  and  $\beta$ . The most likely value of  $\beta$  is  $-1.5$  for all cases, for the reason described above. We note that the current mass estimate is rather *insensitive* to the  $\beta$  prior. As the  $\beta$  prior decreases, the estimated mass generally increases, and the best-fitting  $\beta$  decreases, because the small  $\beta$  prior is biased toward more tangentially anisotropic velocity distributions than the large  $\beta$  prior. However, since most of our sample have high  $\varepsilon$ , the best-fitting  $\beta$  remains  $-1.5$  *regardless* of whether we adopt the uniform prior or the uniform-energy prior for  $\beta$ . This property makes the mass estimate insensitive to the  $\beta$  prior.

Now we apply the Bayesian likelihood method to the subsample of objects with both radial velocities and proper motions available, and consider the derived space motions. In contrast to the above case, where we used radial velocities alone, we find that the maximum probability within the range of  $\beta$  we consider is now bounded (Figure 8a). This may be caused by the characteristic distribution of  $\varepsilon$  for the sample with full space

motions, as shown in Figure 7 (dotted histogram). This figure shows that there exists a larger fraction of low  $\varepsilon$  stars than are found in the sample with radial velocities alone (solid histograms), so a larger  $\beta$  is preferred to achieve a larger  $F$ . The mass estimate in this case is quite insensitive to the presence or absence of Leo I. Figure 8b shows the probabilities as a function of  $M$ , with a fixed value of  $\beta = -1.25$ , for the case of a power-law tracer population with  $\gamma = 3.4$ . Solid and dashed lines denote the probabilities with and without Leo I, respectively. As is evident, the agreement between both probabilities is significantly improved compared to the case of the radial velocities alone (Figure 6b). When Leo I is included, the most likely value of the total mass  $M$  and the scale length  $a$  are  $25.0 \times 10^{11} M_{\odot}$  and 225 kpc, respectively. Excluding Leo I yields  $M = 18.0 \times 10^{11} M_{\odot}$  and  $a = 160$  kpc. Table 4 summarizes the various results obtained when the proper motions of the objects are considered. This Table illustrates that, for all cases, the mass of the Galaxy with Leo I is in good agreement with that obtained without Leo I. Also, the mass estimate depends only weakly on the index  $\gamma$ , unknown prior probabilities for  $a$  and  $\beta$ , as well as on the range of  $r$  for the sample selection, resulting in small changes in the mass estimates over a range of only a few times  $10^{11} M_{\odot}$ .

To estimate the typical errors in this mass determination that are associated with the measurement errors of the 561 tracers we have analyzed, we have conducted Monte Carlo simulations, adopting the assumptions that typical errors in the distances and radial velocities are 10 %, and 10 km s<sup>-1</sup>, respectively, and that the proper-motion errors are 1 mas yr<sup>-1</sup> for globular clusters, 0.3 mas yr<sup>-1</sup> for satellite galaxies, and 5 mas yr<sup>-1</sup> for the FHB stars. We generated 561 data points (including Leo I) drawn from Gaussian distribution functions centered on the observational data, and with dispersions set to the above typical errors. Given a true mass  $M$ , or scale length  $a$  (where we use  $M = 2.3 \times 10^{12} M_{\odot}$  with  $a = 200$  kpc), and prior probabilities for  $a$  and  $\beta$  ( $1/a^2$  and the uniform-energy prior, respectively), we calculate the most likely mass,  $M'$ , and compare it with an input true mass. Figure 9 shows the distribution of the discrepancy between  $M'$  and  $M$ ,  $100 \times (M' - M)/M$ , obtained from 1000 realizations. The error distribution in the current mass estimate has a mean value shifted downward by 20 %, and a dispersion of half-width 20 %. These values suggest that one might adopt an estimate of the systematic error on the order of 20 %, and a random error of  $\pm 20$  %. Exclusion of Leo I does not influence the magnitude of these errors. It is worth noting that WE99 obtained roughly  $\sim 100$  % systematic errors, and  $\sim 90$  % random errors in their mass estimate, which was based on about 30 data points. The significant improvement of our mass estimate is mainly due to our consideration of a much larger data set that includes several hundred FHB stars.

As shown in Table 4, the most likely estimated total mass depends on model assumptions at a level of a few times  $10^{11} M_{\odot}$ . When the model is fixed, the current large

data set allows us to limit both systematic and random errors to a level of about 20 %. If we follow WE99's procedure for the adoption of the most likely total mass, i.e., if we adopt the mass estimate that provides the smallest difference between the masses obtained with Leo I and without Leo I, we obtain  $2.5_{-1.0}^{+0.5} \times 10^{12} M_{\odot}$  (Leo I included) and  $1.8_{-0.7}^{+0.4} \times 10^{12} M_{\odot}$  (Leo I excluded). On the other hand, the mass estimate within the distance of the LMC (50 kpc) is quite robust, covering the narrow range  $5.4$  to  $5.5 \times 10^{11} M_{\odot}$ .

## 5. DISCUSSION AND CONCLUDING REMARKS

We have placed new limits on the mass of the Galaxy, based on a newly assembled set of halo objects with the latest available proper-motion data, using two alternative methods for mass determination. The first method, based on the escape velocity argument, enables us to obtain a lower limit on the total mass of the Galaxy of  $1.3$  to  $1.4 \times 10^{12} M_{\odot}$ . We have shown that this mass estimate depends on *neither* the presence or absence of Leo I, *nor* on the large velocity errors for Draco and Pal 3. The second method, based on a Bayesian likelihood approach that reproduces all of the positions and velocities of the sample, also provides a mass estimate that is insensitive to the presence or absence of Leo I, at least when proper motions are taken into account. Although the best mass estimate obtained from this second approach depends somewhat on model assumptions (prior probabilities for  $a$  and  $\beta$  and possibly the shape of  $F$ , see below), the resultant systematic change of the total mass is confined within a few times  $10^{11} M_{\odot}$ . The most likely total mass of the Galaxy we derive is  $2.5_{-1.0}^{+0.5} \times 10^{12} M_{\odot}$ . This is in good agreement with the total mass obtained by WE99 ( $1.9_{-1.7}^{+3.6} \times 10^{12} M_{\odot}$ ) and that obtained from other methods (e.g., Peebles 1995,  $2 \times 10^{12} M_{\odot}$ ). Since the size of our tracer sample is significantly larger than used in previous studies, both systematic and random errors are reduced to a great extent. We note that consideration of the numerous FHB stars plays a vital role in this mass estimate, as demonstrated in § 3.

It is also worth noting that, if we fix the mass of the Galaxy equal to our most likely mass estimate, there is insufficient matter present to gravitationally bind the LMC, if we adopt the recent proper-motion measurement by Anguita et al. (2000). These authors reported rather high proper motions,  $(\mu_{\alpha} \cos \delta, \mu_{\delta}) = (+1.7 \pm 0.2, +2.9 \pm 0.2)$ , compared to previous measurements,  $(\mu_{\alpha} \cos \delta, \mu_{\delta}) = (+1.94 \pm 0.29, -0.14 \pm 0.36)$  (Kroupa & Bastian 1997). Thus their results need confirmation from other studies.

The current work also implies that the Galactic rotation curve at outer radii,  $R > R_{\odot}$ , does not decline out to at least  $R \sim 20$  kpc (as long as local disturbances to circular motions, such as warping motions and/or non-axisymmetric motions, are ignored). As illustrated in

Figure 2, a declining rotation curve corresponding to  $a = 20$  kpc and  $V_{LSR} = 211 \text{ km s}^{-1}$  fails to bind *many* sample objects to the Galaxy. The smallest possible value for  $a$  to bind all objects in the isothermal-like density distribution (eq. 4) is  $a = 195$  kpc, yielding  $V_{LSR} \simeq 220 \text{ km s}^{-1}$ .

In a more general context, the detailed shape of the rotation curve at and beyond  $R = R_\odot$  reflects the interplay between the disk and halo mass distributions, as this region is located near the boundary of both components. Thus, determining the rotation curve at  $R_\odot \lesssim R \lesssim 15$  kpc will set useful limits on the mass distribution in the inner parts of the Galaxy. Indeed, the Japanese project VLBI Exploration of Radio Astrometry (*VERA*) will be able to determine both inner and outer rotation curves from measurement of trigonometric parallaxes and proper motions of astronomical maser sources that are widely distributed in the Galactic disk (Sasao 1996; Honma, Kawaguchi, & Sasao 2000). *VERA* will reach unprecedented astrometric precision,  $\sim 10 \mu\text{as}$ , and will yield precise determinations of the Galactic constants  $R_\odot$  and  $V_{LSR}$ . We note that whatever results are derived for the rotation curve, the total mass of the Galaxy ought to be larger than  $10^{12} M_\odot$ , in order to bind the more distant stellar objects.

Our estimate for the mass of the Galaxy inside 50 kpc, i.e., within the distance of the LMC, is  $5.5_{-0.3}^{+0.0} \times 10^{11} M_\odot$  (Leo I included) and  $5.3_{-0.4}^{+0.1} \times 10^{11} M_\odot$  (Leo I excluded). The error estimates are calculated from the maximum and minimum values of the total mass. Thus, about 24% of the total mass of the Galaxy resides within  $r \leq 50$  kpc. This implies that the possibility of brown dwarfs as the origin of the microlensing events toward the LMC may be excluded, because it requires a much smaller mass inside 50 kpc,  $\sim 1.3 \times 10^{11} M_\odot$  (Honma & Kan-ya 1998). Our result is also in good agreement with the recent statistics of the microlensing events obtained from analysis of the 5.7-year baseline of photometry for 11.9 million stars in the LMC (Alcock et al. 2000), showing the absence of short-duration lensing events by brown dwarfs. However, the most recent work has suggested that perhaps one of the microlensing events is actually caused by a nearby low-mass star in the Galactic disk (Alcock et al. 2001). More direct observations for identifying lensing objects are required to settle this issue.

Once the total mass of the Galaxy is fixed, it is possible to place a useful constraint on the mass of the Local Group. Most of the mass in the Local Group is concentrated in M31 and the Galaxy. The total mass of M31 can be estimated from the positions and radial velocities of its satellite galaxies, globular clusters, and planetary nebulae (Evans & Wilkinson 2000; Côté et al. 2000; Evans et al. 2000). If we take it to be  $1.2_{-0.6}^{+1.8} \times 10^{12} M_\odot$  (Evans & Wilkinson 2000), the mass of the Local Group is  $\sim 3.7 \times 10^{12} M_\odot$ . This is in good agreement with the estimate by Schmoldt & Saha (1998),  $(4 - 8) \times 10^{12} M_\odot$ , based on



modified variational principles.

To set tighter limits on the total mass of the Galaxy we require more accurate proper-motion measurements for a greater number of objects at large Galactocentric distances. The high-velocity FHB stars in our sample (with apparent magnitudes  $V < 16$ ) that are responsible for setting the minimum mass of the Galaxy have proper-motion errors of  $\sim 5 \text{ mas yr}^{-1}$ , whereas Draco and Pal 3 have much larger *relative errors*, comparable to their proper motions themselves (see Table 2). Indeed, both the Space Interferometry Mission (*SIM*: Unwin, Boden & Shao 1997) and the Global Astrometry Interferometer for Astrophysics (*GAIA*: Lindegren & Perryman 1996) will be able to provide more accurate proper motions for such high-velocity objects, as well as for numerous other distant tracers of the Galaxy’s mass, up to a precision of a few  $\mu\text{as}$  for targets with  $V \leq 15$ . This corresponds to an error of  $\lesssim 10 \text{ km s}^{-1}$  in the tangential velocity components for many distant objects, i.e., comparable to the error of their (presently determined) radial velocities. Furthermore, roughly half of our sample objects lack proper-motion measurements altogether. To a great extent, the lack of proper-motion measurements (at least for southern sources) will be removed with the completion of the recently re-started Southern Proper Motion survey of van Altena and colleagues, as well as other efforts to substantially increase the numbers of stars with reasonably well-measured proper motions (e.g., UCAC1: Zacharias et al. 2000; UCAC2: Zacharias et al. 2001).

Further assembly of radial velocities for FHB stars, especially those at large  $r$  (beyond distances where accurate ground-based proper motions can be obtained), is also of great importance for a number of reasons. First, as Figure 3 demonstrates, large Galactocentric regions are characterized by small escape velocities. The current sample of FHB stars (because of their locations near the sun) explore distances where the corresponding escape velocities are in the range  $500 \lesssim V_{\text{esc}} \lesssim 600 \text{ km s}^{-1}$ . More remote FHB stars, with distances in the range  $10 \lesssim r \lesssim 50 \text{ kpc}$ , will offer a further constraint on the total mass of the Galaxy by covering the range  $400 \lesssim V_{\text{esc}} \lesssim 500 \text{ km s}^{-1}$ . Secondly, the assembly of samples of more distant FHB stars will enable exploration of the suggested change in velocity anisotropy from the inner to the outer halo (e.g., Sommer-Larsen et al. 1997), and better constrain its dependence on Galactocentric distance.

In exploring the Bayesian approach for mass estimates of the Galaxy, we have adopted a specific form of the phase-space distribution function  $F$  (eq. 8) to facilitate comparison with previous studies. This procedure implicitly assumes that the velocity-anisotropy parameter,  $\beta$ , is constant everywhere in the Galactic volume. However, as noted by Sommer-Larsen et al. (1997), there is an indication that the velocity anisotropy of the halo may be mostly radial at  $R \lesssim 20 \text{ kpc}$  and tangential at  $R \gtrsim 20 \text{ kpc}$ . If so, many of

distant FHB stars, especially those at  $R > 20$  kpc, play a crucial role in the determination of the global distribution of velocity anisotropy. Searches for a more realistic form of the phase-space distribution function, combined with a more elaborate likelihood method, are both worthy pursuits.

Fortunately, prospects are excellent for obtaining a rapid increase in the observational database of FHB stars with the required data. There already exists a substantial body of additional spectroscopy for FHB/A stars observed during the course of the HK survey of Beers and colleagues and the Hamburg/ESO Stellar survey (Christlieb et al. 2001), many of which also have available proper motions, or will soon, from completion of the SPM survey and/or other ground-based efforts. However, as was noted by Wilhelm, Gray, & Beers (1999) (foreshadowed by Norris & Hawkins 1991; Rodgers & Roberts 1993, and references therein; Kinman, Suntzeff, & Kraft 1994; Preston, Beers, & Shectman 1994), a substantial fraction (perhaps as high as 50%) of high-latitude A-type stars are *not* FHB, but rather some (as yet undetermined) mixture of binaries and high-gravity stars (see Preston & Sneden 2000). For some applications, such as estimates of the mass of the Galaxy that rely on space motions of tracers (and in turn on reasonably precise distance estimates of individual objects), confident separation of bona-fide members of the FHB population from possible “contaminants” is crucial<sup>1</sup>. In the past, this has required that one obtain either Strömgren photometry and/or spectrophotometry (e.g., Kinman et al. 1994), broad-band *UBV* photometry in combination with medium-resolution spectroscopy (e.g., Wilhelm et al. 1999a), or reasonably high S/N, high-resolution spectroscopy (e.g., Preston & Sneden 2000). All such endeavors are rather time intensive. However, Christlieb et al. (2002, in preparation) have been exploring means by which adequate separation of FHB stars from higher-gravity A-type stars might be accomplished *directly* from objective-prism spectra, such as those in the Hamburg/ESO stellar survey. Such methods, which look promising, would be most helpful in future investigations of this sort. Wide-field stellar surveys, such as those presently being carried out with the 6dF facility at the UK Schmidt Telescope, are capable of providing large numbers of radial velocities for FHB/A candidates, and are expected to contribute 5,000-10,000 suitable data over the course of the next few years.

We are grateful to B. Fuchs, R. B. Hanson, and I. Platais for assistance with the comparison of the Wilhelm et al. sample with the catalogs of STARNET, NPM1, and SPM 2.0, respectively. We also thank the members of the VERA team for several useful comments on this work. T.C.B. acknowledges partial support from grants AST-00 98508

---

<sup>1</sup> For example, if 10 % of our FHB sample is contaminated by blue metal-poor stars, we obtain a  $2 \sim 3 \times 10^{11} M_{\odot}$  decrease in our total mass estimate, based on Monte Carlo experiments.

and AST-00 98549 awarded by the U.S. National Science Foundation. T.C.B also would like to acknowledge the support and hospitality shown him during a sabbatical visit to the National Astronomical Observatory of Japan, funded in part by an international scholar award from the Japanese Ministry of Education, Culture, Sports, Science, and Technology, during which initial discussions of this work took place.

## REFERENCES

- Alcock, C. et al. 2000, *ApJ*, 542, 281
- Alcock, C. et al. 2001, *Nature*, 414, 617
- Anguita, C., Loyola, P., & Pedreros, M. H. 2000, *AJ*, 120, 845
- Beers, T. C., Chiba, M., Yoshii, Y., Platais, I., Hanson, R. B., Fuchs, B., & Rossi, S. 2000, *AJ*, 119, 2866
- Binney, J., & Tremaine, S. 1987, *Galactic Dynamics*, Princeton Univ. Press, Princeton, NJ, 236
- Carney, B.W., Laird, J.B., & Latham, D.W. 1988, *AJ*, 96, 560
- Carretta, E., Gratton, R.G., & Clementini, G. 2000, *MNRAS*, 316, 721
- Carretta, E., Gratton, R. G., Clementini, G., & Pecci, F. F. 2000, *ApJ*, 533, 215
- Chen, L., Geffert, M., Wang, J. J., Reif, K., & Braun, J. M. 2000, *A&AS*, 145, 223
- Chiba, M., & Beers, T. C. 2001, *ApJ*, 549, 325
- Christlieb, N., Wisotzky, L., Reimers, D., Homeier, D., Koester, D., & Heber, U. 2001, *A&A*, 366, 898
- Côté, P., Mateo, M., Sargent, W. L. W., & Olszewski, W. 2000, *ApJ*, 537, L91
- Dauphole, B., & Colin, J. 1995, *A&A*, 300, 117
- Dejonghe, H. 1986, *Phys. Rep.*, 133, 217
- Dinescu, D. I., Girard T. M., & van Altena, W. F. 1999, *AJ*, 117, 1792
- Dinescu, D. I., Majewski, S. R., Girard, T. M., & Cudworth, K. M. 2000, *AJ*, 120, 1892
- Dinescu, D. I., Majewski, S. R., Girard, T. M., & Cudworth, K. M. 2001, *AJ*, 122, 1916
- Evans, N. W., & Wilkinson, M. I. 2000, *MNRAS*, 316, 929
- Evans, N. W., & Wilkinson, M. I., Guhathakurta, P., Grebel, E. K., & Vogt, S. S. 2000, *ApJ*, 540, L9
- Fich, M., & Tremaine, S. 1991, *ARA&A*, 29, 409

- Fricke, V. W. 1949, *Astr.Nachr.*, 278,49
- Harris, W. E. 1996, *AJ*, 112, 1487
- Held, E. V., Clementini, G., Rizzi, L., & Momany, Y., Saviane, I., & Di Fabrizio, L. 2001, *ApJ*, 562, L39
- Hernquist, L. 1990, *ApJ*, 356, 359
- Hog E., Fabricius C., Makarov V. V., Urban S., Corbin T., Wycoff G., Bastian U., Schwekendiek P., Wicenec A. 2000, *A&A*, 355, L27
- Honma, M., & Sofue, Y. 1997, *PASJ*, 49, 453
- Honma, M., & Kan-ya, Y. 1998, *ApJ*, 503, L139
- Honma, M., Kawaguchi, N., & Sasao 2000, in *Proc. SPIE 4015, Radio Telescope*, ed. H. R. Butcher, 624
- Irwin, M., Ibata, M. J., Gilmore, G., Wyse, R., & Suntzeff, N. 1996, in *ASP Conf. Ser. 92 Formation of the Galactic Halo — Inside and Out*, ed. Morrison, H., & Sarajedini, A. (San Francisco: ASP), 841
- Johnston, K. V., Spergel, D. N., & Hernquist, L. 1995, *ApJ*, 451, 598
- Kerr, F. J., & Lynden-Bell, D. 1986, *MNRAS*, 221, 1023
- Kinman, T.D., Suntzeff, N.B., & Kraft, R.P. 1994, *AJ*, 108, 1722
- Klemola, A. R., Hanson, R. B., Jones, B. F. 1994, *Lick Northern Proper Motion Program: NPM1 Catalog (NSSDC/ADC Cat. A1199)* (Greenbelt, MD: GSFC)
- Kochanek, C. S. 1996, *ApJ*, 457, 228
- Kroupa, P., & Bastian, U. 1997, *NewA*, 2, 77
- Leonard, P. J. T., & Tremaine, S. 1990, *ApJ*, 353, 486
- Lindgren, L., & Perryman, M. A. C. 1996, *A&A*, 116, 579
- Little, B., & Tremaine, S. 1987, *ApJ*, 320, 493
- Martin, J. C., & Morrison, H. L. 1998, *AJ*, 116, 1724
- Mateo, M. 1998, *ARA&A*, 36, 435

- Méndez, R. A., Platais, I., Girard, T. M., Kozhurina-Platais, V., & van Altena, W. F. 1999, *ApJ*, 524, L39
- Miyamoto, M., & Nagai, R. 1975, *PASJ*, 27, 533
- Miyamoto, M., & Zhu, Z. 1998, *AJ*, 115, 1483
- Miyamoto, M., Satoh, C., & Ohashi, M. 1980, *A&A*, 90, 215
- Norris, J.E., & Hawkins, M.R.S. 1991, *ApJ*, 380, 104
- Peebles, P. J. E. 1995, *ApJ*, 449, 52
- Platais, I., Girard, T.M., Kozhurina-Platais, V., van Altena, W.F., López, C.E., Méndez, R.A., Ma, W.-Z., Yang, T.-G., MacGillivray, H.T., & Yentis, D.J. 1998 *AJ*, 116, 2556
- Press, W. H., Teukolsky, S. A., Vetterling, W. T., & Flannery, B. P. 1992, *Numerical Recipes in Fortran 77 : the art of scientific computing*, 2nd ed., Cambridge University Press, Cambridge
- Preston, G.W., & Sneden, C. 2000, *AJ*, 120, 1014
- Preston, G.W., Beers, T.C., & Shectman, S.A. 1994, *AJ*, 108, 538
- Preston, G.W., Shectman, S.A., & Beers, T.C. 1991, *ApJ*, 375, 121
- Rodgers, A.W. & Roberts, W.H. 1993, *AJ*, 106, 1839
- Röser, S. 1996, in *IAU Symp. 172*, ed. S. Ferraz-Mello et al. (Dordrecht: Kluwer), 481
- Sasao, T. 1996, in *Proc. 4th Asia-Pacific Telescope Workshop*, ed. E. A. King (Sidney: Australian Telescope National Facility), 94
- Schmoldt, I. M., & Saha, P. 1998, *AJ*, 115, 2231
- Scholz, R. -D., & Irwin, M. J. 1994, in *IAU Symp. 161, Astronomy from Wide-Field Imaging*, ed. MacGillivray, H. T., Thomson, E. B., Lasker, B. M., Reid, I. N., Malin, D. F., West, R. M., Lorenz, H. (Dordrecht : Kluwer), 535
- Siegel, M. H., Majewski, S. R., Cudworth, K. M., & Takamiya, M. 2001, *AJ*, 121, 935
- Sommer-Larsen, J., Beers, T. C., Flynn, C., Wilhelm, R., & Christensen, P. R. 1997, *ApJ*, 481, 775
- Unwin, S., Boden, A., Shao, M. 1997, *Proc. STAIF, AIP Conf. Proc.*, 387, 63

- Wang, J. J., Chen, L. L., Wu, Z. Y., Gupta, A. C., & Geffert M. 2000, A&AS, 142, 373
- Wilkinson, M. I., & Evans, N. W. 1999, MNRAS, 310, 645 (WE99)
- Wilhelm, R., Beers, T.C., & Gray, R.O. 1999a, AJ, 117, 2308
- Wilhelm, R., Beers, T. C., Sommer-Larsen, J., Pier, J. R., Layden, A. C., Flynn, C., Rossi, S., & Christensen, P. R. 1999b, AJ, 117, 2329
- Zacharias, N., Urban, S.E., Zacharias, M.I., Hall, D.M., Wycoff, G.L., Rafferty, T.J., Germain, M.E., Holdenried, E.R., Pohlman, J.W., Gauss, F.S., Monet, D.G., & Winter, L. 2000, AJ, 120, 2131
- Zacharias, N., Zacharias, M.I., Urban, S.E., & Rafferty, T.J. 2001, BAAS, 199, 129.08
- Zaritsky, D., Olszewski, E. W., Schommer, R. A., Peterson, R. C., & Aaronson, M. 1989, ApJ, 345, 759

Figure 1: Spatial distributions of satellite galaxies (squares), globular clusters (circles), and FHB stars (triangles) on the plane perpendicular to the Galactic disk, where the  $X$  axis connects the Galactic center ( $X=0$ ) and the sun ( $X=8.0$  kpc). The filled and open symbols denote the objects with and without available proper motions, respectively. The plus sign in panel (b) denotes the position of the sun,  $(X, Y) = (8.0, 0)$ .

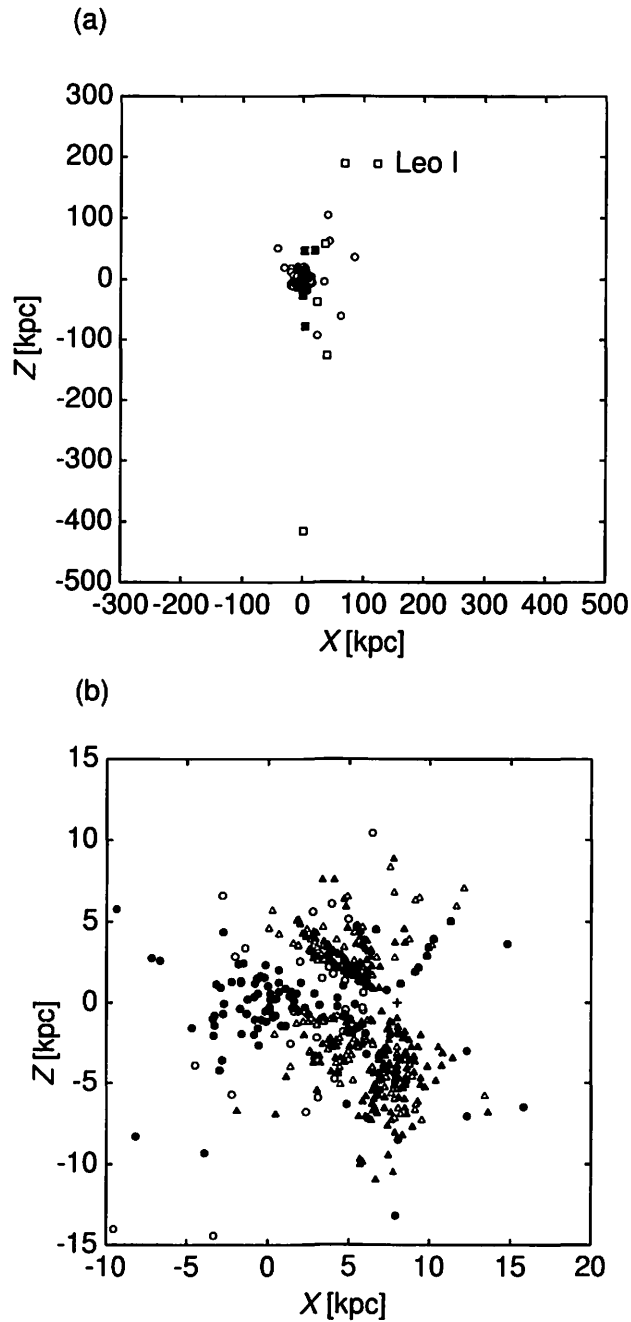




Fig. 2.— Rotation curves for Model A and Model B, parameterizations of the mass distributions considered in this paper. See the text for more information on the nature of these models.

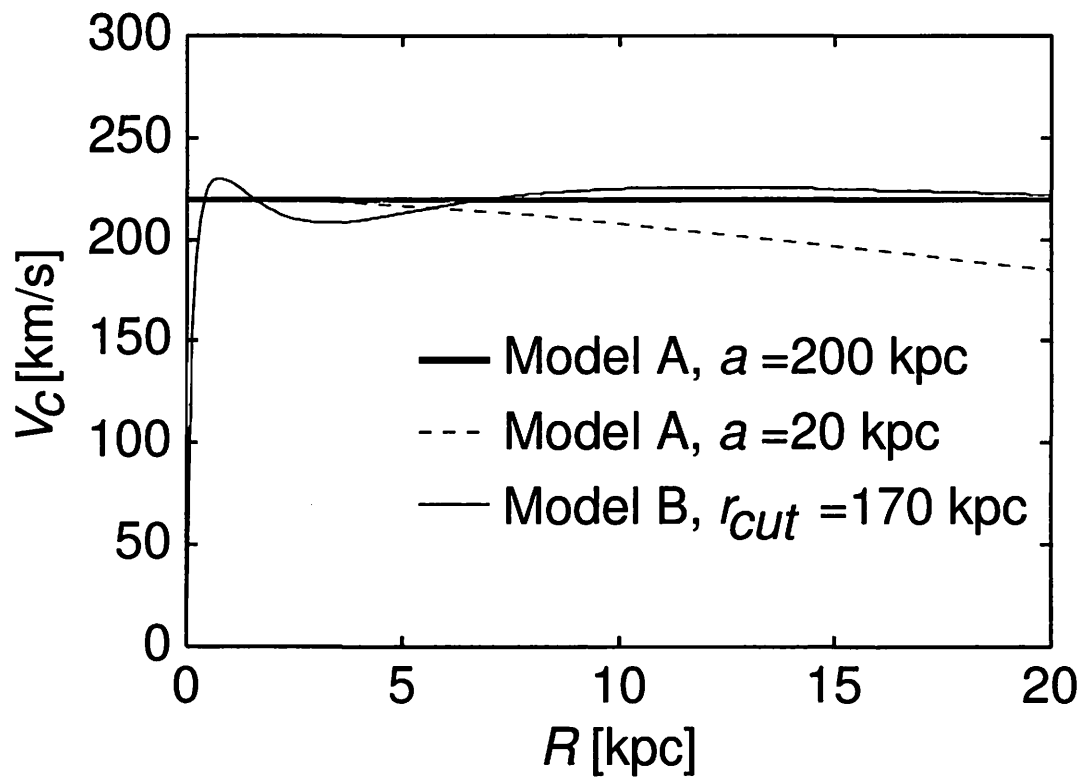


Figure 3: (a) The relation between escape velocities,  $V_{esc}$ , and space velocities,  $V_{RF}$ , for Model A with  $a = 195$  kpc. The symbols are the same as those in Figure 1. The solid line denotes the boundary between the gravitationally bound and unbound objects – those in the region below the line are bound to the Galaxy. For the sake of clarity, velocity errors are plotted for only the high velocity objects relevant to the mass estimate. (b) Same as panel (a) but for Model B with  $r_{cut} = 295$  kpc.

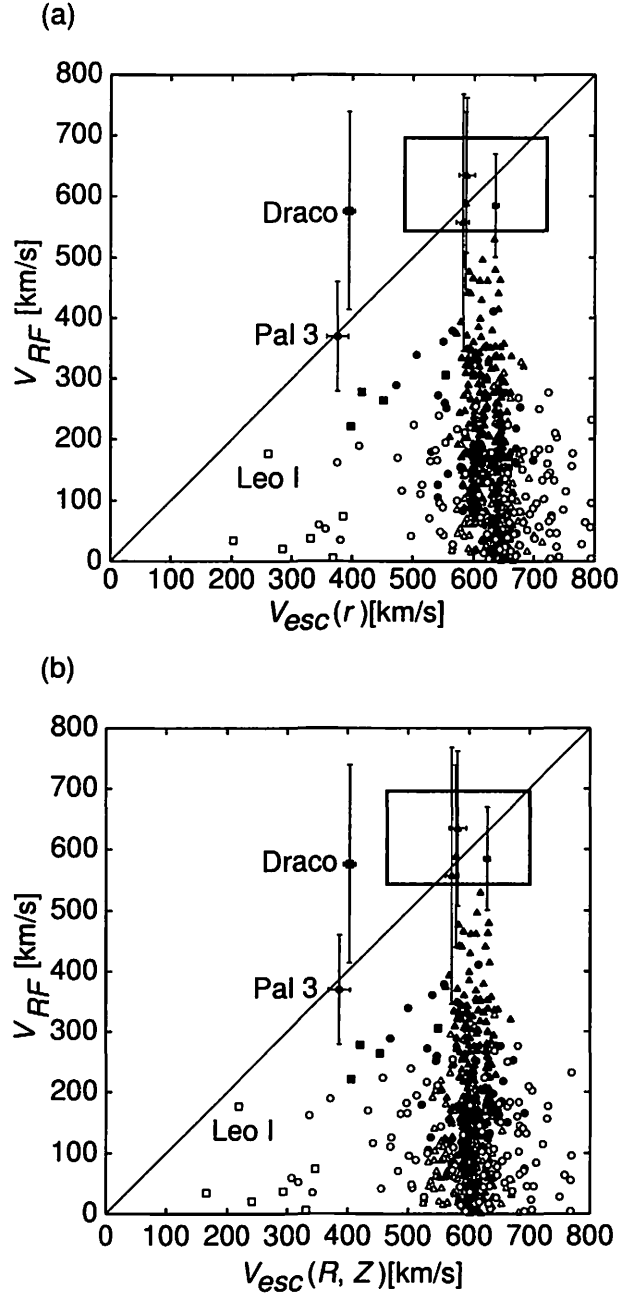


Fig. 4.— The relation between escape velocities,  $V_{esc}$ , and space velocities,  $V_{RF}$ , for Model A with  $a = 20$  kpc. In this case, the rotation curve declines with increasing radii, as shown in Figure 2 (dashed line). Note that, if this situation were to apply, many of the sample objects would be unbound to the Galaxy.

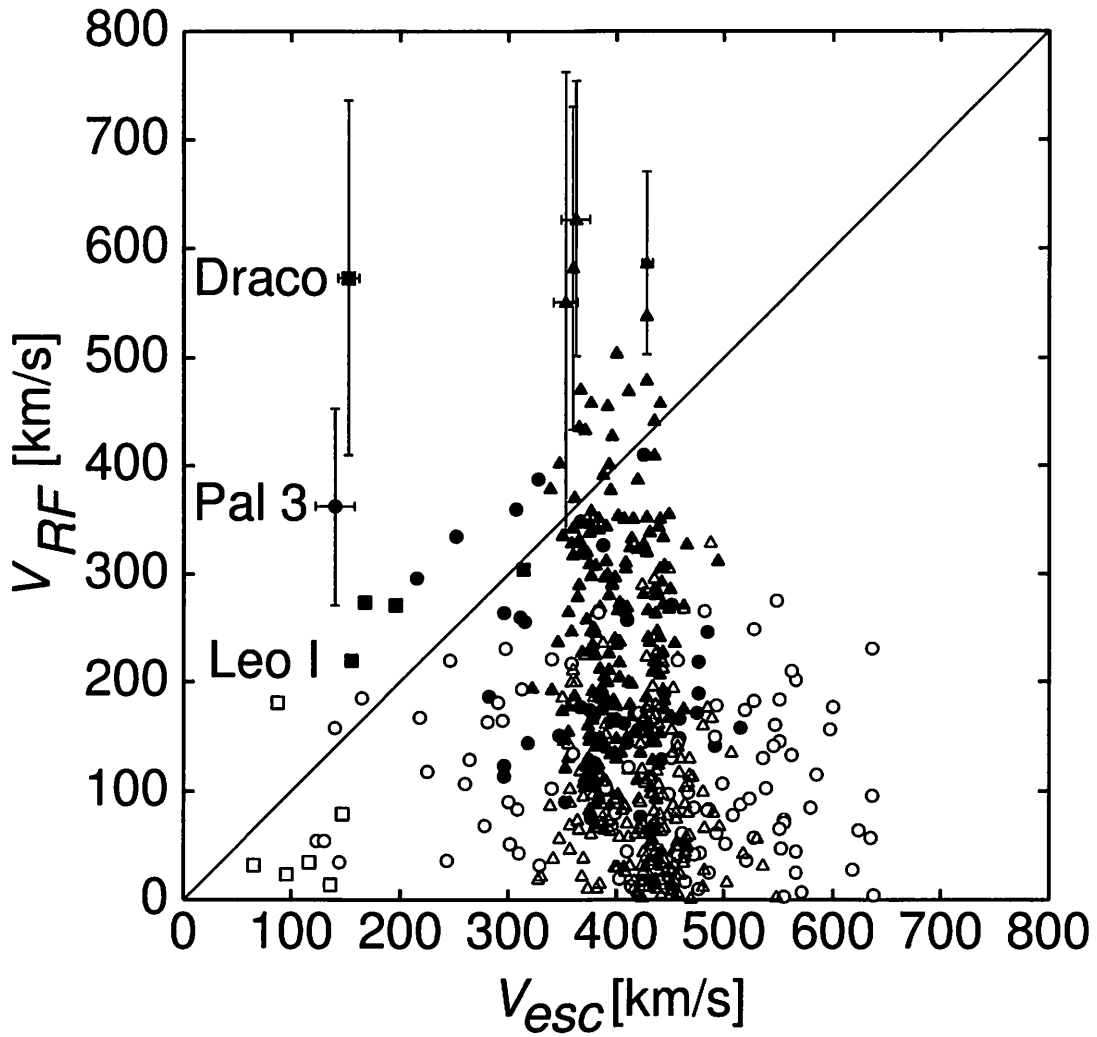


Fig. 5.— Cumulative number distribution,  $N(< r)$ , of the distances of globular clusters and satellite galaxies (solid histogram) in comparison with model distributions (continuous dashed and solid lines). See the text for additional information.

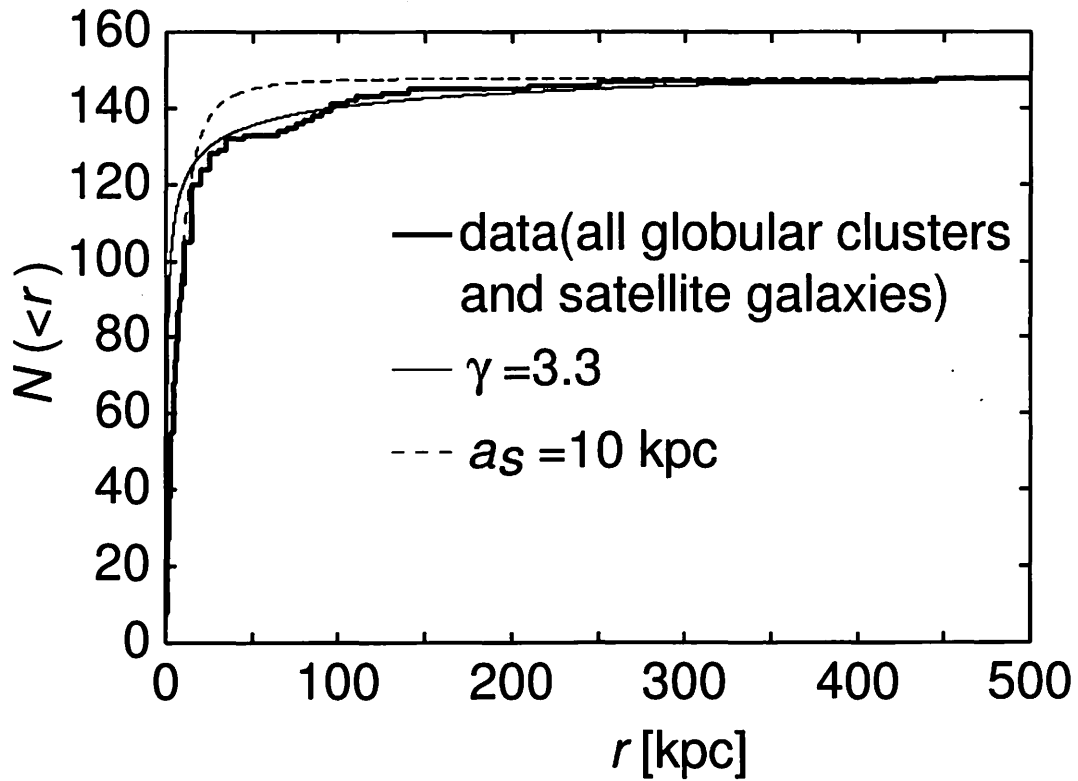


Figure 6: (a) Likelihood contours in the plane of the mass  $M$  and velocity anisotropy  $\beta$ , obtained from an analysis using only radial velocities. The solid and dashed curves show the results including Leo I and excluding Leo I, respectively; the cross and the asterisk show the maxima of the probabilities for each case. Contours are plotted at heights of 0.32, 0.1, 0.045, and 0.01 of the peak height. The spatial distribution of a tracer population is assumed to follow a power-law form with  $\gamma = 3.4$ . (b) Probabilities of the mass  $M$  at  $\beta = -1.5$ , including Leo I (solid line) and excluding Leo I (dashed line).

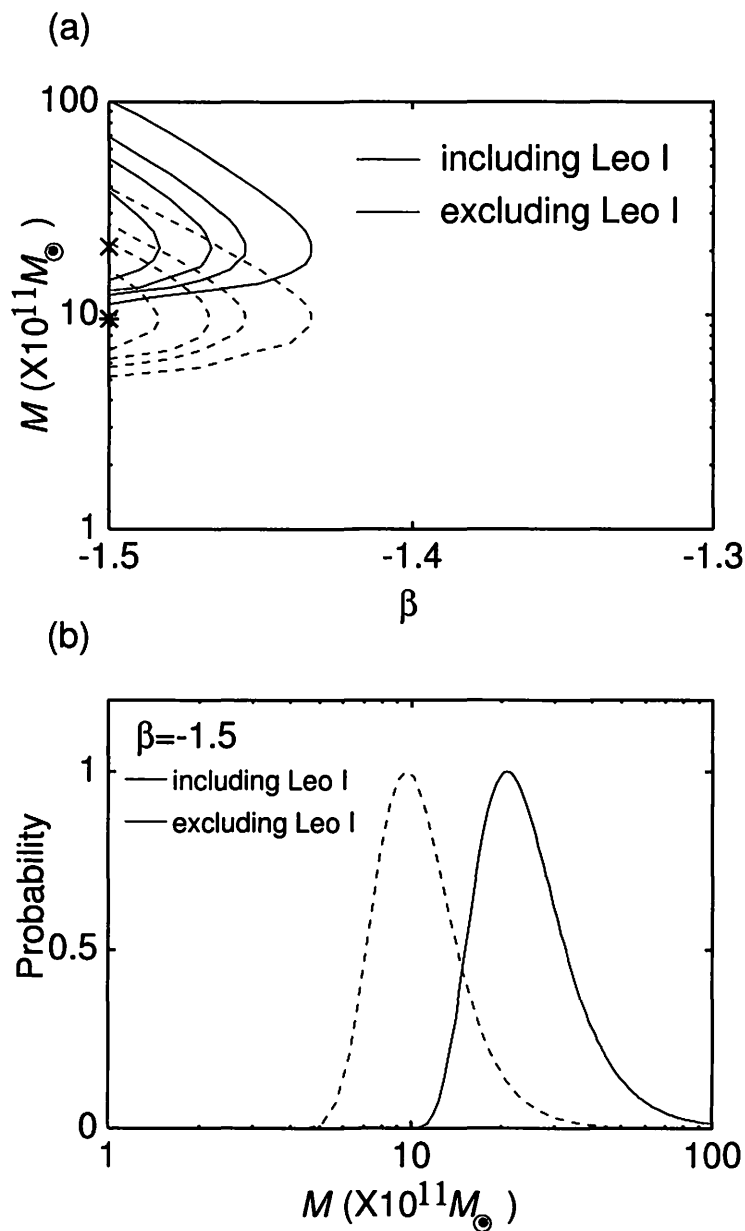


Figure 7: The distribution function,  $F$ , for  $\beta = -1$  (solid lines) and  $\beta = 1$  (dotted lines), at  $r = 10$  kpc (panel a) and  $r = 50$  kpc (panel b). Also plotted are the number distributions  $N(\varepsilon)$  of the stars when  $a = 200$  kpc, where dotted and solid histograms denote the sample with and without available proper motions, respectively. The range of  $r$  for plotting  $N(\varepsilon)$  ( $r < 10$  kpc for panel a and  $20 < r < 80$  kpc for panel (b)) is chosen to approximately match that for  $F$ .

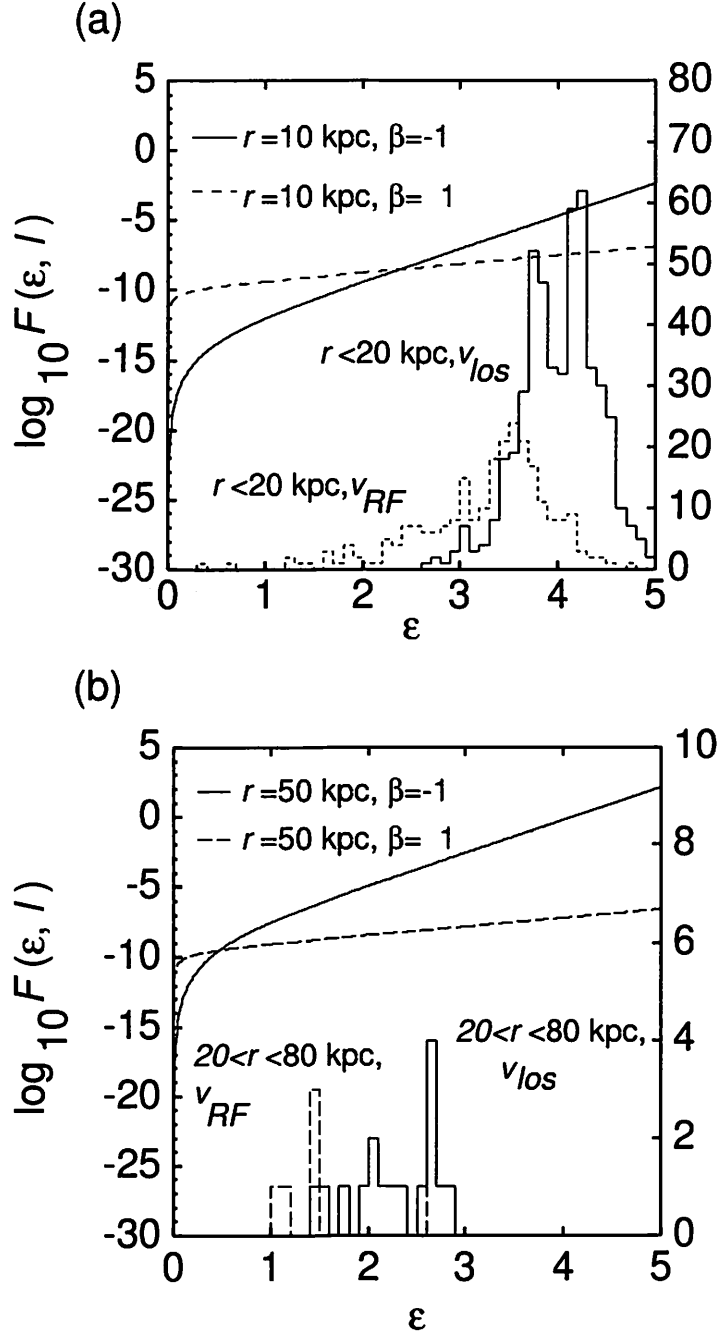


Figure 8: (a) Likelihood contours in the plane of the mass  $M$  and velocity anisotropy  $\beta$ , obtained from an analysis that uses both radial velocities and proper motions. Solid and dashed curves show the results including Leo I and excluding Leo I, respectively; the cross and the asterisk show the maxima of the probabilities for each case. Contours are plotted at heights of 0.32, 0.1, 0.045, and 0.01 of the peak height. The spatial distribution of a tracer population is assumed to follow a power-law form with  $\gamma = 3.4$ . (b) Probabilities of the mass  $M$  at the best-fitting  $\beta$  of  $-1.25$ , including Leo I (solid line) and excluding Leo I (dashed line).

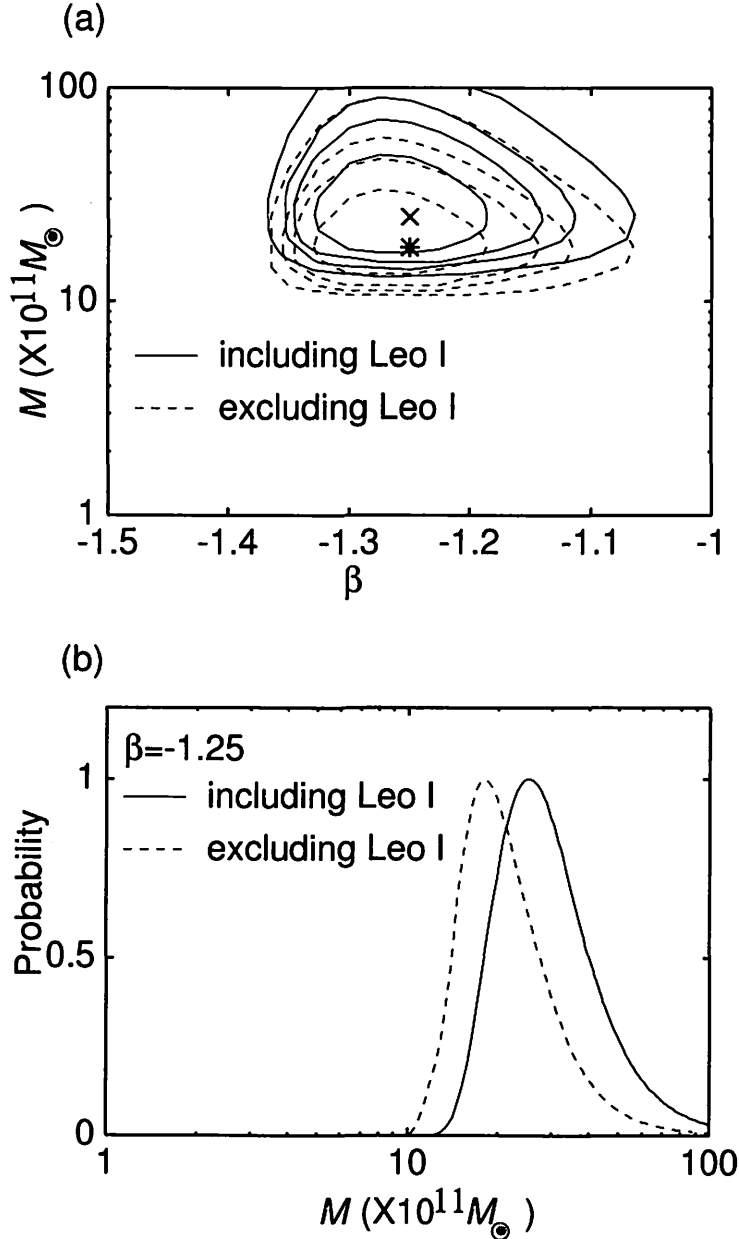


Fig. 9.— An approximate error distribution of the mass estimate caused by the typical measurement errors of the data. The abscissa denotes the relative error in mass,  $100 \times (M' - M)/M$ , where  $M'$  is the mass calculated by a Monte Carlo method and  $M$  is the input true value. See text for more details.

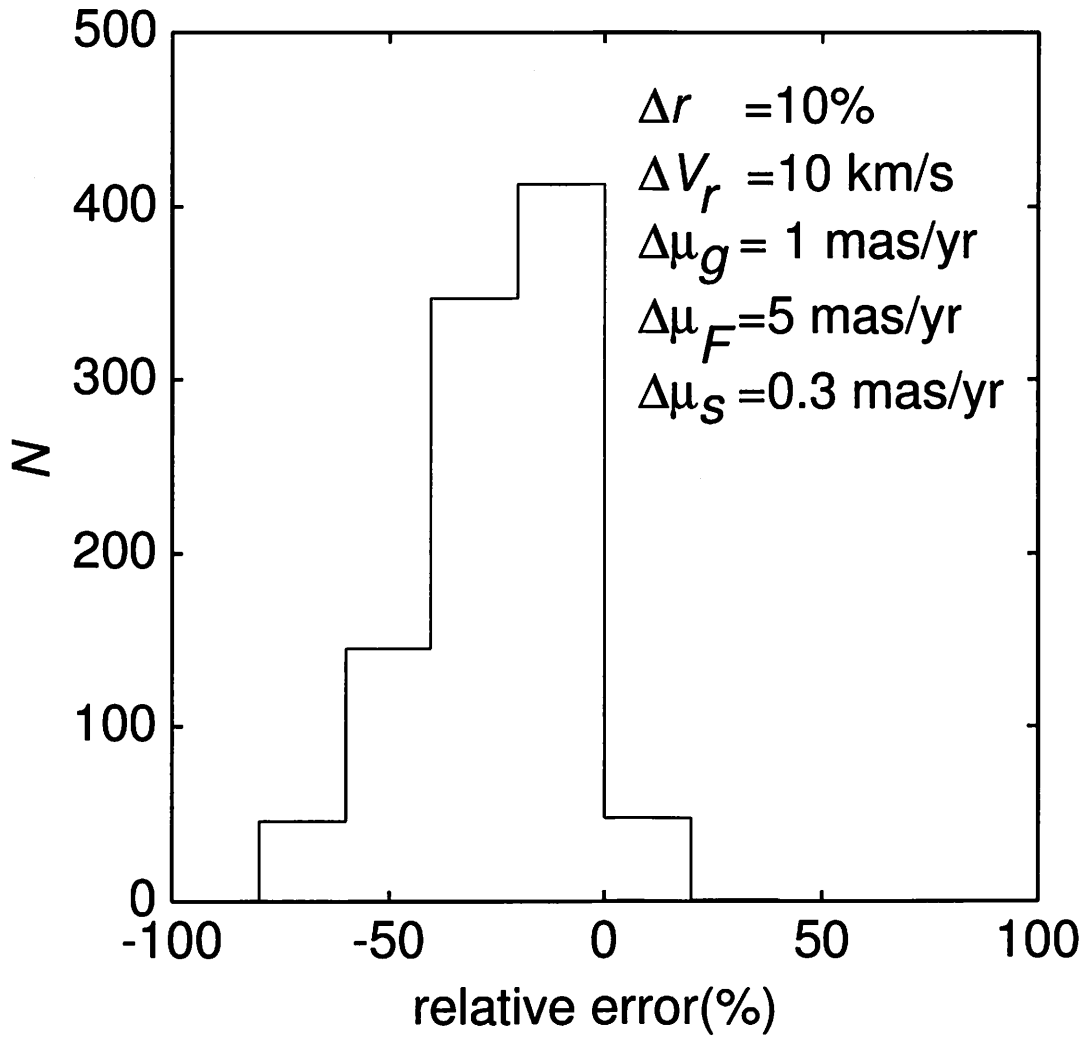




TABLE 1. Distance Estimates and Proper Motions of the FHB stars

NAME <sup>a</sup>	RA (2000.0)	DEC	[Fe/H]	$r^b$ (kpc)	RV <sup>c</sup> (km/s)	$\mu_{\alpha^*}$ (mas/yr)	$\mu_{\delta}$ (mas/yr)	$v$ (km/s)	Source <sup>d</sup>
228760029	0:01:57.6	-36:40:46	-2.8	8.8	48	$-0.8 \pm 3.4$	$-5.6 \pm 2.8$	$148 \pm 90$	S
295170031	0:02:09.8	-14:08:55	-2.7	12.4	-314	$5 \pm 5$	$-10 \pm 5$	$394 \pm 211$	N
295170044	0:03:08.4	-14:24:25	-2.6	10.4	60	$-9 \pm 5$	$-1 \pm 5$	$477 \pm 163$	N
228760030	0:05:36.5	-36:41:28	0	7.9	50	$2.0 \pm 2.4$	$1.7 \pm 2.1$	$244 \pm 28$	A, S
295030008	0:06:02.2	-24:37:09	-2	9.3	45	$1.4 \pm 3.4$	$-3.0 \pm 4.9$	$156 \pm 120$	S
228760034	0:06:20.7	-35:17:14	-1.2	8.2	-94	$3.3 \pm 10.0$	$-6.4 \pm 8.9$	$140 \pm 127$	S
228760031	0:08:25.6	-36:09:15	-2.2	8.7	-54	$-2.2 \pm 2.8$	$-8.7 \pm 2.8$	$193 \pm 62$	S
228760038	0:09:44.7	-34:39:14	-1.9	8.8	-153	$4.4 \pm 2.2$	$-5.4 \pm 2.3$	$177 \pm 31$	S
295030024	0:10:08.5	-25:33:40	-1.4	9.7	143	$13.2 \pm 3.6$	$-3.7 \pm 2.9$	$347 \pm 87$	S
295030029	0:11:19.1	-26:26:39	-2.4	8.6	-49	$12.7 \pm 2.0$	$-16.6 \pm 1.8$	$197 \pm 44$	A, S
295270016	0:27:53.9	-18:57:44	-3	9.5	-43	$11 \pm 5$	$-3 \pm 5$	$193 \pm 125$	N
294970009	0:29:01.7	-23:40:01	-2.1	9.3	-131	$15.5 \pm 3.9$	$-2.4 \pm 5.1$	$308 \pm 100$	S
295270022	0:32:17.1	-19:26:00	-1.9	8.2	-98	$25.0 \pm 2.0$	$-8.1 \pm 2.0$	$168 \pm 18$	A, N, T
295270026	0:33:05.0	-21:11:13	-1.8	9.4	-146	$6 \pm 5$	$-3 \pm 5$	$176 \pm 76$	N
221700002	0:34:37.9	-10:28:50	-2.5	10.0	-270	$-4 \pm 5$	$-10 \pm 5$	$321 \pm 102$	N
295270035	0:35:23.2	-21:00:47	-2.2	8.7	-22	$13.9 \pm 3.8$	$-11.9 \pm 3.8$	$97 \pm 63$	A, N
228820023	0:36:04.3	-30:16:11	-1.8	10.6	-132	$-2.3 \pm 2.7$	$-9.1 \pm 2.8$	$291 \pm 83$	S
295270031	0:36:39.2	-22:25:44	-2.1	9.1	-19	$3.7 \pm 2.2$	$-13.1 \pm 2.9$	$92 \pm 51$	A, N, S
295270039	0:36:53.8	-19:56:51	-1.1	10.0	-26	$2 \pm 5$	$2 \pm 5$	$268 \pm 138$	N
294970031	0:38:41.9	-24:26:56	-2	9.3	-24	$5.8 \pm 2.5$	$-2.0 \pm 2.5$	$148 \pm 57$	S
221790011	0:39:36.1	-04:35:40	-3	9.1	-110	$3.0 \pm 4.4$	$-10.0 \pm 4.4$	$95 \pm 60$	A
221700009	0:40:52.9	-10:19:31	-2.2	9.9	-125	$9 \pm 5$	$-8 \pm 5$	$105 \pm 92$	N
294970033	0:41:08.5	-25:56:31	-1.9	9.2	73	$1.6 \pm 4.0$	$-9.4 \pm 2.7$	$124 \pm 65$	S
221700013	0:41:21.1	-08:23:19	-2	9.1	-7	$16.8 \pm 3.3$	$-30.9 \pm 3.3$	$50 \pm 56$	A, N
221700015	0:43:01.6	-08:15:10	0	9.6	-22	$-11 \pm 5$	$-13 \pm 5$	$362 \pm 100$	N
294970038	0:43:47.1	-26:43:50	-3	13.1	-42	$0.0 \pm 4.8$	$-6.5 \pm 4.1$	$192 \pm 222$	S
221700024	0:43:47.3	-11:22:15	-2.2	12.0	81	$2 \pm 5$	$-6 \pm 5$	$150 \pm 105$	N
295270061	0:43:48.6	-20:45:41	-3	9.5	-67	$7.0 \pm 4.1$	$-3.9 \pm 4.1$	$118 \pm 81$	A, N
221830011	0:52:56.1	-02:52:42	-1.7	8.9	-126	$13.9 \pm 2.1$	$-31.5 \pm 2.1$	$187 \pm 25$	A, T
295090031	0:52:57.0	-29:54:38	-1.6	9.3	131	$13.5 \pm 2.4$	$-11.0 \pm 1.9$	$219 \pm 44$	S
295090039	0:54:41.4	-28:13:54	-1.7	9.0	54	$15.8 \pm 2.1$	$-7.4 \pm 2.2$	$186 \pm 42$	S
221830014	1:00:15.3	-02:17:29	-3	10.9	85	$2 \pm 5$	$16 \pm 5$	$635 \pm 127$	N
221830028	1:02:07.6	-07:02:34	-2.2	10.1	-50	$2 \pm 5$	$12 \pm 5$	$441 \pm 101$	N
221830024	1:02:53.2	-04:44:54	-2.2	8.7	-100	$11.2 \pm 2.1$	$-15.5 \pm 2.1$	$70 \pm 16$	A, T
295140008	1:05:42.6	-23:58:40	-1.3	9.0	113	$4.8 \pm 3.0$	$-8.7 \pm 2.3$	$129 \pm 25$	A, S
295140006	1:06:49.4	-25:01:48	-1	9.8	57	$5.5 \pm 5.7$	$-10.2 \pm 5.7$	$77 \pm 104$	S
221660032	1:07:33.8	-12:11:01	-2.6	9.5	-128	$15.0 \pm 3.1$	$-7.3 \pm 3.1$	$166 \pm 41$	A, N
221660034	1:08:03.3	-12:11:37	-2.2	11.2	128	$4 \pm 5$	$13 \pm 5$	$590 \pm 149$	N
295140013	1:10:31.3	-25:40:47	-2.7	10.3	-50	$0.8 \pm 3.8$	$-12.9 \pm 3.7$	$226 \pm 105$	S
295180028	1:16:44.5	-31:04:58	-2.5	11.7	14	$0.5 \pm 2.5$	$-4.2 \pm 2.6$	$127 \pm 100$	S
295180035	1:16:58.2	-27:45:57	-2.4	10.7	45	$-3.4 \pm 4.3$	$-9.8 \pm 4.3$	$280 \pm 133$	S
295140038	1:22:53.5	-26:17:35	-1.6	9.8	54	$5.9 \pm 4.6$	$-14.3 \pm 3.1$	$154 \pm 87$	S
221740034	1:25:24.5	-09:36:19	-2.5	9.8	37	$1.6 \pm 3.4$	$-10.2 \pm 3.4$	$136 \pm 51$	A, N, T
221740042	1:30:19.3	-09:44:57	-2.5	9.0	7	$4.3 \pm 1.8$	$-18.6 \pm 1.8$	$109 \pm 19$	A, N, T

TABLE 1. (continued)

NAME <sup>a</sup>	RA (2000.0)	DEC	[Fe/H]	$r^b$ (kpc)	RV <sup>c</sup> (km/s)	$\mu_{\alpha^*}$ (mas/yr)	$\mu_{\delta}$ (mas/yr)	$v$ (km/s)	Source <sup>d</sup>
221800003	1:30:46.2	-10:25:54	-2.7	10.6	79	$9 \pm 5$	$11 \pm 5$	$442 \pm 114$	N
295040004	1:31:02.1	-36:38:34	-2.3	10.1	45	$8.6 \pm 2.4$	$-5.6 \pm 2.6$	$100 \pm 68$	S
221800002	1:31:42.6	-10:05:30	-1.1	9.3	-82	$10 \pm 5$	$3 \pm 5$	$232 \pm 60$	N
221800006	1:35:19.5	-10:33:20	-2.2	11.0	47	$5 \pm 5$	$6 \pm 5$	$350 \pm 128$	N
295040028	1:36:12.2	-34:06:08	-2.3	10.0	-56	$6.7 \pm 3.8$	$-8.4 \pm 3.4$	$120 \pm 32$	S
221800017	1:36:49.9	-12:00:53	-1.5	9.6	34	$11.1 \pm 3.1$	$-5.2 \pm 3.1$	$111 \pm 50$	A, N
295040035	1:40:51.5	-33:25:36	-2.2	11.1	43	$10.6 \pm 4.8$	$-12.0 \pm 4.8$	$310 \pm 156$	S
295040045	1:47:30.5	-34:07:23	-2	9.8	189	$9.0 \pm 3.8$	$-6.0 \pm 2.4$	$144 \pm 49$	S
221710022	2:03:11.7	-08:13:10	-2.2	9.7	-201	$-2.4 \pm 2.8$	$-22.7 \pm 2.8$	$310 \pm 38$	A, N
221750001	2:13:30.0	-11:37:36	-1.4	9.7	67	$5.1 \pm 2.9$	$4.2 \pm 2.9$	$256 \pm 41$	A, N
221750003	2:15:32.1	-10:40:28	-1.8	10.7	16	$14 \pm 5$	$6 \pm 5$	$335 \pm 91$	N
221890005	2:32:29.2	-14:31:48	-1.6	11.4	-120	$10 \pm 5$	$-5 \pm 5$	$177 \pm 60$	N
221810032	3:01:34.1	-08:56:03	-1.2	11.5	-126	$5 \pm 5$	$2 \pm 5$	$273 \pm 96$	N
310640037	3:09:30.0	-67:37:08	-1.5	7.9	277	$67.1 \pm 2.1$	$-41.8 \pm 2.1$	$496 \pm 22$	A, T
221670008	3:16:28.3	-07:06:46	0	15.2	14	$1 \pm 5$	$0 \pm 5$	$204 \pm 209$	N
221670017	3:18:45.7	-03:38:50	-2	11.2	198	$14.0 \pm 3.9$	$-21.0 \pm 3.9$	$322 \pm 74$	A
221850020	3:27:47.2	-14:06:46	-1.7	9.7	24	$19.7 \pm 3.1$	$-22.8 \pm 3.1$	$144 \pm 37$	A, N
310750042	3:30:46.9	-66:38:17	-2.1	8.2	145	$13.5 \pm 4.4$	$1.0 \pm 4.4$	$86 \pm 67$	A
221760020	3:45:49.2	-10:23:13	-1.5	10.9	-6	$16 \pm 5$	$-14 \pm 5$	$181 \pm 81$	N
221690025	4:16:05.6	-14:34:02	-1.6	12.1	10	$17 \pm 5$	$2 \pm 5$	$336 \pm 93$	N
310720061	5:27:10.8	-59:05:17	-1.5	8.5	64	$-0.5 \pm 6.6$	$-1.0 \pm 6.6$	$188 \pm 57$	A
310720068	5:29:20.6	-61:16:28	-2.4	8.2	495	$41.8 \pm 3.4$	$-1.8 \pm 3.4$	$369 \pm 20$	A, T
156210043	10:09:27.2	+24:50:05	-3	12.5	67	$4.5 \pm 4.5$	$-2.0 \pm 4.5$	$245 \pm 113$	A
156210039	10:13:50.7	+25:18:24	0	10.3	131	$-4.5 \pm 5.0$	$-8.0 \pm 5.0$	$113 \pm 71$	A
156210070	10:21:56.1	+27:11:19	-1.8	9.6	35	$17.0 \pm 1.8$	$-6.5 \pm 1.8$	$304 \pm 20$	A, T
156210015	10:26:19.1	+23:30:36	-2.2	9.3	181	$-23.7 \pm 1.4$	$-11.4 \pm 1.4$	$246 \pm 17$	A, T
156210009	10:28:01.0	+27:32:38	-1.4	11.0	27	$-17.7 \pm 5.0$	$-2.0 \pm 5.0$	$373 \pm 91$	A
156210010	10:28:06.1	+26:32:51	0	10.5	46	$0.0 \pm 7.2$	$-28.0 \pm 7.2$	$324 \pm 138$	A
156250026	11:51:50.0	+26:21:56	-1.1	8.3	-3	$7.8 \pm 1.6$	$-17.8 \pm 1.6$	$189 \pm 10$	A, N, T
160260011	12:16:50.0	+28:56:03	-1.7	9.3	19	$-9.7 \pm 3.3$	$-8.5 \pm 3.3$	$70 \pm 58$	A, N
160260028	12:23:02.8	+27:27:15	-3	9.6	-81	$5.2 \pm 3.5$	$8.9 \pm 3.5$	$467 \pm 77$	A, N
160270049	13:12:26.9	+30:21:16	-1.9	8.7	57	$-14.4 \pm 3.6$	$-10.3 \pm 3.6$	$143 \pm 49$	A, N
228770008	13:12:49.9	-10:32:31	-2.2	7.8	90	$0 \pm 5$	$-5 \pm 5$	$135 \pm 140$	N
160270051	13:13:01.9	+31:01:28	-2.2	11.8	39	$-6 \pm 5$	$8 \pm 5$	$557 \pm 211$	N
228770012	13:13:32.9	-09:35:18	-2.6	7.4	312	$4.0 \pm 3.7$	$0.6 \pm 3.7$	$355 \pm 57$	A, N
228770005	13:15:28.5	-11:25:31	-3	8.4	31	$2 \pm 5$	$1 \pm 5$	$308 \pm 164$	N
228770031	13:17:08.5	-09:40:25	-2.7	7.6	65	$-6.8 \pm 5.0$	$2.3 \pm 5.0$	$225 \pm 104$	A, N
228770036	13:17:39.4	-11:18:33	-2.3	7.6	202	$-8.7 \pm 5.0$	$1.2 \pm 5.0$	$239 \pm 99$	A, N
228770038	13:18:01.9	-11:58:46	-2.2	7.2	57	$-13.4 \pm 5.0$	$4.3 \pm 5.0$	$270 \pm 82$	A, N
228770027	13:18:10.7	-08:51:05	-1.9	8.6	-11	$-7.4 \pm 5.0$	$-10.6 \pm 5.0$	$292 \pm 147$	A, N
228770030	13:18:26.0	-09:07:12	-3	9.5	136	$-7.8 \pm 5.0$	$-1.7 \pm 5.0$	$227 \pm 204$	A, N
228770020	13:18:32.8	-07:49:17	-2.2	7.5	-18	$-10 \pm 5$	$4 \pm 5$	$277 \pm 99$	N
228770026	13:20:55.2	-08:38:48	-2	7.3	161	$-10.3 \pm 3.3$	$-23.9 \pm 3.3$	$298 \pm 55$	A, N
228770045	13:21:35.2	-11:29:15	-2.1	7.2	42	$4.9 \pm 5.0$	$-16.8 \pm 5.0$	$272 \pm 75$	A, N

TABLE 1. (continued)

NAME <sup>a</sup>	RA (2000.0)	DEC	[Fe/H]	$r^b$ (kpc)	RV <sup>c</sup> (km/s)	$\mu_{\alpha^*}$ (mas/yr)	$\mu_{\delta}$ (mas/yr)	$v$ (km/s)	Source <sup>d</sup>
228770046	13:22:24.0	-11:17:27	-2.3	7.3	60	-5.6 $\pm$ 5.0	-5.8 $\pm$ 5.0	37 $\pm$ 96	A, N
228770049	13:22:52.2	-11:58:37	-2	7.2	-7	-10.4 $\pm$ 5.0	-11.4 $\pm$ 5.0	146 $\pm$ 78	A, N
228890013	13:36:32.1	-08:15:15	-1.9	7.1	32	-5 $\pm$ 5	-3 $\pm$ 5	129 $\pm$ 90	N
228890023	13:38:14.8	-12:09:25	-1.9	7.2	-75	-11 $\pm$ 5	-9 $\pm$ 5	256 $\pm$ 126	N
228890037	13:40:14.5	-09:14:04	-2.5	7.0	174	-6 $\pm$ 5	-14 $\pm$ 5	173 $\pm$ 97	N
228890040	13:42:53.8	-07:49:02	-3	8.9	-80	0 $\pm$ 5	-4 $\pm$ 5	220 $\pm$ 182	N
228890060	13:49:13.8	-10:22:21	-1.1	6.9	-107	-5.9 $\pm$ 3.5	-6.5 $\pm$ 3.5	206 $\pm$ 42	A, N
228890058	13:51:13.4	-11:03:31	0	6.8	65	6.2 $\pm$ 3.2	-24.3 $\pm$ 3.2	351 $\pm$ 46	A, N
228890057	13:52:21.3	-11:07:54	-0.6	6.7	8	-11 $\pm$ 5	-2 $\pm$ 5	158 $\pm$ 101	N
156230001	14:08:30.2	+23:20:51	-2.8	8.1	-159	11.0 $\pm$ 4.7	-12.0 $\pm$ 4.7	431 $\pm$ 89	A
228830007	14:15:34.5	+10:21:43	-2.6	7.3	51	-20.5 $\pm$ 1.8	-39.0 $\pm$ 1.8	304 $\pm$ 35	A, N, T
228830028	14:20:47.1	+09:15:11	-3	7.1	-1	-29.6 $\pm$ 7.3	-14.0 $\pm$ 7.3	462 $\pm$ 145	A
228740025	14:30:24.7	-24:00:11	-2.1	6.0	-27	-41.1 $\pm$ 6.4	9.0 $\pm$ 6.4	585 $\pm$ 85	A
228710008	14:32:03.4	-20:34:17	-2	5.6	83	-1 $\pm$ 5	-8 $\pm$ 5	147 $\pm$ 132	N
228740021	14:33:17.2	-24:32:46	-2.3	5.7	91	-13.7 $\pm$ 6.0	-20.0 $\pm$ 6.0	278 $\pm$ 98	A
228710003	14:33:57.1	-21:12:58	-1.5	5.5	78	-16 $\pm$ 5	-3 $\pm$ 5	281 $\pm$ 113	N
228710013	14:34:14.4	-19:37:26	-3	5.8	-14	-8 $\pm$ 5	-9 $\pm$ 5	204 $\pm$ 157	N
228710005	14:34:30.7	-20:50:56	-2.2	6.1	105	-11 $\pm$ 5	0 $\pm$ 5	318 $\pm$ 168	N
228740036	14:34:57.5	-24:13:04	-1.9	6.0	40	-25.5 $\pm$ 6.1	-13.0 $\pm$ 6.1	260 $\pm$ 79	A
228710009	14:35:30.9	-20:27:53	0	7.0	-179	-15.7 $\pm$ 2.1	7.7 $\pm$ 2.1	333 $\pm$ 15	A, N, T
228740053	14:35:37.4	-26:31:15	-0.6	5.9	-94	-25.1 $\pm$ 5.8	0.0 $\pm$ 5.8	339 $\pm$ 76	A
228710045	14:38:12.2	-21:45:48	-2.3	5.4	-5	-11 $\pm$ 5	-0.2 $\pm$ 5	226 $\pm$ 112	N
228710031	14:38:47.6	-18:54:01	-3	5.5	210	-12 $\pm$ 5	-7 $\pm$ 5	241 $\pm$ 114	N
228710034	14:39:20.1	-19:16:13	-1.5	6.1	84	-14 $\pm$ 5	-13 $\pm$ 5	530 $\pm$ 198	N
228710077	14:41:17.3	-18:03:25	-1.7	6.1	-10	-10 $\pm$ 5	3 $\pm$ 5	354 $\pm$ 148	N
228710062	14:42:05.6	-20:38:44	-2	5.5	35	-5 $\pm$ 5	-5 $\pm$ 5	35 $\pm$ 116	N
228710064	14:42:22.4	-20:31:01	-2	5.4	192	4 $\pm$ 5	-12 $\pm$ 5	352 $\pm$ 121	N
228710052	14:42:53.2	-22:24:30	-2.3	5.5	79	-8 $\pm$ 5	-1 $\pm$ 5	130 $\pm$ 88	N
228710063	14:43:00.6	-20:42:53	-3	5.6	159	-20 $\pm$ 5	-6 $\pm$ 5	259 $\pm$ 83	N
228710057	14:43:10.0	-21:11:28	-1.9	5.8	82	-16 $\pm$ 5	-8 $\pm$ 5	124 $\pm$ 71	N
228710085	14:43:51.6	-18:43:19	-2.2	5.4	209	-16 $\pm$ 5	-2 $\pm$ 5	302 $\pm$ 104	N
228710088	14:44:26.0	-19:20:25	-1.8	5.4	115	-9 $\pm$ 5	3 $\pm$ 5	245 $\pm$ 105	N
228710087	14:45:16.9	-18:58:35	-2.4	5.4	-107	-9 $\pm$ 5	11 $\pm$ 5	464 $\pm$ 103	N
228710092	14:46:08.5	-20:30:13	-2.4	5.3	57	-8 $\pm$ 5	-5 $\pm$ 5	45 $\pm$ 105	N
228710096	14:48:02.9	-21:32:53	-2	5.6	-40	-17 $\pm$ 5	-13 $\pm$ 5	208 $\pm$ 75	N
228710115	14:48:35.2	-17:37:22	-2.7	5.5	17	0 $\pm$ 5	6 $\pm$ 5	353 $\pm$ 99	N
228710109	14:48:42.3	-18:47:59	-2.6	5.3	-54	0 $\pm$ 5	3 $\pm$ 5	302 $\pm$ 111	N
228710099	14:50:24.1	-21:43:10	-2.5	5.3	162	-3 $\pm$ 5	-9 $\pm$ 5	190 $\pm$ 120	N
228710103	14:50:52.9	-20:53:17	-1.2	5.4	75	1 $\pm$ 5	-7 $\pm$ 5	209 $\pm$ 156	N
228710113	14:51:26.5	-19:31:48	-2.8	5.3	97	0 $\pm$ 5	-11 $\pm$ 5	270 $\pm$ 135	N
228900015	15:16:10.2	+02:09:16	-3	6.5	129	-25.0 $\pm$ 4.3	-40.0 $\pm$ 4.3	380 $\pm$ 60	A
228900042	15:18:53.8	+00:45:52	-2.2	6.4	-78	-21.0 $\pm$ 5.0	-13.0 $\pm$ 5.0	130 $\pm$ 55	A
228840015	15:31:46.9	-09:28:32	-2.1	5.9	223	-4.8 $\pm$ 3.3	-9.4 $\pm$ 3.3	241 $\pm$ 32	A, N
228840021	15:34:44.9	-08:16:39	-1.6	5.0	-59	-13 $\pm$ 5	6 $\pm$ 5	358 $\pm$ 92	N

TABLE 1. (continued)

NAME <sup>a</sup>	RA (2000.0)	DEC	[Fe/H]	$r^b$ (kpc)	RV <sup>c</sup> (km/s)	$\mu_{\alpha^*}$ (mas/yr)	$\mu_{\delta}$ (mas/yr)	$v$ (km/s)	Source <sup>d</sup>
228840006	15:35:13.9	-11:40:54	0	4.8	34	-7 ± 5	-6 ± 5	48 ± 106	N
228840036	15:37:38.5	-11:30:38	-2.2	5.6	69	-13.7 ± 3.5	-13.1 ± 3.5	88 ± 49	A, N
228840047	15:41:27.7	-11:28:06	-2.1	5.7	-90	-14.1 ± 3.1	1.6 ± 3.1	225 ± 36	A, N
228720040	16:21:43.7	-03:08:00	-1.2	5.0	-65	2 ± 5	3 ± 5	295 ± 94	N
228720041	16:22:40.4	-03:24:38	-2	6.6	-204	-48.9 ± 3.9	-31.0 ± 3.9	317 ± 33	A
228720071	16:25:40.3	-02:57:18	-3	4.3	74	3.0 ± 8.5	-3.0 ± 8.5	278 ± 214	A
228720067	16:25:51.5	-03:38:59	-2.6	4.3	-16	1.0 ± 5.8	-8.0 ± 5.8	180 ± 134	A
228780105	16:50:04.4	+08:11:15	-2.3	5.7	-179	-4.0 ± 5.2	14.0 ± 5.2	416 ± 72	A
229590022	18:45:25.7	-65:57:31	-1.7	6.1	304	-22.8 ± 4.3	14.0 ± 4.3	481 ± 45	A
229590189	19:10:50.0	-66:33:16	-1.6	6.0	-129	-0.8 ± 4.8	-9.0 ± 4.8	248 ± 54	A
229390058	19:22:28.0	-28:09:09	-2	3.2	17	-2.5 ± 3.3	3.2 ± 3.4	322 ± 83	S
229390167	19:28:43.8	-28:50:26	-2	5.0	-32	-2.7 ± 8.7	3.3 ± 3.6	277 ± 71	S
228960041	19:29:22.3	-52:54:06	-3	5.2	-122	-0.6 ± 5.8	-11.0 ± 5.8	180 ± 76	A
229390164	19:32:31.4	-28:41:21	-1.7	5.3	-264	13.7 ± 2.8	-13.8 ± 2.8	336 ± 25	A, S
228960086	19:36:45.3	-57:18:59	0	5.5	-169	-11.9 ± 5.7	-9.0 ± 5.7	284 ± 83	A
229640125	19:57:17.9	-39:01:18	-2.1	4.2	26	1.3 ± 2.9	3.9 ± 2.9	338 ± 72	S
229640219	20:06:29.4	-39:03:41	-0.3	4.6	19	-5.6 ± 2.9	-1.1 ± 2.9	230 ± 61	S
229500023	20:16:20.7	-14:33:38	-1.3	7.0	-13	19.3 ± 2.0	-14.4 ± 2.0	208 ± 12	T
229500008	20:18:25.5	-16:01:17	-1.2	6.0	-6	-2.9 ± 5.4	-3.0 ± 5.4	196 ± 59	A
228850092	20:26:02.3	-40:44:34	-1.8	4.6	36	-1.1 ± 2.3	-7.0 ± 2.3	60 ± 51	S
228850125	20:26:53.6	-37:57:37	-1.1	4.8	-175	1.6 ± 4.1	-8.5 ± 4.9	171 ± 59	S
229550104	20:31:35.1	-25:51:22	-1.6	5.2	-122	5.0 ± 2.7	-4.3 ± 2.6	230 ± 89	S
229550103	20:31:54.4	-26:06:51	-1.6	7.7	-180	2.5 ± 3.0	-6.9 ± 3.0	268 ± 139	S
229550099	20:32:47.2	-26:33:23	-1.1	6.5	-101	9.4 ± 2.4	-25.9 ± 2.6	109 ± 18	S
228850179	20:33:41.1	-38:42:26	-1.2	4.7	-28	-0.3 ± 2.2	0.4 ± 2.2	246 ± 65	S
229550147	20:38:35.6	-25:48:06	-1.1	5.2	10	-6.8 ± 2.0	-5.8 ± 2.0	188 ± 38	S
228800061	20:41:28.9	-21:47:13	-2.5	5.2	-107	-3.6 ± 2.1	-15.6 ± 2.2	300 ± 68	N, S
228790063	20:45:08.4	-41:14:00	-1.4	5.0	9	-5.7 ± 2.2	-6.1 ± 2.2	170 ± 47	S
304920015	21:03:42.6	-40:23:29	-2	5.4	-39	-5.8 ± 4.7	-10.0 ± 3.8	152 ± 84	S
295010036	21:07:15.8	-36:06:45	-1.2	5.7	-151	2.9 ± 2.3	-8.5 ± 2.3	151 ± 33	S
304920063	21:08:33.1	-40:47:28	-0.4	6.9	15	4.4 ± 3.8	-10.1 ± 3.8	326 ± 171	S
295010057	21:16:13.3	-35:09:10	-1.6	5.7	-98	19.1 ± 2.6	-7.6 ± 2.6	442 ± 46	S
295010102	21:25:59.7	-36:09:49	-1.6	5.8	-37	2.1 ± 1.8	-4.5 ± 1.8	112 ± 47	S
229480002	21:33:50.1	-40:53:38	-1.9	6.2	-56	5.3 ± 2.5	1.7 ± 2.5	336 ± 81	S
294930009	21:39:43.1	-30:53:39	-0.9	6.4	1	-4.5 ± 3.1	-6.4 ± 2.3	168 ± 38	S
294930012	21:41:20.0	-29:53:38	-1.6	6.2	-175	9.8 ± 5.1	-10.8 ± 5.8	278 ± 131	S
295160054	22:19:58.2	+04:52:34	-1.2	7.7	56	-11.3 ± 2.2	-10.1 ± 2.2	275 ± 15	A, T
295160017	22:27:58.7	+03:59:28	-2.2	7.7	-228	-21.0 ± 5.7	-31.0 ± 5.4	351 ± 62	A
295160011	22:28:36.2	+06:21:09	-2.6	7.8	-14	13.6 ± 3.4	-7.3 ± 3.4	186 ± 29	A, T
303320115	22:34:19.2	+09:16:19	-2	8.4	-73	-8.9 ± 5.1	11.0 ± 5.1	463 ± 78	A
303320025	22:46:15.8	+07:51:34	-1.5	8.0	-37	0.0 ± 5.7	-2.0 ± 5.7	194 ± 49	A
303320016	22:48:46.3	+10:51:48	-2.5	8.7	-180	4.9 ± 4.8	-25.0 ± 4.8	386 ± 89	A
295130017	23:19:16.9	-37:40:48	-1.6	7.6	94	9.0 ± 6.3	-9.2 ± 3.1	134 ± 68	S
229410002	23:28:24.9	-35:56:04	-2	8.4	-77	5.7 ± 1.9	-3.9 ± 2.0	143 ± 55	S

TABLE 1. (continued)

NAME <sup>a</sup>	RA (2000.0)	DEC	[Fe/H]	$r^b$ (kpc)	RV <sup>c</sup> (km/s)	$\mu_{\alpha^*}$ (mas/yr)	$\mu_{\delta}$ (mas/yr)	$v$ (km/s)	Source <sup>d</sup>
229410008	23:28:40.2	-33:55:22	-3	7.6	-85	$-5.7 \pm 1.9$	$-11.7 \pm 1.9$	$208 \pm 26$	A, S
229410016	23:31:27.3	-33:57:02	0	7.9	43	$8.7 \pm 3.7$	$3.5 \pm 3.2$	$305 \pm 70$	S
229410022	23:33:11.6	-35:22:03	-1.7	9.0	-25	$11.9 \pm 1.8$	$-2.5 \pm 2.5$	$342 \pm 64$	S
229410028	23:35:19.6	-36:48:25	-2	8.4	-6	$6.1 \pm 3.0$	$-14.3 \pm 3.0$	$218 \pm 89$	S
294990022	23:38:44.8	-22:31:20	-2.2	7.8	101	$11.2 \pm 1.4$	$-16.12 \pm 1.4$	152 12 11	A, N, S, T
229410037	23:38:54.2	-35:16:38	-2	7.7	-69	$2.9 \pm 2.6$	$-9.1 \pm 2.9$	$102 \pm 35$	A, S
294990019	23:39:03.0	-23:12:36	-1.6	9.4	-69	$8.1 \pm 4.5$	$-0.4 \pm 4.4$	$254 \pm 146$	S
294990001	23:39:19.5	-27:24:37	-3	9.5	209	$2.0 \pm 2.5$	$-5.4 \pm 2.5$	$238 \pm 44$	S
229660029	23:41:15.9	-30:24:00	-2.5	8.1	0	$17.9 \pm 5.9$	$-5.9 \pm 4.3$	$291 \pm 111$	S
294990035	23:43:11.9	-26:38:47	-2.5	7.9	-203	$16.3 \pm 2.6$	$-19.6 \pm 3.0$	$235 \pm 38$	A, S
294990026	23:43:31.5	-22:57:50	-1.6	8.8	-73	$10.5 \pm 3.9$	$-6.0 \pm 3.7$	$183 \pm 101$	S
229410054	23:44:03.2	-33:01:35	-2	9.3	-155	$6.3 \pm 4.1$	$-11.1 \pm 4.5$	$259 \pm 134$	S
229660042	23:44:29.1	-28:19:21	-2.1	7.9	-46	$4.2 \pm 3.1$	$-20.3 \pm 2.5$	$129 \pm 44$	A, S
228940048	23:45:29.5	-01:57:27	-2	12.1	-135	$6 \pm 5$	$-4 \pm 5$	$171 \pm 203$	N
294990040	23:46:14.7	-25:45:07	-3	9.3	-66	$2.8 \pm 3.5$	$-8.7 \pm 3.5$	$82 \pm 98$	S
229410050	23:46:28.1	-34:56:01	-1.8	8.4	58	$6.2 \pm 5.5$	$-0.6 \pm 5.5$	$207 \pm 133$	S
229660061	23:46:48.6	-30:00:29	-1.8	8.4	-131	$-0.9 \pm 7.0$	$-10.5 \pm 8.7$	$181 \pm 106$	S
294990036	23:46:51.9	-26:48:08	-2.7	9.7	48	$11.2 \pm 3.3$	$-9.9 \pm 3.3$	$332 \pm 108$	S
295170002	23:47:04.3	-16:41:01	-1.3	8.5	2	$5 \pm 5$	$6 \pm 5$	$321 \pm 95$	N
294960025	23:47:06.2	-30:02:49	-3	9.5	8	$-1.7 \pm 7.4$	$-5.3 \pm 8.0$	$171 \pm 264$	S
294990042	23:47:34.9	-24:55:09	-1.9	8.5	-29	$7.2 \pm 3.6$	$-4.6 \pm 3.8$	$114 \pm 84$	S
295170007	23:47:52.7	-14:47:59	-2.4	13.3	-309	$-2 \pm 5$	$-7 \pm 5$	$376 \pm 220$	N
294990037	23:49:22.6	-26:35:43	-2.8	8.2	-104	$4.5 \pm 3.3$	$-7.5 \pm 4.4$	$104 \pm 47$	S
294990038	23:49:40.0	-26:31:36	-2.4	8.9	-33	$13.1 \pm 3.5$	$-9.5 \pm 3.4$	$264 \pm 95$	S
229660059	23:49:59.0	-29:56:11	-2.8	11.5	-126	$0.6 \pm 3.2$	$-5.5 \pm 2.7$	$153 \pm 81$	S
228760011	23:51:02.7	-34:04:00	-1.7	8.5	117	$15.5 \pm 4.8$	$-3.5 \pm 4.0$	$343 \pm 102$	S
229660065	23:51:16.8	-31:30:34	-2.9	11.3	238	$1.5 \pm 2.9$	$-4.8 \pm 3.7$	$248 \pm 67$	S
229660076	23:51:40.0	-29:10:16	-0.7	8.1	-32	$-2.5 \pm 2.3$	$-7.9 \pm 2.9$	$169 \pm 47$	S
228760006	23:52:19.2	-36:03:54	-2.3	8.3	-28	$7.4 \pm 4.0$	$-0.2 \pm 4.2$	$217 \pm 101$	S
228760010	23:53:52.2	-33:51:51	-2.2	9.0	96	$6.9 \pm 2.4$	$-8.3 \pm 2.5$	$144 \pm 54$	S
229660069	23:54:58.8	-31:05:53	-1.7	8.3	-114	$-8.9 \pm 3.0$	$-20.0 \pm 1.9$	$399 \pm 54$	S
294990064	23:55:32.8	-25:32:52	-2.1	9.7	-48	$7.5 \pm 2.2$	$-4.3 \pm 4.2$	$151 \pm 100$	S
229660071	23:55:48.3	-30:35:43	-3	11.5	35	$2.0 \pm 2.8$	$-6.7 \pm 3.1$	$125 \pm 137$	S
228760022	23:58:06.4	-33:45:11	-2.5	8.3	38	$17.0 \pm 2.7$	$-6.7 \pm 2.6$	$278 \pm 58$	S
228760019	23:59:20.1	-33:17:04	-2.1	8.8	2	$1.5 \pm 2.9$	$-12.9 \pm 2.9$	$172 \pm 79$	S

<sup>a</sup>The names of FHB stars follow Wilhelm et al. (1999).

<sup>b</sup>heliocentric radial velocity

<sup>c</sup>Galactocentric distance

<sup>d</sup>A, N, S, and T denote the STARNET catalogue, the NPM catalogue, the SPM catalogue and the Tycho-2 catalogue, respectively.

TABLE 2. Basic Data for Seven High Velocity Objects

Name	RA (2000.0)	DEC	$l$	$b$	$r$ (kpc)	RV (km/s)	$\mu_{\alpha} \cos \delta$ (mas/yr)	$\mu_{\delta}$ (mas/yr)	Type <sup>a</sup>
Leo I	10:08:27	+12:18.5:00	226	49	250	286	...	...	S
Draco	17:20:19.0	+57:54.8:0.0	86	35	82.0	-293.0	$0.60 \pm 0.50$	$1.10 \pm 0.50$	S
Pal 3	10:05:31.4	+00:04:17.0	240	42	96.8	83.4	$0.33 \pm 0.23$	$0.3 \pm 0.31$	G
160270051	13:13:1.90	+31:01:28.0	74	84	11.8	39.0	$-6 \pm 5$	$8 \pm 5$	F
221660034	01:08:3.30	-12:11:37.0	138	-75	11.2	128	$4 \pm 5$	$13 \pm 5$	F
221830014	01:00:15.3	-02:17:29.0	128	-65	10.9	85.0	$2 \pm 5$	$16 \pm 5$	F
228740025	14:30:24.7	-24:00:11.0	330	34	6.1	-27.0	$-41.1 \pm 6.4$	$9.0 \pm 6.4$	F

<sup>a</sup>S, G, and F denote satellite galaxies, globular clusters and FHB stars, respectively.

TABLE 3. Likelihood Results for Only the Radial Velocities

$a_s$ or $\gamma$	$a$ prior	$\beta$ prior	Leo I	best $\beta$	best $a$ (kpc)	best $M^a$	$M(< 50\text{kpc})^a$	$M(< 100\text{kpc})^a$
Power-law Tracers								
$\gamma = 3.4$	$1/a^2$	Energy	Yes	−1	160	18.0	5.4	9.6
			No	−1	70	7.9	4.6	6.5
$\gamma = 3.4$	$1/a$	Energy	Yes	−1	175	20.0	5.4	9.8
			No	−1	75	8.5	4.7	6.8
$\gamma = 3.4$	$1/a^2$	Uniform	Yes	−1	160	18.0	5.4	9.6
			No	−1	70	7.9	4.6	6.5
$\gamma = 4.0$	$1/a^2$	Energy	Yes	−1	170	19.0	5.4	9.7
			No	−1	80	9.0	4.8	7.1
Shadow Tracers								
$a_s = 100$	$1/a^2$	Energy	Yes	−1	180	20.0	5.4	9.8
			No	−1	65	7.4	4.5	6.2
$a_s = 100$	$1/a$	Energy	Yes	−1	205	23.0	5.5	10.0
			No	−1	70	7.9	4.6	6.5
$a_s = 100$	$1/a^2$	Uniform	Yes	−1	180	20.0	5.4	9.8
			No	−1	65	7.4	4.5	6.2
$a_s = a_{halo}$	$1/a^2$	Energy	Yes	−1	170	19.0	5.4	9.7
			No	−1	70	7.9	4.6	6.5

<sup>a</sup>All masses are in units of  $10^{11} M_\odot$ .

TABLE 4. Likelihood Results for the Full Space Velocities

$a_s$ or $\gamma$	$a$ prior	$\beta$ prior	Leo I	best $\beta$	best $a$ (kpc)	best $M^a$	$M(< 50\text{kpc})^a$	$M(< 100\text{kpc})^a$
Power-law Tracers								
$\gamma = 3.4$	$1/a^2$	Energy	Yes	−1	200	23.0	5.5	10.0
			No	−1	150	17.0	5.3	9.4
$\gamma = 3.4$	$1/a$	Energy	Yes	−1	230	26.0	5.5	10.0
			No	−1	160	18.0	5.4	9.6
$\gamma = 3.4$	$1/a^2$	Uniform	Yes	−1	200	23.0	5.5	10.0
			No	−1	150	17.0	5.3	9.4
$\gamma = 4.0$	$1/a^2$	Energy	Yes	−1	225	25.0	5.5	10.0
			No	−1	165	19.0	5.4	9.6
Shadow Tracers								
$a_s = 100$	$1/a^2$	Energy	Yes	−1	275	31.0	5.5	11.0
			No	−1	190	21.0	5.4	10.0
$a_s = 100$	$1/a$	Energy	Yes	−1	330	37.0	5.6	11.0
			No	−1	220	25.0	5.5	10.0
$a_s = 100$	$1/a^2$	Uniform	Yes	−1	275	31.0	5.5	11.0
			No	−1	190	21.0	5.4	10.0
$a_s = a_{halo}$	$1/a^2$	Energy	Yes	−1	275	31.0	5.5	11.0
			No	−1	190	21.0	5.4	10.0

<sup>a</sup>All masses are in units of  $10^{11} M_\odot$ .



# Gravothermal Catastrophe and Tsallis' Generalized Entropy of Self-Gravitating Systems : Thermodynamic properties of stellar polytrope

Atsushi Taruya<sup>a</sup> and Masa-aki Sakagami<sup>b</sup>

<sup>a</sup>*Research Center for the Early Universe(RESCEU), School of Science, University of Tokyo, Tokyo 113-0033, Japan*

<sup>b</sup>*Department of Fundamental Sciences, FIHS, Kyoto University, Kyoto 606-8501, Japan*

## Abstract

We investigate the thermodynamic properties of stellar self-gravitating system arising from the Tsallis generalized entropy. In particular, physical interpretation of the thermodynamic instability, as has been revealed by previous study(Taruya & Sakagami, Physica A 307 (2002) 185), is discussed in detail based on the framework of non-extensive thermostatics. Examining the Clausius relation in a quasi-static experiment, we obtain the standard result of thermodynamic relation that the physical temperature of the equilibrium non-extensive system is identified with the inverse of the Lagrange multiplier,  $T_{\text{phys}} = 1/\beta$ . Using this relation, the specific heat of total system is computed, and confirm the common feature of self-gravitating system that the presence of negative specific heat leads to the thermodynamic instability. In addition to the gravothermal instability discovered previously, the specific heat shows the curious divergent behavior at the polytrope index  $n > 3$ , suggesting another type of thermodynamic instability in the case of the system surrounded by the thermal bath. Evaluating the second variation of free energy, we check the condition for onset of this instability and find that the zero-eigenvalue problem of the second variation of free energy exactly recovers the marginal stability condition indicated from the specific heat. Thus, the stellar polytropic system is consistently characterized by the non-extensive thermostatics as a plausible thermal equilibrium state. We also clarify the non-trivial scaling behavior appeared in specific heat and address the origin of non-extensive nature in stellar polytrope.

## 1 Introduction

Due to its complexity and peculiarity, stellar self-gravitating system has long attracted much attention in the subject of astronomy and astrophysics, and even statistical physics. For an isolated stellar system, the dynamical equilibrium is rapidly attained after a few crossing time and the thermodynamic description provides useful information in characterizing the late-time behavior of this system. Even in this simplest situation, however, the equilibrium state of self-gravitating system shows various interesting phenomena, which may offer an opportunity to recast the framework of the thermodynamics and/or statistical mechanics.

In earlier study, applying the Tsallis' generalized entropy[1], we have investigated the thermodynamic instability of self-gravitating systems[2]. The self-gravitating stellar system confined in a spherical cavity of radius,  $r_e$ , exhibits an instability, so-called *gravothermal catastrophe*, which has been widely accepted as a fundamental physical process and plays an important role for the long-term evolution of globular clusters [3, 4, 5].

The presence of this instability has been long known since the pioneer work by Antonov[6] and Lynden-Bell & Wood[7]. Historically, the gravitational catastrophe has been studied on the basis of the maximum entropy principle for the phase-space distribution function, with a particular attention to the Boltzmann-Gibbs entropy [8, 9].

In contrast to previous work, we have applied the Tsallis-type generalized entropy to seek the equilibrium criteria for the first time. Then, the distribution function of Vlasov-Poisson system can be reduced to a stellar polytropic system[10, 11]. Evaluating the second variation of entropy around the equilibrium state and solving the zero-eigenvalue problem, the criterion for the onset of gravothermal instability is obtained. The main results of our previous analysis are summarized as follows:

- (i) Local entropy extremum ceases to exist in cases with polytrope index  $n > 5$  for sufficiently larger radius of the wall,  $r_e > \lambda_{\text{crit}} GM^2/(-E)$ , and for highly density contrast,  $\rho_c/\rho_e > D_{\text{crit}}$ , where  $M$  and  $E$  denote the total mass and energy of the system,  $\rho_c$  and  $\rho_e$  mean the density at center and edge, respectively.
- (ii) The critical values  $\lambda_{\text{crit}}$  and  $D_{\text{crit}}$  depend on the polytrope index, both of which respectively approach 0.335 and 709 in the limit of  $n \rightarrow \infty$ , consistent with the well-known result adopting the Boltzmann-Gibbs entropy.
- (iii) The stability/instability criterion obtained from the second variation of Tsallis entropy exactly matches with the result from standard turning-point analysis.

While the successful results suggest that non-extensive generalization of thermodynamics will offer various astrophysical applications involving long-range nature of self-gravitating systems, there still remain some important issues concerning the physical interpretation of thermodynamic instability.

Heuristically, the gravothermal instability is explained by the presence of negative specific heat as follows. In a fully relaxed gravitating system with sufficiently larger radius, negative specific heat arises at the inner part of the system and we have  $C_{V,\text{inner}} < 0$ , while the specific heat at the outer part remains positive,  $C_{V,\text{outer}} > 0$ , since one can safely neglect the effect of self-gravity. In this situation, if a tiny heat flow is momentarily supplied from inner to outer part, both the inner and the outer parts get hotter after the hydrostatic readjustment. Now imagine the case,  $C_{V,\text{outer}} > |C_{V,\text{inner}}|$ . The outer part has so much thermal inertia that it cannot heat up as fast as the inner part, and thereby the temperature difference between inner and outer parts increases. As a consequence, the heat flow never stops, leading to a catastrophe temperature growth.

While the above thought experiment is naive in a sense that we artificially divide the system into the inner and the outer part, the argument turns out to capture an essence of the thermodynamic instability in cases with the Boltzmann-Gibbs entropy. Evaluating the specific heat explicitly, Lynden-Bell and Wood[7] showed that the specific heat of the total system should be greater than zero at the onset of instability, although the central part of this system still has the negative specific heat. Therefore, one can naively expect that the self-gravitating system generally exhibits the thermodynamic instability associated with the negative specific heat and this could even hold in the system characterized by the non-extensive entropy.

To address this issue, however, we should remember the following two remarks that have been never clarified. First note that there exists a subtle point concerning the concept

of temperature in the non-extensive thermodynamics. Framework of the non-extensive formalism is formally constructed keeping the standard result of thermodynamic relations [14, 15, 16], however, the physical temperature,  $T_{\text{phys}}$ , might not be simply related to the usual one, i.e, the inverse of Lagrange multiplier, as has been criticized recently[12, 13]. This point is in particular important in evaluating the specific heat.

Second, as has been mentioned by the pioneer work of Lynden-Bell & Wood[7], self-gravitating system shows various types of thermodynamic instability. While our early study deals with the stellar system confined within an adiabatic wall, one may replace the adiabatic wall with the thermally conducting wall surrounded by a heat bath. In this situation, assuming the Boltzmann-Gibbs entropy, Lynden-Bell & Wood showed that no equilibrium state exists for sufficiently low temperature and high-density contrast. Note that even in this case, the presence of negative specific heat plays an essential role for the appearance of instability.

Keeping the above remarks in mind, in this article, we focus on the thermodynamic property of self-gravitating systems characterized by Tsallis' generalized entropy. For this purpose, we first investigate the thermodynamic temperature of the self-gravitating system from the Clausius relation. To clarify the physical interpretation of thermodynamic instability, the specific heat is computed and a role of negative specific is discussed in detail. Then we turn to focus on the thermodynamic instability in a system surrounded by the heat bath. The stability/instability criterion is derived from the second variation of free energy and a geometrical construction of marginal stability condition is discussed.

This article is organized as follows. in section 2, we recast the problem that finds the most probable state of equilibrium stellar distribution adopting the Tsallis entropy. The main part of this article is section 3, in which the thermodynamic properties of stellar polytrope are investigated in detail. After identification of the thermodynamic temperature, the explicit expression for specific heat is presented and the marginally stability condition for the thermodynamic instability is investigated in both the adiabatic and the isothermal cases. In section 4, thermodynamic instability in a system surrounded by a thermal bath is re-considered by means of the free energy and the marginal stability condition is re-derived from the second variation of free energy. Furthermore, following the preceding results, the origin of the non-extensive nature in stellar polytropic system is discussed in section 5. Finally, section 6 is devoted to the summary and conclusion.

## 2 Stellar polytrope as an extremum state of Tsallis entropy

In this section, we recast the problem finding the most probable state of equilibrium stellar system, based on the maximum entropy principle. In our previous study, the entropy for the phase-space distribution function has been introduced without recourse to the correct dimensions. Although this does not alter the stability/instability criterion for the stellar equilibrium state, for the sake of the completeness and the later analysis, we repeat the same calculation as shown in ref.[2], taking fully account of the correct dimensions.

Suppose a system containing  $N$  particles which are confined within a hard sphere of radius  $r_e$ . For simplicity, each particle is assumed to have the same mass  $m_0$  and interacts via Newton gravity. The problem considered here is to find an equilibrium state in an adiabatic treatment. That is, we investigate the equilibrium particle distribution in which

the particles elastically bounce from the wall, keeping the energy  $E$  and the total mass  $M(= Nm_0)$  constant.

For present purpose, it is better to employ the mean-field treatment that the correlation between particles is smeared out and the system can be fully characterized by the one-particle distribution function,  $f(\mathbf{x}, \mathbf{v})$ , defined in six-dimensional phase-space  $(\mathbf{x}, \mathbf{v})$  [2, 3][6, 7, 8, 9]. Let us denote the phase-space element as  $h^3(= l_0^3 v_0^3)$  with unit length  $l_0$  and unit velocity  $v_0$ . Since the distribution function  $f(\mathbf{x}, \mathbf{v})$  counts the number of particles in a unit cell of phase-space, the energy and the total mass are respectively expressed as follows:

$$E = K + U \equiv m_0 \int \left\{ \frac{1}{2} v^2 + \frac{1}{2} \Phi(\mathbf{x}) \right\} f(\mathbf{x}, \mathbf{v}) d^6 \tau, \quad (1)$$

$$M = m_0 N \equiv m_0 \int f(\mathbf{x}, \mathbf{v}) d^6 \tau, \quad (2)$$

with the quantity  $\Phi$  being the gravitational potential:

$$\Phi(\mathbf{x}) = -G m_0 \int \frac{f(\mathbf{x}', \mathbf{v}')}{|\mathbf{x} - \mathbf{x}'|} d^6 \tau'. \quad (3)$$

In the above expressions, the dimensionless integral measure  $d^6 \tau$  is introduced:

$$d^6 \tau \equiv \frac{d^3 \mathbf{x} d^3 \mathbf{v}}{h^3} ; \quad h = l_0 v_0. \quad (4)$$

Owing to the maximum entropy principle, we explore the most probable state maximizing the entropy. The entropy quoted here is a quantity defined in the phase-space and it counts the number of possible particle state. We are specifically concerned with the equilibrium state for the Tsallis entropy [1]:

$$S_q = -\frac{N}{q-1} \int \left[ \left( \frac{f}{N} \right)^q - \left( \frac{f}{N} \right) \right] d^6 \tau. \quad (5)$$

Maximizing the entropy  $S_q$  under the constraints reduces to the following mathematical problem using Lagrange multipliers  $\alpha$  and  $\beta$ :

$$\delta S_q - \alpha \delta M - \beta \delta E = 0, \quad (6)$$

which leads to [2, 10, 11]:

$$f(\mathbf{x}, \mathbf{v}) = A \left[ \Phi_0 - \Phi(\mathbf{x}) - \frac{1}{2} v^2 \right]^{1/(q-1)}, \quad (7)$$

where the constants  $A$  and  $\Phi_0$  are respectively given by

$$A = N \left\{ \left( \frac{q-1}{q} \right) m_0 \beta \right\}^{1/(q-1)}, \quad \Phi_0 = \frac{1 - (q-1)m_0 \alpha}{(q-1)m_0 \beta}. \quad (8)$$

The one-particle distribution function (7) is often called *stellar polytrope*, which satisfies the polytropic equation of state [3][10]. The density profile  $\rho(r)$  and the isotropic

pressure  $P(r)$  at the radius  $r = |\mathbf{x}|$  are respectively given by

$$\begin{aligned}\rho(r) &\equiv m_0 \int f(\mathbf{x}, \mathbf{v}) \frac{d^3 \mathbf{v}}{h^3} \\ &= 4\sqrt{2}\pi B\left(\frac{3}{2}, \frac{q}{q-1}\right) \frac{m_0 A}{h^3} \{\Phi_0 - \Phi(r)\}^{1/(q-1)+3/2},\end{aligned}\quad (9)$$

and

$$\begin{aligned}P(r) &\equiv m_0 \int \frac{1}{3} v^2 f(\mathbf{x}, \mathbf{v}) \frac{d^3 \mathbf{v}}{h^3} \\ &= \left(\frac{1}{q-1} + \frac{5}{2}\right)^{-1} \rho(r) \{\Phi_0 - \Phi(r)\},\end{aligned}\quad (10)$$

with  $B(a, b)$  being the  $\beta$  function. Thus, these two equations lead to the relation

$$P(r) = K_n \rho^{1+1/n}(r), \quad (11)$$

with the polytrope index given by

$$n = \frac{1}{q-1} + \frac{3}{2}. \quad (12)$$

In equation (11), the dimensional constant  $K_n$  is introduced:

$$K_n \equiv \frac{1}{n+1} \left\{ 4\sqrt{2}\pi B\left(\frac{3}{2}, n - \frac{1}{2}\right) \frac{m_0 A}{h^3} \right\}^{-1/n}. \quad (13)$$

Note that the above quantity is equivalent to the variable  $(n - 3/2)T/(n + 1)$  defined in ref.[2].

Once provided the distribution function, the equilibrium configuration can be completely specified by solving the Poisson equation. Hereafter, we specifically restrict our attention to the spherically symmetric configuration for  $q > 1$  (or  $n > 3/2$ ). From the gravitational potential (3), it reads

$$\frac{1}{r^2} \frac{d}{dr} \left( r^2 \frac{d\Phi(r)}{dr} \right) = 4\pi G \rho(r). \quad (14)$$

Combining (14) with (9), we obtain the ordinary differential equation for  $\Phi$ . Equivalently, a set of equations which represent the hydrostatic equilibrium are derived using (9), (10) and (14):

$$\frac{dP(r)}{dr} = -\frac{Gm(r)}{r^2} \rho(r), \quad (15)$$

$$\frac{dm(r)}{dr} = 4\pi \rho(r) r^2. \quad (16)$$

The quantity  $m(r)$  denotes the mass evaluated at the radius  $r$  inside the wall. We then introduce the dimensionless quantities:

$$\rho = \rho_c [\theta(\xi)]^n, \quad r = \left\{ \frac{(n+1)P_c}{4\pi G \rho_c^2} \right\}^{1/2} \xi, \quad (17)$$

which yields the following ordinary differential equation:

$$\theta'' + \frac{2}{\xi}\theta' + \theta^n = 0, \quad (18)$$

where prime denotes the derivative with respect to  $\xi$ . The quantities  $\rho_c$  and  $P_c$  in (17) are the density and the pressure at  $r = 0$ , respectively. To obtain the physically relevant solution of (18), we put the following boundary condition:

$$\theta(0) = 1, \quad \theta'(0) = 0. \quad (19)$$

A family of solutions satisfying (19) is referred to as the *Emden solution*, which is well-known in the subject of stellar structure (e.g., see Chap.IV of ref.[17]).

Figure 1 shows the numerical solution of equation (18) for various polytrope indices, where the density profile,  $\rho(r)/\rho_c$  is plotted as a function of dimensionless radius,  $\xi$ . Clearly, profiles with index  $n < 5$  rapidly fall off and they abruptly terminate at finite radius(*left-panel*), while the  $n \geq 5$  cases infinitely continue to extend over the outer radius(*right-panel*). As already mentioned in previous study, characteristic feature seen in figure 1 plays an essential role for the thermodynamic instability associated with negative specific heat.

For later analysis, it is convenient to introduce the following set of variables, referred to as homology invariants [17, 18]:

$$u \equiv \frac{d \ln m(r)}{d \ln r} = \frac{4\pi r^3 \rho(r)}{m(r)} = -\frac{\xi \theta^n}{\theta'}, \quad (20)$$

$$v \equiv -\frac{d \ln P(r)}{d \ln r} = \frac{\rho(r)}{P(r)} \frac{Gm(r)}{r} = -(n+1) \frac{\xi \theta'}{\theta}, \quad (21)$$

which reduce the degree of equation (18) from two to one. The derivative of these variables with respect to  $\xi$  becomes

$$\frac{du}{d\xi} = \left(3 - u - \frac{n}{n+1}v\right) \frac{u}{\xi}, \quad \frac{dv}{d\xi} = \left(-1 + u + \frac{1}{n+1}v\right) \frac{v}{\xi}. \quad (22)$$

Equations (18) can thus be re-written with

$$\frac{u}{v} \frac{dv}{du} = \frac{(n+1)(u-1) + v}{(n+1)(3-u) - nv}. \quad (23)$$

The corresponding boundary condition to (19) becomes  $(u, v) = (3, 0)$ . Using these variables, the basic thermodynamic quantities such as the energy and the entropy are evaluated and the results are summed up in Appendix A, which are subsequently used in section 3.

### 3 Thermodynamic properties of stellar polytrope

In this section, we address our main issue, i.e, the physical interpretation of gravothermal instability in stellar polytropes, based on the framework of non-extensive thermodynamics. In section 3.1, we first discuss the thermodynamic temperature of stellar polytrope

calculating both the heat and the entropy changes in a quasi-static treatment. Then we evaluate the specific heat in section 3.2. The connection between the absence of extremum entropy state and the presence of negative specific heat is discussed in detail. Further, we argue that there appears another type of thermodynamic instability, which is subsequently analyzed by means of the free energy.

### 3.1 Thermodynamic temperature from the Clausius relation

As has been mentioned in section 1, the concept of temperature is non-trivial in non-extensive thermostatics. This is because the standard framework of thermodynamics crucially depends on the assumption of extensivity of entropy. According to the recent claim, the definition of physical temperature  $T_{\text{phys}}$  should be altered depending on the choice of energy constraint and is related to the inverse of the Lagrange multiplier,  $1/\beta$ , with *some correction factors* [12, 13]. Note, however, that this discussion heavily relies on the extensivity of the energy as well as the thermodynamic zeroth law. In our present case, the maximum entropy principle was applied subject to the constraints  $E$  and  $M$ , adopting the standard definition of mean values (see eqs.(1)(2)). As a consequence, the resultant energy  $E$  becomes non-extensive and we cannot apply the above definition.

To address the physical temperature in the present case, we therefore consider the relation between the heat transfer and entropy change and seek the most plausible candidate for thermodynamic temperature. That is, we analyze the variation of equilibrium configuration under fixing the total mass. Specifically, we deal with the quasi-static variation along an equilibrium sequence.

Let us first write down the heat change. The thermodynamic first law states that

$$d'Q = dE + P_e dV, \quad (24)$$

where the operation  $d'$  stands for incomplete differentiation. The subscript  $_e$  denotes a quantity evaluated at the edge. In the spherically symmetric configuration, the second term in right-hand side of (24) becomes  $4\pi r_e^2 P_e dr_e$ . As for the first term, the energy of the stellar polytropic system within the radius  $r_e$ , is computed in Appendix A.1. Introducing the dimensionless parameter  $\lambda$ , it is expressed in terms of the homology invariants as follows:

$$\lambda \equiv -\frac{r_e E}{GM^2} = -\frac{1}{n-5} \left[ \frac{3}{2} \left\{ 1 - (n+1) \frac{1}{v_e} \right\} + (n-2) \frac{u_e}{v_e} \right], \quad (25)$$

where the quantity with subscript  $_e$  represents the one evaluated at the boundary  $r = r_e$ . Using (25), the heat change  $d'Q$  is rewritten as follows:

$$\begin{aligned} d'Q &= d \left( -\lambda \frac{GM^2}{r_e} \right) + 4\pi r_e^2 P_e dr_e, \\ &= \frac{GM^2}{r_e} \left\{ \left( \lambda + \frac{u_e}{v_e} \right) \frac{dr_e}{r_e} - \xi_e \frac{d\lambda}{d\xi_e} \frac{d\xi_e}{\xi_e} \right\}, \end{aligned} \quad (26)$$

where the relation  $4\pi r_e^4 P_e / (GM^2) = u_e / v_e$  is used in the last line (see definitions (20)(21)). In the above expression, derivative of  $\lambda$  with respect to  $\xi_e$  can be computed with a help of relation (22) (see eq.(33) of ref.[2]):

$$\xi_e \frac{d\lambda}{d\xi_e} = \frac{n-2}{n-5} \frac{g(u_e, v_e)}{2v_e}, \quad (27)$$

where

$$g(u, v) = 4u^2 + 2uv - \left\{ 8 + 3 \left( \frac{n+1}{n-2} \right) \right\} u - \frac{3}{n-2} v + 3 \left( \frac{n+2}{n-2} \right). \quad (28)$$

Next focus on the change of the entropy. From (71) in Appendix A.2, the entropy of the extremum state is given by

$$S_q = \left( n - \frac{3}{2} \right) \left[ \frac{1}{n-5} \frac{\beta G M^2}{r_e} \left\{ 2 \frac{u_e}{v_e} - (n+1) \frac{1}{v_e} + 1 \right\} + N \right]. \quad (29)$$

Hence, the variation of entropy  $dS_q$  under fixing the total mass can be decomposed into the variation of homology invariants  $(u_e, v_e)$ , radius  $r_e$  and Lagrange multiplier  $\beta$  as follows:

$$dS_q = \frac{n-3/2}{n-5} \frac{\beta G M^2}{r_e} \left[ \left( \frac{d\beta}{\beta} - \frac{dr_e}{r_e} \right) \left\{ 2 \frac{u_e}{v_e} - (n+1) \frac{1}{v_e} + 1 \right\} + \left\{ 2 \frac{u_e}{v_e} \left( \frac{du_e}{u_e} - \frac{dv_e}{v_e} \right) - \frac{n+1}{v_e} \frac{dv_e}{v_e} \right\} \right]. \quad (30)$$

Among these variations, variation of homology invariants is simply rewritten with  $d\xi_e$ , through the relation (22). On the other hand, from the mass conservation, the variation of Lagrange multiplier,  $d\beta$  is related to both the variations of homology invariants and  $dr_e$  as follows. Using the condition of hydrostatic equilibrium at the edge  $r_e$ , one can obtain the following relation (see derivation in Appendix A.3):

$$\eta \equiv \left\{ \frac{(GM)^n (m_0 \beta)^{n-3/2}}{r_e^{n-3} h^3} \right\}^{1/(n-1)} = \alpha_n (u_e v_e^n)^{1/(n-1)}, \quad (31)$$

where the constant  $\alpha_n$  is given by

$$\alpha_n = \left\{ \frac{(n-1/2)^{n-3/2}}{16\sqrt{2}\pi^2 (n+1)^n B(3/2, n-1/2)} \right\}^{1/(n-1)}, \quad (32)$$

which asymptotically approaches unity, in the limit  $n \rightarrow +\infty$ . Keeping the total mass  $M$  constant, variation of (31) yields

$$\frac{n-3/2}{n-1} \frac{d\beta}{\beta} - \frac{n-3}{n-1} \frac{dr_e}{r_e} = \frac{1}{n-1} \left( \frac{du_e}{u_e} + n \frac{dv_e}{v_e} \right). \quad (33)$$

We then rewrite it with

$$\frac{d\beta}{\beta} - \frac{dr_e}{r_e} = \frac{1}{n-3/2} \left( -\frac{3}{2} \frac{dr_e}{r_e} + \frac{du_e}{u_e} + n \frac{dv_e}{v_e} \right). \quad (34)$$

Substituting the relation (34) into equation (30), the dependence of  $d\beta/\beta$  can be eliminated. Thus, using the relation (22), the final form of the entropy change is expressed in terms of the variations  $d\xi_e$  and  $dr_e$ . After some manipulation, we obtain

$$dS_q = \frac{\beta G M^2}{r_e} \left[ -\frac{3/2}{n-5} \left( 2 \frac{u_e}{v_e} - \frac{n+1}{v_e} + 1 \right) \frac{dr_e}{r_e} - \frac{n-2}{n-5} \frac{1}{2v_e} \right. \\ \left. \times \left\{ 4u_e^2 + 2u_e v_e - \left( 8 + 3 \frac{n+1}{n-2} \right) u_e - \frac{3}{n-2} v_e + 3 \left( \frac{n+1}{n-2} \right) \right\} \frac{d\xi_e}{\xi_e} \right]. \quad (35)$$



Now, from the knowledge of the expressions  $\lambda$  and  $\xi_e(d\lambda/d\xi_e)$ , one can easily show that the above equation is just identical to

$$dS_q = \frac{\beta GM^2}{r_e} \left\{ \left( \lambda + \frac{u_e}{v_e} \right) \frac{dr_e}{r_e} - \xi_e \frac{d\lambda}{d\xi_e} \frac{d\xi_e}{\xi_e} \right\}. \quad (36)$$

Therefore, comparison between (36) and (26) immediately leads to the following relation:

$$dS_q = \beta d'Q = \beta (dE + P_e dV), \quad (37)$$

which exactly coincides with the standard result of *Clausius relation* in a quasi-static process.

The relation (37) strongly suggests that the thermodynamic temperature  $T_{\text{phys}}$  is identified with the inverse of Lagrange multiplier,  $T_{\text{phys}} = 1/\beta$ . At first glance, the result seems somewhat trivial, since one can easily expect this relation from the standard thermodynamic relation,  $\partial S_q / \partial E = \beta$ , which generally holds even in the non-extensive Tsallis formalism [14, 15]. As advocated by many author, however, the relation  $\partial S_q / \partial E = \beta$  does not simply imply the thermodynamic temperature  $T_{\text{phys}} = 1/\beta$  and it might even contradict with the thermodynamic temperature defined through the thermodynamic zeroth law [13].

On the other hand, in our case of the self-gravitating system, the thermodynamic temperature  $T_{\text{phys}} = 1/\beta$  is mathematically verified by the integrable condition of the thermodynamic entropy through the Clausius relation. Further, it is remarkably found that the relation  $T_{\text{phys}} = 1/\beta$  holds even in the absence of gravity (the limit  $G \rightarrow 0$ ) and can be proven through an alternative route. In Appendix B, as a pedagogical example, we demonstrate that the relation  $T_{\text{phys}} = 1/\beta$  is indeed obtained in the classical gas model using the Carnot cycle.

### 3.2 Negative specific heat and thermodynamic instability

Once obtained the thermodynamic temperature,  $T_{\text{phys}} = 1/\beta$ , we are in a position to investigate the thermodynamic instability from the straightforward calculation of the specific heat. Let us first discuss the qualitative behavior of the specific heat. By definition, the specific heat at constant volume is given by

$$C_v \equiv \left( \frac{dE}{dT_{\text{phys}}} \right)_e = -\beta^2 \left( \frac{dE}{d\beta} \right)_e = -\beta^2 \frac{\left( \frac{dE}{d\xi} \right)_e}{\left( \frac{d\beta}{d\xi} \right)_e}. \quad (38)$$

Recall that the dimensionless parameters  $\lambda$  and  $\eta$  are respectively proportional to  $-E$  and  $\beta^{(n-1)/(n-3/2)}$  (see eqs.(25)(31)). This implies that for a system of constant mass inside a fixed wall, the qualitative behavior of (38) can be deduced from the relation between  $\eta$  and  $\lambda$ .

Figure 2 depicts the trajectories of the Emden solutions in the  $(\eta, \lambda)$ -plane with various polytrope indices. Each point along the trajectory represents an Emden solution for different value of the radius  $r_e$ . From the boundary condition, all the trajectories start from  $(\eta, \lambda) = (0, -\infty)$ , corresponding to the origin  $r_e = 0$ . As gradually increasing the

radius, the trajectories first move to upper-right direction monotonically, as marked by the arrow. At this stage, the kinematic energy dominates the potential energy and the system lies in a kinematically thermal state ( $\lambda < 0$ ), indicating the positive specific heat. For larger radius, while the curves with index  $n \leq 3$  abruptly terminate, the trajectories with  $n > 3$  suddenly change their direction from upper-right to upper-left. Moreover, in the case of  $n > 5$ , the trajectory progressively changes its direction and it finally spirals around a fixed point.

From these observations, one can roughly infer the existence of the two types of the thermodynamic instability as follows. At first inflection point for  $n > 3$ , the specific heat diverges and the signature of  $C_v$  becomes indefinite. Beyond this point, the specific heat changes from positive to negative. This means that the potential energy conversely dominates the kinetic energy, indicating the system being *gravothermal*. In this case, equilibrium state ceases to exist for a system in contact with a heat bath, but does still exist for a system surrounded by an adiabatic wall. However, for the polytrope index  $n > 5$ , the specific heat of the system turns to increase beyond this inflection point and it next reaches at the point  $d\lambda/d\eta = 0$ , i.e.,  $C_v = 0$ . This means that while the inner part of the system still keeps the specific heat negative, the fraction of the outer normal part grows up as increasing  $r_e$  and it eventually balances with inner gravothermal part. Thus, beyond this critical point, no thermal balance is attainable and the system becomes gravothermally unstable. This is true even in the system surrounded by an adiabatic wall.

Now, let us write down the explicit expression for the specific heat  $C_v$ . In equation (38), the variation of  $\beta$  and  $E$  with  $\xi_e$  can be respectively rewritten with

$$\left(\frac{dE}{d\xi}\right)_e = -\frac{GM^2}{r_e} \frac{d\lambda}{d\xi_e}, \quad (39)$$

and

$$\left(\frac{d\beta}{d\xi}\right)_e = \frac{n-1}{n-3/2} \frac{\beta}{\eta} \frac{d\eta}{d\xi_e}. \quad (40)$$

Here, the variable  $d\lambda/d\xi_e$  has been already given in (27). As for the derivative of  $\eta$  with respect to  $\xi_e$ , we obtain

$$\xi_e \frac{d\eta}{d\xi_e} = \left(u_e - \frac{n-3}{n-1}\right) \eta. \quad (41)$$

Then the quantity  $C_v$  becomes

$$C_v = \frac{(n-3/2)(n-2)}{(n-1)(n-5)} \frac{\beta GM^2}{r_e} \frac{g(u_e, v_e)}{2v_e \left(u_e - \frac{n-3}{n-1}\right)},$$

with the function  $g(u_e, v_e)$  given by (28). Notice that the above expression is still redundant, since there remains the explicit dependence of the variable  $\beta$ . Eliminating the variable  $\beta$  by using the relation (31), one finally obtains

$$\frac{C_v}{N} = \tilde{\alpha}_n \left(\frac{h^2}{GM r_e}\right)^{(3/2)/(n-3/2)} \frac{g(u_e, v_e)}{2 \left(u_e - \frac{n-3}{n-1}\right)} (u_e v_e^{3/2})^{1/(n-3/2)}, \quad (42)$$

where we introduced the new dimensionless constant  $\tilde{\alpha}_n$ :

$$\tilde{\alpha}_n \equiv \frac{(n-3/2)(n-2)}{(n-1)(n-5)} \alpha_n^{1/(n-3/2)}. \quad (43)$$

Note that in the limit  $n \rightarrow +\infty$ , equation (42) consistently recovers the well-known result of isothermal sphere (e.g, eq.(39) of ref.[20]):

$$\frac{C_v}{N} \xrightarrow{n \rightarrow +\infty} \frac{4u_e^2 + 2u_e v_e - 11u_e + 3}{2(u_e - 1)}. \quad (44)$$

Comparing (42) with the isothermal limit, the resultant expression contains a residual dimensional parameter  $h$ , as well as the quantities  $M$  and  $r_e$ . While the residual dependence can be regarded as a natural consequence of the non-extensive generalization of the entropy, it would be helpful to understand the origin of this scaling in more simplified manner. This will be discussed in section 5.

Apart from the residual factor, the expression of specific heat (42) clearly reveals the two types of thermodynamic instability seen in Figure 2. The inflection point with the infinite specific heat,  $C_v \rightarrow \pm\infty$  leads to the condition

$$u_e - \frac{n-3}{n-1} = 0, \quad (45)$$

which immediately yields the conclusion that this is only possible for the polytrope index  $n > 3$ , consistent with Figure 2. On the other hand, critical point with the vanishing specific heat,  $C_v = 0$  corresponds to the following condition:

$$g(u_e, v_e) = 0. \quad (46)$$

This is exactly the same condition as obtained from the second variation of entropy (see eq.(33) or (53) in ref.[2]). According to the previous analysis, the condition (46) represents the marginal stability at which the extremum state of the entropy  $S_q$  is neither maximum nor minimum. This situation turns out to appear when the polytrope index  $n > 5$ .

Therefore, we reach a fully satisfactory conclusion that the thermodynamic instability found from the second variation of entropy is intimately related to the presence of negative specific heat and the stability/instability criterion can be exactly recovered from the critical point of the thermal balance,  $C_v = 0$ , which is also consistent with the analysis in the Boltzmann-Gibbs limit,  $n \rightarrow \infty$  [7]. The successful result can be regarded as an outcome of the correct definition of  $T_{\text{phys}}$ . As for the transition point with  $C_v \rightarrow \pm\infty$ , it clearly indicates the thermodynamic instability of a system in contact with a thermal bath. In next section, by means of the free energy, we confirm that the condition (45) indeed represents the marginal stability of the system surrounded by a thermal wall and beyond this point the system will be unstable.

In Figure 3, by varying the radius  $r_e$ , the normalized specific heat per particle  $C_v^*/N$  is plotted as a function of density contrast,  $\rho_c/\rho_e$  around the critical polytrope indices  $n = 3$  (*upper-panels*) and  $n = 5$  (*middle-panels*). Here, the normalized specific heat  $C_v^*$  is defined by the specific heat  $C_v$  divided by the redundant factor  $(h^2/GMr_e)^{(3/2)/(n-3/2)}$ . Obviously, the transition point  $C_v \rightarrow \pm\infty$  appears when  $n > 3$  (*crosses*), while the existence of critical point  $C_v = 0$  is allowed for higher density contrast of  $n > 5$  cases (*arrows*). The critical

values  $D_{\text{crit}} \equiv (\rho_c/\rho_e)_{\text{crit}}$  indicated by arrows exactly coincide with those obtained from the previous analysis (see Table 1 of ref.[2]). Lower-panels of Figure 3 show the specific heat with large polytrope indices  $n = 10$  and  $30$ , together with the Boltzmann-Gibbs limit ( $n \rightarrow +\infty$ , labeled by *iso*). As increasing the polytrope index  $n$ , the critical/transition points tend to shift to the lower density contrast, while the successive divergent and zero-crossing points appear at the higher density contrast, corresponding to the behavior seen in Figure 2.

## 4 Thermodynamic instability from the second variation of free energy

Previous section reveals that there exists another type of thermodynamic instability in which the marginal stability is deduced from the condition (45). In this section, to check the consistency of the non-extensive thermostatics, we reconsider this issue by means of the Helmholtz free energy:

$$F_q = E - T_{\text{phys}} S_q. \quad (47)$$

Adopting the relation  $T_{\text{phys}} = 1/\beta$ , we re-derive the marginal stability condition (45) from the second variation of  $F_q$ .

Consider a system surrounded by the thermally conducting wall in contact with a heat bath. Usually, the stable equilibrium state should keep the free energy  $F_q$  minimum. Thus the presence of thermodynamic instability implies the absence of minimum free energy, which can be deduced from the signature of the second variation  $\delta^2 F_q$  around the extremum state of free energy. Since the non-extensive formalism still verifies the Legendre transform structure leading to the standard result of thermodynamic relation[14, 15], the extremum state of the free energy exactly coincides with that of the entropy. One thus skips to find the extremum state of  $F_q$  and proceeds to evaluate the second order variation.

In contrast to the adiabatic treatment, we here deal with the density perturbation  $\rho \rightarrow \rho + \delta\rho$ , surrounded by a thermal wall. To be specific, we evaluate the second variation under keeping the radius  $r_e$ , the total mass  $M$  and the temperature  $T_{\text{phys}}$  constant. Then the variation of energy up to the second order leads to

$$\begin{aligned} \delta E &= \delta \left[ \int \left\{ \frac{3}{2} P(x) + \frac{1}{2} \rho(x) \Phi(x) \right\} d^3 \mathbf{x} \right], \\ &= \int \left\{ \frac{3}{2} \delta P + \frac{1}{2} (\delta \rho \Phi + \rho \delta \Phi) + \frac{1}{2} \delta \rho \delta \Phi \right\} d^3 \mathbf{x}. \end{aligned} \quad (48)$$

Similarly, using the expression (70) in Appendix A.2, the variation of Tsallis entropy becomes

$$\begin{aligned} \delta S_q &= \delta \left[ \left( n - \frac{3}{2} \right) \left\{ N - \beta \int P(x) d^3 \mathbf{x} \right\} \right], \\ &= - \left( n - \frac{3}{2} \right) \beta \int \delta P(x) d^3 \mathbf{x}. \end{aligned} \quad (49)$$

The above expressions include the variation of pressure  $\delta P$ , which can be expanded with

a help of the polytropic equation of state (11):

$$\delta P = \left(1 + \frac{1}{n}\right) \frac{P}{\rho} \delta \rho + \frac{1}{2} \left(1 + \frac{1}{n}\right) \frac{1}{n} \frac{P}{\rho^2} (\delta \rho)^2. \quad (50)$$

Combining the above result with equations (48) and (49) and collecting the second order terms only, the second variation of free energy becomes

$$\delta^2 F_q = \delta^2 E - T_{\text{phys}} \delta^2 S_q = \frac{1}{2} \int \left\{ \frac{n+1}{n} \frac{P}{\rho^2} (\delta \rho)^2 + \delta \rho \delta \Phi \right\} d^3 \mathbf{x}, \quad (51)$$

where the relation  $T_{\text{phys}} = 1/\beta$  is used in the last line. Now, restricting our attention to the spherical symmetric perturbation, we introduce the following perturbed quantity (see refs. [2][8]):

$$\delta \rho(r) = \frac{1}{4\pi r^2} \frac{dQ(r)}{dr}. \quad (52)$$

Then the mass conservation  $\delta M = 0$  implies the boundary condition  $Q(0) = Q(r_e) = 0$ . Substituting (52) into (51) and repeating the integration by part, one finally reaches the following quadratic form:

$$\delta^2 F_q = -\frac{1}{2} \int_0^{r_e} dr Q(r) \left[ \frac{n+1}{n} \frac{d}{dr} \left\{ \frac{1}{4\pi r^2 \rho} \left( \frac{P}{\rho} \right) \frac{d}{dr} \right\} + \frac{G}{r^2} \right] Q(r). \quad (53)$$

Thus, the problem just reduces to the eigenvalue problem and the stability of the system can be deduced from the signature of the eigenvalue. More specifically, the onset of instability corresponds to the marginally stability condition,  $\delta^2 F_q = 0$ , and it is sufficient to analyze the zero-eigenvalue equation:

$$\hat{L} Q(r) \equiv \left[ \frac{d}{dr} \left\{ \frac{1}{4\pi r^2 \rho} \left( \frac{P}{\rho} \right) \frac{d}{dr} \right\} + \frac{n}{n+1} \frac{G}{r^2} \right] Q(r) = 0, \quad (54)$$

with the boundary condition,  $Q(0) = Q(r_e) = 0$ . Equation (54) has quite similar form to the zero-eigenvalue equation found in the adiabatic treatment (see eq.(46) of ref.[2]). Except for the non-local term, one can utilize the previous knowledge to solve the equation (54):

$$\hat{L} (4\pi r^3 \rho) = \frac{n-3}{n+1} \frac{G m(r)}{r^2}, \quad \hat{L} m(r) = \frac{n-1}{n+1} \frac{G m(r)}{r^2}. \quad (55)$$

These two equation leads to the ansatz of the solution:

$$Q(r) = c \left\{ 4\pi r^3 \rho(r) - \frac{n-3}{n-1} m(r) \right\}. \quad (56)$$

Here, the variable  $c$  is an arbitrary constant. The above equation (56) automatically satisfies the boundary condition  $Q(0) = 0$ , while the remaining condition  $Q(r_e) = 0$  puts the following constraint:

$$Q(r_e) = c \left( 4\pi r_e^3 \rho_e - \frac{n-3}{n-1} M \right) = c \left( u_e - \frac{n-3}{n-1} \right) M = 0. \quad (57)$$

Again, we arrive at the satisfactory result that the solution of zero-eigenvalue equation exactly recovers the condition (45).

Now, remaining task is to show that the second variation  $\delta^2 F_q$  becomes negative beyond the transition point of  $C_v \rightarrow \pm\infty$ . One can rewrite the expression (53) with

$$\delta^2 F_q = \frac{1}{2} (H - 1) \int_0^{r_e} \frac{GQ^2}{r^2} dr,$$

with the constant  $H$  given by

$$H \equiv \frac{\frac{n+1}{n} \int_0^{r_e} \frac{1}{4\pi r^2 \rho} \left(\frac{P}{\rho}\right) \left(\frac{dQ}{dr}\right)^2 dr}{\int_0^{r_e} \frac{GQ^2}{r^2} dr}. \quad (58)$$

That is, the condition  $H > 1$  implies stable local minimum state of free energy, while the inequality  $H < 1$  represents unstable local maximum state. Integrating by part, equation (58) can be regarded as an eigenvalue equation with eigenvalue,  $H$ :

$$-\frac{d}{dr} \left\{ \frac{1}{4\pi r^2 \rho} \left(\frac{P}{\rho}\right) \frac{dQ}{dr} \right\} = H \frac{n}{n+1} \frac{GQ}{r^2}. \quad (59)$$

Obviously, equation (56) becomes the solution of above equation with the minimum eigenvalue,  $H_{\min} = 1$ , if the condition (57) is fulfilled. In this case, solution (56) can be regarded as the ground state of the eigensystem (59), since the function (56) does not possess any nodes between  $[0, r_e]$ . Therefore, for a suitably smaller radius  $r_e$  or a smaller density contrast  $\rho_e/\rho_c$  below the transition point, the eigenvalue  $H$  should be larger than unity. Conversely, from continuity, the condition  $H < 1$  must be satisfied beyond the critical radius.

Finally, using the  $(u, v)$ -variables, the geometrical meaning of onset of thermodynamic instability is briefly discussed in similar manner to the adiabatic case. In Figure 4, the thick solid lines show the Emden trajectories with various polytrope indices in  $(u, v)$ -plane. The thin-solid lines in Figure 4 represents the straight lines,  $u - (n-3)/(n-1) = 0$ . Since the equilibrium state only exists along the Emden trajectory, the condition (57) is satisfied at the intersection of these two solid lines, which is only possible for  $n > 3$ . On the other hand, as seen in previous section, the equilibrium system surrounded by a thermal wall is characterized by the three parameters,  $r_e$ ,  $M$  and  $\beta$  (or  $T_{\text{phys}}$ ), through the relation (31). In other words, the system must lie on the curve:

$$v = \left( \frac{\eta}{\alpha_n} \right)^{(n-1)/n} u^{-1/n}, \quad (60)$$

with some constant value  $\eta$ . We have seen in Figure 2 that the constant value  $\eta$  is bounded from above,  $\eta \leq \eta_{\text{crit}}$ . Thus, the critical curve (60) with  $\eta = \eta_{\text{crit}}$  must intersect with both the Emden trajectory and the straight line  $u - (n-3)/(n-1) = 0$  simultaneously. This is clearly shown in Figure 4, where the critical curve is plotted as dashed lines. Since the critical curves tangentially intersect with Emden solutions, it always satisfies the condition  $d\eta/d\xi = 0$  at the contact point, leading to the condition (45) consistently.

Table 1 summarizes the dimensionless quantities  $\eta_{\text{crit}}$  and  $D_{\text{crit}} \equiv (\rho_c/\rho_e)_{\text{crit}}$  evaluated at the contact point. As increasing the polytrope index  $n$ , these values asymptotically approach the well-known results of Boltzmann-Gibbs limit,  $\eta_{\text{crit}} \rightarrow 2.52$  and  $D_{\text{crit}} \rightarrow 32.1$ .

## 5 Origin of non-extensive nature in stellar polytrope

As has been mentioned in section 3.2, specific heat of the stellar polytropic system explicitly depends on the residual dimensional parameter  $h$ , in contrast to the isothermal limit (44). In this section, to contact the physical meaning of the non-extensivity in stellar polytrope, we discuss the origin of this residual dependence. Indeed, the appearance of the residual factor can be recognized as the breakdown of both the intensivity of temperature and the extensivity of energy and entropy as follows. From equation (18), the asymptotic behavior of the Emden solution becomes

$$\theta \sim \xi^{-2/(n-1)}, \quad \rho \sim r^{-2n/(n-1)}, \quad (\xi, r \rightarrow \infty)$$

so that the mass within a sphere of radius  $r$  is given by

$$M \sim \rho r^3 \propto r_e^{(n-3)/(n-1)}. \quad (61)$$

Then the energy of a virialized stellar system is roughly estimated as

$$E \sim \frac{GM^2}{r_e} \propto r_e^{(n-5)/(n-1)} \propto M^{(n-5)/(n-3)},$$

and the relation (31) tells

$$\beta \propto r_e^{-(n-3)/(n-1)/(n-3/2)} \propto M^{-1/(n-3/2)}.$$

These relations clearly show the breakdown of the intensivity of temperature and the extensivity of energy, which lead to the scaling of the specific heat per mass:

$$\frac{C_V}{N} = \frac{1}{M} \frac{dE}{dT_{\text{phys}}} \sim \frac{\beta E}{M} \propto M^{-3(n-2)/(n-3)/(n-3/2)}. \quad (62)$$

On the other hand, the dimensionless combination  $h^2/(GM r_e)$  represents the ratio of a typical scale of the stellar system,  $GM r \sim (GM/r)r^2 \sim v^2 r^2$ , to that of the reference cell,  $h = v_0 l_0$ . This behaves as

$$\frac{h^2}{GM r_e} \propto \frac{1}{M r_e} \propto M^{2(n-2)/(n-3)}. \quad (63)$$

Thus, these two equations (62) and (63) lead to the scaling relation of (42):

$$\frac{C_V}{N} \sim \left( \frac{h^2}{GM r_e} \right)^{(3/2)/(n-3/2)}. \quad (64)$$

Notice that the Clausius relation (37) suggests that the entropy per unit mass has the same scaling relation:

$$\frac{S_q}{M} \sim \frac{\beta E}{M} \sim \frac{C_V}{N},$$

Therefore, resultant dependence (64) for the stellar polytrope can be a natural outcome of the non-extensivity of the entropy.

In fact, framework of the thermostatics generally requires an introduction of the scale of the unit cell in order to count the available number of states in phase spaces. This is even true in the case of the isothermal stellar system ( $n \rightarrow +\infty$  or  $q \rightarrow 1$ ), but, the thermodynamic quantities show somewhat peculiar dependence of the scale  $h$ . A typical example is the entropy:

$$S_{\text{BG}} = \frac{M}{m_0} \left\{ \left( 2u_e + v_e - \frac{9}{2} \right) - \ln \left( \frac{u_e v_e^{3/2}}{4\pi} \right) - \frac{3}{2} \ln \left( \frac{h^2}{2\pi G M r_e} \right) \right\},$$

where  $u_e$  and  $v_e$  are the homology invariants for the isothermal system. The above equation shows that in the Boltzmann-Gibbs limit,  $h$ -dependence of the entropy can be recognized as a matter of choice of an additive constant, so that its derivatives, e.g., specific heat, is free from the residual dependence.

It should be emphasized that the stellar equilibrium system recovers the extensivity in the limit  $n \rightarrow \infty$  and it behaves as

$$E \sim M \sim r, \quad C_V \sim M. \quad (65)$$

Also, the temperature becomes intensive in this limit. Thus, we readily understand that the scaling behavior shown in (42) or (64) has nothing to do with the long-range nature of the gravity. Even in the free polytropic gas model in Appendix B, the residual dependence emerges as

$$\frac{C_V}{N} \sim \left\{ \frac{h^2}{(P/\rho)V^{2/3}} \right\}^{(3/2)/(n-3/2)}.$$

It follows that the explicit dependence of the specific heat on the reference cell scale  $h$  just originates from the non-extensive nature of Tsallis entropy.

## 6 Summary

In this article, thermodynamic properties of the stellar self-gravitating system arising from Tsallis' non-extensive entropy have been studied in detail. In particular, physical interpretation of the thermodynamic instability previously found from the second variation of entropy is discussed in detail within a framework of the non-extensive thermostatics. After briefly reviewing the equilibrium state of Tsallis entropy, we first address the issues on thermodynamic temperature in the case of equilibrium stellar polytrope. Analyzing the heat transfer and the entropy change in a quasi-static process, standard form of the Clausius relation is derived, irrespective of the non-extensivity of entropy. According to this result, we explicitly calculate the specific heat and confirm the presence of negative specific heat. The onset of instability found in previous work just corresponds to the zero-crossing point,  $C_V = 0$ , supporting the fact that the heuristic explanation of gravothermal catastrophe holds even in the non-extensive thermostatics.

Further, the analysis of specific heat shows divergent behavior at  $n > 3$ , suggesting another type of thermodynamic instability, which occurs when the system is surrounded by a thermal wall. We then turn to the stability analysis by means of the Helmholtz free energy. Similar to the previous early work, the stability/instability criterion just reduces



to the solution of the zero-eigenvalue problem and solving the eigenvalue equation, we recover the marginal stability condition derived from the divergence of specific heat (45).

In addition to the thermostatic treatment, we have also discussed the origin of non-extensivity in stellar polytrope. The residual dependence of the reference scale  $h$  appeared in the specific heat (42) naturally arises from the non-extensivity of the entropy and the resultant scaling dependence can be simply deduced from the asymptotic behavior of the Emden solutions.

The stability analysis using the free energy in section 4 is consistent with recent claim by Chavanis [21], who has investigated the dynamical instability of polytropic gas sphere. According to his early paper [20], the thermodynamic stability of stellar system is intimately related to the dynamical stability of gaseous system, which has been clearly shown in the case of the isothermal distribution. Thus, the correspondence between Chavanis' recent result [21] and a part of our present analysis can be regarded as a generalization of his early work to the polytropic system. Note, however, that starting from the Tsallis entropy, we extensively discuss the thermodynamic temperature and the specific heat of stellar polytrope. Therefore, at least, from the thermodynamic point of view, our present analysis provides a valuable insight to the stellar equilibrium systems.

At present, the results shown in this article seems fully consistent with the general framework of the thermostatics. Apart from the thermodynamic instability, the stellar polytropic system can be a plausible thermodynamic equilibrium state, as well as the isothermal stellar distribution. In the isothermal case, existence of the thermodynamic limit has been discussed by de Vega and Sánchez [19]:

$$M, V \rightarrow \infty, \quad \frac{M}{V^{1/3}} = \text{fixed},$$

where  $V \sim r^3$  is a volume of the system. Recalling the discussion in section 5, the above condition merely reflects the extensivity of the isothermal system (65). Thus, similar argument can hold for the non-extensive system. According to the scaling relation (61), the existence of the thermodynamic limit in stellar polytrope yields the condition:

$$M, V \rightarrow \infty, \quad \frac{M}{V^{(n-3)/(3n-3)}} = \text{fixed}.$$

Note, however, that this discussion relies on the non-uniqueness of the Boltzmann-Gibbs theory, which can be proven only mathematically[22]. Indeed, framework of the thermostatics cannot answer the question whether the stellar polytropic distribution is really achieved as a thermodynamic equilibrium. To address this issue, we must study the detailed process of the long-term stellar dynamical evolution. In the light of this, the analysis using Fokker-Planck model or direct N-body simulation can provide an invaluable insight to the non-extensive nature of stellar gravitating systems. This issue is now in progress and will be presented elsewhere.

## Appendix A: Thermodynamic variables in a stellar polytropic system

In this appendix, using the equilibrium state of stellar polytrope described in section 2, we explicitly evaluate the thermodynamic variables, which have been used in section 3 and 4.

### A.1 Energy

Recall that the equilibrium system confined in a spherical container satisfies the following virial theorem (e.g, p.502 of Ref.[3]):

$$2K + U = 4\pi r_e^3 P_e.$$

The energy (1) is then expressed as

$$E = K + U = 4\pi r_e^3 P_e - K = 4\pi r_e^3 P_e - \frac{3}{2} \int_0^{r_e} P(r) 4\pi r^2 dr. \quad (66)$$

To evaluate the above integral in the spherically symmetric case, we use the following integral formula:

$$\int_0^{r_e} P(r) 4\pi r^2 dr = -\frac{1}{n-5} \left\{ 8\pi r_e^3 P_e - (n+1) \frac{M P_e}{\rho_e} + \frac{G M^2}{r_e} \right\}, \quad (67)$$

which can be derived from the conditions of hydrostatic equilibrium, (15) and (16) (see Appendix A of ref.[2]). Thus, the energy of extremum state becomes

$$E = \frac{1}{n-5} \left[ \frac{3}{2} \left\{ \frac{G M^2}{r_e} - (n+1) \frac{M P_e}{\rho_e} \right\} + (n-2) 4\pi r_e^3 P_e \right]. \quad (68)$$

In terms of the homology invariants, we obtain

$$E = \frac{1}{n-5} \frac{G M^2}{r_e} \left[ \frac{3}{2} \left\{ 1 - (n+1) \frac{1}{v_e} \right\} + (n-2) \frac{u_e}{v_e} \right]. \quad (69)$$

### A.2 Entropy

First note the definition of Tsallis entropy (5):

$$S_q = - \left( n - \frac{3}{2} \right) \left\{ \int N \left( \frac{f}{N} \right)^{(n-1/2)/(n-3/2)} d^6 \tau - N \right\}.$$

Substituting the distribution function (7) into the above equation, after some manipulation, we obtain

$$S_q = - \left( n - \frac{3}{2} \right) \left\{ \beta \int P(x) d^3 \mathbf{x} - N \right\}. \quad (70)$$

Thus, the substitution of integral formula (67) immediately leads to

$$S_q = \left(n - \frac{3}{2}\right) \left[ \frac{1}{n-5} \left\{ 8\pi r_e^3 P_e - (n+1) \frac{MP_e}{\rho_e} + \frac{GM^2}{r_e} \right\} \beta + N \right],$$

which can be expressed in terms of the homology invariants:

$$S_q = \left(n - \frac{3}{2}\right) \left[ \frac{1}{n-5} \frac{\beta GM^2}{r_e} \left\{ 2 \frac{u_e}{v_e} - (n+1) \frac{1}{v_e} + 1 \right\} + N \right]. \quad (71)$$

### A.3 Radius-mass-temperature relation

The mass-radius-temperature relation (31) is derived from the equilibrium stellar polytropic configuration. Using (15), we first write down the condition of hydrostatic equilibrium at the boundary  $r_e$ :

$$\frac{GM}{r_e^2} = -\frac{1}{\rho_e} \left( \frac{dP}{dr} \right)_e.$$

The right-hand-side of this equation is rewritten with dimensionless quantities in (17):

$$\frac{GM}{r_e^2} = -(n+1) K_n \rho_c^{1/n} \left( \frac{\xi_e}{r_e} \right) \theta'_e. \quad (72)$$

We wish to express the above equation only in terms of the variables at the edge. To do this, we eliminate the residual dependences,  $\rho_c$  and  $K_n$  from (72). The definition (17) leads to

$$\frac{\xi_e}{r_e} = \left\{ \frac{4\pi G \rho_c^2}{(n+1)P_c} \right\}^{1/2} = \left\{ \frac{4\pi G}{(n+1)K_n} \right\}^{1/2} \rho_c^{(n-1)/(2n)},$$

which can be rewritten with

$$\rho_c^{1/n} = \left\{ \frac{4\pi G}{(n+1)K_n} \right\}^{1/(n-1)} \left( \frac{\xi_e}{r_e} \right)^{2/(n-1)}.$$

Substituting the above relation into (72), the  $\rho_c$ -dependence is first eliminated and we obtain

$$\frac{G^{n/(n-1)} M}{r_e^{(n-3)/(n-1)}} = - \left[ \frac{\{(n+1)K_n\}^n}{4\pi} \right]^{1/(n-1)} \xi_e^{(n+1)/(n-1)} \theta'_e. \quad (73)$$

As for  $K_n$ -dependence, the definition (13) together with (8) yields

$$(n+1) K_n = \left\{ 4\sqrt{2}\pi \frac{B(3/2, n-1/2)}{(n-1)^{n-3/2}} \frac{M}{h^3} \right\}^{-1/n} (m_0 \beta)^{-(n-3/2)/n}. \quad (74)$$

Hence, substituting the above expression into (73), the relation between mass  $M$ , radius  $r_e$  and Lagrange multiplier  $\beta$  can be finally obtained. In terms of the homology invariants, it follows that

$$\left\{ \frac{(GM)^n (m_0 \beta)^{n-3/2}}{r_e^{n-3} h^3} \right\}^{1/(n-1)} = \alpha_n (u_e v_e^n)^{1/(n-1)}, \quad (75)$$

where the constant  $\alpha_n$  is given by

$$\alpha_n \equiv \left\{ \frac{(n-1/2)^{n-3/2}}{16\sqrt{2}\pi^2 (n+1)^n B(3/2, n-1/2)} \right\}^{1/(n-1)},$$

which asymptotically approaches unity in the limit  $n \rightarrow \infty$ .

## Appendix B: Thermodynamic temperature of classical gas model from the Carnot cycle

In a standard framework of thermodynamics, the temperature is defined by means of an efficiency of the Carnot cycle. Here we apply the standard procedure to seek the physical temperature  $T_{\text{phys}}$  for so-called polytropic system of which distribution function is given by the extremization of the Tsallis entropy (see eqs.(5)(6)). For simplicity, we discuss a case of the free classical gas without gravity, which corresponds to the  $G \rightarrow 0$  limit of the *stellar polytropic system*.

From the  $G \rightarrow 0$  limit of the formula (68), free polytropic system of the volume  $V$  with *homogeneous* pressure  $P$  and density  $\rho$  has an (internal) energy:

$$E = K = \frac{3}{2} PV = \frac{3}{2} \frac{MP}{\rho}. \quad (76)$$

Here we drop the subscript  $e$  for the pressure and density, since both are constant within the system in absence of gravity. And equation of state (11) becomes

$$P = K_n \rho^{1+1/n} = K_n \left( \frac{M}{V} \right)^{1+1/n}. \quad (77)$$

From equations (8) and (13), the constant  $K_n$  is related to the Lagrange multiplier  $\beta$  as

$$K_n \propto \beta^{-(n-3/2)/n}, \quad (78)$$

so that this constant can be used as a parameter which characterizes the temperature of the system. However, it is not sure whether  $K_n$  itself has a role of the physical temperature, which should be determined through the efficiency of the Carnot cycle.

The internal energy (76) and the equation of state (77) give the thermodynamic first law:

$$\begin{aligned} d'Q &= dE + P dV \\ &= M^{1+1/n} \left\{ \frac{3}{2} \frac{dK_n}{V^{1/n}} + \left( \frac{n-3/2}{n} \right) K_n \frac{dV}{V^{1+1/n}} \right\}, \end{aligned} \quad (79)$$

from which adiabatic changes  $d'Q = 0$  is expressed as

$$K_n V^{(2/3-1/n)} = \text{constant}, \quad P V^{5/3} = \text{constant}'. \quad (80)$$

Note that adiabatic lines in a  $P$ - $V$  plane become steeper than isothermal ones when  $n > 3/2$ .

Now, let us consider the Carnot cycle shown in Figure 5. As usual, quasi-static changes  $B \rightarrow C$  and  $D \rightarrow A$  are adiabatic. As for the process  $A \rightarrow B$ , the system is in a thermal contact with a heat bath which has a higher temperature  $K_n^H$ . Similarly, during the change  $C \rightarrow D$ , the system lies in a thermal equilibrium with another heat bath that has a lower temperature  $K_n^L$ . The system absorbs amount of heat  $Q^H$  from the higher temperature bath and disposes  $Q^L$  to the lower one during the isothermal processes  $A \rightarrow B$  and  $C \rightarrow D$ , respectively. They are easily evaluated from (79):

$$\begin{aligned} Q^H &= (n - \frac{3}{2}) M^{1+1/n} K_n^H \left( V_A^{-1/n} - V_B^{-1/n} \right), \\ Q^L &= (n - \frac{3}{2}) M^{1+1/n} K_n^L \left( V_D^{-1/n} - V_C^{-1/n} \right). \end{aligned} \quad (81)$$

On the other hand, a relation between the parameters of the cycle can be obtained from the equation of state (77) and the adiabatic changes (80):

$$\left(\frac{K_n^H}{K_n^L}\right)^\gamma = \frac{V_C}{V_B} = \frac{V_D}{V_A}; \quad \gamma = \frac{3}{2} \frac{n}{n-3/2}. \quad (82)$$

Thus, equations (81) and (82) lead to the following efficiency of the Carnot cycle:

$$\eta \equiv 1 - \frac{Q^L}{Q^H} = 1 - \left(\frac{K_n^L}{K_n^H}\right)^{n/(n-3/2)} = 1 - \frac{\beta^H}{\beta^L}, \quad (83)$$

where we used the relation (78) in the last line. This clearly shows that the inverse of the Lagrange multiplier  $\beta$  has a role of the physical temperature.

## References

- [1] C. Tsallis, J.Stat.Phys. 52 (1988) 479.
- [2] A. Taruya, M. Sakagami, Physica A 307 (2002) 185.
- [3] J. Binney, S. Tremaine, *Galactic Dynamics* (Princeton Univ. Press, Princeton, 1987).
- [4] R. Elson, P. Hut and S. Inagaki, Ann. Rev. Astron. Astrophys. 25 (1987) 565.
- [5] G. Meylan, D.C. Heggie, Astron.Astrophys.Rev. 8 (1997) 1.
- [6] V.A. Antonov, *Vest. Leningrad Gros. Univ.*, 7 (1962) 135 (English transl. in *IAU Symposium 113, Dynamics of Globular Clusters*, ed. J. Goodman and P. Hut [Dordrecht: Reidel], pp. 525–540 [1985])
- [7] D. Lynden-Bell, R. Wood, Mon.Not.R.Astr.Soc. 138 (1968) 495.
- [8] T. Padmanabhan, Astrophys.J.Suppl. 71 (1989) 651.
- [9] T. Padmanabhan, Phys.Rep. 188 (1990) 285. 651.
- [10] A.R. Plastino, A. Plastino, Phys.Lett. A 174 (1993) 384.
- [11] A.R. Plastino, A. Plastino, Braz. J. Phys. 29 (1999) 79.
- [12] S. Martínez, F. Nicolás, F. Pennini, A. Plastino, Physica A 286 (2000) 489.
- [13] S. Abe, S. Martínez, F. Pennini, A. Plastino, Phys.Lett. A 281 (2001) 126.
- [14] E.M.F. Curado, C. Tsallis, J.Phys.A 24 (1991) L69.
- [15] A. Plastino, A.R. Plastino, Phys.Lett. A 226 (1997) 257.
- [16] C. Tsallis, R.S. Mendes, A.R. Plastino, Physica A 261 (1998) 534.
- [17] S. Chandrasekhar, *Introduction to the Study of Stellar Structure* (New York, Dover, 1939)
- [18] R. Kippenhahn, A. Weigert, *Stellar Structure and Evolution* (Springer, Berlin, 1990)
- [19] H.J. de Vega, N. Sánchez, Nucl.Phys. B 625 (2002) 409.
- [20] P.H. Chavanis, Astron. & Astrophys. 381 (2002) 340.
- [21] P.H. Chavanis, astro-ph/0108378.
- [22] S. Abe, A.K. Rajagopal, Phys.Lett. A 272 (2000) 345; J. Phys. A 33 (2000) 8733; Europhys. Lett. 52 (2000) 610.

Table 1: Critical values of the radius-mass-temperature relation,  $\eta_{\text{crit}}$  and the density contrast between center and edge,  $D_{\text{crit}} = (\rho_c/\rho_e)_{\text{crit}}$  in the case of a system in contact with a heat bath for given polytrope index  $n$  or  $q$ .

n	q	$\eta_{\text{crit}}$	$D_{\text{crit}}$
3	$\frac{5}{3}$	—	—
4	$\frac{7}{5}$	0.9421	153.5
5	$\frac{9}{7}$	1.193	88.15
6	1.22	1.379	68.38
7	1.18	1.520	58.86
8	1.15	1.631	53.28
9	1.13	1.720	49.62
10	1.12	1.793	47.04
30	1.04	2.263	35.89
50	1.02	2.363	34.28
100	1.01	2.440	33.17
$\infty$	1	2.518	32.13

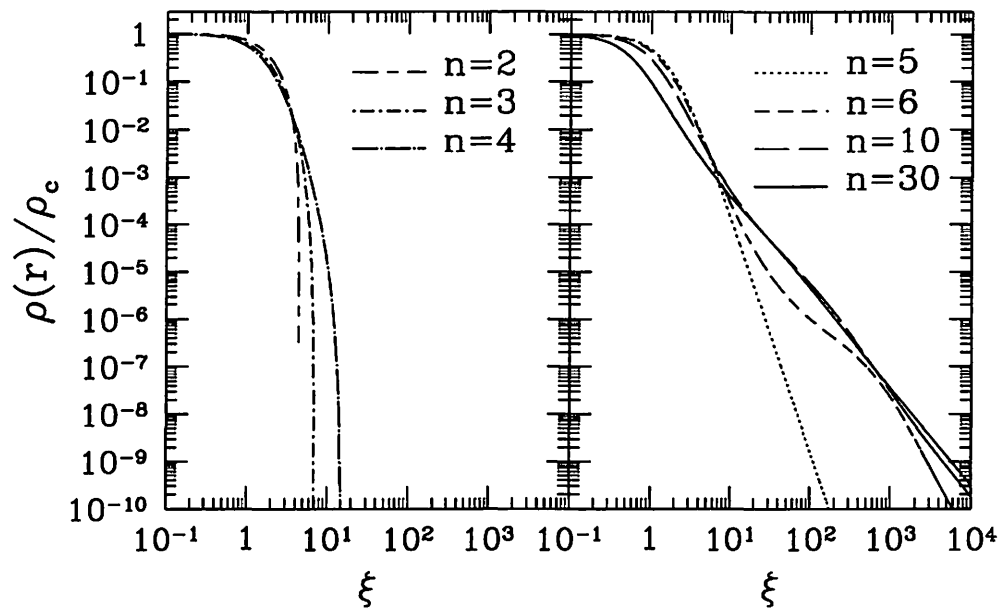


Figure 1: Density profiles of stellar polytrope for  $n < 5$  (left) and  $n \geq 5$  (right).

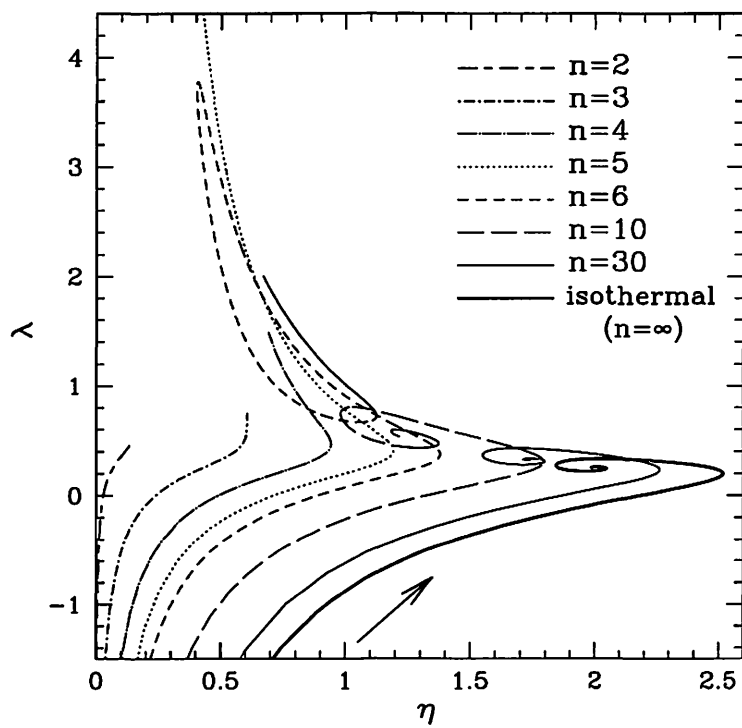


Figure 2: Trajectory of Emden solutions in  $(\eta, \lambda)$ -plane.



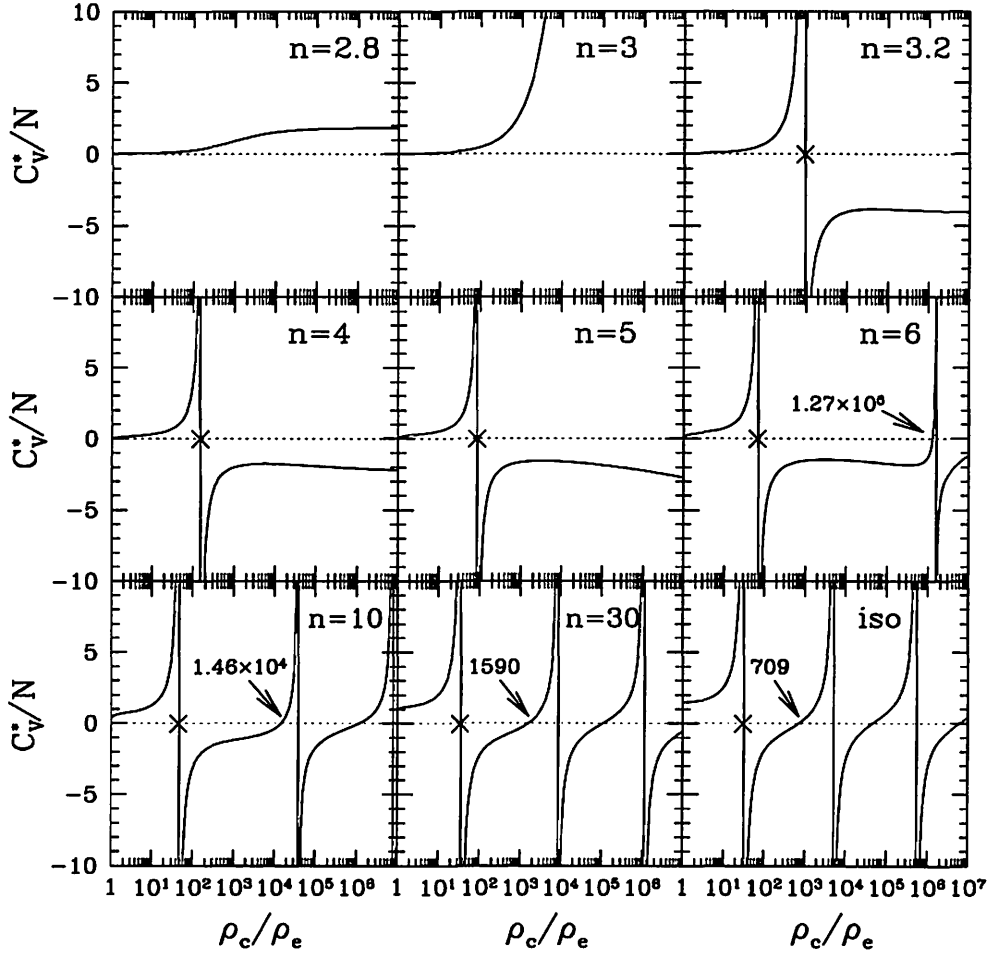


Figure 3: Normalized specific heat per particle  $C_v^*/N$  as a function of density contrast  $\rho_c/\rho_e$  near the critical polytrope indices  $n = 3$ (upper) and  $n = 5$ (middle), and large  $n$  cases(lower). Here, the normalized specific heat  $C_v^*$  is defined by  $C_v/(\hbar^2/GMr_e)^{(3/2)/(n-3/2)}$ .

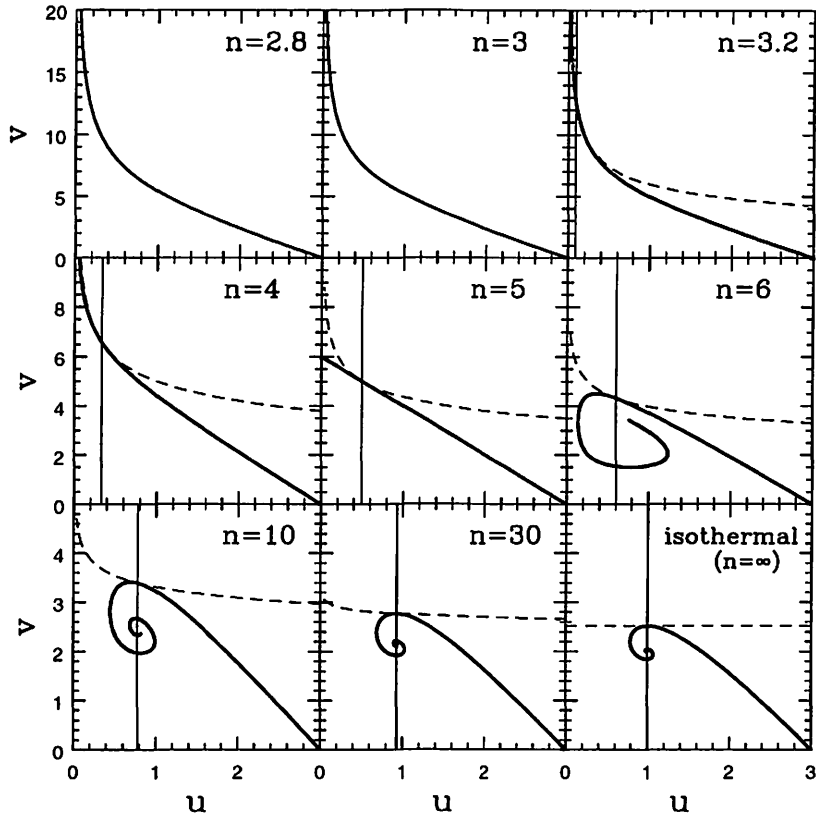


Figure 4: Stability/instability criterion for a system in contact with a thermal bath in  $(u, v)$ -plane. The thick solid lines represent the trajectories of Emden solutions, while the thin-solid and dashed lines respectively denote the conditions (45) and (60).

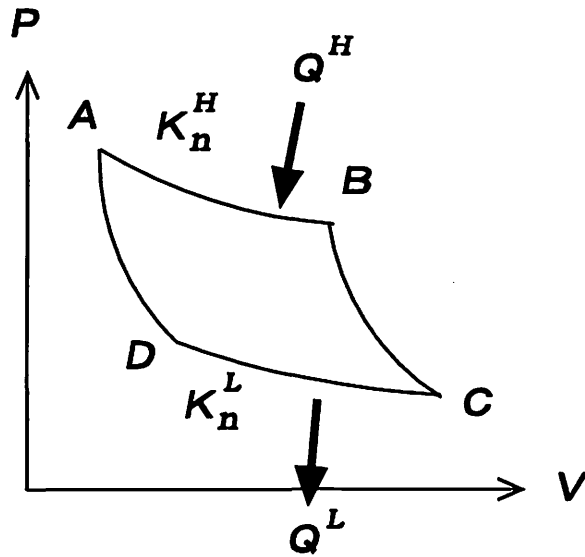


Figure 5: A schematic description of Carnot cycle.

# Stationary state in N-body System with power law interaction

Osamu Iguchi\*

Department of Physics, Ochanomizu University, 2-1-1 Ohtuka, Bunkyo, Tokyo, 112-8610 Japan

Observations and simulations show many scaling properties in self-gravitating system. In order to study the origin of these scaling properties, we consider the stationary state in N-body system with inverse power law interaction. As a simple case, we consider the self-similar stationary solution in the collisionless Boltzmann equation with power law potential and investigate its stability in the term of a linear symplectic perturbation. The stable scaling solution obtained are expressed by the power of the potential and the virial ratio of the initial state. The nonextensive system has many various stable scaling solution compared with the extensive one.

## I. INTRODUCTION

There are many self-gravitating system which are characterized by some scaling properties. For example, the inter stellar medium show that its velocity dispersion  $\sigma$  has scaling relation with the system size  $L$  or mass  $M$  [1] ( $\sigma \sim L^{0.38} \sim M^{0.2}$ ) and isothermal contour are characterized by the fractal dimension  $D \sim 1.36$  [2]. The observations by the Hubble Space Telescope show elliptical galaxies has a power law density distribution  $\rho \sim r^{-n}$  (At outer region,  $n \sim 4$  and at inner region,  $n \sim 0.5 - 1.0$  for the bright elliptical galaxies and  $n \sim 2$  for the faint ones [3]). The distribution of the galaxies and the cluster of galaxies can be characterized by the fractal dimension  $D \sim 2$  [4]. In cosmological simulation based on the standard cold dark matter scenario, the density profile is a power law distribution (At outer region,  $n \sim 3$  and at inner region,  $n \sim 1.0 - 1.5$  [5, 6]).

Recently, in order to study the statistical properties of self-gravitating system, we proposed self-gravitating ring model [7], where each particles are moving, on a circular ring fixed in three-dimensional space, with mutual interaction of gravity in three-dimensional space. The numerical simulation shows that the system at the intermediate energy scale, where the specific heat becomes negative, has some peculiar properties such as non-Gaussian and power law velocity distribution ( $f(v) \sim v^{-2}$ ), the scaling mass distribution, and the self-similar recurrent motion. In this model, the *halo* particles which belong to the intermediate energy scale, play an important role in such specific characters.

We are interested in the origin of these scaling properties. In order to study the statistical properties of long range interaction such as gravity, Ising model, and spin glass, the model with power law potential has been used and revealed anomalous properties [8–10]. For example, a gravitational-like phase transition [8], reduction of mixing [9], and long relaxation [10] are observed. Using a model with an attractive  $1/r^\alpha$  potential in general D-dimensional space, we can control the extensivity of the system and the sign of the specific heat by the spatial

dimension  $D$  and the exponent of inverse power of the potential  $\alpha$ .

In this paper, we study the quasi-equilibrium state of N-body system with a power law potential. As first step, we consider the collisionless Boltzmann equation (CBE) in replace of N-body system and derive the self-similar stationary solution of CBE which has a scaling property appeared in the quasi-equilibrium state and discuss the linear stability by use of energy functional analysis [11–13].

In section II, we show some property in N-body system with power law potential. In section III, we derive the self-similar stationary solution of CBE with an attractive  $1/r^\alpha$  potential assuming spherical symmetry and isotropic orbit case in D-dimensional space. The stability for the linear perturbation around the self-similar stationary solution is investigated in section IV. The discussion is devoted to section V.

## II. N-BODY SYSTEM WITH POWER LAW POTENTIAL

In this section, we show some character in N-body system with power law potential.

We consider the Hamiltonian in N-body system with power law potential can be written in the form.

$$H = \sum_{i=1}^N \frac{P_i^2}{2m} - \sum_{i<j}^N \frac{Gm^2}{r_{ij}^\alpha}, \quad (1)$$

where  $r_{ij} := |r_i - r_j|$  and  $\alpha$  is a parameta characterized potential energy.

In this system, the virial condition is

$$2 \langle K \rangle + \alpha \langle \Phi \rangle = 0, \quad (2)$$

where  $\langle K \rangle$  is a time averaged kinetic energy and  $\langle \Phi \rangle$  is a time averaged potential energy. From the equation  $H = K + \Phi$ , we have

$$H = -\frac{2-\alpha}{\alpha} \langle K \rangle = \frac{2-\alpha}{2} \langle \Phi \rangle. \quad (3)$$

From Eq.(3), the sign of the specific heat is determined by the sign of the term  $-(2-\alpha)/\alpha$ .

\*Electronic address: osamu@phys.ocha.ac.jp

TABLE I: In the system with an attractive  $1/r^\alpha$  potential in  $D$ -dimensional space, the property of the specific heat and the extensivity are shown.

	$0 < \alpha < 2$	$\alpha > 2$
specific heat	negative	positive
	$\alpha < D$	$\alpha > D$
extensivity	nonextensive	extensive

In order to show the extensivity of the system, we use the  $N$  dependence of the potential energy  $\Phi$  under the fixing the number density  $N/L^D$  which can be calculated as follows[15].

$$\frac{\Phi}{N} \sim \int^{N^{1/D}} dr r^{D-1} r^{-\alpha} \sim N^{1-\alpha/D}. \quad (4)$$

If the  $N$  dependence of the potential energy per one particle disappears when  $N$  goes to infinity, we call the system is extensive. Otherwise, we call the system is nonextensive. In the case of the gravity in  $D$ -dimensional space, since  $\alpha = D - 2$ , the system is always nonextensive. We summarize the sign of the specific heat of the system and the extensivity in Table.I.

### III. SELF-SIMILAR STATIONARY SOLUTION IN COLLISIONLESS BOLTZMANN EQUATION (CBE)

In this section, we derive a self-similar stationary solution in collisionless Boltzmann equation (CBE);

$$\frac{df}{dt} = \frac{\partial f}{\partial t} + [f, H] = 0. \quad (5)$$

where  $f = f(\mathbf{x}, \mathbf{v}, t)$  is a mass distribution function and  $[A, B]$  denotes a Poisson bracket.

The stationary solution  $f_0$  satisfies following equation.

$$[f_0, H] = \sum_{i=1}^{2D} \frac{\partial w^i}{\partial t} \frac{\partial f_0}{\partial w^i} = 0, \quad (6)$$

where  $w^i = \{\mathbf{x}, \mathbf{v}\}$ .

For the coupled CBE and Poisson equation, R.N. Henriksen and L.M. Widrow[16] studied the self-similar stationary solution in CBE with spherical symmetry case in three-dimensional space by the systematic method which is based on the work of B. Carter and R.N. Henriksen[17].

Following R.N. Henriksen and L.M. Widrow[16], we study self-similar stationary solution for spherical symmetry and isotropic orbit case in  $D$ -dimensional space. The case that  $D = 3$  and  $\alpha = 1$  corresponds to the work by R.N. Henriksen and L.M. Widrow[16]. By the extension of the spatial dimension  $D$  and the exponent of power of potential  $\alpha$ , it is possible to investigate the

relation between the extensivity of the system and the self-similarity.

From Eq.(6), mass distribution function  $f(r, v)$  obeys

$$v \partial_r f - \partial_r \Phi \partial_v f = 0, \quad (7)$$

where  $v := \sqrt{v_r^2 + v_\theta^2 + v_\phi^2}$  and  $\Phi$  is a potential. The potential  $\Phi$  satisfies a following equation,

$$\frac{1}{r^{D-1}} \partial_r (r^{\alpha+1} \partial_r \Phi) = S_D^2 G \int v^{D-1} f dv, \quad (8)$$

where  $S_D := 2\pi^{D/2}/\Gamma(D/2)$ . In the case of  $\alpha = D - 2$ , the above equation corresponds to Poisson equation.

A self-similar stationary solution satisfies the following equation.

$$\mathcal{L}_k f = 0, \quad (9)$$

where

$$\mathcal{L}_k := k^i \partial_i = \delta r \partial_r + \nu v \partial_v + \mu m \partial_m \quad (10)$$

is a Lie derivative with respect to the vector  $k$  in phase space, and  $\delta, \nu$ , and  $\mu$  are arbitrary constants.

In a dimensional space of length, velocity, and mass, we introduce these vector  $\mathbf{a} = (\delta, \nu, \mu)$  and  $\mathbf{d}_f$ . The vector  $\mathbf{a} = (\delta, \nu, \mu)$  describes changes in the logarithms of dimensional quantities. Each dimensional quantity  $f$  in the problem has its dimension represented by the vector  $\mathbf{d}_f$ . Using these vector  $\mathbf{a}$  and  $\mathbf{d}_f$ , the action of  $k$  reads

$$\mathcal{L}_k f = (\mathbf{d}_f \cdot \mathbf{a}) f. \quad (11)$$

The dimensional quantities in current problem  $f, \Phi$ , and  $G$  have the following dimensional covectors,

$$\begin{aligned} \mathbf{d}_f &= (-D, -D, 1), \\ \mathbf{d}_\Phi &= (0, 2, 0), \\ \mathbf{d}_G &= (\alpha, 2, -1). \end{aligned} \quad (12)$$

The requirement of the invariance of  $G$  under rescaling group (10) implies  $\mathbf{d}_G \cdot \mathbf{a} = 0$ ,

$$\mu = \alpha \delta + 2\nu. \quad (13)$$

The dimensional space is reduced to the subspace of (length, velocity), wherein the rescaling group element  $\mathbf{a} = (\delta, \nu)$  and

$$\begin{aligned} \mathbf{d}_f &= (\alpha - D, 2 - D), \\ \mathbf{d}_\Phi &= (0, 2). \end{aligned} \quad (14)$$

Here we define the new coordinate  $R(r)$  and  $X$  in place of the original coordinate  $r$  and  $v$  such that

$$\mathcal{L}_k R = 1, \quad (15)$$

$$\mathcal{L}_k X = 0. \quad (16)$$

From Eqs.(15) and (16), we choose

$$r|\delta| = e^{\delta R}, \quad (17)$$

$$v = X e^{\nu R}. \quad (18)$$

Under the new coordinate, these physical quantities  $f$  and  $\Phi$  can be written in the form.

$$f(X, R) = \bar{f}(X) e^{-[(D-\alpha)\delta + (D-2)\nu]R}, \quad (19)$$

$$\Phi(X, R) = \bar{\Phi}(X) e^{2\nu R}. \quad (20)$$

Substituting Eqs.(19) and (20) into Eqs.(7) and (8), these equations for a bounded solution yield

$$\frac{d \ln \bar{f}}{d \ln X} = -\frac{[D-2 + (D-\alpha)\frac{\delta}{\nu}] X^2}{X^2 + 2\bar{\Phi}}, \quad (21)$$

$$2\nu^2 |\delta|^{D-\alpha-2} \left[ 2 + \frac{\alpha\delta}{\nu} \right] \bar{\Phi} = S_D^2 G \int_0^{\sqrt{-2\bar{\Phi}}} X^{D-1} \bar{f} dX. \quad (22)$$

Without loss of generality, we can set  $\nu = 1$ .

Solving Eqs.(21) and (22), we have a following solution,

$$\bar{f} = C |X^2 + 2\bar{\Phi}|^{-[(D-\alpha)\delta + (D-2)]/2}, \quad (23)$$

where

$$C = \frac{|2 + \alpha\delta| |\delta|^{D-\alpha-2} \Gamma(D/2) \Gamma(2 + (\alpha-D)\delta/2)}{2\pi^D G | -2\bar{\Phi}|^{(\alpha-D)\delta/2} \Gamma([4-D + (\alpha-D)\delta - D]/2)}, \quad (24)$$

if the following condition satisfies

$$(D-\alpha)\delta < 4-D. \quad (25)$$

Since  $\bar{\Phi} < 0$ , from Eq.(22) we obtain the additional condition,

$$\alpha\delta < -2. \quad (26)$$

If these condition Eqs.(25) and (26) satisfy, we have the bounded self-similar stationary solution (23) and (24). The mass distribution function  $f$ , the mass density  $\rho$ , and the velocity distribution  $f(v)$  become respectively

$$f(r, v) = C |2E|^{-[(D-\alpha)\delta + (D-2)]/2}, \quad (27)$$

$$\rho := S_D \int dv v^{D-1} f(r, v) \sim r^{\alpha-D+2/\delta}, \quad (28)$$

$$f(v) := S_D \int dr r^{D-1} f(r, v) \sim v^{\alpha\delta+2-D}, \quad (29)$$

where  $E$  denotes the mean field energy;

$$E := \frac{1}{2} v^2 + \Phi_0. \quad (30)$$

Since the solution (27) we obtained is a bounded solution, the specific heat of the self-similar stationary solution is always negative. The ratio of the average of the kinetic energy to the potential energy is as follows.

$$\frac{\langle \Phi_0 \rangle}{\langle K \rangle} = \frac{(D-\alpha)\delta - 4}{D}. \quad (31)$$

The  $\delta_*$  where the virial condition satisfies is

$$\delta_* = \frac{2(2\alpha - D)}{\alpha(D - \alpha)}. \quad (32)$$

If  $(D-\alpha)(\delta - \delta_*) < 0$ , the potential energy is dominant compared with the virial state.

The relation between pressure  $P$  and mass density  $\rho$  can be written in the form.

$$P \sim \rho^{1 + \frac{1}{1 + (\alpha-D)\delta/2}}. \quad (33)$$

The above equation of state corresponds to one of Polytropes gas when Polytropes index  $n$  equals  $1 + (\alpha-D)\delta/2$ . Note that for  $\alpha = D$ , there is no self-similar stationary solution that corresponds to an isothermal state.

As for gravity case ( $\alpha = D-2$ ), the above solution (24) and (27) in  $D = 3$  corresponds to the solution derived by R.N. Henriksen and L.M. Widrow[16]. For  $D = 1$  and  $D = 2$  where  $\alpha = D-2 \leq 0$ , we show the self-similar stationary solution in Appendix.A.

#### IV. LINEAR PERTURBATION ANALYSIS

In this section, we investigate the stability of the self-similar stationary solution derived in previous section for a symplectic linear perturbation by the energy functional analysis[11–14].

As for the linear stability of the stationary solution in CBE of the gravity in three-dimensional space, there are many works [11–14, 18–24]. For the stationary state, assuming spherical symmetry, characterized by the mass distribution function  $f_0$  specified as a function of the mean field energy  $E$  and the squared angular momentum  $J^2$ , if  $\partial f_0 / \partial E < 0$  and  $\partial f_0 / \partial J^2 < 0$ , then the system is stable to the linear perturbation.

Following the work by J. Perez and J.J. Aly[13] where the stability of stationary solution in the coupled CBE and Poisson equation with a spherical symmetry in three-dimensional space is studied, we study stability of solution obtained in previous section.

At first, we explain a symplectic linear perturbation by the energy functional analysis[11–14]. In term of the mass distribution function  $f(\mathbf{x}, v, t)$ , Hamiltonian  $H$  is written as follows,

$$H = \int d\Gamma \frac{v^2}{2} f(\mathbf{x}, v, t) + \frac{G}{2} \int d\Gamma \int d\Gamma' \mathcal{G}(|\mathbf{x} - \mathbf{x}'|) f(\mathbf{x}, v, t) f(\mathbf{x}', v', t), \quad (34)$$

where  $d\Gamma := d^D \mathbf{x} d^D v$  is a  $2D$ -dimensional phase volume element and the kernel  $\mathcal{G}$  satisfies

$$\Phi(\mathbf{x}) = G \int d\Gamma' \mathcal{G}(|\mathbf{x} - \mathbf{x}'|) f(\mathbf{x}', v', t). \quad (35)$$

We consider the small perturbation around some stationary solution  $f_0$ . The distribution function and Hamiltonian can be expanded around the stationary solution as follows.

$$f(\mathbf{x}, \mathbf{v}, t) = f_0 + \delta^{(1)}f + \delta^{(2)}f + \dots, \quad (36)$$

$$H = H_0 + \delta^{(1)}H + \delta^{(2)}H + \dots. \quad (37)$$

Here we consider any symplectic perturbation, which can be generated from the stationary solution  $f_0$  by use of a canonical transformation. By using some generating function  $K$ , any symplectic deformation can be expressed in the form

$$f = e^{[K, \cdot]} f_0. \quad (38)$$

From the above definition (38),  $f$  can be also expressed as follows.

$$\begin{aligned} f = & f_0 + [K, f_0] + \frac{1}{2!}[K, [K, f_0]] \\ & + \frac{1}{3!}[K, [K, [K, f_0]]] + \dots \end{aligned} \quad (39)$$

In the term of the parameta  $\epsilon$  which represents the amplitude of the perturbation, the generating function  $K$  is expanded in the form

$$K = \epsilon K^{(1)} + \epsilon^2 K^{(2)} + \epsilon^3 K^{(3)} + \dots, \quad (40)$$

and identifying  $g^{(n)} = \epsilon^n K^{(n)}$ , the perturbed quantities in the Eq.(36) are written as follows.

$$\delta^{(1)}f = [g^{(1)}, f_0], \quad (41)$$

$$\delta^{(2)}f = [g^{(2)}, f_0] + \frac{1}{2}[g^{(1)}, [g^{(1)}, f_0]]. \quad (42)$$

The first order term in Eq.(37) yields

$$\delta^{(1)}H = \int d\Gamma E[g^{(1)}, f_0], \quad (43)$$

where  $E$  is the energy of a particle,

$$E := \frac{v^2}{2} + \Phi_0, \quad (44)$$

where  $\Phi_0$  is the potential energy generated by  $f_0$ . Since  $E$  and  $f_0$  are conserved quantities,  $\delta^{(1)}H = 0$ .

The next order term in Eq.(37) yields

$$\begin{aligned} \delta^{(2)}H = & \int d\Gamma E[g^{(2)}, f_0] + \frac{1}{2} \int d\Gamma E[g^{(1)}, [g^{(1)}, f_0]] \\ & + \frac{G}{2} \int d\Gamma \int d\Gamma' \mathcal{G}(|\mathbf{x} - \mathbf{x}'|) [g^{(1)}, f_0] [g^{(1)'}', f_0']. \end{aligned} \quad (45)$$

The first term in Eq.(45) also vanishes and by an integration by parts, Eq.(45) is rewritten in the form.

$$\begin{aligned} \delta^{(2)}H = & -\frac{1}{2} \int d\Gamma [g^{(1)}, f_0] [g^{(1)}, E] \\ & + \frac{G}{2} \int d\Gamma \int d\Gamma' \mathcal{G}(|\mathbf{x} - \mathbf{x}'|) [g^{(1)}, f_0] [g^{(1)'}', f_0']. \end{aligned} \quad (46)$$

Hereafter we consider the case that the stationary solution  $f_0$  is a function of only the energy  $E$ . In this case, we obtain

$$[g^{(1)}, f_0] = F_E[g^{(1)}, E], \quad (47)$$

$$\int d^D \mathbf{v} [g^{(1)}, f_0] = \int d^D \mathbf{v} \partial_{\mathbf{v}} (F_E \mathbf{v} g^{(1)}), \quad (48)$$

where  $F_E := \partial_E f_0$ .

Integrating by parts and using Eqs.(47) and (48), we have

$$\begin{aligned} \delta^{(2)}H = & \frac{1}{2} \int d\Gamma (-F_E) |[g^{(1)}, E]|^2 + \frac{G}{2} \int d\Gamma \partial_{\mathbf{v}} (-F_E \mathbf{v} g^{(1)}) \\ & \times \int d\Gamma' \partial_{\mathbf{v}'} (-F_E' \mathbf{v}' g^{(1)'}) \mathcal{G}(|\mathbf{x} - \mathbf{x}'|). \end{aligned} \quad (49)$$

The linear perturbation  $g^{(1)}$  has two kind of gauge mode. One is the case that  $g^{(1)} = g^{(1)}(E)$ . In this case, the linear perturbation of the mass distribution  $\delta^{(1)}f$  is trivially zero. The other is the case that  $g^{(1)} = \mathbf{a} \mathbf{v}$  where  $\mathbf{a}$  is a constant. This perturbation means the translation of the center of mass. In order to consider the physical perturbation, we investigate the linear perturbation except the above gauge mode.

The stability for the linear perturbation[24, 25] reads that

$$\text{If } \delta^{(2)}H > 0, \text{ then the system is stable.} \quad (50)$$

### A. spherical mode

Since the first order perturbed potential  $\delta^{(1)}\Phi(r)$  satisfies

$$\begin{aligned} \frac{1}{r^{D-1}} \partial_r \left( r^{\alpha+1} \partial_r \delta^{(1)}\Phi(r) \right) &= S_D G \int d^D \mathbf{v}' [g^{(1)'}, f_0'] \\ &= S_D G \frac{1}{r^{D-1}} \partial_r \left( r^{D-1} \int d^D \mathbf{v}' F_E' \mathbf{v}' g^{(1)'} \right), \end{aligned} \quad (51)$$

the spatial derivative of  $\delta^{(1)}\Phi(r)$  becomes

$$\partial_r \delta^{(1)}\Phi(r) = \frac{S_D G}{r^{\alpha-D+2}} \int d^D \mathbf{v}' F_E' \mathbf{v}' g^{(1)'}. \quad (52)$$

From Eqs.(49) and (52), we have

$$\begin{aligned} 2\delta^{(2)}H &= \int d\Gamma (-F_E) |[g^{(1)}, E]|^2 \\ &\quad - \int d\Gamma \partial_r (-F_E \mathbf{v} g^{(1)}) \delta^{(1)}\Phi(r) \\ &= \int d\Gamma (-F_E) |[g^{(1)}, E]|^2 - S_D G \int \frac{d^D \mathbf{x}}{r^{\alpha-D+2}} \\ &\quad \times \int d^D \mathbf{v} (-F_E \mathbf{v} g^{(1)}) \int d^D \mathbf{v}' (-F_E' \mathbf{v}' g^{(1)'}). \end{aligned} \quad (53)$$

Introducing new variables,

$$g^{(1)} =: rv^r \mu(r, v, t), \quad (54)$$

and using Schwartz's inequality, we have

$$\begin{aligned} 2\delta^{(2)}H &= \int d\Gamma(-F_E)[|\mu rv^r, E]|^2 - GS_D \int \frac{d^D x}{r^{\alpha-D+2}} \\ &\quad \times \int d^D v [-F_E r(v^r)^2 \mu] \int d^D v' [-F'_E r(v'^r)^2 \mu'] \\ &\geq \int d\Gamma(-F_E)[|\mu rv^r, E]|^2 - GS_D \int \frac{d^D x}{r^{\alpha-D+2}} \\ &\quad \times \int d^D v [-F_E r(v^r)^2 \mu^2] \int d^D v' [-F'_E r(v'^r)^2] \\ &= \int d\Gamma(-F_E) \left\{ |\mu rv^r, E]|^2 - \frac{GS_D (rv^r)^2 \mu^2 \rho_0}{r^{\alpha-D+2}} \right\}, \end{aligned} \quad (55)$$

where  $\rho_0$  is non-perturbed mass density:

$$\rho_0 := \int d^D v f_0 = \int d^D v (-F_E)(v^r)^2. \quad (56)$$

Using the property of the Poisson bracket, and the fact that the integral of Poisson bracket over the phase space vanishes, the equation (55) can be rewritten in the form.

$$\begin{aligned} 2\delta^{(2)}H &\geq \int d\Gamma(-F_E) \left\{ |\mu rv^r, E]|^2 - \frac{GS_D (rv^r)^2 \mu^2 \rho_0}{r^{\alpha-D+2}} \right\} \\ &= \int d\Gamma(-F_E) \left\{ (rv^r)^2 |\mu, E]|^2 + |\mu|^2 |[rv^r, E]|^2 \right. \\ &\quad \left. + rv^r [\mu^2, E][rv^r, E] - \frac{GS_D (rv^r)^2 \mu^2 \rho_0}{r^{\alpha-D+2}} \right\} \\ &= \int d\Gamma(-F_E) \left\{ (rv^r)^2 |\mu, E]|^2 + |\mu|^2 |[rv^r, E]|^2 \right. \\ &\quad \left. + [\mu^2 rv^r [rv^r, E], E] - |\mu|^2 rv^r [[rv^r, E], E] \right. \\ &\quad \left. - |\mu|^2 |[rv^r, E]|^2 - \frac{GS_D (rv^r)^2 \mu^2 \rho_0}{r^{\alpha-D+2}} \right\} \\ &= \int d\Gamma(-F_E) \left\{ (rv^r)^2 |\mu, E]|^2 \right. \\ &\quad \left. - |\mu|^2 rv^r [[rv^r, E], E] - \frac{GS_D (rv^r)^2 \mu^2 \rho_0}{r^{\alpha-D+2}} \right\}. \end{aligned} \quad (57)$$

Using the following relation,

$$\begin{aligned} [[rv^r, E], E] &= -rv^r \left( \frac{d^2 \Phi_0}{dr^2} + \frac{3}{r} \frac{d\Phi_0}{dr} \right) \\ &= -rv^r \left( \frac{GS_D \rho_0}{r^{\alpha-D+2}} + \frac{2-\alpha}{r} \frac{d\Phi_0}{dr} \right), \end{aligned} \quad (58)$$

we obtain the final expression in the form.

$$\begin{aligned} \delta^{(2)}H &\geq \\ &\frac{1}{2} \int d\Gamma(-F_E)(rv^r)^2 \left( |\mu, E]|^2 + |\mu|^2 \frac{2-\alpha}{r} \frac{d\Phi_0}{dr} \right). \end{aligned} \quad (59)$$

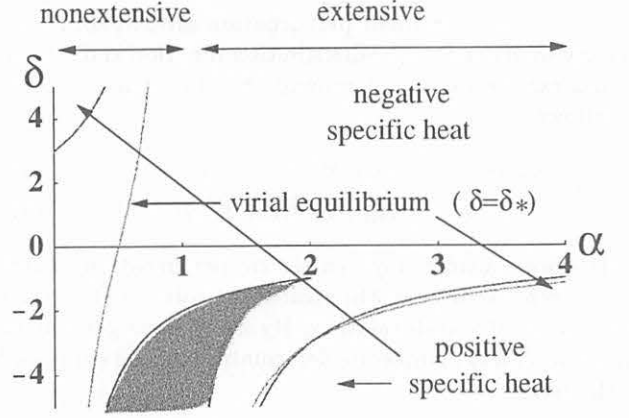


FIG. 1: One-dimensional case ( $D = 1$ ). The dark region corresponds to the stable region (62) in the parameta space  $(\delta, \alpha)$ . The property of the specific heat and the extensivity also are shown.

From Eq.(59), if  $\delta^{(2)}H = 0$  when  $F_E < 0$  and  $\alpha \leq 2$ ,  $\delta^{(1)}f = 0$ . Since such a perturbation is a gauge mode, we conclude that

$$\text{If } F_E < 0 \text{ and } \alpha \leq 2, \text{ then } \delta^{(2)}H > 0. \quad (60)$$

From the self-similar stationary solution Eq.(23), we have

$$F_E = \text{sgn}(E)[(\alpha - D)\delta + 2 - D]C|2E|^{[(\alpha-D)\delta-D]/2}. \quad (61)$$

As an explicit example, we consider  $D = 1$  case. From Eqs.(25), (26), and (60), if the following condition satisfies, the self-similar stationary solution Eq.(27) is stable.

$$\begin{aligned} -\frac{1}{\alpha-1} &< \delta < -\frac{2}{\alpha} & (1 < \alpha \leq 2) \\ \delta &< -\frac{2}{\alpha} & (0 < \alpha \leq 1) \end{aligned} \quad (62)$$

In Fig.1, the region of existence of the stable self-similar stationary solution in the parameta space  $(\alpha, \delta)$  is shown.

Note that in the above calculation, we use the integration by parts and neglect the surface term. Since the self-similar stationary solution obtained in this paper is singular at the boundary, the surface term can not be neglected in general. However we hope that the self-similar stationary solution can be connected with some regular solution near the boundary and the boundary term can be neglected.

## B. aspherical mode

Next, we consider an aspherical mode. Since it is difficult to analyze a general case, we study the gravity case in  $D$ -dimensional space ( $\alpha = D - 2$ ).

By the integral of Poisson equation over the configuration space and integration by parts, we have

$$-\frac{1}{S_D} \int d^D \mathbf{x} |\nabla \delta^{(1)} \Phi|^2 = \int d\Gamma \delta^{(1)} f \delta^{(1)} \Phi, \quad (63)$$

where

$$\delta^{(1)} \Phi := \int d^D \mathbf{v} \mathcal{G} \delta^{(1)} f. \quad (64)$$

From Eqs.(46), (63), and (64), we have

$$\delta^{(2)} H = \frac{1}{2} \int d\Gamma \frac{(\delta^{(1)} f)^2}{-F_E} - \frac{G}{2S_D} \int d^D \mathbf{x} |\nabla \delta^{(1)} \Phi|^2. \quad (65)$$

Here we introduce  $\delta^{(1)} \tilde{f}$  as follows.

$$\delta^{(1)} f =: F_E \delta^{(1)} \Phi + \delta^{(1)} \tilde{f}. \quad (66)$$

Substituting Eq.(66) into Eq.(65), we get

$$\begin{aligned} \delta^{(2)} H &= \frac{1}{2} \int d\Gamma \left[ \frac{(\delta^{(1)} \tilde{f})^2}{-F_E} + F_E (\delta^{(1)} \Phi)^2 - 2\delta^{(1)} f \delta^{(1)} \Phi \right] \\ &\quad - \frac{G}{2S_D} \int d^D \mathbf{x} |\nabla \delta^{(1)} \Phi|^2, \\ &= \frac{1}{2} \int d\Gamma \frac{(\delta^{(1)} \tilde{f})^2}{-F_E} + \frac{1}{2S_D} \int d^D \mathbf{x} \left\{ |\nabla \delta^{(1)} \Phi|^2 \right. \\ &\quad \left. - S_D \left[ \int d^D \mathbf{v} (-F_E) |\delta^{(1)} \Phi|^2 \right] \right\}. \end{aligned} \quad (67)$$

Moreover, using the new variable  $w$  which is defined by

$$\delta^{(1)} \Phi =: w(\mathbf{x}, t) \partial_r \Phi_0, \quad (68)$$

we can rewrite the equation (67) in the form.

$$\begin{aligned} \delta^{(2)} H &= \frac{1}{2} \int d\Gamma \frac{(\delta^{(1)} \tilde{f})^2}{-F_E} \\ &\quad + \frac{1}{2S_D} \int d^D \mathbf{x} \left\{ (\partial_r \Phi_0)^2 \left[ |\nabla w|^2 - S_D \int d^D \mathbf{v} (-F_E) |w|^2 \right] \right. \\ &\quad \left. - |w|^2 \partial_r \Phi_0 \nabla^2 \partial_r \Phi_0 \right\}. \end{aligned} \quad (69)$$

By the straightforward calculation, we obtain

$$\nabla^2 \partial_r \Phi_0 = S_D \int d^D \mathbf{v} F_E \partial_r \Phi_0 + (D-1) \frac{\partial_r \Phi_0}{r^2}. \quad (70)$$

By using the Wirtinger's inequality (see Appendix B), we have

$$\int \left[ |\nabla_s w|^2 - \frac{D-1}{r^2} |w|^2 \right] d\Omega \geq 0, \quad (71)$$

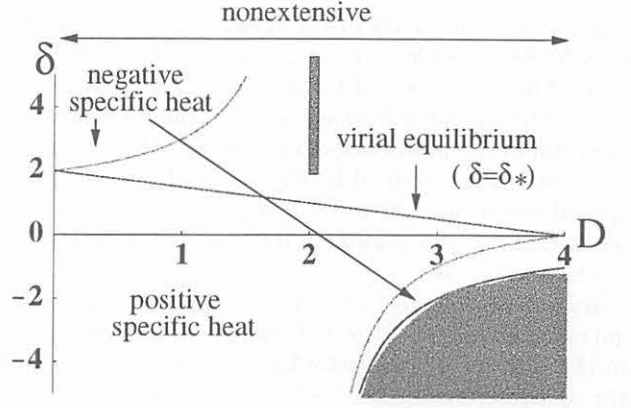


FIG. 2: Gravity case ( $\alpha = D - 2$ ). The dark region corresponds to the stable region (74) in the parameta space  $(\delta, D)$ . The property of the specific heat and the extensivity also are shown.

where  $\nabla_s := \nabla - \frac{\mathbf{r}}{|\mathbf{r}|} \frac{\partial}{\partial r}$ .

From Eqs.(69), (70), and (71), we get the final expression in the form.

$$\begin{aligned} \delta^{(2)} H &= \frac{1}{2} \int d\Gamma \frac{(\delta^{(1)} \tilde{f})^2}{-F_E} \\ &\quad + \frac{1}{2S_D} \int d^D \mathbf{x} (\partial_r \Phi_0)^2 \left[ |\partial_r w|^2 + |\nabla_s w|^2 - \frac{D-1}{r^2} |w|^2 \right] \\ &\geq \frac{1}{2} \int d\Gamma \frac{(\delta^{(1)} \tilde{f})^2}{-F_E} + \frac{1}{2S_D} \int d^D \mathbf{x} (\partial_r \Phi_0)^2 |\partial_r w|^2. \end{aligned} \quad (72)$$

From Eq.(72), if  $\delta^{(2)} H = 0$  when  $F_E < 0$ ,  $\delta^{(1)} f = 0$  or  $g^{(1)} = av$ . Since such a perturbation is a gauge mode, we conclude that

$$\text{If } F_E < 0, \text{ then } \delta^{(2)} H > 0. \quad (73)$$

This condition is weaker than the condition (60). In the gravity case ( $\alpha = D - 2$ ), from Eqs.(25), (26), and (60), if the following condition satisfies, the self-similar stationary solution Eq.(27) is stable.

$$\delta < -\frac{2}{D-2} \quad (2 < D \leq 4) \quad (74)$$

In Fig.2, the region of existence of the stable self-similar stationary solution in the parameta space  $(D, \delta)$  is shown. This stability condition (60) is consistent with the work by J. Perez and J.J. Aly[13]( $\alpha = 1, D = 3$ ).

## V. DISCUSSION

We study the self-similar stationary solution in the collisionless Boltzmann equation with an attractive  $1/r^\alpha$  potential assuming the spherical symmetric and isotropic orbit case in D-dimensional space and investigate the linear stability of its solution. In the above model, we can



control the extensivity of the system and the sign of the specific heat by the spatial dimension  $D$  and the exponent of inverse power of the potential  $\alpha$ .

The self-similar stationary solution can be expressed in the form of the power law of the energy. The exponent of the power is determined by the power of the potential  $\alpha$ , spatial dimension  $D$ , and the scaling parameter  $\delta$ . Here we interpret  $\delta$  as a parameter which denotes a virial ratio of the initial state.

By use of the energy functional approach, we investigate the stability of the self-similar stationary solution in the term of a symplectic linear perturbation. As for the spherical symmetric and isotropic orbit case of the gravity in  $D$ -dimensional space ( $\alpha = D - 2$ ), if the mass distribution function decrease monotonically and the spatial dimension is less than 4, then the system is stable. As for the power-law potential case in one-dimensional space ( $D = 1$ ), if the mass distribution function decrease monotonically and the inverse power of the potential is less than 2, then the system is stable. The self-gravitating ring model[7] is similar to the case of  $\alpha = 1$  in one-dimensional space. From the form of the velocity distribution obtained by a numerical simulation,  $\delta \sim -3$ . This case belongs to the stable self-similar stationary solution.

The stable self-similar stationary solution we obtained is a state where the potential energy is dominate compared with the virial equilibrium state. As for the extensivity of the system, the nonextensive system has many various stable scaling solution compared with the extensive one in the parameter space ( $\delta, \alpha, D$ ).

In the time evolution of the collisionless system assuming the spherical symmetry and isothermal case, Larson-Penston solution which shows self-similar collapse is attractor[26]. By such a self-similar time evolution of system, we hope that the class of the stable self-similar stationary solution obtained in this paper plays an important role as a quasi-equilibrium state of system with a long range interaction such as gravity. In a realistic situation, since the anisotropic velocity space is important, we will extend this analysis to the anisotropic case in the future work.

## APPENDIX A: STABILITY CONDITION FOR GRAVITY CASE IN $D = 1$ AND $2$ ( $\alpha = D - 2$ )

For the case that the potential  $\bar{\Phi}$  is positive, Eq.(22) is modified as follows.

$$2\nu^2 \left[ 2 + (D-2)\frac{\delta}{\nu} \right] \bar{\Phi} = S_D^2 G \int_0^\infty X^{D-1} \bar{f} dX. \quad (A1)$$

### 1. $D = 1$ case

Since the potential  $\bar{\Phi}$  is positive, from Eq.(A1),

$$\delta < 2. \quad (A2)$$

By integrating Eq.(A1), we have self-similar solution (23) and

$$C = \frac{|2 - \delta| \Gamma(\delta - 1/2) |2\bar{\Phi}|^\delta}{\sqrt{\pi} G \Gamma(\delta - 1)}, \quad (A3)$$

if the following condition satisfies

$$\delta > 1. \quad (A4)$$

From Eqs. (A2), (A4), (60) and (73), the stability condition for linear perturbation yields

$$1 < \delta < 2. \quad (A5)$$

The ratio of the average of the kinetic energy to the potential energy is same as Eq.(31) in  $D = 1$ . However if  $\delta \leq 2$ , the integral of the kinetic energy over the velocity space diverges. By this reason, there does not exist the stable self-similar stationary solution in  $D = 1$ .

### 2. $D = 2$ case

If the potential  $\bar{\Phi}$  is negative, the condition that the bounded self-similar solution exists is the same one (25) and (26). Since this case does not satisfies the condition (26), the only case that  $\bar{\Phi} > 0$  is possible.

In this case, from Eq.(A1), the condition that a bounded self-similar solution exists yields

$$\delta > \frac{1}{2}, \quad (A6)$$

and the integral constant of (23) is

$$C = \frac{2\Gamma(\delta) |2\bar{\Phi}|^\delta}{\pi^{5/2} G \Gamma(\delta - 1/2)}. \quad (A7)$$

From Eqs. (60), (73), and (A6), the stability condition for linear perturbation yields

$$\delta > \frac{1}{2}. \quad (A8)$$

The ratio of the average of the kinetic energy to the potential energy is same as Eq.(31) in  $D = 2$ . However, same as  $D = 1$  case, if  $\delta \leq 2$ , the integral of the kinetic energy over the velocity space diverges. Finally, if  $\delta > 2$ , the self-similar stationary solution in  $D = 2$  is stable. In this case, the specific heat is always positive.

## APPENDIX B: WIRTINGER'S INEQUALITY

Following J.J.Aly and J.P'erez[27], we show simple proof of Wirtinger's inequality (71). Let us consider a function  $f(r, a^i)$ , where  $a^i$  denote an angular coordinate in  $(D-1)$ -sphere and  $i = 1, 2, \dots, D-1$ , with zero average value over the  $(D-1)$ -sphere  $S_r$ . Using spherical

harmonics  $Q_{m^i}^l$  in  $S^{D-1}$ , the function  $f$  can be written as follows.

$$f(r, a^i) = \sum_{l=1}^{\infty} \prod_{i=1}^{D-2} \sum_{m^i} c_{m^i}^l(r) Q_{m^i}^l(a^i), \quad (\text{B1})$$

where  $m^i$  denotes the proper mode corresponds to each angular coordinate. From the completeness of the spherical harmonics and the following relation

$$\nabla_s^2 Q_{m^i}^l = -\frac{l(l+D-2)}{r^2} Q_{m^i}^l, \quad (\text{B2})$$

where  $\nabla_s := \nabla - \left(\frac{r}{r}\right) \partial_r$ , we have

$$\begin{aligned} & \int_{S_r} \nabla_s Q_{m^i}^l \nabla_s Q_{m'^i}^{l'} d\Omega_{D-1} \\ &= \frac{l(l+D-2)}{r^2} \int_{S_r} Q_{m^i}^l Q_{m'^i}^{l'} d\Omega_{D-1} \\ &= \frac{l(l+D-2)}{r^2} \delta_{l'}^l \delta_{m'^i}^{m^i}, \end{aligned} \quad (\text{B3})$$

where  $d\Omega_{D-1}$  denotes a volume element in (D-1)-sphere.

From Eqs.(B1) and (B3), we obtain

$$\begin{aligned} \int_{S_r} |\nabla_s f|^2 d\Omega_{D-1} &= \sum_{l \geq 1; m^i} |c_{m^i}^l|^2 \frac{l(l+D-2)}{r^2} \\ &\geq \frac{D-1}{r^2} \sum_{l \geq 1; m^i} |c_{m^i}^l|^2 \\ &= \frac{D-1}{r^2} \int_{S_r} |f|^2 d\Omega_{D-1}, \end{aligned} \quad (\text{B4})$$

with equality satisfying if  $f = f_0(r) \left(\frac{r}{r}\right) \alpha$  for some constant vector  $\alpha$ .

- 
- [1] R.B. Larson Mon. R. astr. Soc. **194** (1981), 809.
  - [2] E. Falgarone, T.G. Phillips, and C.K. Walker, ApJ. **378** (1991), 186.
  - [3] D. Merritt and T. Fridman, ApJ. **460** (1996), 136.
  - [4] F.S. Labini, M. Montuori, and L. Pietronero, Phys. Rep. **61** (1998), 293.
  - [5] J.F. Navarro, C.S. Frenk, and S.D.M. White, ApJ. **490** (1997), 493.
  - [6] T. Fukushige and J. Makino, ApJ. **557** (2001), 533.
  - [7] Y. Sota, O. Iguchi, M. Morikawa, T. Tatekawa, and K. Maeda, Phys. Rev. E **64** (2001), 056133.
  - [8] I. Ispolatov and E.D.G. Cohen Phys. Rev. Lett. **87** (2001), 210601.
  - [9] A. Campa, A. Giansanti, D. Moroni, and C. Tsallis, Phys. Lett. A **286** (2001), 251.
  - [10] A. Campa, A. Giansanti, and D. Moroni, Physica A **305** (2002), 137.
  - [11] H.E. Kandrup, ApJ. **351** (1990), 104.
  - [12] H.E. Kandrup, ApJ. **370** (1991), 312.
  - [13] J. Perez and J.J. Aly, Mon. Not. R. Astron. Soc. **280** (1996), 689.
  - [14] J. Perez and M. Lachieze-Rey, ApJ. **465** (1996), 54.
  - [15] P. Jund, S.G. Kim, and C. Tsallis, Phys. Rev. B **52**

- (1995), 50.
- [16] R.H. Henriksen and L.M. Widrow, Mon. Not. R. Astron. Soc. **276** (1995), 679.
- [17] B. Carter and R.H. Henriksen, J. Math. Phys. **32** (1991), 2580.
- [18] J. Binney and S. Tremaine, *Galactic Dynamics*. ( Princeton Univ. Press, Princeton ) (1993).
- [19] V.A. Antonov, Sov. Astron. **4** (1961), 859.
- [20] V.A. Antonov, Vestnik Leningrad Univ. **7** (1962), 135.
- [21] D. Lynden-bell and N. Sannit, Mon. Not. R. Astron. Soc. **143** (1969), 167.
- [22] J.R. Ipser, ApJ. **232** (1974), 863.
- [23] J.F. Sygnet, G. Des Forets, M. Lachieze-Rey, and R. Pellet, ApJ. **276** (1984), 737.
- [24] P. Bartholomew, Mon. Not. R. Astron. Soc. **151** (1971), 333.
- [25] D.D. Holm, J.E. Marden, T. Ratiu, and A. Weinstein, Phys. Rep. **123** (1985), 1.
- [26] T. Hanawa and K. Nakamura, ApJ. **484** (1997), 238.
- [27] J.J. Aly and J. Perez, Mon. Not. R. Astron. Soc. **259** (1992), 95.

# The effect of tidal force and mass loss in star clusters sinking process

Tatsushi Matsubayashi<sup>1</sup> and Toshikazu Ebisuzaki<sup>2</sup>

<sup>1</sup> Department of Earth and Planetary Sciences  
Tokyo Institute of Technology, Japan  
tatsushi@geo.titech.ac.jp

<sup>2</sup> Computational Science Division, RIKEN, Japan

## Abstract

Compact star clusters in the star burst galaxies sink toward the galactic center through dynamical friction. If they survive well against the mass loss by the stellar evolution and the tidal stripping from the parent galaxy, then they convey intermediate mass black holes ( $\sim 1000M_{\odot}$ ), produced in the cluster through runaway mergings of the massive stars, into the galactic center to form a supermassive black hole (Ebisuzaki et al. 2001). In the present paper, we investigated the condition for surviving of the cluster by means of numerical simulations which include stellar evolution and tidal stripping. As stars evolve, they eject their mass, which is lost away from the cluster. Furthermore, the tidal force of the parent galaxy stripped the stars in the outside of the cluster. Through these processes, both the number of the members and the gravitational binding energy of the cluster become smaller as they sink. Finally they totally disrupted, when their identity are lost. We performed a series of the gravitational N-body simulations for the star clusters of the star burst galaxy, M82.

We found that a cluster with an initial cluster mass of  $M_c \geq 3 \times 10^6 M_{\odot}$  and the lower limit mass of IMF  $M_{min} \leq 0.5M_{\odot}$ . the surviving condition of the cluster well survive until it sinks down to the center of the parent galaxy. We also found that the results depend strongly on the total mass of cluster and initial mass function, IMF. These results support the formation scenario of supermassive black hole described above.

We adopted Hernquist spherical galaxy model for M82 Galaxy, which is truncated at galactic radius  $r_b = 0.5\text{kpc}$  and whose mass is  $M_b = 2.08 \times 10^9 M_{\odot}$ . and King model for the compact star cluster, whose central potential is  $W_0 = 5.0$  and core radius is  $r_0 \sim 1\text{pc}$ . We used special-purpose computer, MDGRAPE-2.

## 1 Introduction

### 1.1 SMBHs in Galaxies

There is rapidly growing evidence for supermassive black hole (SMBHs) in the centers of many galaxies (Kormendy & Richstone 1995).

Many authors have pointed out that the mass of the central BH,  $m_{BH}$ , correlates with the mass of the bulge,  $M_b$ . The ratio of  $m_{BH}$  to  $M_b$  is almost constant at 0.002 (Kormendy & Richstone 1995), 0.006 (Magorrian et al. 1998) as figure 1. This suggests that the formation of the central BH is somehow related to that of the bulge.

The formation mechanism of SMBHs is not well understood. In the famous diagram by Rees (1978, 1984), there were basically two paths from gas clouds to SMBHs. The first is direct monolithic collapse; the second is via the formation of a star cluster, with subsequent runaway collisions leading to BH formation. Previous numerical studies, however, have demonstrated that neither path is likely. In the first, a massive gas cloud is much more likely to fragment into many small clumps in which stars then form, so direct formation of a massive BH from a gas cloud seems impossible. In the second, stellar dynamics in star clusters does not easily lead to the formation of SMBHs. A number of low-mass BHs (masses around  $10M_{\odot}$ ) are formed via the evolution of massive stars, and these BHs do indeed sink to the center of the cluster through dynamical friction and form binaries by three-body encounters. Taniguchi et al. (2000) argued that intermediate-mass BH (IMBHs) could be formed through successive merging of compact objects. However, recent

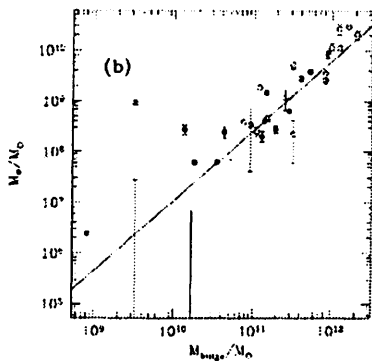


Figure 1:  $M_{BH}$  and  $M_{bulge}$  produced by Magorrian et al. (1998) models. The filled and open circles plot power-law and core galaxies, respectively. The solid line plot  $M_{BH,fit}$  as described  $M_{BH}/M_b = 0.006$ .

N-body simulations (Portegies Zwart & McMillian 2000) have demonstrated that practically all of these BH binaries are ejected from the cluster by recoil of interactions with other BHs (or BH binaries) before they merge through gravitational radiation.

## 1.2 Discovery of IMBH in M82 and New formation scenario for SMBHs

Matsumoto et al. (2001) have identified nine bright compact X-ray source in the central region of M82 using recent *Chandra* data. The brightest source, which exist some 200pc away from dynamical center of M82, (number 7 in their Table 1) had a luminosity of  $9 \times 10^{40} \text{ ergs}^{-1}$  in 2000 January, corresponding to a BH with a minimum mass of  $700M_{\odot}$ . Assuming the Eddington luminosity, the relation  $M_{BH}$  and the luminosity of the brightest source is given as

$$M_{BH} > 770M_{\odot} \left( \frac{L}{9 \times 10^{40} \text{ ergs}^{-1}} \right). \quad (1)$$

It probably consists of a single compact object, as its X-ray flux shows rapid time variation (Matsumoto et al. 2001). This is first detection of a BH candidate with a mass much greater than  $100M_{\odot}$  but much less than  $10^6M_{\odot}$ . Among the eight other sources, at least three (5, 8, and 9) have Eddington masses greater than  $30M_{\odot}$ .

Matsushita et al. (2000) observed the same region with the Nobeyama Millimeter Array and found a huge expanding shell of the molecular gas. They estimated the age and kinetic energy of the shell to be around 1 Myr and  $10^{55} \text{ ergs}^{-1}$ , suggesting that a strong star-burst took place a few milliyears ago. We show more details about the cluster in §2.3

We now have two important observational results. The first is that a BH with intermediate mass ( $770 < M_{BH}/M_{\odot} < 10^6$ ) may have been found. The second is that it coincides with a young compact star cluster.

Based on these findings, Ebisuzaki et al.(2001) suggest a new formation scenario for SNBHs. In this scenario, IMBHs first form in young compact star clusters through runaway merging of massive stars. While these IMBHs are forming, the host star clusters sink toward the galactic nucleus through dynamical friction and upon evaporation deposit their IMBHs near the galactic center. The IMBHs then form binaries and eventually merge via gravitational radiation, forming an SMBH.

In the following, we discuss how IMBHs can be formed in young compact star clusters in §1.3, then IMBHs might grow into SMBHs in §1.4.

### 1.3 IMBH formation through runaway growth

Ebisuzaki et al. (2001) proposed that IMBHs form and grow through successive mergings of massive star (and IMBHs) in dense star clusters (Fig.1.4).

A gas cloud fragments to form many less massive clouds as it cools by radiation. Many stars are formed through this fragmentation, and a star cluster comes into being. There are two possible evolutionary paths for this cluster depending on its stellar density.

If the cluster is not dense enough for mass segregation to occur in 10 Myr, massive stars evolve into compact stellar remnants such as neutron stars and stellar mass BHs ( $10M_{\odot}$ ). Those stellar remnants slowly sink to the cluster center since they are heavier than other stars in the system and eventually form binaries. Successive three-body interactions make these binaries more tightly bound, and eventually they are ejected from the cluster by the slingshot mechanism.

If the star cluster is so dense that stellar mass segregation is faster than stellar evolution for the most massive stars (time-scale  $10^6 yr$ ), those stars sink to the cluster core by dynamical friction and form a dense inner core of massive stars at the cluster center. In this inner core, the massive stars undergo a runaway stellar merging and a very massive star forms, with mass exceeding  $100M_{\odot}$ . This very massive star eventually collapses into a BH, which continues to grow by swallowing nearby massive stars. More massive stars in star clusters have higher merging rates than less massive cluster members (or field stars) because of their larger geometrical cross sections, a stronger gravitational focusing and concentration to the central region by mass segregation in the cluster. In fact, Portegies Zwart et al. (1999) demonstrate N-body simulations that runaway merging can take place in systems containing  $\sim 12,000$  stars before stellar evolution eliminates the most massive stars.

Portegies Zwart et al. (1999) found that in one case, the most massive star experienced more than 10 collisions and reached a mass of around  $200M_{\odot}$  before evolving into a supernova. There is considerable uncertainty as to how much mass would remain as a star approached within its tidal radius, leading to a relatively large merger cross section.

In order for runaway merging to occur, the dynamical friction time-scale for the most massive stars must be short enough that they can sink to the center during their lifetimes of several milliyers. The dynamical friction time-scale can be expressed as follows (§A.3):

$$\begin{aligned} t_{fric} &= \frac{0.519}{\ln \Lambda} \frac{r^2 v_c}{Gm} \\ &\simeq 1.17 \times 10^7 \times \left( \frac{r}{1pc} \right)^2 \times \left( \frac{v_c}{10km/s} \right) \times \left( \frac{10M_{\odot}}{m} \right) yr, \end{aligned} \quad (2)$$

where  $\ln \Lambda$  is the Coulomb logarithm,  $G$  is the gravitational constant,  $v_c$  is the circular velocity, which value is same order with random velocity,  $r$  is the distance from the center of the cluster, and  $m$  is the mass of the star. Here it is assumed that the background stellar distribution is that of the singular isothermal sphere. Equation(2) is a useful approximation at  $r \sim r_0$  ( $r_0$  is the core radius  $\simeq 1pc$  for compact star cluster).

In the following, we consider how dynamical friction works in the cluster found in M82. If the total mass of the cluster has  $3 \times 10^6 M_{\odot}$ , about 50% of the total mass is included within a radius of  $r = 1pc \sim r_c$ . Then, the dynamical friction time-scale of stars about 25% of the total mass, is about 10 Myr. Marchant & Shapiro (1980) performed Monte Carlo simulations of this stage for a simplified cluster containing  $3 \times 10^5 M_{\odot}$  stars and one  $50M_{\odot}$  seed BH. They found that the BH mass jumped to over  $10^3 M_{\odot}$  (0.3% of the cluster mass) almost immediately after they put the BH into the system. Their result should be regarded as a lower limit on the BH growth rate since realistic effects, in particular the presence of a mass spectrum, would greatly enhance the accretion rate. Taking these effects into account, that it seems safe (even conservative) to suppose that 0.1% of the total cluster mass accretes to form a  $\sim 1000M_{\odot}$  central BH in a 10 Myr.

Presently there are more than 100 star clusters discovered in M82 galaxy, some of them apparently hosting small BHs. Their age is around 10 Myr (T. Harashima et al. 2001, in preparation). Also the starburst in M82 is a long-duration event, having started at least 200 Myr ago. A close encounter with a large galaxy, M81, in the last 100 million years is thought to be the cause of the starburst activity (see §2.2). As stated above, we conclude that around 100 clusters similar to our host cluster have formed in total and that a considerable fraction of them host IMBHs.

## 1.4 Building up the central SMBH

We now describe how IMBHs formed in star clusters combine to form a central SMBH (Fig.1.4).

The growth rate of the IMBH in a star cluster slows down once all the massive stars are swallowed (after  $\sim 100$  Myr). Subsequently, the cluster is subject to two evolutionary processes: evaporation through two-body relaxation and orbital decay (sinking) via dynamical friction. Evaporation is driven partly by thermal relaxation and partly by stellar mass loss. Portegies Zwart et al. (2001) estimated that the evaporation time-scale for a tidally limited compact star cluster is around 2~3 half-mass relaxation times, which is of the order of a few gigayears for our star clusters. Rewriting equation (1) using appropriate scaling for this case  $r \gg r_0$  (§A.3), we find that the time-scale on which the cluster sinks to the galactic center via dynamical friction is

$$t_{fric} \simeq 1.0 \times 10^9 \times \left( \frac{r}{1 \text{ kpc}} \right)^2 \times \left( \frac{v_c}{100 \text{ km/s}} \right) \times \left( \frac{3 \times 10^6 M_\odot}{m} \right) yr. \quad (3)$$

Clusters initially within 1 kpc of the galactic center can therefore reach the center within a few Gyr.

According to my estimate in §1.3, around 100 compact clusters have formed close to the center of M82 in the last 200 Myr. If we assume that half of these clusters contain  $1000 M_\odot$  IMBHs and that these IMBHs actually merge, then the total BH mass at the center of the galaxy will be at least  $1.0 \times 10^5 M_\odot$ . Successive mergings of IMBHs form an SMBH with a mass of  $10^6 M_\odot$ .

We have demonstrated that  $1000 M_\odot$  IMBHs can form and reach the galactic center in a reasonable time-scale. We now turn to the question of whether the multiple IMBHs at the center can merge. Begelman, Blandford, & Rees (1980) discussed the evolution of an SMBH binary at the center of a galaxy, taking dynamical friction from field stars and energy loss via gravitational radiation into account. They found that the merging time-scale depends strongly on mass, and for a very massive BH with a mass of  $10^8 M_\odot$  in which they were interested, merging took much longer than a Hubble time.

For the IMBHs, however, the time-scale for merging through gravitational radiation is many orders of magnitude shorter than that for the SMBHs. Recent extensive numerical simulations (Makino et al. 1993; Makino 1997) have shown that the hardening of the BH binary through dynamical friction is in fact several orders of magnitude faster than the prediction from loss cone arguments. Although the number of particles employed (up to 256,000) was not large enough to model SMBH binaries, it was certainly large enough to model evolution of IMBH binaries.

Once one BH has become more massive than typical infalling BHs, it becomes extremely unlikely that it will be ejected since the recoil velocity from three-body interactions is inversely proportional to the mass (because of momentum conservation). Thus, even though some of the infalling BHs might be ejected by the slingshot mechanism, the central BH will continue to grow.

Since we now have the first candidate for IMBHs, it seems natural to expect that SMBHs might be formed from them. Thus, IMBHs are created and transported to the center of the galaxy, where they eventually merge to form SMBHs.

Ebisuzaki et al.(2001) propose that IMBHs are formed in the cores of young compact star clusters through mergings of massive stars and BHs formed from them. These compact young clusters sink to the galactic center by dynamical friction. The cluster is dissolved through Stellar mass loss, the tidal stripping of the parent galaxies, and the thermal relaxation of stars.

First, Fukushige and Heggie (1995) investigated the effect of the Stellar evolution and galactic tide of globular clusters. Their study included the following realistic effects: the spectrum of stellar masses; mass loss arising from stellar evolution; and a tidal cut-off to model the effect of the galactic tidal field. They performed an extensive survey of models that differ with regard to the initial mass function, the central potential of the cluster, and the distance from the galactic center. For example, they obtained the result that the cluster having sufficiently deep central potential survived during over  $2 \times 10^9 yr$ . They chose the dimensionless central potential of King's model,  $W_0 = 5.0$  for the cluster model. They adopt the parent galaxy model in which the cluster is assumed to move in a spherically symmetric galactic potential, taken to be that of a distant point mass. In the region of galactic center, the life time of cluster may be shorter than  $2 \times 10^9 yr$ .

Second, Binney and Tremaine(1987) discussed that the evaporation time-scale for the local globular cluster is  $t_{life} \sim 100 \times t_{relax}$ , where  $t_{relax} \sim 10^{10} yr$  for the typical globular cluster, and  $t_{life}$  for the globular cluster  $\sim 10^{12} yr$ . However, a globular cluster move under the influence of the

mean potential generated by all the other particles. In the core of a globular cluster, therefore, life time may be shorter than  $10^9 \text{yr}$ , and play a key role.

We estimated that the time-scale of Stellar evolution, dissolving by tidal force from the parent galaxy, and evaporation is  $10^6 \sim 10^9 \text{yr}$ . In this paper, We had taking into account these effect except an evaporation effect, because the evaporation time-scale is longer than other two effects, Stellar evolution and galactic tidal force. We performed these effects by the simulations using N-body integrator assumed in the region of a galactic center.

The organization of the paper is as follows: The numerical method and more detail the initial condition for my simulations is discussed in §2. In §3 the results are presented. Finally, §4 notes conclusions.

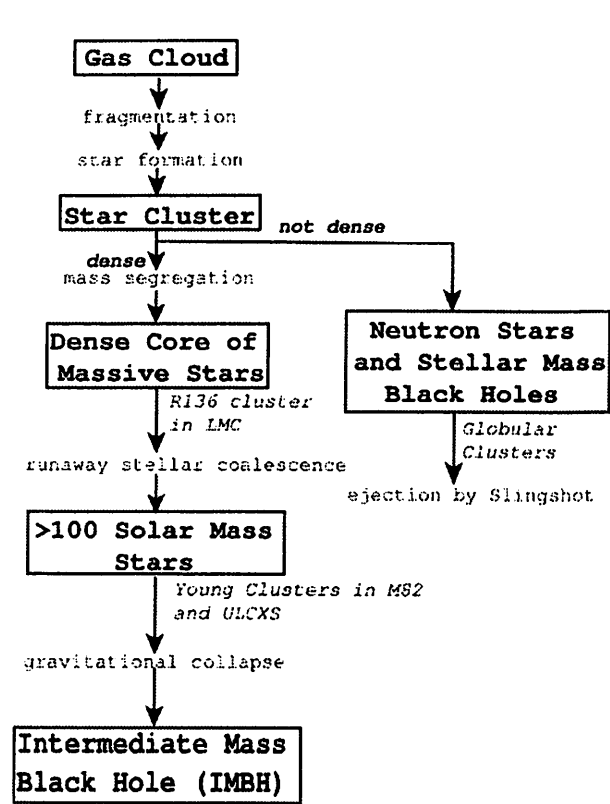


Figure 2: Schematic diagram of the formation process of an IMBH.\*

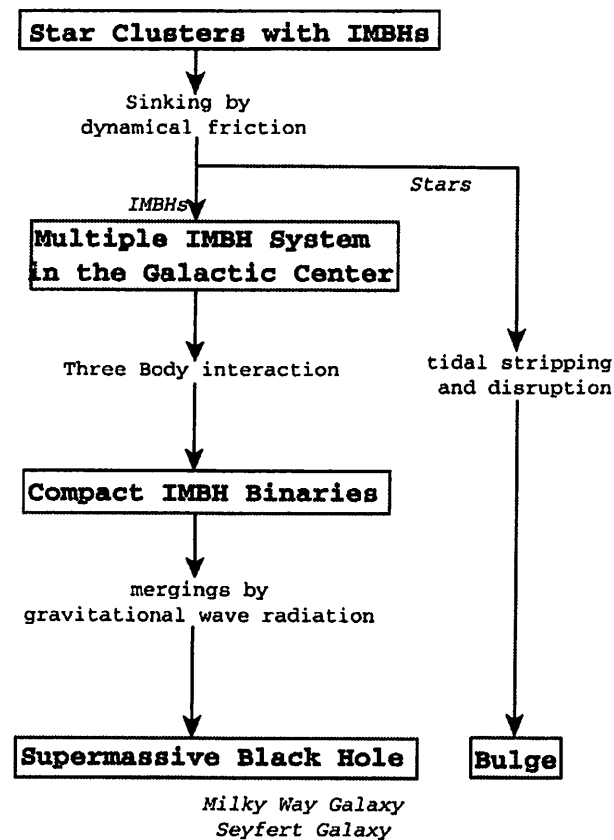


Figure 3: Schematic diagram of the formation of SMBHs from star clusters containing IMBHs.\*

\*:Ebisuzaki et al. (2001)

## 2 Numerical methods and Initial conditions

The  $N$ -body integration algorithm, used in this paper, is described in §2.1. In §2.3.1, we described how the evolution of stars is calculated; the tidal boundary of cluster is described in §2.3.3.

## 2.1 The $N$ -body integrator

The  $N$ -body portion of the simulation is carried out using the tree code (Barnes & Hut 1986). The Burns-Hut tree code is widely used algorithm that reduces the cost of the force calculation. In this tree code, forces on a particle from distant particles are replaced by multipole expansions of groups of particles. More distant particles are organized into larger groups, so that the truncation error of the expansion is similar everywhere. A hierarchical tree structure is used to form groups efficiently. The calculation cost is reduced from  $O(N^2)$  to  $O(N \log N)$ .

Even with this tree code, the cost of the force calculation is still high, and it dominates the total calculation cost. In order to accelerate the tree code further, we can use GRAPE (GRAvity PipE; Sugimoto et al. 1990, Makino and Taji 1998). GRAPE is special-purpose hardware for the calculation of the gravitational force between particles. In this paper, we performed all simulations with MDGRAPE-2 (Narumi et al. 1999, Susukita et al. 2002). For the implementation of the tree code on GRAPE, see Makino 1991.

We chose the system of units in which the total mass of each galaxy, the typical velocity dispersion, and gravitational constant are both 1 and in which the initial maximum radius of galaxy is  $1/2$ . In §??, we described detail. We integrated this system up to  $T_{end} = 30.0$  with constant time-step. In models with  $N = 114000$ , the time-step,  $\Delta t$ , was  $1/4000 = 0.00025$ , and 120000 steps took for about 5 days. while for point mass cluster model  $N = 104000$  the time-step,  $\Delta t$ , was  $1/200 = 0.005$ , and 6000 steps took for a few hour. We took same value for softening parameter with the value of  $\Delta t$ , because of velocity dispersion  $\sigma \approx 1.0$ .

## 2.2 M82 Galaxy component

The distance from our Galaxy to M82 is assumed to be  $3.25 Mpc$ . The rotation curve of small-mass starburst galaxy M82 has a steep nuclear rise, peaking at 500 pc radius, which then declines in a Keplerian fashion. This rotation curve mimics that for a central bulge of spiral galaxies with a high concentration of stellar mass. The declining rotation indicates that its extended disk mass is missing.

Sofue (1998) propose that M82 is a surviving central bulge of a much larger disk galaxy, whose outer disk was truncated during a close encounter with M81. Through the close encounter with M81, when M82 penetrated the disk of M81, the outer disk of M82 was tidally truncated, but the bulge and nuclear disk have survived the tidal disruption. The truncated disk may have become the HI envelope and tails in M81-M82 system. The central gas disk of M82 was dense enough. This close encounter has caused the high-density molecular disk in M82, and starburst.

### 2.2.1 Rotation curve and mass distribution

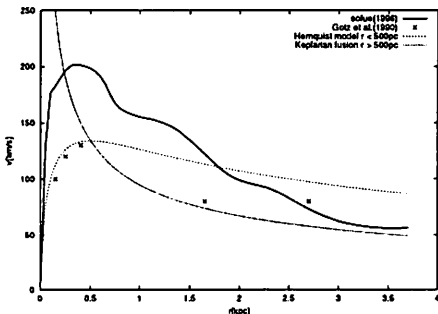


Figure 4: Rotation curve of M82. Full line shows observation of HII line by Sofue (1998). Dot point shows observation of CO and Na line Gotz et al (1990). Dashed and dotted line shows Hernquist model, using our simulations

In figure 4, we compare the rotation curve of Sofue (1998), Gotz et al (1990), and the simulations model. We used Hernquist model for Galaxy model in all of the simulations, which models



detail is discussed in below section §2.2.2.

Gotz et al.(1990) propose a total dynamical mass model within the starburst region (radius= $15'' \sim 230pc$ ) is close to  $6 \times 10^8 M_\odot$ . They propose mass profile giving  $M_{(500pc)} = 2.08 \times 10^9 M_\odot$ .

### 2.2.2 Galaxy model (Hernquist model)

We adopted galaxy model are spheroid Herunquist model(Herunquist 1993) proposed for spherical galaxies and bulges. The bulge density profile is

$$\rho_b(r) = \frac{M_b}{2\pi} \frac{r_b}{r(r+r_b)^3}, \quad (4)$$

where  $r_b$  is scale length for bulge,  $M_b$  is defined as the mass within a infinite radius,  $M_b = M_{b(\infty)}$ . The cumulative mass profile and potential corresponding to  $\rho_b(r)$  can be written

$$M_b(r) = M_b \frac{r^2}{(r+r_b)^2}, \quad \Phi_b(r) = -\frac{GM_b}{r+r_b}. \quad (5)$$

### 2.2.3 Effect of dynamical friction

We checked the effect of dynamical friction by means of test calculations with different particle number,  $N$ . Figure 5 shows the evolution of distance from the center of galaxy for several different particle numbers. The calculations were performed with  $N = 10400, 20800, 52000$ , and  $104000$  (We defined as 10K, 20K, 50K, and 100K below) for all cases for  $(M_{galaxy}, M_{cluster}) = (2.08 \times 10^9 M_\odot, 2.0 \times 10^6 M_\odot \text{ or } 3.0 \times 10^6 M_\odot)$ .

As shown in Figure 5, the clusters with smallest particle number ( $N \leq 20K$ ) are substantially affected by underestimate of the influence of dynamical friction. The results with large particle number ( $N \geq 50K$ ) differ little from each other. Therefore, we adopted 100K for the number of particle simulations.

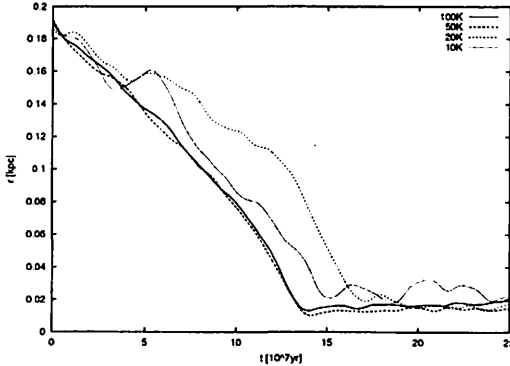


Figure 5: Dependence of the evolution on the number of particles  $N$  used in the simulation. The coordinate is the periodic average of an orbital radius.

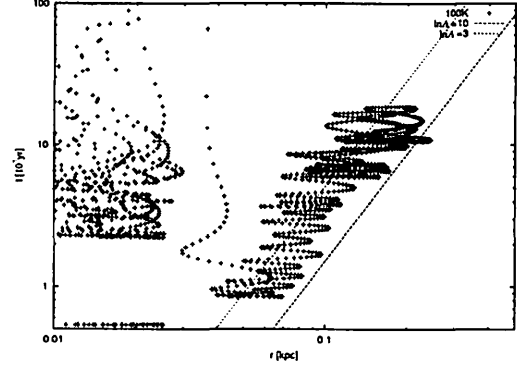


Figure 6: Relation between  $r$  and  $t_{fric} \equiv r/\dot{r}$  in the case of 100K. The both dashed line gives the theoritcal line as equation (7) cased  $\ln \Lambda \simeq 10$  and  $\ln \Lambda \simeq 3$ .

Second, we estimated the effect of dynamical friction between theoretical curve and simulation result. Dynamical friction formula definition as bellow( equation(17) ),

$$t_{fric} = \frac{r}{\dot{r}} = \frac{M_c v_c}{F} = -\frac{2.34}{4\pi \ln \Lambda} \frac{v_c \sigma^2}{G^2 M_c \rho(r)}, \quad (6)$$

where,  $M_c$  is mass of the cluster,  $\sigma$  is velocity dispersion. Using equation (6), and eliminating from the relation,  $\sigma = v_c/\sqrt{2}$ , we have

$$t_{fric} = -\frac{2.34}{8\pi \ln \Lambda} \times \frac{1}{G^2 M_c} \frac{v_c^3}{\rho(r)} = -\frac{0.5841}{\ln \Lambda} \times \frac{M_b}{M_c} \left( \frac{r}{1 \text{ kpc}} \right)^{5/2}. \quad (7)$$

Unfortunately, equation(7) has the unfixed value due to  $1/\ln \Lambda$ . Theoretically, this value is  $\ln \Lambda \simeq 10$  at a typical for spheroid galaxy. However, numerical integration couldn't realize the real galaxy, because the number of stars in a typical galaxy is over  $10^{10}$ . Therefore, we should modeling a galaxy using softening parameter  $\varepsilon$  as a galactic dynamics studies. Consequently,  $\ln(b_{max}/b_{min})$  is used for  $\ln \Lambda$  in a numerical simulation, where  $b_{max}$  is the largest impact parameter,  $b_{max} \approx$  the system size,  $b_{min}$  is the right-angle impact parameter,  $b_{min} \approx \varepsilon$ . In our simulations, therefore,  $\ln \Lambda \sim 3$  is good estimates.

## 2.3 Cluster component

Matsushita et al.(2000) observed that the luminosity of the  $2.2\mu m$  secondary peak is equivalent to  $\sim 1500 M_2$  supergiants. From their observations, using an extended IMF of  $dN \propto M^{-2.5} dM$  with lower and upper mass limits of 1 and 30, respectively, and assuming that there are 1500 stars whose masses are  $25 \sim 30 M_\odot$ , the total mass formed would be about  $2 \times 10^6 M_\odot$ , and with lower mass limits se of 0.5,  $3 \times 10^6 M_\odot$ , respectively. The stars of  $\geq 30 M_\odot$  would have already exploded in this cluster. These massive stras lifetimes of less than  $2 \times 10^6 \text{ yr}$ ; see Table1. We adopted mainly their estimates for the total mass of the compact cluster model,  $2 \times 10^6 M_\odot$  to  $3 \times 10^6 M_\odot$ .

### 2.3.1 Stellar evolution

We modeled effects of stellar evolution by changing the mass of each star. At the last stage of stellar evolution, stars lose a significant fraction of their mass in stellar wind and supernova explosions. The potential well of a cluster typically  $10 \text{ km s}^{-1}$  is not deep enough to retain the gas ejects from stars mass, since the escape velocity is only a few times  $10 \text{ km s}^{-1}$ . At birth, a neutron star or black hole recieves a high velocity kick in a random direction. This distribution is flat at velocities below  $250 \text{ km s}^{-1}$ . We assume, therefore that lost mass disappears abruptly from the cluster. We give the mass,  $m(t)$ , of each particle at each time step;

$$m(t) = \begin{cases} m_{ini} & : t + \Delta t < t_{seq} \\ m_{rm} + [m_{ini} - m_{rm}] \frac{(t_{seq} - t)}{\Delta t} & : t < t_{seq} < t + \Delta t \\ m_{rm} & : t_{seq} < t, \end{cases} \quad (8)$$

where  $m_{ini}$  is the initial mass,  $m_{rm}$  is the mass of any remnant after mass loss, and  $t_{seq}$  is the main-sequence time scale (Table 1). We obtain values between the points listed in Table 1 by linear interpolation. The remnant mass,  $m_{rm}$ , is summarized in Table 2. These tables are due to Iben & Renzini (1983) from which these tables have been copied.

We assume that a star with a mass larger than  $40 M_\odot$  leaves a black hole after ejecting its envelope during the main-sequence and Wolf-Rayet phase. The mass of the black hole is assumed as  $0.35 m_{ini} - 12 M_\odot$  (Table 2). Stars with masses between  $8 M_\odot$  and  $40 M_\odot$  are assumed to become neutron stars, through the Type II supernova explosion. Stars with masses between  $4.7 M_\odot$  and  $8 M_\odot$  are assumed to become no remnant, through the Type we supernova explosion. The mass of the white dwarf is taken to be equal to the core mass of its progenitor at the tip of the asymptotic giant branch. Stars with masses less than  $4.7 M_\odot$  are assumed to become white dwarfs. Iben and Renzini 1983 give the final white dwarf mass  $0.53\eta^{-0.082} + 0.15\eta^{-0.35} \times (m/M_\odot - 1.0)$ , and we take the mass-loss rate given by  $\eta = \frac{1}{3}$  for all models. This formula is accurate for initial masses  $m \geq 0.8 M_\odot$ .

### 2.3.2 King model

We used king's model (King 1966) to generate the initial conditons for globular cluster. Thus the distribution function,  $g(E)$ , is a lowered Maxwellian, given by

$$g(E) = K [\exp(-(E - E_t))], \quad (9)$$

Initial Mass $m$ $\log_{10}[m/M_{\odot}]$	Main Sequence Time $\log_{10}[t_{se}/yr]$
-0.08	10.18
-0.01	9.93
0.07	9.63
0.16	9.28
0.27	8.90
0.40	8.50
0.54	8.11
0.72	7.68
0.91	7.33
1.11	7.02
1.33	6.76
1.55	6.57
1.79	6.50

Table 1: Stelalr evolution time\*

Initial Mass ( $M_{\odot}$ )	Remnant Mass ( $M_{\odot}$ )	Comments
$< 4.7$	$0.58 + 0.22 \times (m_{ini} - 1)$	White dwarf
$[4.7, 8.0]$	0	No remnant
$[8.0, 40.0]$	1.4	Neutron star
$[40.0 \sim]$	$0.35m_{ini} - 12$	Black hole

Table 2: Mass evolution\*

\*:Iben&Renzini(1983)

for  $E < E_t = \phi(r_t)$ , where  $E = v^2/2 + \phi(r_t)$ ,  $K$  is constant, and  $r_t$  is the (tidal) radius of the edge of the cluster. (We used  $\phi$  for the Newtonian potential to distinguish it from the the softened potential  $\Phi_c$ .) The King model is determined by the dimensionless center potential,  $W_0 = \beta[\phi(r_t) - \phi(0)]$ . For this paper we performed a survey of models defined by combinations of the value of the dimensionless central potential of the King model,  $W_0 = 5$ .

### 2.3.3 Tidal boundary

During the course of a simulation stars escape from the cluster. The precise dynamical definition of escape is not easy if there is a tidal field, and there we adopt a simple geometric definition: escape are defined to be those stars beyond the tidal radius. All stars within the tidal radius are taken to be members, even though the tidal field is not spherically symmetric. More precisely, we use the distance between the center of the cluster (defined below equation (11)) and the Lagrangian point in the direction of the galactic center as the tidal radius. If, as before, the galaxy is represented by a point mass, it follows that the tidal radius is given by

$$r_t = \left( \frac{M_c}{3M_g} \right)^{1/3} R_g, \quad (10)$$

approximately, where  $M_c$  is the mass of the cluster,  $R_g$  is the distance to the galaxy,  $M_{g(R_g)}$  is the mass of the galaxy within  $R_g$ . If we assumed a spherically galactic potential, taken to be that of a distant point mass  $M_{g(R_g)}$ . Indeed galactic potential is not possible to be point mass model, but the force which a cluster  $R_g$  away from the center of the Galaxy receives can be considered as received from a point mass  $M_{g(R_g)}$ . Here,  $M_c$  is taken to be the total mass of the 'members', and since this depends on  $r_t$  itself, some iteration is usually required. We define the (mass-weighted) center of the cluster by

$$\mathbf{r}_c = \frac{\sum_i^N m_i \mathbf{r}_i}{\sum_i^N m_i}, \quad (11)$$

where  $r_i$  and  $m_i$  is the distance to the galaxy and the mass of  $i$ th particle which fulfills the conditions of ( $r_t \gg |\mathbf{r}_c - \mathbf{r}_i|$ ),  $N$  is the total particle number of cluster.

### 3 Results

We present the results of our simulations, in which we performed a survey of models differing in the slope,  $\alpha$ , of the initial power-law mass function(IMF), and in the dimensionless central potential of King's model,  $W_0$ .

All clusters sank in the center of the Galaxy in the simulation of Point mass. In the point mass model simulations, all clusters sank toward galactic center.

Case	Stellar evolution	Tidal
Point mass	×	×
Point mass	○	×
N particles	×	○

Table 3: Stellar evolution and Tidal force

model	$M_{min}$	$M_{max}$	total mass	$\alpha$
A	$1M_{\odot}$	$100M_{\odot}$	$2.0 \times 10^6 M_{\odot}$	-2.5
B	$1M_{\odot}$	$30M_{\odot}$	$2.0 \times 10^6 M_{\odot}$	-2.5
C	$1M_{\odot}$	$30M_{\odot}$	$2.0 \times 10^6 M_{\odot}$	-3.0
D	$0.5M_{\odot}$	$30M_{\odot}$	$3.0 \times 10^6 M_{\odot}$	-2.5

Table 4: Cluster model

Figure 3 shows the evolution of the total mass of the cluster for each model. Detail discussion for Stellar evolution is in §2.3.1. A massive star is shorter life time and larger difference mass before and after supernova than that of light star component. Therefore, Stellar evolution strongly depend on the slope of IMF,  $\alpha$ .

As can be seen from figure 8, the star cluster of model-A, model-B and model-C ( $M_c = 2 \times 10^6 M_{\odot}$ ) broke, the star cluster of model-D ( $M_c = 3 \times 10^6 M_{\odot}$ ) did not broke. Compared with model-A, model-B, model-C and model-D, we can see that the effect of Stellar mass loss in model-D is little from the other. The total remnant mass in the case of model-D is about twice larger than another model, and about 75% of initial mass of it self. Therefore, we could understand that the minimum mass of IMF is very important.

The life time of the cluster is listed in table.5 for each cluster model. The second column means the life time of the cluster defined as the time while the total mass of cluster was left over 1%, and model D was survived. In the third column, the distance from galactic center at the cluster collapse time (life time), and model-D perfectly sank.

model	life time	Nearest distance [kpc]
A	$1.59 \times 10^8 yr$	0.144
B	$1.71 \times 10^8 yr$	0.102
C	$1.94 \times 10^8 yr$	0.054
D	—survive—	0.005

Table 5: Simulation results. Every case takes into accounted the both effects Stellar evolution and tidal force from galaxy.

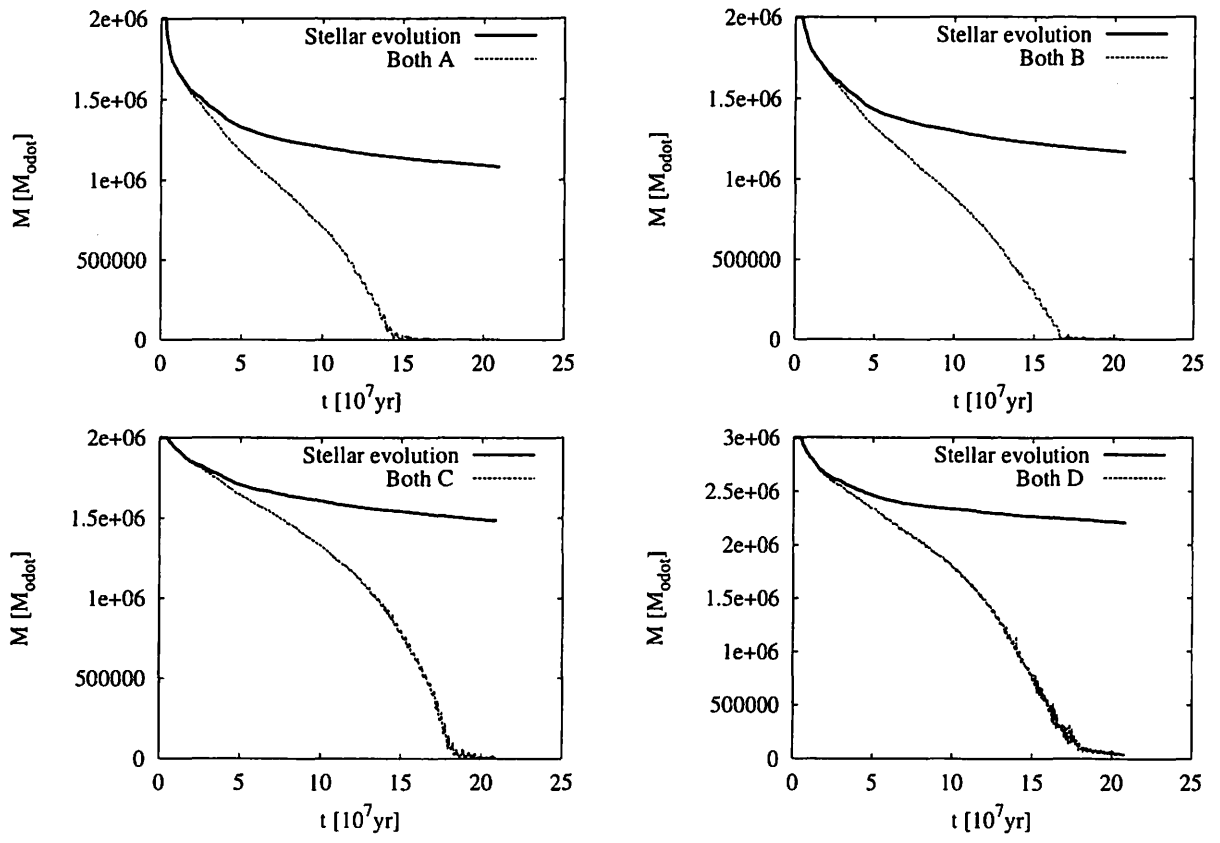


Figure 7: Evolution of mass of cluster. Bold line shows the effect of stellar evolution only, thin dashed line shows the both effects included tidal force.

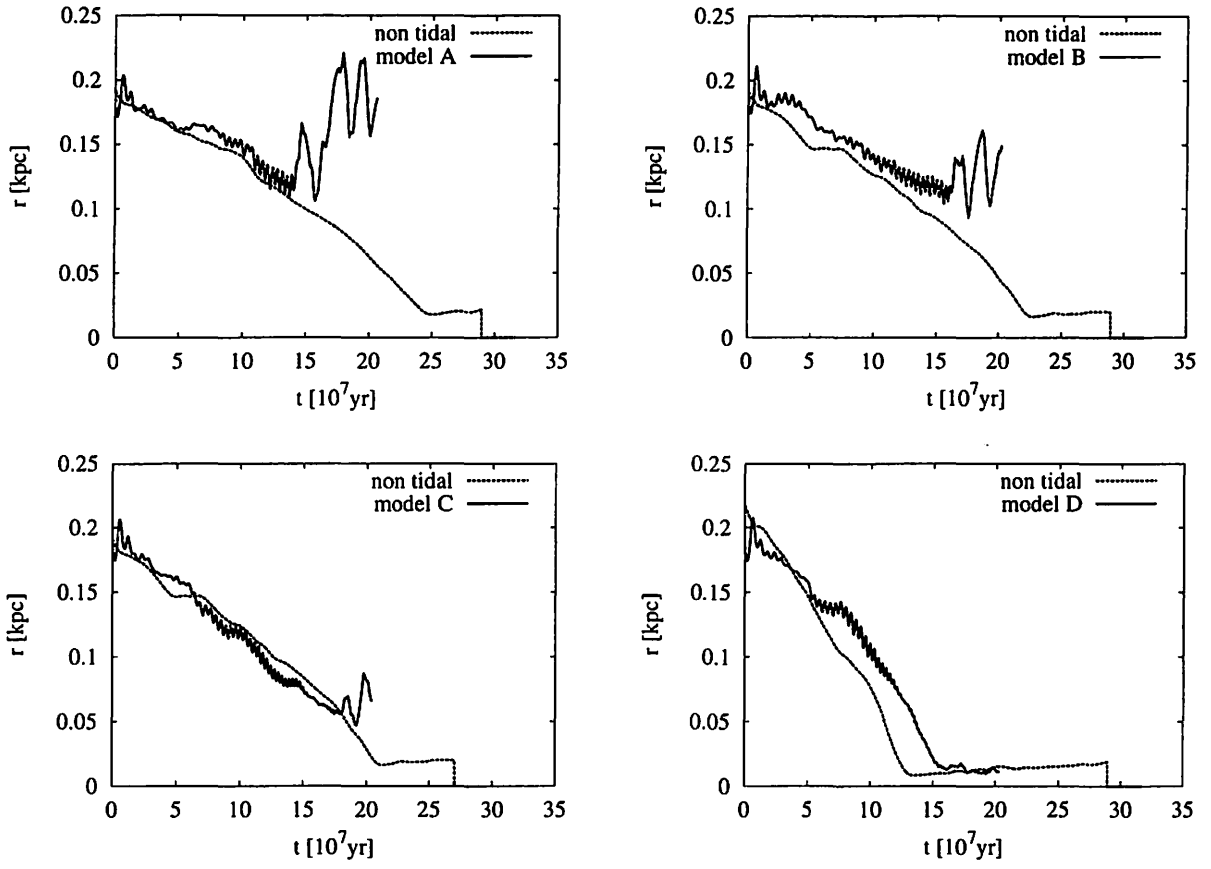


Figure 8: The evolution of total cluster mass. In each figure, the dashed line shows the results of point mass simulations taking account into the effect of stellar evolution, and the bold thick line shows the results of simulations included the both effects, Stellar evolution and tidal force.

## 4 Conclusion and Discussion

We investigated the surviving condition of clusters by means of gravitational N-body simulations which include stellar evolution and tidal stripping. We put a cluster initially at a distance of 200pc from the center of the parent galaxy. We adopted Henquits spherical galaxy model of the truncated galactic radius  $r_b = 0.5\text{kpc}$ , and galactic mass  $M_b = 2.08 \times 10^9 M_\odot$ . We found that the surviving condition strongly depends on the total mass,  $M_c$  of the cluster and initial mass function (IMF). IMF is characterized by three parameters; those are the slope index,  $\alpha$ , the upper limit mass,  $M_{max}$ , and lower limit mass  $M_{min}$ . We found one cluster with  $M_c = 3 \times 10^6 M_\odot$  and  $M_{min} = 0.5 M_\odot$  falls into galactic center for  $1.5 \times 10^8 \text{yr}$ , while the other clusters with  $M_c \leq 2 \times 10^6 M_\odot$  are collapsed before they reach the galactic center. Since the recent observation suggests the compact cluster in M82 galaxy has  $M_c \gg 3 \times 10^6 M_\odot$  (§2.3, it will well survive until it sink toward galactic center within  $1.5 \times 10^8 \text{yr}$ . This supports the formation scenario for SMBHs proposed by Ebisuzaki et al.(2001).

It is well known that there is a linear correlation between the mass,  $M_{BH}$  of central black hole and those of the bulge. The ration between then is about 0.006 (§1.1). This correlation is well explained if the stars, which are born in the clusters but stripped in the precess of the sinking of the cluster, forms the bulge of the galaxies. For example, a thousand of cluster with a mass of  $\sim 3 \times 10^6$  and an IMBH (1000 solar mass) can evolve into a system consist of one bulge with a mass of  $3 \times 10^9$  solar mass and a central black hole with a mass of  $10^6$  solar mass. We will study this connection between the bulge and the central black hole further and report in near future.

I am grateful to many other people who have supported my work: Junichiro Makino and Toshiyuki Fukushima at University of Tokyo, who gave me invaluable advice on the initial condition and the simulation method. Makoto Ideta at Kyoto University, who has instructed me the way to generate initial condition of the simulation.

## A Appendix

### A.1 Relaxation time

From Binney and Tremaine(1987), individual stellar encounters will perturb a star from the course it would take if the other matter of the system were perfectly smoothly distributed only over of order  $0.1N/\ln N$  crossing times. The relaxation time defined as

$$t_{relax} = \frac{0.1N}{\ln N} t_{cross}. \quad (12)$$

Consequently, even if  $N$  is as small as 50, each stars is deflected from its mean trajectory only after several crossing times, and it is possible to obtain some understanding of the dynamics of even small system by investigating the orbits of the stars in a suitable mean potential.

In a globular cluster, on the other hand,  $N \approx 10^5$  and the crossing time  $t_{cross} \approx 10^5 \text{yr}$ , so that stellar encounters may be important over the cluster lifetime of  $10^{10} \text{yr}$ . Indeed, in the core of a globular cluster, where  $t_{cross}$  is very short and  $N \approx 10^4$ , encounters play a key role. But in the case of a cluster of galaxies, or of a globular cluster, as in the case of a galaxy, the fundamental dynamics is that of a collision system in which the constituent particles (galaxies or stars) move under the influence of the mean potential generated by all the other particles.

### A.2 Evaporation time

From time to time an encounter gives enough energy to a star that it can escape from the system. Thus there is a slow but irreversible leakage of stars from system, and in a sense the only permanent equilibrium state of a stellar system consist of two stars in a Kepler orbit, with all the other having escaped to infinity. The timescale over which the stars "evaporate" in this way can be directly related to the relaxation timescale by the following simple argument. The escape speed  $v_e$  at  $\mathbf{x}$  is given by  $v_e^2 = -2\Psi(\mathbf{x})$ . The mean-square escape speed in a system whose density is  $\rho(\mathbf{x})$  is therefore

$$\langle v_e^2 \rangle = \frac{\int \rho(\mathbf{x}) v_e^2 d^3\mathbf{x}}{\int \rho(\mathbf{x}) d^3\mathbf{x}} = -2 \frac{\int \rho(\mathbf{x}) \Psi(\mathbf{x}) d^3\mathbf{x}}{M} = \frac{4W}{M}, \quad (13)$$

where  $M$  and  $W$  are the total mass and potential energy of the system. According to the virial theorem,  $-W = 2K$ , where  $K = \frac{1}{2}M \langle v_e^2 \rangle$  is the total kinetic energy. Here  $\langle v_e^2 \rangle = 4 \langle v^2 \rangle$ .

Thus the root mean square (rms) escape speed is just twice the rms speed. The fraction of particles in a Maxwellian distribution that have speeds exceeding twice the rms speed is  $\gamma = 7.38 \times 10^{-3}$  (Binney and Tremaine(1987)). We can crudely represent the evaporation process as simply removing a fraction  $\gamma$  of the stars every relaxation time. Thus the rate of loss is  $dN/dt = -\gamma N/t_{relax} \equiv -N/t_{evap}$ , where the **evaporation time**, the characteristic time in which the system's stars evaporate, is  $t_{evap} = 136t_{relax}$ . Thus we expect that evaporation sets an upper limit to the lifetime of any bound stellar system of about  $10^2 t_{relax}$ .

### A.3 Dynamical friction

#### Decay of globular cluster orbits

As a globular cluster orbits through a galaxy, it is subject to dynamical friction. This drag causes the cluster to lose energy and spiral in toward the galaxy center. We now estimate the time  $t_{fric}(r_i)$  required for a cluster that is initially on a circular orbit of radius  $r_i$  to reach the center.

The flatness of many observed rotation curves suggests that we approximate the density interior to  $r_i$  ( $r_i > r_0$ : *coreradius*) with the density distribution,

$$\rho(r) = \frac{v_c^2}{4\pi G r^2}, \quad (14)$$

of the singular isothermal sphere with circular speed  $v_c$  and velocity dispersion  $\sigma = v_c/\sqrt{2}$ . Binney and Tremaine (1987) gives the friction force on a cluster of mass  $M$  moving at speed  $v_c$  at radius  $r$  as

$$F = -0.428 \ln \Lambda \frac{GM^2}{r^2}, \quad (15)$$

where  $\Lambda = b_{max} V_0^2 / G(M + m)$ , where  $b_{max}$  is the largest impact parameter that need the system size,  $V_0$  is field stars relative velocity that need the velocity dispersion.

The force eq(15) is tangential and thus causes the cluster to lose angular momentum per unit mass  $L$  at a rate

$$\frac{dL}{dt} = \frac{Fr}{M} \sim -0.428 \frac{GM}{r} \ln \Lambda. \quad (16)$$

Since the cluster continues to orbit at speed  $v_c$  as it spirals to the center, its angular momentum per unit mass at radius  $r$  is at all times  $L = rv_c$ . Substituting the time derivative of this expression into eq(16), we obtain

$$r \frac{dr}{dt} = -0.428 \frac{GM}{v_c} \ln \Lambda. \quad (17)$$

Solving this differential equation subject to the initial condition  $r(0) = r_i$ , we find that the cluster reaches the center after a time

$$\begin{aligned} t_{fric} &= \frac{1.17 r_i^2 v_c}{\ln \Lambda Gm} \\ &\simeq 1.5 \times 10^9 \left( \frac{r}{1kpc} \right)^2 \left( \frac{v_c}{100km/s} \right) \left( \frac{2 \times 10^6 M_\odot}{m} \right) yr. \end{aligned} \quad (18)$$

In reality, some mass will be stripped from the cluster by the galaxy's tidal field. However, for most globular cluster ( $r_i > 1kpc$ ) this process will not greatly lengthen  $t_{fric}$ .



## References

- [1] Barnes J. E., Hut P. 1986, *Nature*, 324, 446
- [2] Begelman, M. C., Blandford, R. D., & Rees, M. J. 1980, *Nature*, 287, 307
- [3] Binney J. and Tremaine S. 1980, *Galactic Dynamics* (Princeton University Press ,Princeton )
- [4] Ebisuzaki T., Makino J., Tsuru T., Funato Y., Portegies Zwart S. F., Hut P., McMillan S., Matsushita S., Matsumoto H. and Kawabe R. 2001, *ApJ*. 562L, 19E
- [5] Fukushige T., and Heggie D. C. 1995, *MNRAS*. 276, 206
- [6] Fukushige T., and Heggie D. C. 2000, *MNRAS*. 318, 753
- [7] Götz M., McKeith C. D., Downes D., and Greve A. 1990, *A&A*. 240, 52
- [8] Harashima 2002, in preparation
- [9] Hernquist L., 1993, *ApJS*. 86, 389
- [10] Iben I., and Renzini A. 1983, *ARA&AJ*. 21, 271
- [11] King I. R. 1966, *AJ*, 71, 64
- [12] Kormendy, J., & Richstone, D. 1995, *ARA&A*, 33, 581
- [13] Magorrian J., Tremaine S., Richstone D., Bender R., Bower G., Dressler A., Faber S. M., Gebhardt K., Green R., Grillmair C., Kormendy J., Lauer, T. 1998, *AJ*, 115, 2285
- [14] Makino J. 1991, *PASJ* 43, 621
- [15] Makino, J., Fukushige, T., Okumura, S. K., & Ebisuzaki, T. 1993, *PASJ*, 45, 303
- [16] Makino, J. 1997, *ApJ*, 478, 58
- [17] Makino J., Taiji M. 1998, *Scientific Simulations with Special-Purpose Computers — The GRAPE Systems*, John Wiley and Sons, Chichester
- [18] Marchant, A. B., & Shapiro, S. L. 1980, *ApJ*, 239, 685
- [19] Matsumoto, H., Tsuru, T.G., Koyama, K., Awaki, H., Canizares, C. R., Kawai, N., Matsushita, S., & Kawabe, R. 2001, *ApJ*, 547, L25
- [20] Matsushita S., Kawabe R., Matsumoto H., Tsuru T., Kohno K., Morita K., Okumura S., Vila-Vialo B. 2000, *Astrophys.J.* 545, L107
- [21] Narumi T., Kawai A., Koishi T. 2001, *Proceedings of SC2001 (ACM, in CD-ROM)*
- [22] Rees, M. J. 1978, *Observatory*, 98, 210
- [23] Rees, M. J. 1984, *ARA&A*, 22, 471
- [24] Sofue Y. 1998, *PASJ*, 50, 227
- [25] Susukita R., Ebisuzaki T., Elmegreen Bruce G., Furusawa H., Kato K., Kawai A., Kobayashi Y., Koishi T., McNiven Geoffrey D., Narumi T., Yasuoka K. 2002, submitted to *J. Comp. Phys.*
- [26] Taniguchi, Y., Shioya, Y., Tsuru, T. G., & Ikeuchi, S. 2000, *PASJ*, 52, 533
- [27] Portegies Zwart, S. F., & McMillan, S. L. W. 2000, *ApJ*, 528, L17
- [28] Portegies Zwart, S. F., Makino, J., McMillan, S. L. W., & Hut, P. 1999 , *A&A*, 348, 117
- [29] Portegies Zwart, S. F., Makino, J., McMillan, S. L. W., & Hut, P. 2001, *ApJ*, 546, L101
- [30] Sugimoto D., Chikada Y., Makino J., Ito T., Ebisuzaki T., Umemura M. 1990, *Nature* 345, 33

# 星団沈降プロセスにおける潮汐力と質量損失の効果

東京工業大学地球惑星 松林達史  
理化学研究所 戎崎俊一

## Abstract

銀河の中で生成された星団は力学的摩擦により銀河の中心に沈んでいく。その際、星団は母銀河からの潮汐力を受け、星を剥がされていく。また星団内の星は超新星爆発により、重たい星から爆発を起こして質量を失っていく。この時、星団内でも力学的摩擦により、重たい星が星団の中心に集まり、それらの重たい星が星団の中心ポテンシャルを支配する。しかし超新星爆発は重たい星から起きていくため中心の密度は大きく影響を受ける。結果として銀河の中心に向かって星団がどこまで落ちれるかというのは、星団モデルの中心ポテンシャルに強く依存し、つまりは星団内の大質量星に依存すると言える。

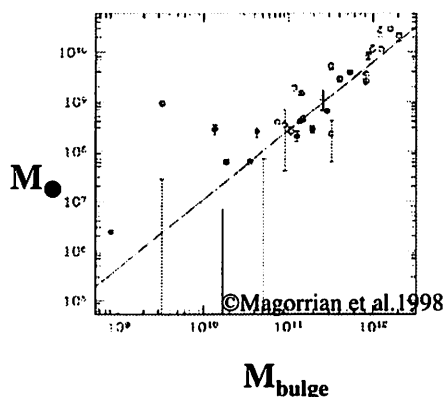
## 巨大ブラックホールの「形成」と「謎」

◆現在観測により、巨大BHと銀河バルジの質量には正の比例関係があることが分かっている。（右図）  
ならば、バルジの形成は巨大BHの形成過程に関係があるのではないだろうか？

◆しかし、近年まで見つかったBHは、  
恒星質量サイズ  $1 \sim 10 M_{\odot}$   
巨大質量サイズ  $10^6 \sim 10^9 M_{\odot}$

大きくわけてこの二つ。

巨大BHの形成過程は1970年代から議論されたが、「ちょっと無理そうだ」ということで謎に包まれたままだった。



●分子雲から直接作る→ 先にclumpができてしまい、BHよりも先に星ができるから ダメ

●星同士の合体から直接作る→ Relaxation time を考えれば宇宙年齢以内に作ることは不可能なので ダメ

# 新種のブラックホール

## 中質量ブラックホールの発見

◆2000年、日本の天文グループによってM82の中心から200pc離れた場所に中質量ブラックホール(1000 $M_{\odot}$ )を含んだ星団が発見された。(右図)

(M82銀河はStar burst galaxy の中では最も我々の銀河に近いものである。)



©Turu(京都大学)

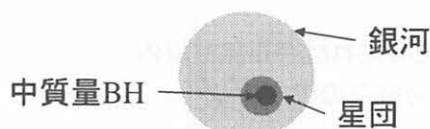
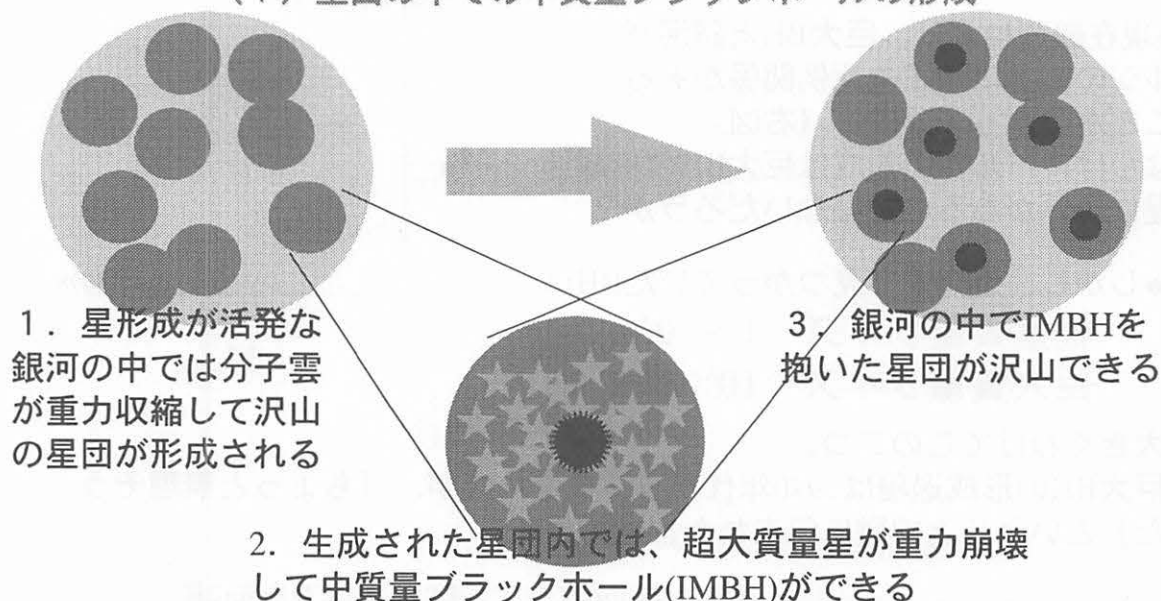


「この発見は巨大ブラックホール形成の足掛かりとなるのではないか！」  
新しいシナリオの提唱(Ebisuzaki et al. 2001)



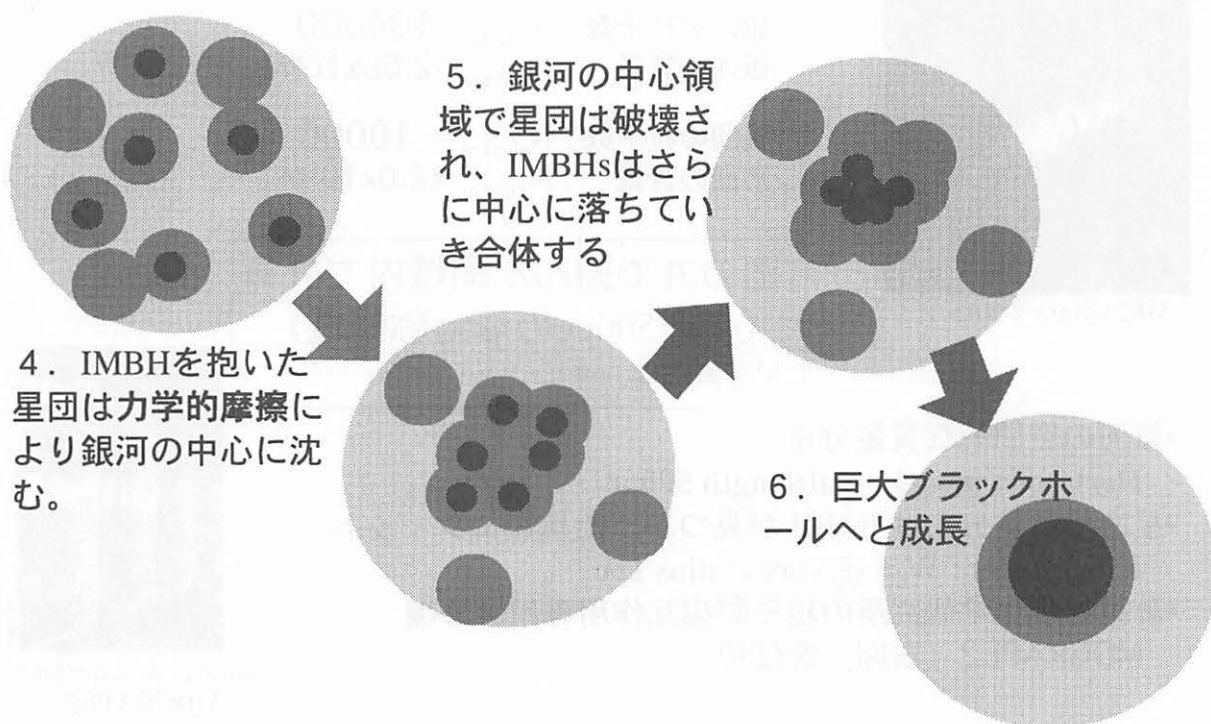
## 巨大ブラックホール形成シナリオ

### (1) 星団の中での中質量ブラックホールの形成



# 巨大ブラックホール形成シナリオ

(2) 星団の銀河中心への落下 → IMBHの合体



## シナリオの問題点と本研究

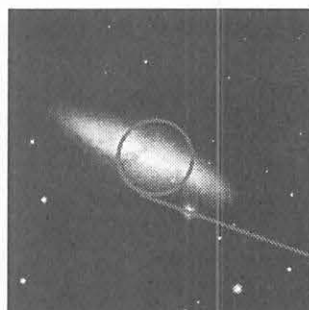
- I. 星団内でIMBHは作れるのか？
- II. 星団は銀河の中心に落ちていけるか？
- III. IMBHが合体して巨大BHが作れるのか？

- I. 星団内でIMBHを作るためには巨大質量星が重力崩壊してしまう前に超大質量星ができる必要がある。約10万年のイベント。
- II. 星団は母銀河からの潮汐力を受けて銀河の中心領域では壊れやすくなる。また、星団内の星の超新星爆発により質量を損失していくため星団自身が質量を失っていく。本研究で扱う問題である。
- III. 銀河の中心で、IMBH同士の合体は重力波放出によるものだが、重力波が効くためにはIMBHが1pcよりも近付かなければいけない。



左図は星団を質点（赤い点）として計算を行ったもの。円盤銀河からの力学的摩擦を受け、約一億年ほどで中心へと沈んでいく。しかしこの計算では星団は潮汐力を受けない。。。

# 銀河と星団のモデル



M82 Galaxy ©NASA

銀河の粒子数 :  $N_{\text{galaxy}} = 104000$   
 銀河の質量 :  $M_{\text{galaxy}} = 2.08 \times 10^9 M_{\odot}$

星団の粒子数 :  $N_{\text{cluster}} = 10000$   
 星団の質量 :  $M_{\text{cluster}} = 2.0 \times 10^6 M_{\odot}$

図の丸で囲んだ領域内で計算  
 (半径500pcの球対称領域)

M82銀河は non-halo disk galaxy

•銀河の空間的な質量分布

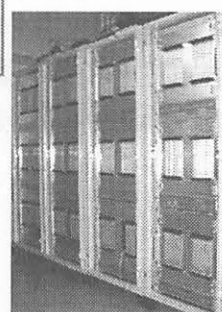
Hernquist model : scalelength 500pc

•星団は質量分布は IMBH が見つかった星団にフィット

King model :  $W_0 = -5$ , core radius 1pc

•計算は理化学研究所の粒子間相互作用専用計算機

MDGRAPE2 (右図) を使用。



Special-purpose computer  
MDGRAPE2

## 星団の質量分布モデル

中質量BHが観測された星団は  $25 \sim 30 M_{\odot}$  の赤色巨星をおよそ 1500個持つことが分っている。これに質量分布をあわせる。

星団の初期の総質量  $2.0 \times 10^6 M_{\odot}$

星の分布のベキ  $\alpha$  -2.5

星の分布の最大質量  $30 M_{\odot}$

星の分布の最小質量  $1 M_{\odot}$

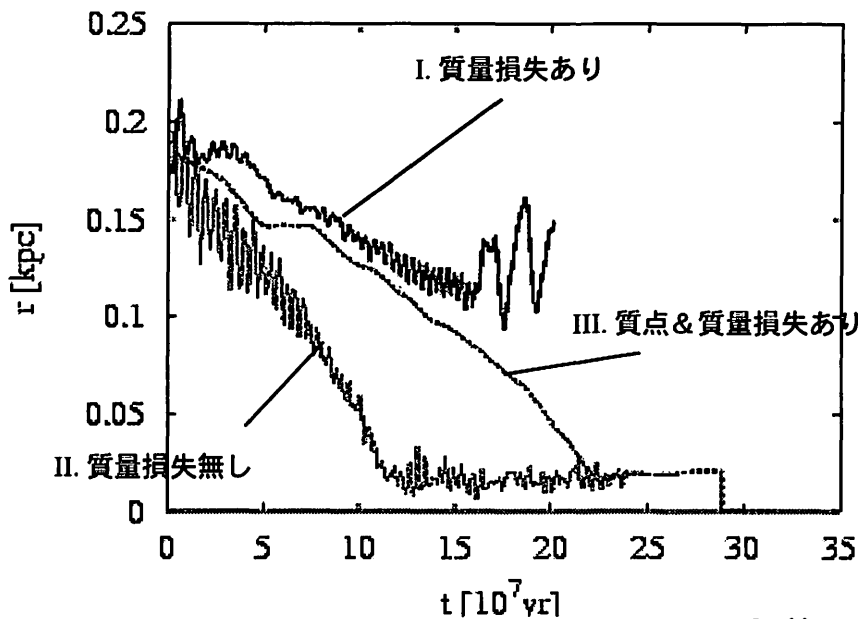
$$dN \propto M^{\alpha} dM$$

Initial Mass Function (IMF)

### モデルの対比

I.	質量損失「あり」	・・・	mass loss & tidal
II.	質量損失「なし」	・・・	----- tidal
III.	質量損失「あり」 & 質点	・・・	mass loss -----

# 計算結果



●質量損失を入れた計算では銀河の中心から100pcで星団は崩壊。  
(Fig1 参照)

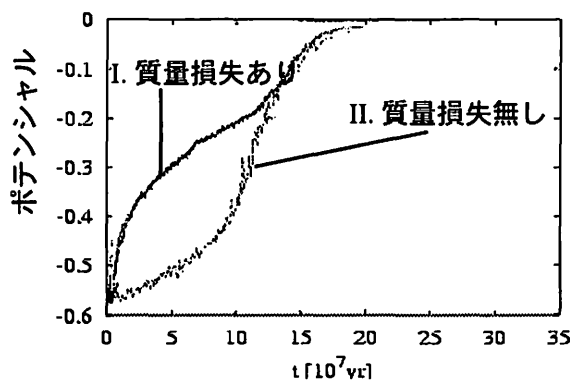
●質量損失無しの計算では銀河中心まで落下。  
(Fig2 参照)

縦軸は銀河の中心からの距離  
横軸は時間

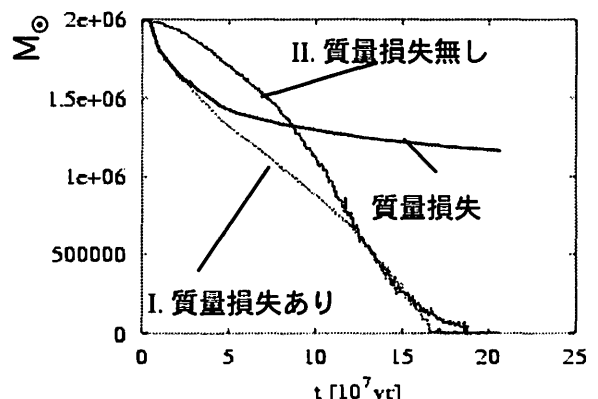
## 定義：星団の崩壊

星団のコア密度が、銀河のその場所での密度と等しくなった時

## 星団中心のポテンシャルの時間進化



## 星団質量の時間進化

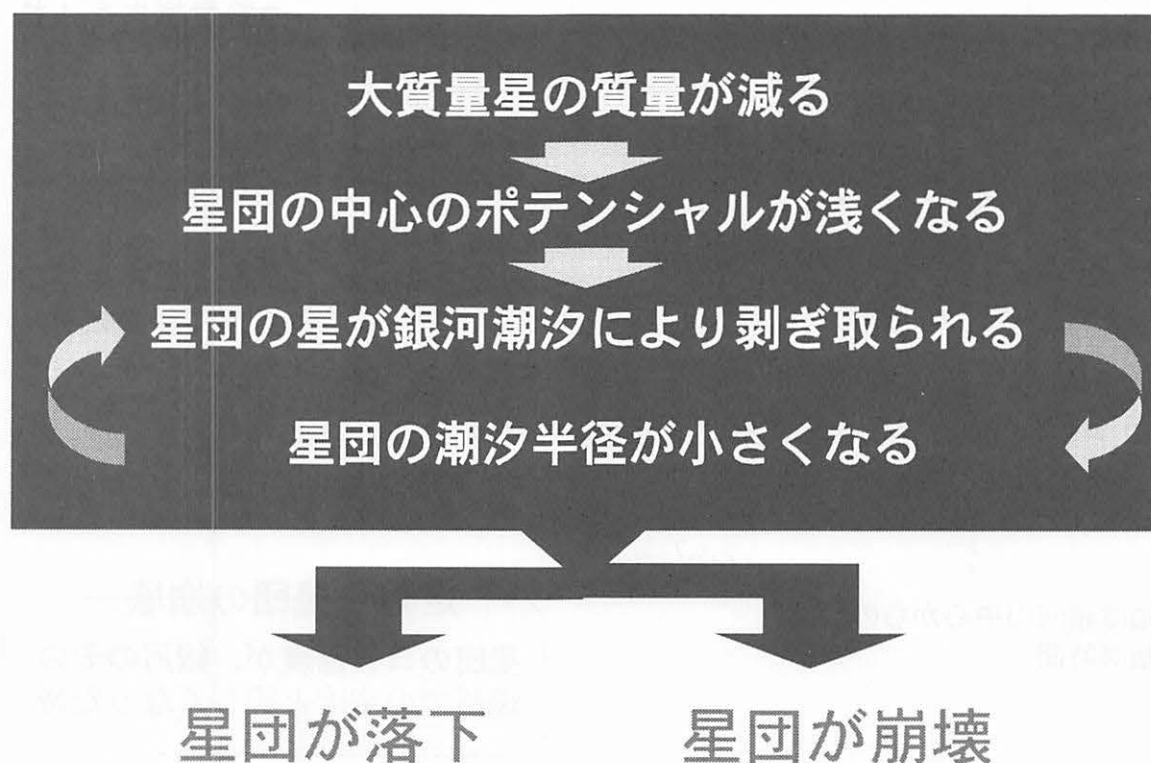


大質量星は力学的摩擦により星団の中心に沈む。  
重たい星は軽い星よりも早く超新星爆発を起こすため、  
超新星爆発は星団の中心で頻繁に起こる。  
星団中心のポテンシャルは徐々に浅くなり、それに伴い  
多数の星が星団から抜けていく。

大質量星の質量損失が大きく影響



# 星の進化の効果



## 結論

- 星団は力学的摩擦により銀河の中心に向かって落ちる。
- 質量損失が無ければ星団は銀河の中心まで落ちる。
- 星団が崩壊するかどうかは、大質量星の質量損失が重要である。

## 今後の計画

- 星団の質量分布を変えて、統計的に見積もる
  - 星の質量分布は、星団内で中質量ブラックホールを形成する重要なパラメーター
  - 巨大ブラックホールを作るパラメーターが決まる？
- 星団から引き剥がされた星の空間分布を計算
  - バルジの形成の関連性
  - バルジ形成の新しいシナリオ？

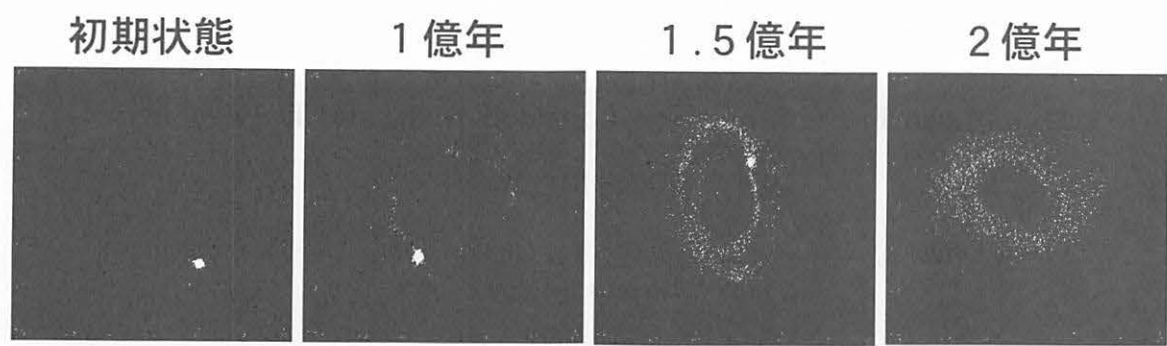


Fig1. I. 質量損失あり

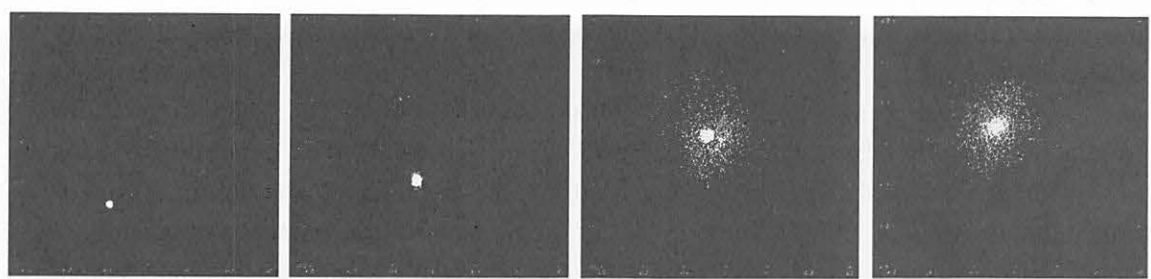
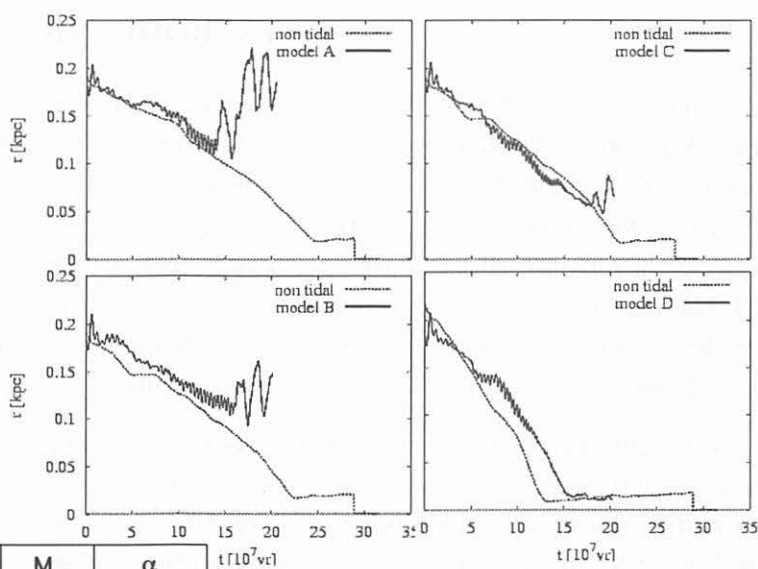


Fig2. II. 質量損失無し

# パラメーターを変えた計算例

最接近距離 & Life time  
 A- 144pc  $1.59 \times 10^8$ yr  
 B- 102pc  $1.71 \times 10^8$ yr  
 C- 54pc  $1.94 \times 10^8$ yr  
 D- 5pc  $1.69 \times 10^8$ yr



## 質量分布

	$M_{\text{cluster}}$	$M_{\text{max}}$	$M_{\text{min}}$	$\alpha$
Model-A	$2.0 \times 10^6 M_{\odot}$	100	1	-2.5
Model-B	$2.0 \times 10^6 M_{\odot}$	30	1	-2.5
Model-C	$2.0 \times 10^6 M_{\odot}$	30	1	-3.0
Model-D	$3.0 \times 10^6 M_{\odot}$	30	0.5	-2.5



# Dynamical Friction between Lopsided Disks and Dark Halos

Makoto IDETA

*Department of Astronomy, University of Tokyo, Tokyo 113-0033, Japan;  
ideta@astron.s.u-tokyo.ac.jp*

## ABSTRACT

The effect of dynamical friction on time evolution of lopsided disks is examined by using a linear perturbation theory. The friction is caused by the gravitational interaction of a rotating lopsided pattern with a density wake induced in halos. The density wake is determined by solving the linearized collisionless Boltzmann and Poisson equations by means of the Fourier-Laplace transform. Then, it is found that dynamical friction always damps a lopsided pattern in our halo model. In addition, the damping time is much shorter than a Hubble time, typically 1 Gyr, unless the pattern speed is quite slow. Considering such a short damping time scale and the observed large fraction of lopsided disks in spirals, say  $\simeq 30$  per cent, it will be unlikely that all of the lopsided disks are recently excited. Thus, it is suggested that most of the observed lopsided disks are very slowly rotating pattern. Significance of weakly damped modes that have a slow pattern speed is discussed.

*Subject headings:* celestial mechanics, stellar dynamics — galaxies: halos — galaxies: kinematics and dynamics — galaxies: structure — method: analytical

## 1. INTRODUCTION

It has long been known that some spiral galaxies have a large-scale lopsided structure (e.g. M101, Arp 1966). Although such a structure is often found at the wavelength of 21 cm (Baldwin, Lynden-Bell, & Sancisi 1980), it is also observed at optical and near-infrared wavelengths (Rix & Zaritsky 1995). Thus, some galactic disks will have a lopsided mass distribution. Moreover, the frequency of lopsided disks in spiral galaxies reaches to half of the H I disks (Richter & Sancisi 1994; Haynes et al. 1998) and one third of the stellar disks (Zaritsky & Rix 1997; Rudnick & Rix 1998; Kornreich, Haynes, & Lovelace 1998). This large fraction of lopsided disks indicates that the lopsidedness would be a repeatedly excited structure or a long-sustained one.

Although the fraction of lopsided disks is large, their origin is not understood well. Theoretically, there may exist a stationary lopsided disk that is responding to the asymmetry of the surrounding dark matter halo potential (Jog 1997, 1999; see also Syer & Tremaine 1996). In addition, Levine & Sparke (1998) considered off-center disks embedded in a flat-cored halo and found that lopsided disks would be maintained for a long time when the disk is orbiting in a retrograde

manner around the halo center. These findings suggest the longevity of lopsided disks. However, in these theoretical studies, the halo is treated as a static potential, and so, the effect of dynamical friction on lopsided disks is not taken into account. Hence, it is necessary to investigate the effect of dynamical friction on dynamical evolution of a lopsided pattern.

One direct way to handle dynamical friction is the Chandrasekhar dynamical friction formula (Chandrasekhar 1943). However, the formula is restricted to a point mass embedded in a uniform, infinite, and non-self-gravitating background, and thus, it cannot be applied to lopsided disks embedded in spherical halos. Another approach to take into account dynamical friction is an  $N$ -body simulation. As an example, a numerical simulation made by Walker, Mihos, & Hernquist (1996) demonstrates that the lopsided structure caused by a minor merger would last for up to 1 Gyr (see also Zaritsky & Rix 1997). However, significant disk thickening, which may affect the evolution of lopsided disks, is also reported in their simulation. Moreover,  $N$ -body simulations with an insufficient number of particles have the problem of discreteness noise. In fact, to achieve a sufficient signal-to-noise ratio, the simulations with a huge number of particles  $N \gtrsim 10^7$  would be required (Weinberg 1998a). However, it is hard to simulate a galaxy with such a huge number of particles. Then, an alternative way that is completely free from discreteness noise is to solve the linearized collisionless Boltzmann and Poisson equations by means of the Fourier-Laplace transform, which is known as the matrix method. The matrix method was first applied to problems in stellar dynamics by Kalnajs (1977) to find the unstable modes of galactic disks. Subsequently, this method was employed by, for example, Palmer & Papaloizou (1987) in the study of the radial orbit instability and was adopted by Weinberg (1989) to study the satellite decay in a spherical halo. A similar approach was used to estimate the bar deceleration rate (Weinberg 1985) and the damping/excitation time scale of galactic warps (Nelson & Tremaine 1995) due to dynamical friction with surrounding dark matter halos.

In this paper, a lopsided pattern rotating in spherical dark matter halos is considered. To examine the lifetime of such a pattern, the effect of dynamical friction on lopsided disks is investigated by using the matrix method. Then, it is found that a lopsided pattern induces a significant density wake in halos. This density wake interacts with the original lopsided pattern through dynamical friction, and then, the friction damps a lopsided pattern within a time scale shorter than a Hubble time unless the rotational period of the pattern is very slow.

This paper is organized as follows. In §2, the method to solve the collisionless Boltzmann and Poisson equations is described. Numerical models and assumptions are also given. Results are shown in §3. In §4, some implications of the results and possible effects of the assumptions on lopsided disks are discussed, and the results are summarized.

## 2. NUMERICAL METHOD AND MODELS

In this paper, dynamical friction is treated as a drag force due to the gravitational interaction of a lopsided disk with a density wake induced in a primary system. Such a density wake is determined by solving the collisionless Boltzmann and Poisson equations by means of the matrix method, which was developed by Kalnajs (1977) (see also Weinberg 1989). To solve these equations, two assumptions are made. First, the amplitude of a lopsided pattern is sufficiently small to adopt a linear perturbation theory. Then, the collisionless Boltzmann equation is linearized. Second, an unperturbed potential is spherical in shape. Such an assumption is employed because in any spherical potential there exist three independent isolated integrals, so that the orbits are analytically solvable. On the other hand, a perturbed potential need not be spherical in shape. Possible effects of these assumptions on the estimate of dynamical friction are discussed in §4.

### 2.1. Matrix Method

In this section, the method to calculate the density wake induced by a perturbed density is mentioned. The method used in this paper is the same as that in Weinberg (1989), who described the method in detail. Then, details should be referred to Weinberg (1989), although the principal formulae are summarized in Appendix A.

The density wake in a primary system will be determined by coupled-solutions of the linearized collisionless Boltzmann and Poisson equations,

$$\frac{\partial f_1}{\partial t} + \frac{\partial f_1}{\partial \mathbf{w}} \frac{\partial H_0}{\partial \mathbf{I}} - \frac{\partial f_0}{\partial \mathbf{I}} \frac{\partial H_1}{\partial \mathbf{w}} = 0, \quad (1)$$

$$\nabla^2 \Phi_1 = 4\pi G \rho_1, \quad (2)$$

where the subscript 0 denotes the equilibrium quantities of the collisionless Boltzmann equation and the subscript 1 denotes the first order perturbation of a six-dimensional distribution function  $f$ , a Hamiltonian  $H$ , a potential  $\Phi$ , and a density  $\rho$ . The collisionless Boltzmann equation is described by action-angle variables,  $(\mathbf{I}, \mathbf{w})$ .

Let  $\Phi_1^{\text{res}}$  be the response potential to an external potential  $\Phi_1^{\text{ext}}$ . Then, the perturbed potential  $\Phi_1$  will be written as the sum of  $\Phi_1^{\text{ext}}$  and  $\Phi_1^{\text{res}}$ . Similarly, the perturbed density will be expressed as  $\rho_1 = \rho_1^{\text{ext}} + \rho_1^{\text{res}}$ . Here, the response density  $\rho_1^{\text{res}}$  to an external density  $\rho_1^{\text{ext}}$  is related to the perturbed distribution function,  $\rho_1^{\text{res}} = \int d^3\mathbf{v} f_1$ . Thus, the linearized Boltzmann-Poisson equation is an integrodifferential equation. Such an equation would be simplified by means of the Fourier-Laplace transform. Then, in the frequency-domain, the linearized Boltzmann-Poisson equation becomes a simple algebraic equation for a particular harmonic  $(l, m)$ , or the matrix equation,

$$\tilde{\mathbf{A}}^{lm}(\omega) = \mathcal{D}^{-1lm} \cdot \mathcal{R}^{lm}(\omega) \cdot \tilde{\mathbf{B}}^{lm}(\omega), \quad (3)$$

where the dispersion matrix  $\mathcal{D}^{lm}$  is

$$\mathcal{D}^{lm} = \mathcal{I} - \mathcal{R}^{lm}. \quad (4)$$

Here,  $\mathcal{I}$  is a unit matrix with the same rank as  $\mathcal{R}^{lm}$ , the response matrix  $\mathcal{R}^{lm}$  describes the information on a primary system,  $A^{lm}$  and  $B^{lm}$  are the expansion coefficients of response and external potentials in biorthonormal basis sets, respectively (Clutton-Brock 1972; Hernquist & Ostriker 1992), and the tilde denotes the Laplace transform in the time variable.

To find the time dependence of expansion coefficients, the inverse Laplace transform for equation (3) is required. In this paper, to avoid the transient wave that originates from an initial condition, the time asymptotic approximation ( $t \rightarrow \infty$ ) is adopted. Then, in the time-domain, the matrix equation is

$$A^{lm}(t) = \mathcal{D}^{-1lm}(m\Omega_p) \cdot \mathcal{R}^{lm}(m\Omega_p) \cdot B^{lm}(t), \quad (5)$$

where  $\Omega_p$  is the pattern speed of a lopsided pattern. Here, there may exist weakly damped modes (see Weinberg 1994) that satisfy the relation

$$\tilde{A}^{lm}(\omega_d) = \mathcal{R}^{lm}(\omega_d) \cdot \tilde{A}^{lm}(\omega_d), \quad (6)$$

or the dispersion relation

$$\det \mathcal{D}^{lm}(\omega_d) = 0, \quad (7)$$

where  $\omega_d$  is the complex frequency of each weakly damped mode. When weakly damped modes should be included, e.g., the damping time of the mode,  $\Im(\omega_d)^{-1}$ , is longer than a Hubble time, one must use equation (A21) instead of equation (5).

Once the expansion coefficients of the response density  $A^{lm}$  are calculated via the matrix equation for each harmonic ( $l, m$ ), the response density and potential, which are both real functions, can be found straightforwardly,

$$\rho_1^{\text{res}}(\mathbf{r}, t) = \frac{1}{2} \sum_{n,l,m} \left[ A_n^{lm}(t) d_n^{lm}(r) Y_{lm}(\theta, \phi) + A_n^{lm*}(t) d_n^{lm*}(r) Y_{lm}^*(\theta, \phi) \right], \quad (8)$$

$$\Phi_1^{\text{res}}(\mathbf{r}, t) = \frac{1}{2} \sum_{n,l,m} \left[ A_n^{lm}(t) u_n^{lm}(r) Y_{lm}(\theta, \phi) + A_n^{lm*}(t) u_n^{lm*}(r) Y_{lm}^*(\theta, \phi) \right]. \quad (9)$$

Here, the asterisk denotes a complex conjugate and  $Y_{lm}(\theta, \phi)$  are the spherical harmonics.  $u_n^{lm}(r)$  and  $d_n^{lm}(r)$  are potential and density basis functions, respectively, which are normalized by

$$-\frac{1}{4\pi G} \int dr r^2 u_n^{lm*}(r) d_{n'}^{lm}(r) = \delta_{nn'}. \quad (10)$$

Then, the gravitational torque felt by a lopsided pattern,  $\tau_z$ , can be written

$$\tau_z = \int d^3r \rho_1^{\text{ext}}(\mathbf{r}, t) \left[ -\frac{\partial \Phi_1^{\text{res}}(\mathbf{r}, t)}{\partial \phi} \right] = -8\pi G \sum_{l=1}^{\infty} \sum_{m=1}^l m \Im \left[ A^{lm}(t) \cdot B^{lm*}(t) \right], \quad (11)$$

where  $\Im$  denotes the imaginary part.

Here, it will be useful to consider the non-self-gravity case, which corresponds to  $\Phi_1 = \Phi_1^{\text{ext}}$  and  $\rho_1 = \rho_1^{\text{ext}}$ , for examining a general property of dynamical friction. Then, the gravitational torque can be written by a further simple formula, or the Lynden-Bell & Kalnajs (1972) formula

$$\tau_z = -8\pi G \sum_{l=1}^{\infty} \sum_{m=1}^l m \left[ \Im(\mathcal{R}^{lm}) \cdot B^{lm} \right] \cdot B^{lm}. \quad (12)$$

Hence, it is not necessary to know the real part of a response matrix  $\mathcal{R}^{lm}$  for calculating the gravitational torque due to dynamical friction in this case. Furthermore, when  $f_0$  depends only on the energy  $E$ , using the explicit formula for a response matrix (eq. [A19]), the gravitational torque can be rewritten

$$\begin{aligned} \tau_z = & \sum_{l=1}^{\infty} \sum_{m=1}^l \frac{32\pi^4 m^2 \Omega_p}{2l+1} \iint \frac{dE dJ}{\Omega_1(E, J)} \sum_l \frac{df_0}{dE} |Y_{l2}(\pi/2, 0)|^2 \left| \sum_n B_n^{lm} W_{l2m}^{l_1 n}(E, J) \right|^2 \\ & \times \delta(m\Omega_p - l_1\Omega_1 - l_2\Omega_2), \end{aligned} \quad (13)$$

where  $\delta$  denotes the Dirac delta function, both  $l_1$  and  $l_2$  are integers, and  $\Omega_{1(2)}$  is an angular frequency of the angle variable  $w_{1(2)}$ . Here,  $w_1$  and  $w_2$  are the conjugate variables with respect to the radial action  $I_r$  and angular momentum  $J$ , respectively. The potential transform  $W_{l2m}^{l_1 n}$  is defined by equation (A8). Meantime, in most galactic models whose distribution function depends only on the energy, the relation  $df_0/dE < 0$  is satisfied for any energy  $E$ ; the halos with such a distribution function are sometimes called IDDF (isotropic decreasing distribution function) halos (e.g., Goodman 1988). Then, as seen from equation (13), in IDDF halos, the gravitational torque is negative (positive) when a lopsided pattern is orbiting in a prograde (retrograde) manner. In addition, the above equation contains the Dirac delta function. Thus, the gravitational torque is caused by the resonance stars that satisfy the resonance condition

$$l_1\Omega_1 + l_2\Omega_2 = m\Omega_p. \quad (14)$$

In the epicyclic approximation,  $\Omega_1$  and  $\Omega_2$  are equal to the epicyclic frequency  $\kappa$  and the orbital frequency  $\Omega$ , respectively.

## 2.2. Model Description

Here, halo and disk models are described. In this paper, it is assumed that an unperturbed potential is spherical in shape and that the orbits of halo particles are not affected by the existence of flat disks. In addition, galactic disks are assumed to be stable against  $m = 1$  distortions. This means that the amplification of lopsided perturbation due to the interaction with disk stars is ignored. Lopsided instabilities are discussed in §4.

First, the halo model is an isotropic King model. The explicit formula of a six-dimensional distribution function is

$$f_0(\mathcal{E}) = \begin{cases} \frac{\rho_1}{2\sqrt{2}\pi^{3/2}\sigma^3} (e^{\mathcal{E}/\sigma^2} - 1) & \mathcal{E} \geq 0; \\ 0 & \mathcal{E} < 0. \end{cases} \quad (15)$$

Here,  $\mathcal{E} = -E + E_0$ . Parameters  $E_0$ ,  $\rho_1$ , and  $\sigma$  are chosen such that a normalized central potential  $W_0 = 3.0$ , a total mass  $M = 6.0 \times 10^{11} M_\odot$ , and a tidal radius  $R_t = 200$  kpc. These choices of parameters are appropriate for the Milky Way (e.g., Kochanek 1996) by combining with a standard exponential disk. This halo model is the same model as used in Weinberg (1998b) and Vesperini & Weinberg (2000).

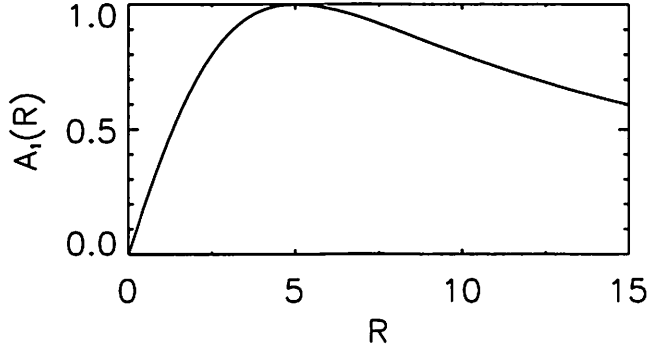


Fig. 1.— Amplitude of the  $m = 1$  component of disks as a function of radius for  $d = 10.0$ .

Second, the disk model is a lopsided exponential disk model. Since the lopsided pattern is assumed to be sufficiently small as compared with the background density, the model disk could be written

$$\Sigma_d(R, \phi, t) \simeq \Sigma_0(R) + \Sigma_1(R, \phi, t) \quad (16)$$

$$= \frac{M_d}{2\pi} \exp(-R) [1 + A_1(R) \cos(\phi - \Omega_p t)], \quad (17)$$

where  $M_d$  is the disk mass,  $A_1(R)$  is the amplitude of the  $m = 1$  Fourier component of disks at a particular radius  $R$ , and  $\Omega_p$  is a pattern speed of lopsided disks. Clearly, the surface density of a lopsided pattern  $\Sigma_1$  is

$$\Sigma_1(R, \phi, t) = \frac{M_d}{2\pi} \exp(-R) A_1(R) \cos(\phi - \Omega_p t). \quad (18)$$

In this paper, the functional form of  $A_1$  that determines the shape of the lopsidedness is chosen to express the lopsidedness induced by a fly-by encounter (Vesperini & Weinberg 2000), which yields

$$A_1(R) \equiv \frac{Rd}{25 + R^2}, \quad (19)$$

where  $d$  is the parameter that determines the amplitude. Here, the  $m = 1$  Fourier amplitude  $A_1(R)$  is shown in Figure 1. Then,  $A_1(R)$  reaches to the maximum value  $d/10.0$  at  $R = 5$ . Thus, for  $d = 2.0$ , the maximum amplitude is 0.2. This value is related to the observational fact that about one third of field spirals have a lopsided mass distribution with the  $m = 1$  Fourier amplitude larger than 0.2 at 1.5 to 2.5 disk scale lengths (Rix & Zaritsky 1995; Zaritsky & Rix 1997; Rudnick & Rix 1998).

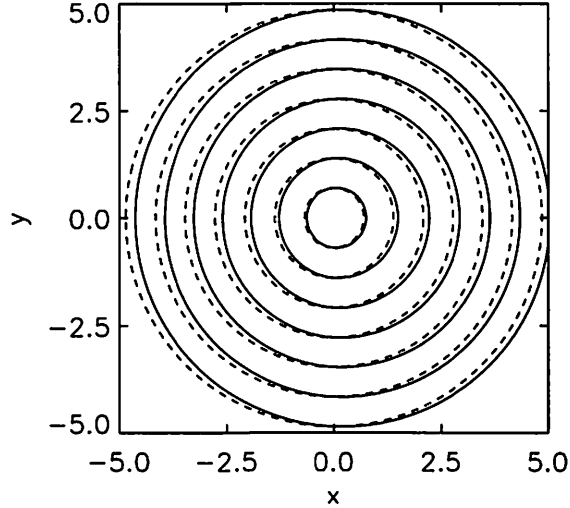


Fig. 2.— Surface density profiles of the lopsided exponential disk model. The solid contours show a lopsided model with  $d = 2.0$  in which the maximum density occurs at  $(0, 0)$ , while the dashed contours represent the corresponding on-center model as a marked reference to a lopsided structure. The contour levels are 1.0, 1/2, ..., 1/128 of the maximum density.

Finally, the external density  $\rho_1^{\text{ext}}(\mathbf{r}, t)$  is set to be

$$\rho_1^{\text{ext}}(\mathbf{r}, t) = \Sigma_1(R, \phi, t) \delta(z). \quad (20)$$

The surface density profiles of  $\Sigma_d$  and  $\Sigma_1$  are represented in Figures 2 and 3, respectively.

When each annulus with a radius  $R$  is displaced at  $\Delta R(R)$  from the disk center in the  $x$ -direction, the surface density profile can be written

$$\Sigma_d(R, \phi) \simeq \Sigma_0(R) + \frac{d\Sigma_0}{dx} \Delta R(R) \quad (21)$$

$$= \Sigma_0(R) + \frac{d\Sigma_0}{dR} \Delta R(R) \cos(\phi - \Omega_p t). \quad (22)$$

As found from equations (17) and (22),  $A_1(R)$  is also considered the displacement of an annulus with the radius  $R$ . Therefore, the angular momentum associated with a lopsided motion,  $L_{z,\text{ext}}$ , can be written

$$L_{z,\text{ext}} = \int_0^\infty dR 2\pi R \Sigma_0(R) A_1(R)^2 \Omega_p \quad (23)$$

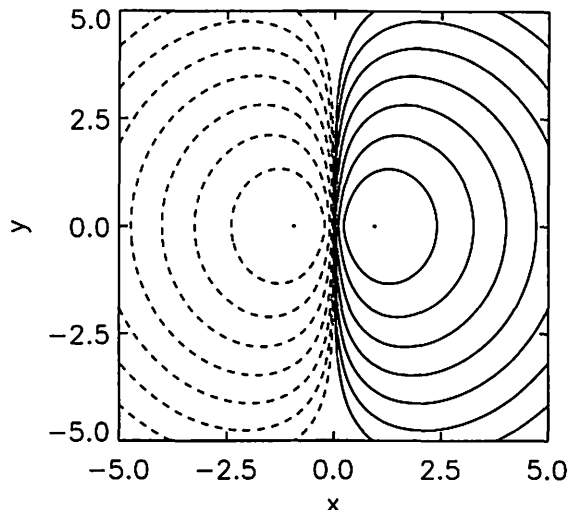


Fig. 3.— Surface density contours of the lopsided pattern on the disk plane. The solid (dashed) contours show overdensity (underdensity). The contour levels are  $\pm 1.0, \pm 1/2, \dots, \pm 1/128$  of the maximum density, which occurs at  $(0.933, 0)$ .

### 2.3. Numerical Procedure

Units are chosen such that the disk mass  $M_d = 1$  and the exponential scale length  $R_d = 1$ . The gravitational constant is set to be 1.35 so that the circular velocity at the solar radius, i.e., 8.5 kpc, is equal to  $220 \text{ km s}^{-1}$ . If these units are scaled to physical values appropriate for the Milky Way, i.e.,  $R_d = 3.5 \text{ kpc}$  and  $M_d = 6.0 \times 10^{10} M_\odot$ , unit time and velocity are  $1.09 \times 10^7 \text{ yr}$  and  $321 \text{ km s}^{-1}$ , respectively.

A biorthogonal basis set to expand the density-potential pair is numerically obtained by solving the Sturm-Liouville problem according to the method described in Weinberg (1999). The angle variables,  $w_1$  and  $w_2$ , and the potential transforms,  $W_{l_2 m}^{l_1 i}$ , are calculated on a  $1000 \times 100$  grid in  $E$  and  $\kappa \equiv J/J_{\max}(E)$  by using the Romberg method with the error tolerance parameter being  $10^{-4}$ . Here, the summations over  $l$  and  $l_1$  in calculating the response matrix (see eq. [A19]) must be truncated at  $l = l_{\max}$  and  $|l_1| = l_{1,\max}$ , respectively. These truncation parameters are chosen such that  $l_{\max} = 5$  and  $l_{1,\max} = 10$ . The expansions of the potential and density in radial basis functions are also truncated at  $n_{\max} = 30$ . With varying parameters, these choices of truncations seem to have an accuracy of 1 per cent for calculating the gravitational torque.

Since the numerical techniques used here are rather complex, it is useful to run a test calculation for checking the validity of our numerical implementation. Here, the test introduced by Weinberg (1989) is done. This test is to calculate the response density and potential when the constant force field that rotates at a constant rate  $\Omega_p$  is imposed to halos as a perturbation. Since the perturbation is a constant force field, no torque acts on the perturbation, and the response of



halos is a barycentric shift against the perturbation (see Weinberg 1989, for a detail). Then, the calculated response potential agrees with a predicted potential to an accuracy of 1 per cent at the range of the pattern speed  $\Omega_p = [0.003 : 0.1]$  for the calculated range  $\Omega_p = [0.001 : 0.1]$ . For  $\Omega_p = 0.002$ , the difference between calculated and predicted values is about 4 per cent, for  $\Omega_p = 0.001$ , the difference reaches to 15 per cent.

### 3. RESULTS

Using the unperturbed model  $f_0(E)$  and perturbed model  $\rho_1^{\text{ext}}(\mathbf{r}, t)$  described in the previous section, the gravitational torque can be calculated by equation (11). Here, since the external density  $\rho_1^{\text{ext}}$  is proportional to  $M_d d$ , its expansion coefficients  $\mathbf{B}^{lm}(t)$  depend also on  $M_d d$ . In addition, the response matrix  $\mathcal{R}^{lm}$  contains no information on the strength of the perturbation. Hence, the gravitational torque  $\tau_z$  is dependent on  $M_d^2 d^2$ .

To quantify the effect of dynamical friction on lopsided disks, it is useful to define the rate of angular momentum change,

$$T = -\frac{\dot{L}_z}{L_z} = -\frac{\tau_z}{L_z}, \quad (24)$$

where  $L_z$  is the angular momentum associated with the lopsided motion, which can be written

$$L_z = L_{z,\text{ext}} + L_{z,\text{res}} \quad (25)$$

Here,  $L_{z,\text{ext}}$  is defined by equation (23) and  $L_{z,\text{res}}$  is defined in the same manner as  $L_{z,\text{ext}}$ ,

$$L_{z,\text{res}} = \int_0^\infty dr 4\pi r^2 \rho_0(r) \Delta r(r)^2 \Omega_p, \quad (26)$$

where  $\rho_0$  is the unperturbed halo density and  $\Delta r$  is the displacement of a shell with the radius  $R$  from the halo center. Clearly, both  $L_{z,\text{ext}}$  and  $L_{z,\text{res}}$  depend on  $d^2$ , and then, the rate of angular momentum change  $T$  is independent of the amplitude  $d$ . This is an expected result since the calculations are linear.

If  $T$  is positive (negative), it will be natural to guess that dynamical friction damps (excites) the lopsided pattern. Here, in IDDF halos and without the self-gravity of a wake,  $\bar{\tau}_z$  is negative (positive) when  $\Omega_p$  is positive (negative) as seen in §2.1. This is found to be valid even if the self-gravity is included. Then,  $T$  is always positive, and so dynamical friction always damps the lopsided pattern in our halo models. In addition, since  $\tau_z$  and  $L_z$  are both odd functions of the pattern speed, the damping rate is an even function  $T(\Omega_p) = T(-\Omega_p)$ , which is expected by the symmetry of the unperturbed system.

In Figure 4, the rate of angular momentum change  $T$  is shown in the unit of Gyr as a function of the pattern speed  $\Omega_p$ . The units are scaled to the values suitable for the Milky Way. The solid line shows the friction due to the self-gravitating response and the dashed line shows the friction

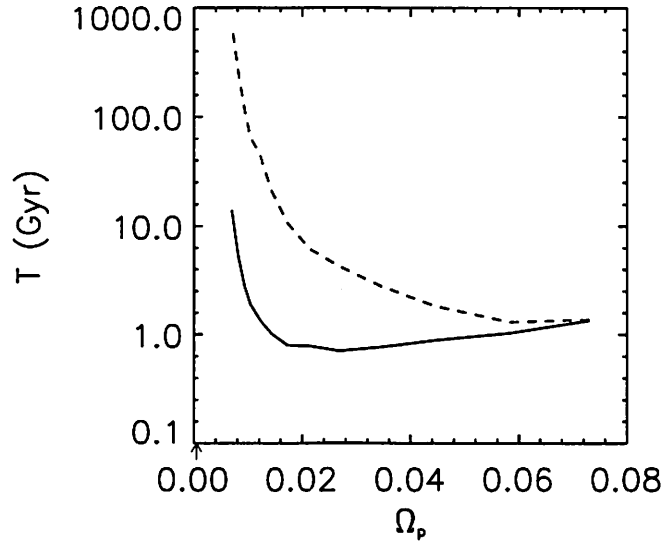


Fig. 4.— Damping time defined by the rate of angular momentum change as a function of the pattern speed for a self-gravitating response (solid line) and non-self-gravitating response (dashed line). The units are scaled to the values suitable for the Milky way. The arrow at the left side shows the real part of the frequency of weakly damped modes.

due to the non-self-gravitating response for a reference. Clearly, there exists a large difference between the self-gravitating and non-self-gravitating responses. In addition, to see the corotation radius corresponding to each pattern speed, orbital frequencies at a particular radius  $R$  are shown in Table 1. As seen in Figure 4, the damping time  $T$  reaches to the minimum value  $\simeq 1.0$  Gyr at  $\Omega_p \simeq 0.03$ . Moreover, the damping rate is shorter than a Hubble time  $\simeq 10$  Gyr when the corotation radius is smaller than the truncation radius,  $R_t = 56$  (see Table 1).

#### 4. DISCUSSION AND CONCLUSIONS

In this paper, lopsided perturbation to stable galactic disks is considered. To estimate the lifetime of the lopsided perturbation, the effect of dynamical friction on lopsided disks is examined. Then, it is found that dynamical friction always damps the lopsided perturbation. The damping time scale of lopsided disks is obtained as a function of the pattern speed. Figure 4 shows the damping time for the self-gravitating and non-self-gravitating cases. Then, the difference between self-gravitating and non-self-gravitating responses is found to be quite large, especially for the slow pattern speed. This could be understood as follows. When the pattern speed of the lopsidedness becomes slower, since the amplitude of the lopsidedness is constant, the barycentric shift of halos becomes larger. In addition, the non-self-gravitating response does not contain the proper information for the barycentric shift, and so, the difference between self-gravitating and non-self-gravitating

Table 1: Orbital frequencies at a particular radius.

Radius	Freq.
5.0	$7.30 \times 10^{-2}$
10.0	$5.84 \times 10^{-2}$
15.0	$4.47 \times 10^{-2}$
20.0	$3.44 \times 10^{-2}$
25.0	$2.69 \times 10^{-2}$
30.0	$2.14 \times 10^{-2}$
35.0	$1.72 \times 10^{-2}$
40.0	$1.44 \times 10^{-2}$
45.0	$1.22 \times 10^{-2}$
50.0	$1.04 \times 10^{-2}$
56.0	$8.78 \times 10^{-3}$

responses will become large with decreasing the pattern speed. Furthermore, Figure 4 clearly shows that the damping time is shorter than a Hubble time unless the rotational period of lopsided disks is quite long.

Here, let us consider the case that the pattern speed is so fast that the corotation radius is within a few optical radius, i.e.,  $\Omega_p \gtrsim 0.04$ . Then, the damping time scale is typically  $\simeq 1$  Gyr. Since the fraction of the lopsidedness in spirals reaches to about 30 per cent, an excitation should occur per  $\simeq 3$  Gyr. One possible recurrent excitation mechanism is a gravitational interaction between host and satellite galaxies. However, since such interaction can easily thicken galactic disks (e.g., see Tóth & Ostriker 1992; Walker et al. 1996), it will be unlikely that most of lopsided disks are excited repeatedly or recently. Even though some of the lopsidedness may be recently excited (Rudnick, Rix, & Kennicutt 2000), it is suggested that most of lopsided disks will be slowly rotating pattern. Unfortunately, the pattern speed of lopsided disks remains unknown observationally. Then, some implications of the results to theoretical models are discussed here.

First, Baldwin et al. (1980) proposed the scenario that a lopsided disk is the pattern that consists of elongated orbits, which is a similar idea to the Lindblad's idea (1963, references therein) of kinematic spiral arms. Then, such a pattern rotates at an angular frequency  $(\Omega - \kappa)(r)$ . Here,  $\Omega$  and  $\kappa$  are the orbital and epicyclic frequencies, respectively. Since the relation  $\Omega(r) \leq \kappa(r) \leq 2\Omega(r)$  is satisfied, the pattern rotates in a retrograde manner. In the King model used here, the pattern speed  $(\Omega - \kappa)(r)$  is almost comparable with the orbital frequency  $-\Omega(r)$ . If  $r$  is set to be an optical radius  $R_{\text{opt}} \simeq 3R_d$ , the pattern speed is about  $-0.08$ . Therefore, such a pattern may damp owing to dynamical friction within a Hubble time.

Second, during a fly-by encounter, a lopsided pattern that rotates at an angular frequency  $V_{\text{rel}}/R_p$  would be induced. Here,  $R_p$  is the pericentric radius and  $V_{\text{rel}}$  is the velocity of a perturber

relative to the primary system at  $R_p$ . The relative velocity will be comparable to the rotational velocity,  $\simeq 200 \text{ km s}^{-1}$ , and hence, the pattern speed would be close to the orbital frequency at  $R_p$ . To induce significant lopsidedness, a small  $R_p$  is favorable. Then, the rotational speed of the significant lopsided pattern would be large. Hence, such a pattern may also damp.

On the other hand, Vesperini & Weinberg (2000) found that weakly damped modes were induced by interactions and that they would play an important role in the lopsidedness. To see the effect of such damped modes, the frequencies of such modes are calculated according to equation (7). Here, the numerical procedure is in the same manner as Weinberg (1994). Then, weakly damped modes with complex frequencies  $\omega_d = (\pm 5.0 \times 10^{-4}, -1.7 \times 10^{-6})$  are found only for  $(l, m) = (1, \pm 1)$ . This value is consistent with the result of Weinberg (1994) within the error. Here, since  $l_1\Omega_1 + l_2\Omega_2$  reaches to the minimum value  $5.0 \times 10^{-3}$  for  $(l_1, l_2) = (1, -1)$  (see equation [14]), there would exist no resonance star with such slowly rotating weakly damped modes. Then, the possibility of  $\Im(\omega_d) = 0$  cannot be excluded. In addition, the weakly damped mode is found to have a very slow pattern speed, which is indicated by an arrow in Figure 4. Then, such modes could survive for a long time against the friction. Furthermore, the excitation of such weakly damped modes can be calculated by using equation (A21). Then, it is found that the perturbation used here can excite strong modes and that the peak density of such modes is one order of magnitude larger than that of the self-gravitating response. This is consistent with the finding of Vesperini & Weinberg (2000). In addition, such modes can be also easily excited by the fly-by encounters (Murali 1999; Vesperini & Weinberg 2000). Thus, the weakly damped modes would play an important role in the lopsidedness, as suggested in Vesperini & Weinberg (2000). However, as mentioned in the end of §2.3, in the slow pattern speed region, the linear perturbation techniques would cause a large error. This is due to the large barycentric shift to very slowly rotating perturbation, and so, the density response may be too large to use the linear perturbation theory. Thus, to confirm the significance of weakly damped modes, it might be required to include the non-linear effects, e.g., by using  $N$ -body simulations with a huge number of particles.

To examine whether such weakly damped modes in reality play a role in lopsided disks, it is useful to observe the pattern speed of the lopsidedness directly, e.g., by using the Tremaine-Weinberg method (Tremaine & Weinberg 1984a; see also Sambhus & Sridhar 2000). In the Milky Way, the next generation astrometric satellites such as SIM and GAIA will help us to measure the dynamics of the lopsidedness directly. In addition, since our results suggest that the lopsidedness will be slowly rotating pattern, and so, they will be long-sustained structure, a large survey of the lopsidedness in the sample of isolated galaxies, which have not any companion galaxies that could cause the significant lopsidedness, would also show the large fraction of the lopsidedness in such sample galaxies.

In this paper, the shape of lopsided disks,  $A_1(R)$ , is chosen to be  $Rd/(25 + R^2)$ . However, since each element of the response matrix,  $R_{nn'}^{lm}$ , is already obtained, dynamical friction acting on another type of lopsided pattern can be readily calculated. Then, the dependence of the functional form of  $A_1$  on the estimate of dynamical friction is found to be weak. This is because the change

of the functional form of  $A_1$  does not change significantly the shape of  $\Sigma_1$  or  $\rho_1^{\text{ext}}$ , because the term of  $\exp(-R)$  is dominant in  $\Sigma_1$ .

Before closing this section, possible effects of the assumptions made in this paper on the estimate of dynamical friction are discussed.

First, the effect of flat disks on the orbits of halo stars is ignored in this work. Taking into account the existence of flat-disks, the halo density near disks would increase owing to the additional flat-disk potential as mentioned in Nelson & Tremaine (1995). For the fast rotating lopsided disks, since the friction will be caused by the halo stars with high orbital frequencies, the damping time scale may be shortened by adding the effect of flat disks. However, for the slowly rotating lopsided disks, the friction will be caused mainly by the stars with low orbital frequencies, since such stars will not be significantly affected by including flat disks, the friction would not change significantly. As a result, the conclusion, the damping time scale is short unless the pattern speed is slow, may be strengthened.

Second, galactic disks are assumed to be stable against lopsided ( $m = 1$ ) perturbation. When galactic disks are unstable to lopsided distortions, the amplification of lopsided modes due to the interaction with disk stars will be important. Then, dynamical friction might not reduce the amplitude of unstable modes ( $d$ ) but might make the pattern speed of the modes ( $\Omega_p$ ) slow down as similar to the case of galactic bars (Weinberg 1985; Debattista & Sellwood 2000). Such  $m = 1$  instabilities will exist in stellar disks in which the fraction of retrograde stars is large (Zang & Hohl 1978; Sawamura 1988; Hozumi & Fujiwara 1989; Sellwood & Merritt 1994). Although there exist some galaxies that have counter-rotating components (e.g., Merrifield & Kuijken 1994), the counter-rotation is a rare phenomenon in stellar disks (Kuijken, Fisher, & Merrifield 1996; Kannappan & Fabricant 2001). Hence, it would be difficult to explain the large fraction of lopsided disks in spiral galaxies. In the mean time, some galactic disks in which the contribution of halos to galactic rotation is small will also have a lopsided instability (e.g., Sellwood 1985; Athanassoula, Bosma, & Papaioannou 1987; Lovelace et al. 1999). Then, dynamical friction may not affect dynamical evolution of such unstable modes owing to the small halo contribution. However, since galactic disks with massive halos will be stable to lopsided distortions (e.g., Athanassoula et al. 1987), this kind of instability would not be major cause of the lopsidedness in real spiral galaxies, especially in late-type spirals.

To confirm the results in this paper, and to avoid the effect of discreteness noise, numerical simulations with a huge number of particles, say  $\simeq 10^7$ , would be valuable. Such a calculation can be done, e.g., with a hierarchical tree algorithm using the GRAPE-5, a special-purpose computer for gravitationally interacting particles (Sugimoto et al. 1990; Kawai et al. 2000). This line of investigation is in progress.

The author is grateful to Shogo Inagaki, Shunsuke Hozumi, Junichiro Makino and Martin Weinberg for their helpful discussions. The author also acknowledges the anonymous referee for his

useful comments and suggestions, which helps me very much to improve the paper. The numerical package 'LAPACK' is used for calculating the inverse and singular value decomposition of a general complex matrix. Numerical computations are carried out on workstations at the Astronomical Data Analysis Center of the National Astronomical Observatory, Japan (ADAC/NAOJ), which is an inter-university research institute of astronomy operated by Ministry of Education, Culture, Sports, Science and Technology. The author acknowledges the Japan Society for the Promotion of Science (JSPS) for financial support.

## A. Matrix Equation

### A.1. Expansion of Density-Potential Pairs

Following Tremaine & Weinberg (1984b) and Weinberg (1989), the response density  $\rho_1^{\text{res}}$ , the response potential  $\Phi_1^{\text{res}}$ , and the external potential  $\Phi_1^{\text{ext}}$  are expanded in biorthonormal basis sets (Clutton-Brock 1972; Hernquist & Ostriker 1992),

$$\rho_1^{\text{res}}(\mathbf{r}, t) = \sum_{n,l,m} A_n^{lm}(t) d_n^{lm}(\mathbf{r}) Y_{lm}(\theta, \phi), \quad (\text{A1})$$

$$\Phi_1^{\text{res}}(\mathbf{r}, t) = \sum_{n,l,m} A_n^{lm}(t) u_n^{lm}(\mathbf{r}) Y_{lm}(\theta, \phi), \quad (\text{A2})$$

$$\Phi_1^{\text{ext}}(\mathbf{r}, t) = \sum_{n,l,m} B_n^{lm}(t) u_n^{lm}(\mathbf{r}) Y_{lm}(\theta, \phi), \quad (\text{A3})$$

where  $Y_{lm}(\theta, \phi)$  are the spherical harmonics. Here, the radial basis functions  $u_n^{lm}(\mathbf{r})$  and  $d_n^{lm}(\mathbf{r})$  are normalized by

$$-\frac{1}{4\pi G} \int dr r^2 u_n^{lm*}(\mathbf{r}) d_{n'}^{lm}(\mathbf{r}) = \delta_{nn'}, \quad (\text{A4})$$

and satisfy the Poisson equation  $\nabla^2 u_n^{lm}(\mathbf{r}) = 4\pi G d_n^{lm}(\mathbf{r})$ .

Since the orbits are periodic with respect to the angle variables  $\mathbf{w}$ , the perturbed potential,  $\Phi_1 = \Phi_1^{\text{res}} + \Phi_1^{\text{ext}}$ , can be also expanded in the Fourier series,

$$\Phi_1(\mathbf{I}, \mathbf{w}, t) = \sum_{l,m} \sum_l \Psi_{1l}(\mathbf{I}, t) e^{i\mathbf{l} \cdot \mathbf{w}}, \quad \text{where} \quad \sum_l \equiv \sum_{l_1=-\infty}^{\infty} \sum_{l_2=-l_1}^{l_1}, \quad (\text{A5})$$

where  $\mathbf{l} = (l_1, l_2, l_3 = m)$  is a triple of integers and  $\mathbf{I} = (I_r, J, J_z \equiv J \cos \beta)$  is a vector of the action variables. Here,  $I_r$  is the radial action,  $J$  is the total angular momentum, and  $J_z$  is its z-component. The Fourier coefficients  $\Psi_{1l}$  are

$$\Psi_{1l}(\mathbf{I}, t) = V_{l2m}(\beta) \sum_n W_{l2m}^{l_1 n}(\mathbf{I}) \left[ A_n^{lm}(t) + B_n^{lm}(t) \right], \quad (\text{A6})$$

$$V_{l2m}(\beta) = i^{m-l_2} Y_{l2}(\pi/2, 0) r_{l2m}^l(\beta), \quad (\text{A7})$$

$$W_{l2m}^{l_1 n}(\mathbf{I}) = \frac{1}{\pi} \int_0^\pi dw_1 \cos[l_1 w_1 + l_2(w_2 - \psi)] u_n^{lm}(\mathbf{r}), \quad (\text{A8})$$

where  $\psi$  is the angle in the orbital plane measured from the ascending node and  $r_{l_2m}^l(\beta)$  is the rotation matrix (see e.g., Edmonds 1960). The explicit formulae for the angle variables  $w_1$  and  $w_2$  are described in Tremaine & Weinberg (1984b).

## A.2. Matrix Equation

The response density to an external density will be determined by coupled-solutions of the linearized Boltzmann-Poisson equation (eqs. [1] and [2]). To find such solutions, it is convenient to use the Fourier-Laplace transform. Then, the linearized collisionless Boltzmann equation becomes

$$\tilde{f}_{1l} = -l \cdot \frac{\partial f_0}{\partial \mathbf{I}} \tilde{\Psi}_{1l} \frac{1}{\omega - l \cdot \boldsymbol{\Omega}}, \quad (\text{A9})$$

where  $\boldsymbol{\Omega}(\mathbf{I}) = \partial H_0 / \partial \mathbf{I}$  are the frequencies of the angle variables and the assumption that the perturbed distribution function vanishes at  $t = 0$  is made. Here, the subscript  $l$  denotes the Fourier transform in the angle variables, and the tilde denotes the Laplace transform in the time variable. Using equation (A9), the response density can be written

$$\tilde{\rho}_1^{\text{res}}(\mathbf{r}, \omega) = \int d^3v \tilde{f}_1 = - \int d^3v \sum_l l \cdot \frac{\partial f_0}{\partial \mathbf{I}} \tilde{\Psi}_{1l} \frac{1}{\omega - l \cdot \boldsymbol{\Omega}} e^{il \cdot \mathbf{w}}. \quad (\text{A10})$$

Then, the expansion coefficients of the response density  $A_n^{lm}$  are

$$\tilde{A}_n^{lm}(\omega) = -\frac{1}{4\pi G} \int d^3r u_n^{lm*}(\mathbf{r}) Y_{lm}^*(\theta, \phi) \tilde{\rho}_1^{\text{res}}(\mathbf{r}, \omega) = \sum_{n'} \mathcal{R}_{nn'}^{lm}(\omega) [\tilde{A}_{n'}^{lm}(\omega) + \tilde{B}_{n'}^{lm}(\omega)], \quad (\text{A11})$$

where the response matrix  $\mathcal{R}_{nn'}^{lm}$  is

$$\begin{aligned} \mathcal{R}_{nn'}^{lm}(\omega) &\equiv \frac{(2\pi)^3}{4\pi G} \frac{2}{2l+1} \iint \frac{dE dJ}{\Omega_1(E, J)} \sum_l l \cdot \frac{\partial f_0(E, J)}{\partial \mathbf{I}} \frac{1}{\omega - l \cdot \boldsymbol{\Omega}} |Y_{l_2}(\pi/2, 0)|^2 \\ &\times W_{l_2m}^{l_1n*}(E, J) W_{l_2m}^{l_1n'}(E, J). \end{aligned} \quad (\text{A12})$$

Here, the relation for canonical variables,  $d^3r d^3v = d^3w d^3I$ , is used, and the variables are changed from  $\mathbf{I} = (I_r, J, J_z)$  to  $(E, J, \beta)$ . The integration over  $\mathbf{w}$  and  $\beta$  is done with the assumption that an unperturbed distribution function does not depend on  $\beta$ . In addition, the following expression for  $u_n^{lm*} Y_{lm}^*$  is used.

$$u_n^{lm*}(\mathbf{r}) Y_{lm}^*(\theta, \phi) = \sum_l V_{l_2m}^*(\beta) W_{l_2m}^{l_1n*}(\mathbf{I}) e^{-il \cdot \mathbf{w}}. \quad (\text{A13})$$

The matrix equation (A11) can be rewritten in a symbolic form,

$$\tilde{A}^{lm}(\omega) = \mathcal{R}^{lm}(\omega) \cdot [\tilde{A}^{lm}(\omega) + \tilde{B}^{lm}(\omega)] = \mathcal{D}^{-1lm}(\omega) \cdot \mathcal{R}^{lm}(\omega) \cdot \tilde{B}^{lm}(\omega), \quad (\text{A14})$$

where

$$\mathcal{D}^{lm} = \mathcal{I} - \mathcal{R}^{lm}. \quad (\text{A15})$$

Here,  $\mathcal{I}$  is a unit matrix with the same rank as  $\mathcal{R}^{lm}$ . The inverse Laplace transform of the matrix equation (A14) is

$$A^{lm}(t) = \int_0^t d\tau \mathcal{D}^{-1lm}(m\Omega_p) \cdot \mathcal{K}^{lm}(t-\tau) \cdot B^{lm}(\tau), \quad (\text{A16})$$

where the kernel matrix  $\mathcal{K}_{nn'}^{lm}$  is the inverse Laplace transform of a response matrix  $\mathcal{R}_{nn'}^{lm}$  which yields

$$\begin{aligned} \mathcal{K}_{nn'}^{lm}(t-\tau) = & -i \frac{(2\pi)^3}{4\pi G} \frac{2}{2l+1} \iint \frac{dEdJJ}{\Omega_1(E, J)} \sum_l l \cdot \frac{\partial f_0}{\partial \mathbf{I}} \exp[-il \cdot \boldsymbol{\Omega} \times (t-\tau)] \\ & \times |Y_{l2}(\pi/2, 0)|^2 W_{l2m}^{l_1 n*}(E, J) W_{l2m}^{l_1 n'}(E, J). \end{aligned} \quad (\text{A17})$$

In principle, the density and potential responses can be calculated by integrating equation (A16) as an initial value problem. In this work, the time asymptotic assumption ( $t \rightarrow \infty$ ) is adopted to equation (A16). Then, all transient waves originating from an initial condition will disperse and the forced harmonic will dominate. When the perturbed density is assumed to rotate at a constant rate  $\Omega_p$ , both  $A_n^{lm}(t)$  and  $B_n^{lm}(t)$  are proportional to  $\exp(-im\Omega_p t)$ . Then, the integration over  $\tau$  in equation (A16) can be readily done,

$$A^{lm}(t) = \mathcal{D}^{-1lm}(m\Omega_p) \cdot \mathcal{R}^{lm}(m\Omega_p) \cdot B^{lm}(t). \quad (\text{A18})$$

Here, each element of a response matrix,  $\mathcal{R}_{nn'}^{lm}$ , is

$$\begin{aligned} \mathcal{R}_{nn'}^{lm} = & \frac{8\pi^3}{4\pi G} \frac{2}{2l+1} \iint \frac{dEdJJ}{\Omega_1(E, J)} \sum_l l \cdot \frac{\partial f_0}{\partial \mathbf{I}} |Y_{l2}(\pi/2, 0)|^2 W_{l2m}^{l_1 n*}(E, J) W_{l2m}^{l_1 n'}(E, J) \\ & \times \left[ \mathcal{P} \left( \frac{1}{m\Omega_p - l \cdot \boldsymbol{\Omega}} \right) - \pi i \delta(m\Omega_p - l \cdot \boldsymbol{\Omega}) \right] \end{aligned} \quad (\text{A19})$$

where  $\delta$  denotes the Dirac delta function and  $\mathcal{P}$  denotes the Cauchy principal value. These equations are identical to equations (48) and (49) in Weinberg (1989). In addition, when  $f_0$  depends only on the energy  $E$ ,

$$l \cdot \frac{\partial f_0}{\partial \mathbf{I}} = (l \cdot \boldsymbol{\Omega}) \frac{df_0}{dE}. \quad (\text{A20})$$

Thus, the real (imaginary) part of a response matrix is an even (odd) function of  $\Omega_p$ .

The equation (A18) is valid for  $t \rightarrow \infty$ , however, if weakly damped modes should be included, one must use the following equation instead of equation (A18)

$$\begin{aligned} A^{lm}(t) = & \mathcal{D}^{-1lm}(m\Omega_p) \cdot \mathcal{R}^{lm}(m\Omega_p) \cdot B^{lm}(t) \\ & - i \sum_d \text{Res} \left( \mathcal{D}^{-1lm}; \omega_d \right) \cdot \mathcal{R}^{lm}(\omega_d) \cdot \int_0^\infty d\tau B^{lm}(\tau) \exp[-i\omega_d \times (t-\tau)]. \end{aligned} \quad (\text{A21})$$

Here,  $\text{Res}$  denotes the residue,  $\omega_d$  is the complex frequency of damped modes,  $\det \mathcal{D}(\omega_d) = 0$ , and the summation is over the number of damped modes.



## REFERENCES

- Arp, H. 1966, *ApJS*, 14, 1
- Athanassoula, E., Bosma, A., & Papaioannou, S. 1987, *A&A*, 179, 23
- Baldwin, J. E., Lynden-Bell, D., & Sancisi, R. 1980, *MNRAS*, 193, 313
- Chandrasekhar, S. 1943, *ApJ*, 97, 251
- Clutton-Brock, M. 1972, *Ap&SS*, 16, 101
- Debattista, V. P., & Sellwood, J. A. 2000, *ApJ*, 543, 704
- Edmonds, A. R. 1960, *Angular Momentum in Quantum Mechanics* (Princeton: Princeton University Press)
- Goodman, J. 1988, *ApJ*, 329, 612
- Haynes, M. P., Hogg, D. E., Maddalena, R. J., Roberts, M. S., Zee, L. V. 1998, *AJ*, 115, 62
- Hernquist, L., & Ostriker, J. P. 1992, *ApJ*, 386, 375
- Hozumi, S., & Fujiwara, T. 1989, *PASJ*, 41, 841
- Jog, C. J. 1997, *ApJ*, 488, 642
- Jog, C. J. 1999, *ApJ*, 522, 661
- Kalnajs, A. J. 1977, *ApJ*, 212, 637
- Kannappan, S. J., & Fabricant, D. G. 2001, *AJ*, 121, 140
- Kawai, A., Fukushige, T., Makino, J., & Taiji, M. 2000, *PASJ*, 52, 659
- Kochanek, C. S. 1996, *ApJ*, 457, 228
- Kornreich, D. A., Haynes, M. P., & Lovelace, R. V. E. 1998, *AJ*, 116, 2154
- Kuijken, K., Fisher, D., Merrifield, M. R. 1996, *MNRAS*, 283, 543
- Levine, S. E., & Sparke, L. S. 1998, *ApJ*, 496, L13
- Lindblad, B. 1963, *Stockholm Obs. Ann.*, 22, 3
- Lovelace, R. V. E., Zhang, L., Kornreich, D. A., & Haynes, M. P. 1999, *ApJ*, 524, 634
- Lynden-Bell, D., & Kalnajs, A. J. 1972, *MNRAS*, 157, 1
- Merrifield, M. R., & Kuijken, K. 1994, *ApJ*, 432, 575

Murali, C. 1999, *ApJ*, 519, 580  
 Nelson, R. W., & Tremaine, S. 1995, *MNRAS*, 275, 897  
 Palmer, P. L., & Papaloizou, J. 1987, *MNRAS*, 224, 1043  
 Richter, O.-G., & Sancisi, R. 1994, *A&A*, 290, L9  
 Rix, H.-W., & Zaritsky, D. 1995, *ApJ*, 447, 82  
 Rudnick, G., Rix, H.-W. 1998, *AJ*, 116, 1163  
 Rudnick, G., Rix, H.-W., & Kennicutt, R. C., Jr. 2000, *ApJ*, 538, 569  
 Sambhus, N., & Sridhar, S. 2000, *ApJ*, 539, L17  
 Sawamura, M. 1988, *PASJ*, 40, 279  
 Sellwood, J. A. 1985, *MNRAS*, 217, 127  
 Sellwood, J. A., & Merritt, D. 1994, *ApJ*, 425, 530  
 Sugimoto, D., Chikada, Y., Makino, J., Ito, T., Ebisuzaki, T., Umemura, M. 1990, *Nature*, 345, 33  
 Syer, D., & Tremaine, S. 1996, *MNRAS*, 281, 925  
 Tóth, G., & Ostriker, J. P. 1992, *ApJ*, 389, 5  
 Tremaine, S., & Weinberg, M. D. 1984a, *ApJ*, 282, L5  
 Tremaine, S., & Weinberg, M. D. 1984b, *MNRAS*, 209, 729  
 Vesperini, E., & Weinberg, M. D. 2000, *ApJ*, 534, 598  
 Walker, I. R., Mihos, J. C., & Hernquist, L. 1996, *ApJ*, 460, 121  
 Weinberg, M. D. 1985, *MNRAS*, 213, 451  
 Weinberg, M. D. 1989, *MNRAS*, 239, 549  
 Weinberg, M. D. 1994, *ApJ*, 421, 481  
 Weinberg, M. D. 1998a, *MNRAS*, 297, 101  
 Weinberg, M. D. 1998b, *MNRAS*, 299, 499  
 Weinberg, M. D. 1999, *AJ*, 117, 629  
 Zang, T. A., & Hohl, F. 1978, *ApJ*, 226, 521  
 Zaritsky, D., & Rix, H.-W. 1997, *ApJ*, 477, 118

# Formation of Terrestrial Planets in a Dissipating Gas Disk

Junko Kominami

E-mail: kominami@geo.titech.ac.jp

and

Shigeru Ida

Department of Earth and Planetary Sciences,  
Tokyo Institute of Technology,  
Ookayama, Meguro-ku, Tokyo, 152-8551, Japan

May 13, 2002

## Abstract

We performed N-body simulation on formation of terrestrial planets from protoplanets including damping of velocity dispersion caused by gravitational interaction with a dissipating gas disk. In a gas-free case, the resulting planets have relatively high eccentricities compared to those of Earth and Venus. These high eccentricities are the remnant of orbital crossings; collisional damping is not strong enough. The damping due to almost dissipated disk gas is strong enough to damp their eccentricities down to the present values of Earth and Venus, while it allows the protoplanets to grow to the size of the Earth (Kominami & Ida 2002, *Icarus*, in press). In this paper, decay of disk gas is considered as an additional factor. Exponential decay is assumed for time evolution of disk gas. We investigate how planetary formation and depletion timescale are related. We found out that if we assume the minimum mass disk model, Earth like planets are formed when gas depletion timescale is  $10^6 - 10^7$  years.

## 1 Introduction

Terrestrial planets are formed through the collisions of planetesimals that are  $\sim$  km size. Oligarchic growth predicts formation of about twenty Mars-sized protoplanets (about one tenth of Earth mass) on circular orbits (Kokubo & Ida 1998, 2000). Mutual gravitational interaction between the protoplanets (Chambers et al. 1996) and/or effect of the giant planets increase the eccentricities of the protoplanets on a time scale of  $10^6 - 10^7$  years. (Nagasawa et al. 2000, Ito & Tanikawa 1999) to cause orbital instabilities. The protoplanets start to collide and grow till about the size of the Earth and Venus. If only the mutual gravitational interaction between the

protoplanets are considered, the resulting planets have relatively high eccentricities compared to those of Earth and Venus (Chambers and Wetherill 1998, Agnor et al. 1999). These high eccentricities are the remnant of orbital crossings. Collisional damping is not strong enough to reproduce relatively small ( $\lesssim 0.03$ ) eccentricities, which are comparable to time-averaged eccentricities of Venus and Earth. However, it is reasonable to assume that leftover planetesimals and remnant disk gas still exist during this stage. Gravitational interaction with the disk causes damping of eccentricities (and inclinations). If this interaction is also taken into account, the protoplanets can grow to the size of the Earth and acquire low eccentricities (Kominami & Ida 2002).

Kominami & Ida (2002) shows that if orbital crossing and growth of the planets occur when gas with surface density  $\sim 10^{-3} - 10^{-4} \times$  the minimum mass model (Hayashi 1981) is left in the disk, planets with mass of  $\sim M_{\oplus}$  and eccentricities  $\lesssim 0.03$  are formed. This result was acquired from the calculations with time independent amount of disk gas. They assumed a constant gas model in order to study the effect of remnant gas on the planet formation more clearly. Although they also did several calculations with disk gas dissipating exponentially and showed in the case with disk decaying time scale of  $\sim 10^6 - 10^7$  years planets with  $m \sim M_{\oplus}$  and  $e \sim 0.01$  are formed. However, the effects of depleting gas has not been fully understood yet. In this paper, performing much more calculations with decaying disk gas, we investigated the relation between the depletion timescale of the disk gas and the planetary formation. The relation would impose constraints upon evolution of disk gas.

## 2 Calculation Model

### 2.1 Gravitational Gas Drag

We consider the damping due to disk-planet interaction, which we call “gravitational drag”, as in Kominami & Ida (2002). Protoplanets are large enough to ignore the aerodynamic gas drag force.

The effects of the disk-planet interaction can be expressed by the drag force ( $\mathbf{f}_{\text{GD}}$ ) as (Kominami & Ida, 2002)

$$\mathbf{f}_{\text{GD}} = -\frac{\mathbf{v} - \mathbf{v}_{\text{gas}}}{\tau_{\text{damp}}}, \quad (2.1)$$

where  $\mathbf{v}$  and  $\mathbf{v}_{\text{gas}}$  are the velocity of a protoplanet and the gas. We here assume the gas motion is non-inclined circular Keplerian motion. Damping timescale of gravitational gas drag is

$$\tau_{\text{damp}} \simeq \left(\frac{M_{\odot}}{M}\right) \left(\frac{M_{\odot}}{\Sigma_{\text{gas}} r^2}\right) \left(\frac{c_s}{v_{\text{kep}}}\right)^4 \Omega_{\text{kep}}^{-1}, \quad (2.2)$$

where  $\Sigma$  is the surface density of a gas disk,  $c_s$  is sound velocity of disk gas and  $\Omega_{\text{kep}}$  is Keplerian frequency (Ward 1989, 1993, Artymowicz 1993). Supposing the minimum mass disk model with gas surface density given by  $\Sigma^{\text{min}} = 1700 (r/1\text{AU})^{-3/2} \text{gcm}^{-2}$  (Hayashi 1981),

$$\tau_{\text{damp}} \simeq 0.5 \times 10^3 \left(\frac{M}{M_{\oplus}}\right)^{-1} \left(\frac{r}{1\text{AU}}\right)^2 \left(\frac{\Sigma}{\Sigma^{\text{min}}}\right)^{-1} \text{years}. \quad (2.3)$$

We assume exponential decay of the surface density as

$$\Sigma(t) = \Sigma_0 \exp\left(-\frac{t}{\tau_{\text{gas}}}\right). \quad (2.4)$$

As a consequence, damping time scale lengthens as

$$\tau_{\text{damp}}(\Sigma(t)) = \tau_{\text{damp}}(\Sigma_0) \exp\left(\frac{t}{\tau_{\text{gas}}}\right). \quad (2.5)$$

In our calculations, we assume time-independent  $\tau_{\text{gas}}$ . However, since orbital evolution of protoplanets occurs when  $\Sigma \sim 10^{-2} - 10^{-4} \Sigma^{\text{min}}$ , the evolution is regulated by  $\tau_{\text{gas}}$  at those stages. Different values of  $\tau_{\text{gas}}$  at the other stages do not affect the result.

## 2.2 Orbital Integration

We integrate orbits with 4th order Hermite scheme (Makino & Aarseth 1992) and hierarchical individual timestep (Makino 1991) as in Kominami & Ida (2002). The equation of motion of particle  $k$  is

$$\frac{d\mathbf{v}_k}{dt} = -\frac{GM_\odot}{|\mathbf{r}_k|^3} \mathbf{r}_k - \sum_{j \neq k} \frac{GM_j}{|\mathbf{r}_j - \mathbf{r}_k|^3} (\mathbf{r}_j - \mathbf{r}_k) - \frac{\mathbf{v}_k - \mathbf{v}_{\text{gas}}}{\tau_{\text{grav}}}. \quad (2.6)$$

The first term is the gravity from the sun. The second term is the mutual gravity between the protoplanets. And the last term is the gravitational drag from disk gas. The drag force has time dependence as in Eq. (2.5). When protoplanets collide, perfect accretion is assumed. The physical radius of a protoplanet is determined by its mass and internal density as

$$r_p = \left(\frac{3}{4\pi} \frac{M}{\rho_p}\right)^{1/3}. \quad (2.7)$$

The internal density  $\rho_p$  is set to be  $3 \text{ g cm}^{-3}$ .

## 2.3 Initial Condition

The initial protoplanet distribution is also the same as Kominami & Ida (2002). The number of the protoplanets is fifteen, and separation of semimajor axis is 6 - 10  $r_H$ .  $r_H$  is Hill radius which is defined as

$$r_H = \left(\frac{2M}{3M_\odot}\right)^{1/3} r \simeq 0.007 \left(\frac{M}{0.2M_\oplus}\right)^{1/3} r. \quad (2.8)$$

Angular distribution is random. Each protoplanet has mass of  $0.2M_\oplus$ . The initial eccentricities of the planets are the order of  $10^{-3} - 10^{-4}$ , which means the orbits are almost circular.

Initial amount of gas ranges from 1% to 100% of minimum mass model. It is determined as following. Timescale for orbital instability to occur in a gas free case ( $\tau_{\text{cross}}$ ) is a function of orbital separation and masses of the protoplanets (Chambers et al., 1996). Considering gravitational gas drag, orbital instability is suppressed if

$$\tau_{\text{damp}}(\Sigma) \lesssim \tau_{\text{cross}} \quad (2.9)$$

(Iwasaki et al. 2002). We define  $\Sigma_1$  such that orbital instability is allowed when  $\Sigma < \Sigma_1$ .  $\tau_{\text{damp}}$  increases with time as Eq.(2.5). The system becomes ready for the instability to take place when the damping time scale becomes (Eq.(2.5)).

$$\tau_{\text{orb}} = \tau_{\text{gas}} \log\left(\frac{\tau_{\text{cross}}}{\tau_{\text{damp}}(\Sigma_0)}\right). \quad (2.10)$$

If  $\tau_{\text{orb}} \ll \tau_{\text{cross}}$ , disk gas dissipated almost completely by the time when orbital crossing starts. The following orbital evolution will be as the same as gas free case. On the other hand, if  $\tau_{\text{orb}} \gtrsim \tau_{\text{cross}}$ , orbital crossing starts when  $t \gtrsim \tau_{\text{orb}}$ , which means orbital crossing is triggered by dissipation of disk gas and remnant disk gas may affect the following orbital evolution. We consider the latter case, because the former case would result in too large eccentricities. Hence, the initial amount of gas ( $\Sigma_0$ ) must satisfy

$$\tau_{\text{damp}}(\Sigma_0) \exp\left(\frac{\tau_{\text{cross}}(\Delta a_0)}{\tau_{\text{gas}}}\right) < \tau_{\text{cross}}(\Delta a_0). \quad (2.11)$$

If  $\Sigma^{\text{min}}$  is adopted as  $\Sigma_0$ , this condition is usually satisfied. However, in order to reduce computation time, we often start with  $\Sigma_0 < \Sigma^{\text{min}}$ . As long as  $\Sigma_0$  satisfies Eq.(2.11), the results would not change. Initial separation  $\Delta a_0$  and the corresponding  $\tau_{\text{cross}}$  is shown in Table I. We did several runs with no gas on each  $\Delta a_0$  and took the average time for the crossing to happen.

**Table I**  
**Initial Orbital Separation and Instability Time Scale**

$\Delta a_0(r_{\text{H}})$	$\tau_{\text{cross}}(\text{yrs})$
7	$\sim 1 \times 10^5$
8	$\sim 1 \times 10^6$
9	$\sim 2 \times 10^6$
10	$\sim 3 \times 10^6$
12	$\sim 4 \times 10^6$

### 3 Result

We did 15 runs with decaying gas. Initial conditons and the final planets' mass and eccntricityies are listed in Table II.

**Table II**  
**List of Simulations with Initial Conditions, and Final Planets**

Simulation	$\Delta a(r_H)$	$\Sigma_0(\Sigma_H)$	$\tau_{\text{gas}}(\text{yr})$	$n_{\text{final}}$	$M_{\text{max}}(M_{\oplus})$	$e_{\text{max}}$	$M_{\text{max}} \geq 0.8M_{\oplus}$	$e_{\text{max}} \leq 0.04$
Run <sub>4</sub> <sup>3</sup>	7	0.01	$3 \times 10^7$	6	0.6	0.0001	No	Yes
Run <sub>8</sub> <sup>2</sup>	7	0.2	$3 \times 10^6$	8	0.6	0.001	No	Yes
Run <sub>1</sub> <sup>1</sup>	9	0.01	$1 \times 10^7$	8	0.6	0.006	No	Yes
Run <sub>6</sub> <sup>2</sup>	8	0.1	$3.3 \times 10^6$	8	0.6	0.0001	No	Yes
Run <sub>7</sub> <sup>2</sup>	8	0.1	$3.3 \times 10^6$	10	0.4	0.001	No	Yes
Run <sub>4</sub> <sup>2</sup>	8	0.03	$3 \times 10^6$	7	1.0	0.0067	Yes	Yes
Run <sub>5</sub> <sup>2</sup>	8	0.05	$2 \times 10^6$	5	1.0	0.0138	Yes	Yes
Run <sub>2</sub> <sup>1</sup>	12	0.05	$6.7 \times 10^6$	6	0.8	0.021	Yes	Yes
Run <sub>2</sub> <sup>2</sup>	9	0.01	$3 \times 10^6$	4	1.2	0.042	Yes	No/Yes
Run <sub>3</sub> <sup>2</sup>	9	0.1	$3 \times 10^6$	5	1.0	0.142	Yes	No
Run <sub>9</sub> <sup>2</sup>	10	0.05	$2.4 \times 10^6$	4	0.2	0.14	Yes	No
Run <sub>3</sub> <sup>3</sup>	8	0.3	$3 \times 10^5$	4	1.8	0.079	Yes	No
Run <sub>1</sub> <sup>3</sup>	9	1	$3 \times 10^5$	4	1.0	0.053	Yes	No

(\*) a,b,c,labeled in the simulation number indicates the initial angular distribution type. Angular distribution is given randomly. Each distribution is made of different set of random number.  $M_{\text{max}}$  is the largest final planets in each run.  $e_{\text{max}}$  are their eccentricities, respectively. If there are more than one largest planets, the average is taken within the same mass.

The typical result of orbital evolution when the gas dessipates too quickly is shown in fig.1(a). This is a figure of Run<sub>2</sub><sup>3</sup>. Since  $\tau_{\text{gas}}$  is short ( $\tau_{\text{gas}} = 3 \times 10^5$  yrs), gas is being depleted so much during the accretion, that the situation becomes equivalent to gas free case. At  $t = 6.2 \times 10^6$  years, when an Earth-sized planet is formed,  $\tau_{\text{damp}}$  for  $m = 1M_{\oplus}$  is  $2.4 \times 10^{12}$  years, which is much longer than  $\tau_{\text{gas}}$ . Planets can grow but the eccentricity cannot be damped. Final planets are showm in fig.1(b). Filled circles represent the terrestrial planets. Area of the circles are proportional to the mass of the planets. When the gas depletion is too slow, as in Run<sub>4</sub><sup>3</sup>, the protoplanets cannot grow to the size of the Earth. Figure 1(c) is the typical result. When orbital crossing starts at  $t = 1 \times 10^5$  yrs, however, there is too much gas to allow the accretion to continue:  $\Sigma \sim 10^{-2}\Sigma^{\text{min}}$ . The eccentricities are quickly damped and the planets cannot grow. The final planets are shown in fig.1(d).

When  $\tau_{\text{gas}} \sim \tau_{\text{cross}}$ , orbital evolution is like fig.1(e), which shows the orbital evolution of Run<sub>1</sub><sup>1</sup>. Initially, the gas supresses excitation of the eccentricities and prevents the orbits from becoming unstable. The gas is gradually depleted and when  $t \sim 3.5 \times 10^7$  yrs, orbits start crossing. The amount of gas then is  $\sim 2.7 \times 10^{-4}\Sigma^{\text{min}}$ . This corresponeds to  $\tau_{\text{damp}} \sim 9.3 \times 10^6$  years. During this time, the amount of gas is  $\sim 10^{-3} - 10^{-4}\Sigma^{\text{min}}$ , It is already known that if there is gas of amount of  $\sim 10^{-3} - 10^{-4}\Sigma^{\text{min}}$ , during the accretion, planets with large masses ( $M_{\oplus}$ ) and low eccentricities can be formed (Kominami & Ida 2002). Remnant gas is enough to damp the eccentricities of surviving planet with mass  $\sim M_{\oplus}$  (Kominami & Ida 2002; Agnor & Ward 2002). The final planets are shown in fig.1(f). The largest mass has  $0.8M_{\oplus}$  and its eccentricity is  $\sim 0.02$ .

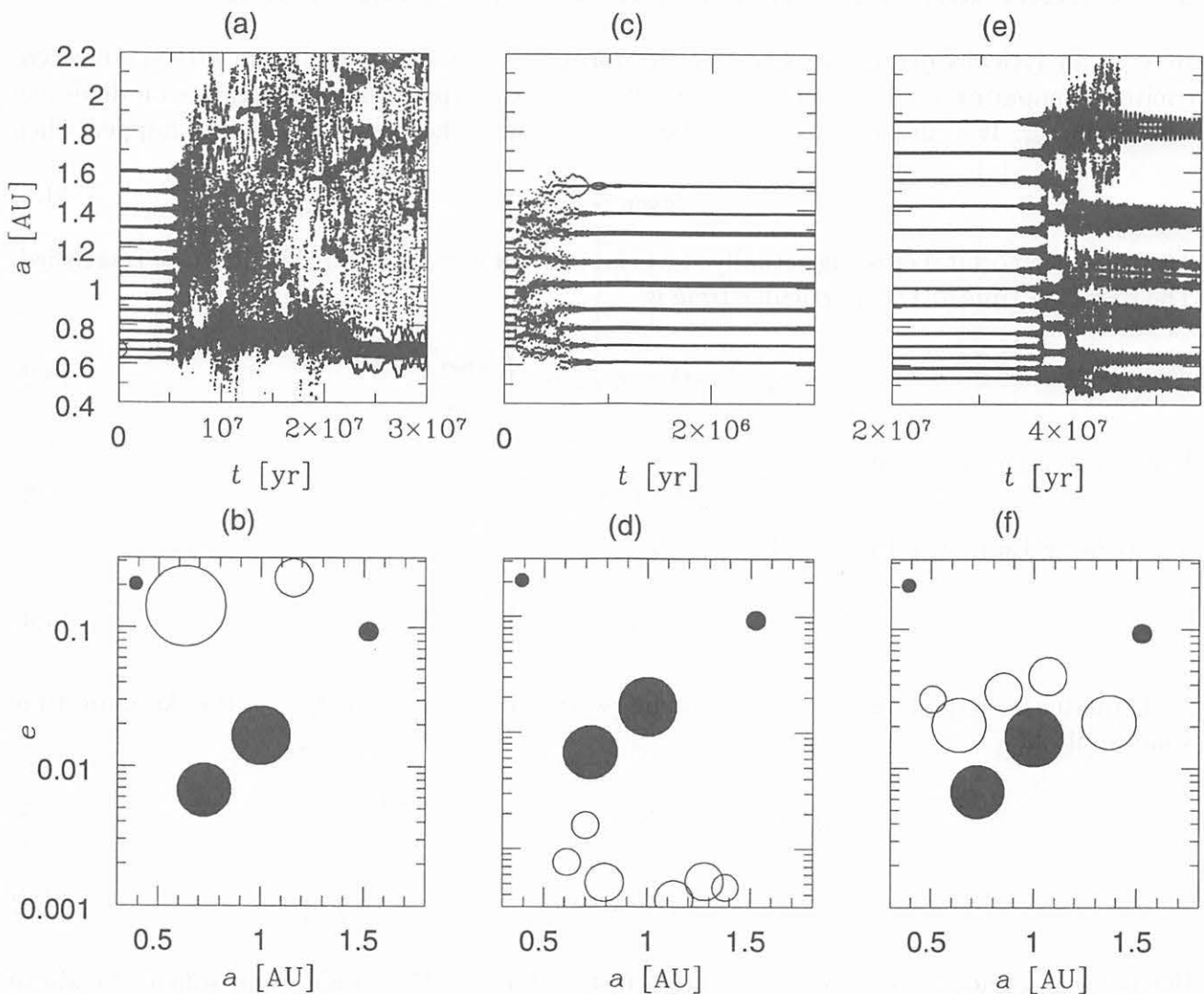


Figure 1: (a) Typical orbital evolution when  $\tau_{\text{gas}} = 3 \times 10^5$  yrs. Semimajor axis, pericenters and apocenters are plotted. (b) Final planets of (a). Area of the circle is proportional to the planets' mass. Filled circles are terrestrial planets. (c) Same as (a) except  $\tau_{\text{gas}} = 3 \times 10^7$  yrs. (d) Final planets of (c). (e) Same as (a) except  $\tau_{\text{gas}} = 3 \times 10^6$  yrs.



The accretion timescale( $\tau_{\text{growth}}$ ), which is from when the orbital crossing starts till the mass becomes  $\sim M_{\oplus}$ , is  $\sim 1 - 3 \times 10^6$  years, although  $\tau_{\text{growth}} \sim 5 \times 10^6$  years in Fig.1(e).

## 4 Condition for Formation of Earth-like Planets

Here we analytically derive the conditions for formation of Earth-like planets with small eccentricities, comparing with the numerical results. We can impose some constraints on depletion time scale  $\tau_{\text{gas}}$ . It is already shown that the system allows the orbital crossing to happen when

$$\tau_{\text{damp}} \gtrsim \tau_{\text{cross}} \quad (4.1)$$

The time when orbital crossing actually starts is some time ( $\tau_{\text{cross}}$ ) after this situation is satisfied. The damping time at this particular time is

$$\tau_{\text{damp}}(\text{inst}) = \tau_{\text{cross}} \exp\left(\frac{\tau_{\text{cross}}}{\tau_{\text{gas}}}\right). \quad (4.2)$$

In order to allow accretion,

$$\tau_{\text{damp}}(\text{inst}) \gtrsim \tau_{\text{growth}}. \quad (4.3)$$

Substituting Eq.(4.2) into Eq.(4.3), we find

$$\frac{\tau_{\text{cross}}}{\tau_{\text{gas}}} \gtrsim \ln\left(\frac{\tau_{\text{growth}}}{\tau_{\text{cross}}}\right) \sim 1. \quad (4.4)$$

If the mass of the largest planets become  $\sim M_{\oplus}$  on a time scale  $\tau_{\text{growth}}$ , the damping time scale would be

$$\tau_{\text{damp}}(\text{end}) = \frac{1}{5} \tau_{\text{damp}}(\text{inst}) \exp\left(\frac{\tau_{\text{growth}}}{\tau_{\text{gas}}}\right) \quad (4.5)$$

$$= \frac{1}{5} \tau_{\text{cross}} \exp\left(\frac{\tau_{\text{cross}} + \tau_{\text{growth}}}{\tau_{\text{gas}}}\right). \quad (4.6)$$

Because  $\tau_{\text{damp}}(\text{end})$  is for  $\sim M_{\oplus}$  while  $\tau_{\text{damp}}(\text{inst})$  is for  $\sim 0.2M_{\oplus}$ , right hand side of the above equations are divided by 5 (Eq.(2.3)). In order to damp the eccentricities of the planets,

$$\tau_{\text{damp}}(\text{end}) \lesssim \tau_{\text{gas}}. \quad (4.7)$$

Substituting Eq.(4.6), we acquire

$$\frac{\tau_{\text{cross}} + \tau_{\text{growth}}}{\tau_{\text{gas}}} \lesssim \ln\left(\frac{5\tau_{\text{gas}}}{\tau_{\text{cross}}}\right) \sim \text{afew}. \quad (4.8)$$

Taking 2 as the above numerical factor, we obtain the conditions for formation of Earth-like planets as

$$\frac{\tau_{\text{cross}} + \tau_{\text{growth}}}{2} \lesssim \tau_{\text{gas}} \lesssim \tau_{\text{cross}}. \quad (4.9)$$

**Table III**  
**Comparison of the Time Scales and the Final Planets**

Simulation	$A = \frac{\tau_{\text{cross}}}{\tau_{\text{gas}}}$	$B = \frac{2\tau_{\text{gas}}}{\tau_{\text{cross}} + \tau_{\text{growth}}}$	$A \geq 0.5$	$M_{\text{max}} \geq 0.8M_{\oplus}$	$B \geq 0.5$	$e_{\text{max}} \leq 0.04$
Run <sub>4</sub> <sup>3</sup>	0.003	30	no	no	yes	yes
Run <sub>8</sub> <sup>2</sup>	0.03	3.0	no	no	yes	yes
Run <sub>1</sub> <sup>1</sup>	0.20	5.0	no	no	yes	yes
Run <sub>6</sub> <sup>2</sup>	0.30	2.2	no	no	yes	yes
Run <sub>7</sub> <sup>2</sup>	0.30	2.2	no	no	yes	yes
Run <sub>4</sub> <sup>2</sup>	0.33	2.0	no	yes	yes	yes
Run <sub>5</sub> <sup>2</sup>	0.50	1.33	yes	yes	yes	yes
Run <sub>2</sub> <sup>1</sup>	0.60	2.2	yes	yes	yes	yes
Run <sub>2</sub> <sup>2</sup>	0.67	1.5	yes	yes	yes	yes/no
Run <sub>3</sub> <sup>2</sup>	0.67	1.5	yes	yes	yes	no
Run <sub>9</sub> <sup>2</sup>	1.25	1.0	yes	yes	yes	no
Run <sub>3</sub> <sup>3</sup>	3.3	0.20	yes	yes	no	no
Run <sub>1</sub> <sup>3</sup>	6.7	0.15	yes	yes	no	no

$A$  is the ratio of time scales  $\tau_{\text{cross}}$  to  $\tau_{\text{gas}}$ . If this value is larger than 0.5, mass of the largest planet is likely to be  $\sim M_{\oplus}$ .  $B$  is the ratio of  $2 \times \tau_{\text{gas}}$  to  $\tau_{\text{cross}} + \tau_{\text{growth}}$ . If this value is larger than 0.5, the eccentricities of the largest planets tend to be smaller than 0.04.

In Table III,  $\tau_{\text{cross}}/\tau_{\text{gas}}$  and  $2\tau_{\text{gas}}/(\tau_{\text{cross}} + \tau_{\text{growth}})$  with  $\tau_{\text{growth}} = 2 \times 10^6$  yrs are shown for each run. If both quantities  $\gtrsim 1$ , the conditions (4.9) are satisfied. Table III shoes the condition (4.9) is consistent with the numerical results.

Note that the case with  $\tau_{\text{cross}} \ll \tau_{\text{gas}}$  does not necessarily mean that Earth-sized planets cannot be formed. In this case, the amount of gas then is too much to allow the planets to grow enough. The eccentricities are damped rapidly. The separations  $\Delta a$  become wider resulting in a longer  $\tau_{\text{cross}}$ . Thus, the condition  $\tau_{\text{damp}} \gtrsim \tau_{\text{growth}}$  may be satisfied after the first orbital crossing stage. However, since  $\tau_{\text{cross}}$  increases so much with the expansion of  $\Delta a$ , the other condition  $(\tau_{\text{cross}} + \tau_{\text{growth}})/2 \lesssim \tau_{\text{gas}}$  would not be satisfied. In the case of Fig.1(c),  $\Delta a$  becomes  $\sim 15r_{\text{H}}$  after the orbital crossing, so that  $\tau_{\text{cross}}$  in the new orbital separation becomes  $10^9 - 10^{10}$  years. Therefore, we can say that the condition  $\tau_{\text{cross}} \gtrsim \tau_{\text{gas}}$  is required.

Oligarchic growth model (Kokubo and Ida 1998, 2000) shows that the protoplanets are formed with separations about  $10 r_{\text{H}}$ .  $\tau_{\text{cross}}$  for  $\Delta a \sim 10r_{\text{H}}$  is  $\sim 10^6 - 10^7$  years near 1 AU (Table I; Chambers et al., 1996). Considering the terrestrial planets, if the gas dissipates on a time scale of  $10^6 - 10^7$  years, planets with large mass and low eccentricities can be formed. This time scale is consistent with observation.

## 5 Summary

Here we summarize what we have done. We performed N-body simulation on formation of terrestrial planets from protoplanets including damping of velocity dispersion caused by gravitational interaction with a dissipating gas disk. We have investigated the effect of dissipating disk gas on orbital evolution of terrestrial planets. Calculations are started with 15 protoplanets and initial amount of gas ranges from 1 - 100 % of minimum mass model. The eccentricities are held low until the gas is depleted to the point when its damping time scale become comparable

to instability time scale (Eq.(2.9)). If we consider the terrestrial planets, damping time scale of the gas when the orbital crossing starts is  $10^6 - 10^7$  years. If the depletion time scale is several million years, the damping time scale after the accretion would be  $10^6 - 10^7$  years and damps the eccentricities of survived planets. As a result, planets with  $m \sim M_{\oplus}$  and  $e \sim 0.01$  are formed.

We have analytically derived the constraints on depletion time scale of the disk gas. We acquired a relation among the time scales,  $\tau_{\text{cross}}$ ,  $\tau_{\text{growth}}$ ,  $\tau_{\text{damp}}$  in order to form planets with  $m \sim M_{\oplus}$  and  $e \sim 0.01$ , as Eq.(4.9).

## References

- [1] Agnor, C. B., R. M. Canup, and H. F. Levison 1999. On the character and consequences of large impacts in the late stage of terrestrial planet formation. *Icarus*,**142**,219-237.
- [2] Agnor, C. B., and Ward, W. R. 2002. Damping of terrestrial-planet eccentricities by density-wave interactions with a remnant gas disk. *APJ*,**567**,579-586.
- [3] Artymowicz, P. 1993. Disk-satellite interaction via density wave and the eccentricity evolution of bodies embedded in disks. *Astron. J.*,**419**,166-180.
- [4] Chambers, J. E., G. W. Wetherill, and A. P. Boss 1996 The stability of multi-planet systems. *Icarus*,**119**,261-268.
- [5] Chambers, J. E., and G. W. Wetherill 1998. Making the terrestrial planets: N-body integrations of planetary embryos in three dimensions. *Icarus*,**136**,304-327.
- [6] Hayashi, C. 1981. Structure of the solar nebula, growth and decay of magnetic fields and effects of magnetic and turbulent viscosities on the nebula. *Prog. Theor. Phys. Suppl.*,**70**,35-53.
- [7] Ito, T., and K. Tanikawa 1999. Stability and instability of the terrestrial protoplanet system and their possible roles in the final stage of planet formation. *Icarus*,**139**,336-349.
- [8] Iwasaki, K., H. Emori, K. Nakazawa, and H. Tanaka, 2002. Orbital stability of a protoplanet system under the drag force proportional to the random velocity. *PASJ*,in press.
- [9] Kokubo, E., S. Ida 1998. Oligarchic growth of protoplanets. *Icarus*,**131**,171-178.
- [10] Kokubo, E., and S. Ida 2000. Formation of protoplanets from planetesimals in the solar nebula. *Icarus*,**143**,15-27.
- [11] Kominami, J., and S. Ida 2002. The effect of tidal interaction with a gas disk on formation of terrestrial planets *Icarus*,in press
- [12] Makino, J. 1991. Optimal order and time-step criterion for Aarseth-type N-body integrators. *A.P.J.*,**369**,200-212.
- [13] Makino, J., and S. J. Aarseth 1992. On a hermite integrator with Ahmad-Cohen scheme for gravitational many-body problems. *Publ. Astron. Soc. Jpn.*,**44**,141-151.
- [14] Nagasawa, M., H. Tanaka, and S. Ida 2000. Orbital evolution of asteroids during depletion of solar nebula. *Astron. J.*,**119**,1480-1497.
- [15] Ward, W. R. 1989. On the rapid formation of giant planet cores. *Astrophys.J.Lett*,**345**,L99-L102.

- [16] Ward, W. R. 1993. Density Wave In The Solar Nebula: Planetesimal Velocities. *ICARUS*, **106**, 274-287.

# The Evidence of a Stellar Encounter on the distribution Edgeworth-Kuiper Belt Object.

Hiroshi Kobayashi, Shigeru Ida, and Hidekazu Tanaka

*Department of Earth and Planetary Sciences,  
Tokyo Institute of Technology*

hkobayas@geo.titech.ac.jp

## ABSTRACT

We show that a stellar encounter may explain high eccentricity ( $e$ ) and inclination ( $i$ ) [rad] in the outer part ( $\geq 40$  AU) Edgeworth-Kuiper belt objects and considering the effect of gas drag after the stellar encounter, we may explain the observed bimodal orbits of outer Edgeworth-Kuiper belt objects, that is,  $i \sim e$  or  $i \gg e$ . We investigated the  $e$  and  $i$  of planetesimals pumped-up by a passing star and, then the change in  $e$  and  $i$  for nebula gas drag and Neptune scattering. We model a protoplanetary system as a disk of massless particles circularly orbiting a host star. The massless particles represent planetesimals. A single star as massive as the host star encounters the protoplanetary system. Numerical orbital simulations show that in the inner region at semimajor axis  $a \lesssim 0.3D$  where  $D$  is pericenter distance of the encounter, the disk is intact, and that in the outer region  $a \gtrsim 0.2D$ ,  $e$  and  $i$  are highly pumped up. If  $D \sim 120$  AU, the pumped-up  $e$  and  $i$  in outer region ( $\gtrsim 40$  AU) are as large as  $e$  and  $i$  of Edgeworth-Kuiper belt objects. However, Their  $e$  is  $\sim i$ . To investigate the effect of gas drag for a long time (the lifetime of nebula gas  $\sim 5 \times 10^7$  yr), we derive analytical equation of change in  $a$ ,  $e$ , and  $i$  of planetesimal caused by the gas drag. We calculate the orbital changing of planetesimals pumped-up by the passing star due to gas drag, using the equations. The gas drag affects on only a planetesimal with high  $e$ , and then they migrate to the inner region. We also consider Neptune scattering of the planetesimal, so that the planetesimals with high  $e$  are ejected from Edgeworth-Kuiper belt. As the result, the stellar encounter in some parameters and the effect of gas drag can explain the observed bimodal orbits of Edgeworth-Kuiper belt objects.

## 1. INTRODUCTION

More than 400 small objects are observed in Edgeworth-Kuiper belt region. The observed EKBO are on very eccentric and inclined orbits (Fig. 1). That is the evidence that Edgeworth-Kuiper belt objects (EKBO) are strongly excited. EKBO are dynamically divided into three classes: bodies locking in the 3:2 resonance with Neptune (Plutinos), bodies with semi-major axis between 42 AU and 50 AU (classical disk), bodies with perihelion distance 30-35 AU scattered by Neptune (scattered disk). Plutinos may be caused by the sweeping of the 3:2 resonance of Neptune during outward migrating (Malhotra 1995). The classical belt do not locate in strong mean motion resonances or secure resonances. There is possibility that classical EKBO experienced dynamical excitation in past.

Some theoretical mechanism are proposed as the excitation. petit *et al.* (1999) proposed that the hypothetical Earth-size planets pumped up eccentricity ( $e$ ) and inclination ( $i$ ) of small planetesimals, before they are ejected by Neptune. Nagasawa & Ida (2000) suggested that the sweeping secular resonances during the primitive solar nebula depletion caused the excitations of  $e$  and  $i$  of EKBO. These models can roughly explain high  $e$  and  $i$  of the classical EKBO. However, inclination distribution of the classical EKBO can not fit single Gaussian but well fit two Gaussians of widths about 0.04 and 0.3 Brown (2001). The theoretical model never explain the bimodal inclination distribution, because  $e$  is pumped-up as highly as  $i$  (see Fig. 2).

Ida *et al.* (2000) shows that a early stellar encounter with the pericenter distance  $D = 100\text{-}200$  AU can explain high  $e$  and  $i$  of EKBO and that the outward migration can capture objects after a stellar encounter. After a stellar encounter, the bodies are excited in outer region and ones are intact in inner region. The stellar encounter result in deciding the boundary radius of inner Edgeworth-Kuiper belt, that is about  $D/3$  (Kobayashi & Ida 2001). It is also difficult for the stellar encounter model to explain the bimodal distribution of EKBO, because  $i$  is pumped-up as highly as  $e$ . However, there are many encounter parameters and we search for the successive sets of parameters. We also consider the effects of Neptune and the gas drag on EKBO after the early stellar encounter. If the orbit of the body is across the Neptune orbit, it is scattered by Neptune. As the small-body which is excited by a stellar encounter would have lived in Solar nebula, they are under the effect of gas drag. The gas drag does not affect on the bodies in outer region well, because the gas nebula is not so dense there. However, the bodies with high  $e$  are made to change on the orbits by the gas drag, because they pass the dense gas nebula at the near perihelion. The gas drag effect may change the distribution of  $e$  and  $i$ .

We investigate the effect of the stellar encounter on a planetesimal disk, and that the orbital changing of planetesimals under the effects of Neptune scatter and the gas drag in

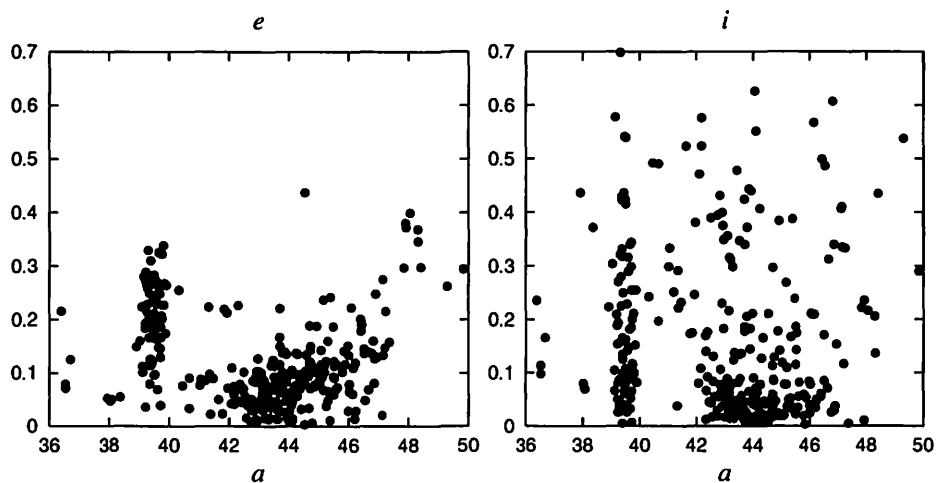


Fig. 1.—  $e$  and  $i$  of the classical EKBO and Prutinos as a function of  $a$ .

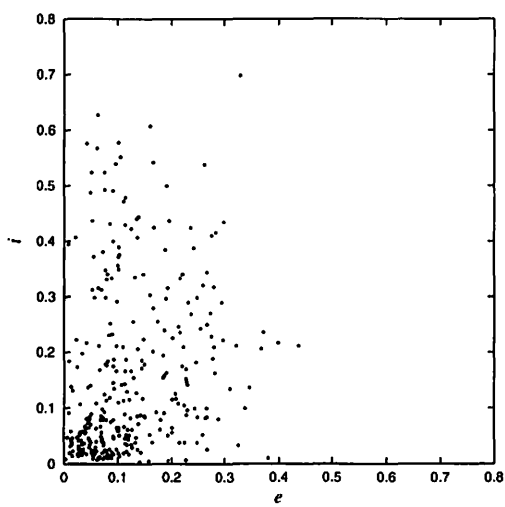


Fig. 2.—  $i$  of the classical EKBO as a function of  $e$ .

EKBO region.

## 2. STELLAR ENCOUNTER

### 2.1. MODEL

To study the dynamical effect of a stellar encounter on a planetesimal disk, we follow Kobayashi & Ida (2001). We model a planetesimal disk as non-self-gravitating, collision-less particles that initially have coplanar circular orbits around a primary (host) star, because two-body relaxation time and mean collision time of planetesimals are much longer than encounter time scale (Kobayashi & Ida 2001). The particulate disk encounters a hypothetical passing star.

We integrated orbits of 10,000 particles with surface number density  $n_s \propto a^{-3/2}$ . The particles are distributed in the region  $a = 40\text{-}80$  AU. We took the scale length of encounter the pericenter distance  $D$  of the passing star, and the parameters for our modeling stellar encounter are the inclination ( $i_*$ ) relative to the initial planetesimal disk, eccentricity ( $e_*$ ), and argument of perihelion ( $\omega_*$ ) of orbit of the passing star, and the scaled passing star mass ( $M_* = M_2/M_1$ ). The encounter geometry is illustrated in Fig. 3. We approximate  $e_*$  and  $M_*$  are the unit (Kobayashi & Ida 2001).

### 2.2. RESULTS

Figures 4 and 5 show  $e$  and  $i$  pumped-up by the stellar encounter with some sets of parameters as a function with  $a/D$ , where  $a$  is semi-major axis of a planetesimal. In inner region ( $a \ll 0.2D$ ) of the disk, the planetesimals are intact. On the other hand,  $e$  and  $i$  are highly pumped-up and have steep radius gradient in outer region. The distributions of  $e$  and  $i$  pumped-up by the stellar encounters strongly depend on the encounter parameters in outer region.

We choose 1,000 bodies in the initial disk from 40 AU to 80 AU to decide  $D$ . We investigate the planetesimals whose semi-major axis  $a$  are distribute from 40 AU to 50 AU, that is, Edgeworth-Kuiper belt region. Fig. 6 and 7 shows  $i$  of these planetesimals as a function of  $e$  of themselves, in the case with  $D = 120$  AU and 140 AU. In the case of some sets of parameters,  $e$  and  $i$  are pumped-up enough highly. In the case of  $D = 140$  AU,  $i_* = 60^\circ$ , and  $\omega_* = 90^\circ$ , the distribution of  $e$  and  $i$  is the best fit to one of EKBO. However, some planetesimal's  $e$  ( $\gg 0.1$ ) is as large as  $i$ . In section 4, we investigate the effects after a



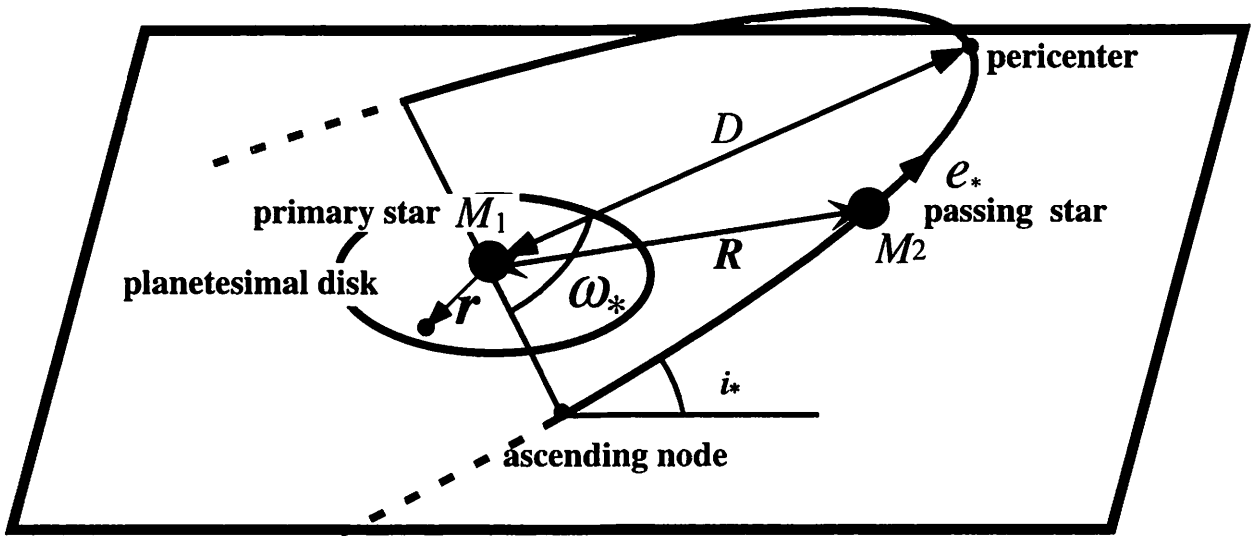


Fig. 3.— Encounter configuration in the frame centered at the primary star with mass  $M_1$ . The orbit of the passing star with mass  $M_2$  is characterized by pericenter distance  $D$ , eccentricity  $e_*$ , inclination  $i_*$  and argument of perihelion  $\omega_*$ . If length and mass are scaled by  $D$  and  $M_1$ , the encounter parameters are  $M_*$  ( $= M_2/M_1$ ),  $e_*$ ,  $i_*$  and  $\omega_*$ .

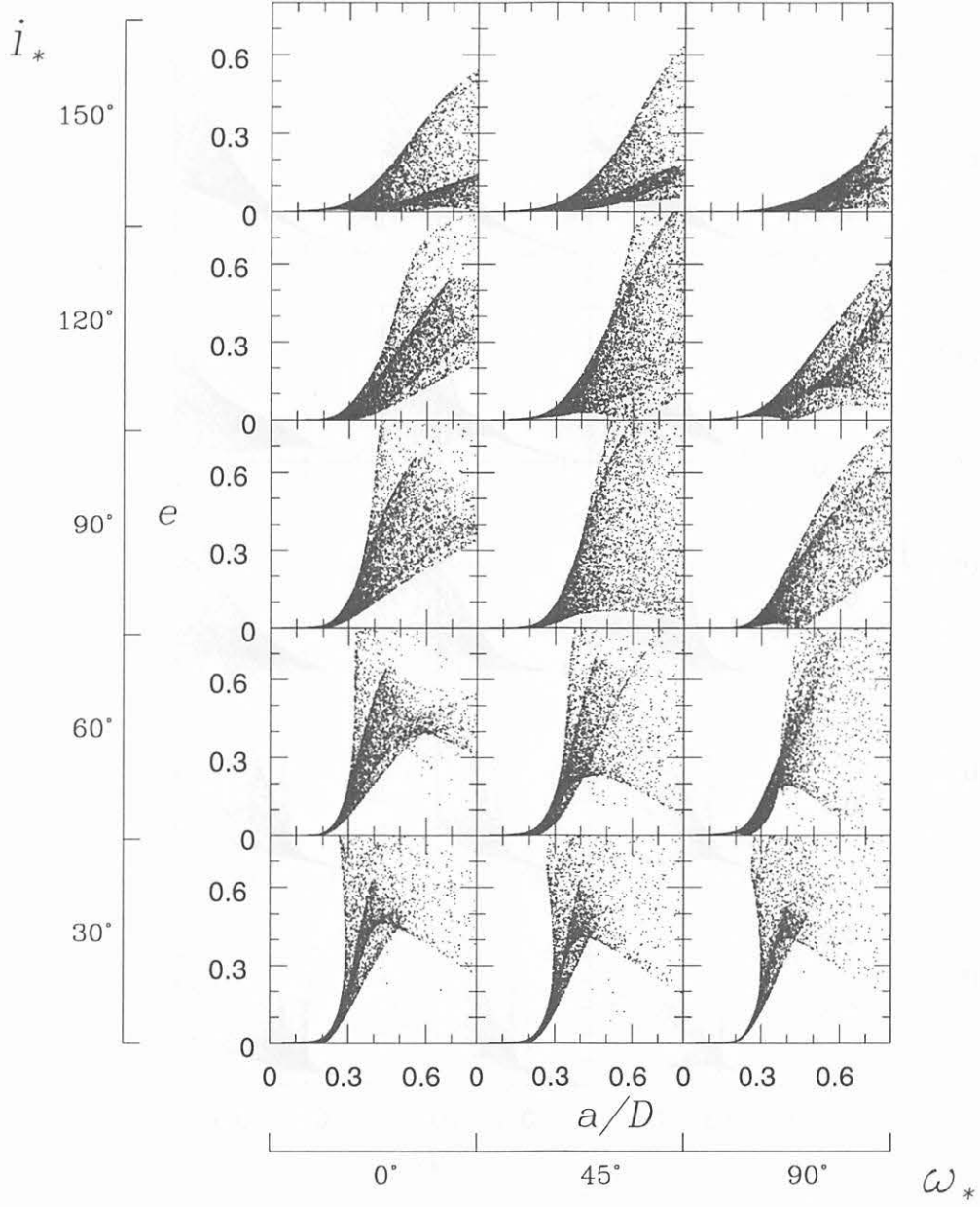


Fig. 4.—  $e$  of the planetsimals pumped-up by the passing star as a function of  $a/D$ . Top to bottom,  $i_*$  are  $i_* = 150^\circ$ ,  $i_* = 120^\circ$ ,  $i_* = 90^\circ$ ,  $i_* = 60^\circ$ , and  $i_* = 30^\circ$ . Left to right,  $\omega_*$  are  $\omega_* = 0^\circ$ ,  $\omega_* = 45^\circ$ , and  $\omega_* = 90^\circ$ .

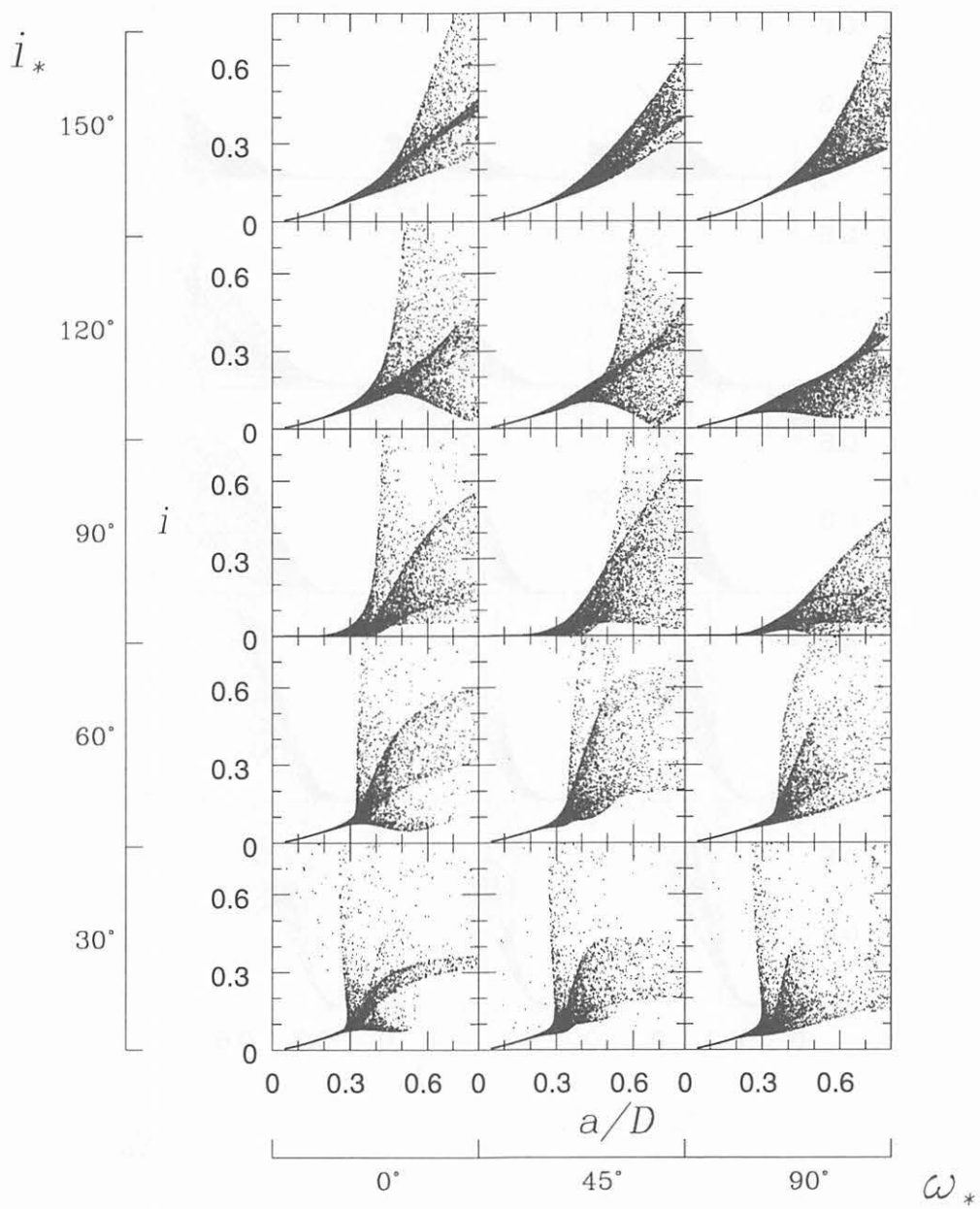


Fig. 5.—  $i$  of the planetsimals pumped-up by the passing star as a function of  $a/D$ .

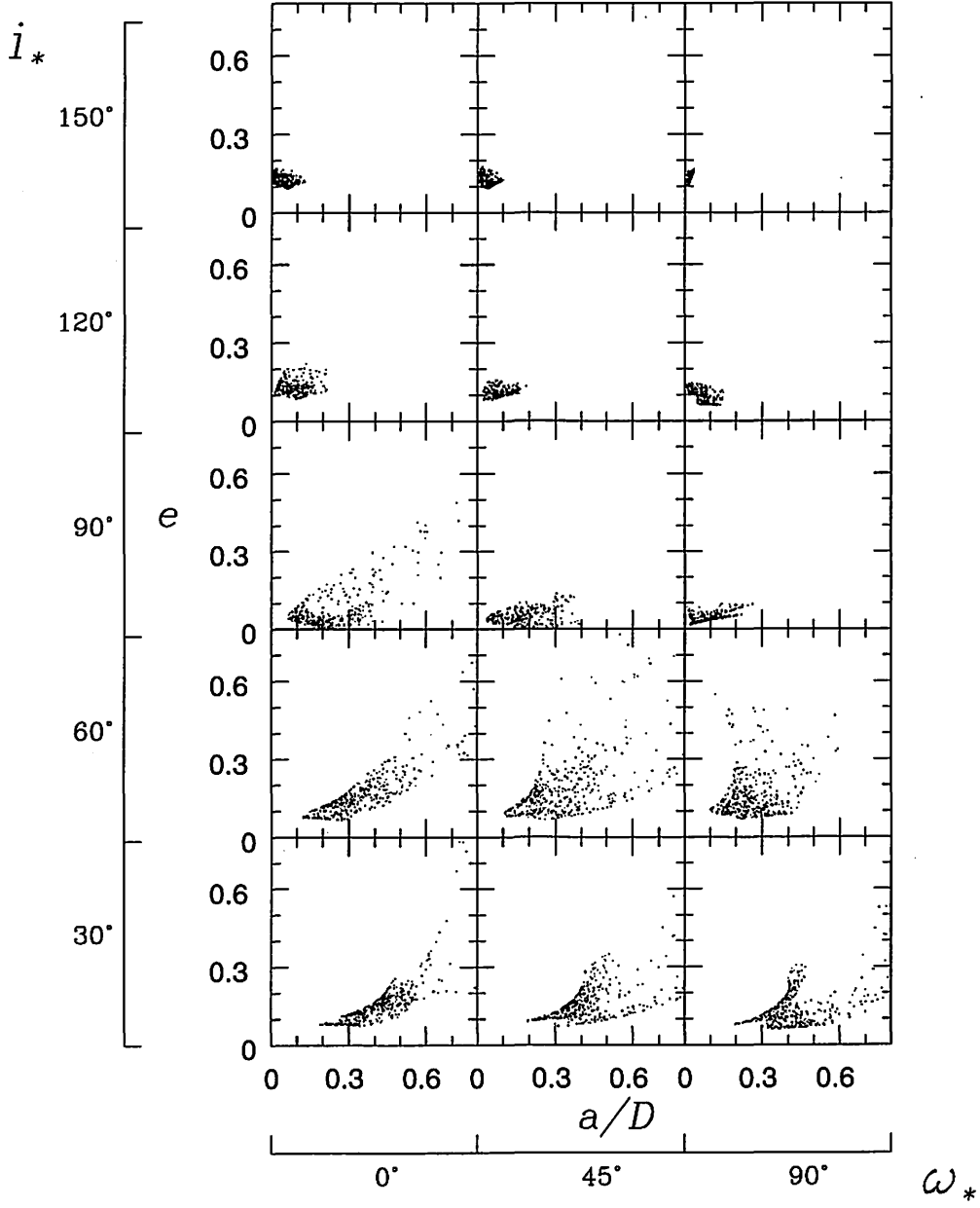


Fig. 6.—  $i$  of the planetesimals with  $a$  between 40 AU and 50 AU as a function of  $e$ , after a stellar encounter with  $D = 120$  AU. We select 1000 bodies with initial  $a$  between 40 AU to 80 AU.

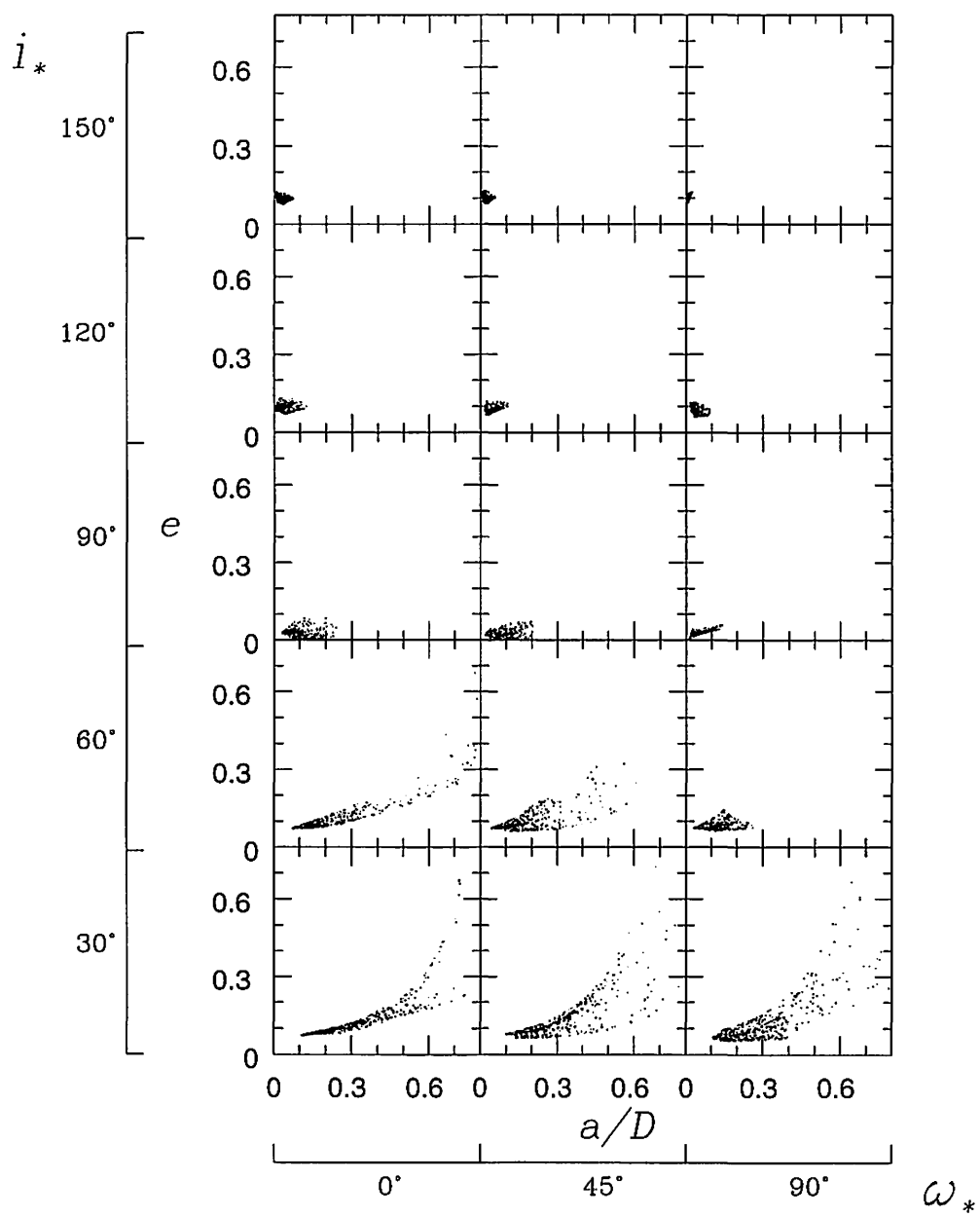


Fig. 7.— Same as Fig. 6, except for  $D = 140$  AU.

stellar encounter.

### 3. GAS DRAG EFFECT

We investigate the orbital evolution of planetesimals exited by a stellar encounter due to gas drag. Adachi *et al.* (1976) derive the mean variation of  $a$ ,  $e$ , and  $i$  for gas drag, in the case that  $e$  and  $i$  are much smaller than unity. However,  $e$  and  $i$  pumped-up by the stellar encounter is as large as unity. We newly derive the mean variation of  $a$ ,  $e$ , and  $i$ , in the case with  $e \sim 1$  and  $i \ll h/a$ , or with  $i \gg h/a$ , where  $h$  is the scale height of gas nebula.

#### 3.1. NEBULA DISK AND GAS DRAG LAW

We consider a gaseous nebula disk rotating around a central star with mass  $M$ . We assume that the nebula disk is axisymmetric and in a steady state. To describe the disk, we use a cylindrical coordinate system  $(r, \theta, z)$ . The  $z$  axis is coincide with the rotation axis of the disk. Then, due to the above assumption, the gas velocity has only the  $\theta$ -component:  $(0, r\Omega_{\text{gas}}, 0)$ , where  $\Omega_{\text{gas}}$  is the angular velocity of the gas. The angular velocity  $\Omega_{\text{gas}}$  and the density distribution of the nebular disk is described by Euler's equation

$$(\mathbf{v}_{\text{gas}} \cdot \nabla) \mathbf{v}_{\text{gas}} = -\frac{1}{\rho} \nabla p - \nabla \left( -\frac{GM}{\sqrt{r^2 + z^2}} \right) \quad (1)$$

where  $G$  is the gravitational constant and  $\rho$ ,  $p$  and  $\mathbf{v}_{\text{gas}}$  are the density, the pressure, and the velocity of gas, respectively. We further assume the disk is geometrically thin and isothermal, that is,  $p = c^2 \rho$ , where  $c$  is the isothermal sound speed. Using these assumptions and the  $z$ -component of Eq. (1), we obtain the  $z$  dependence of  $\rho$  as

$$\rho = \frac{\sigma}{\sqrt{2\pi}h} \exp \left( -\frac{z^2}{2h^2} \right), \quad (2)$$

where  $\sigma (= \int_{-\infty}^{\infty} \rho dz)$  is the surface density of the nebula disk;  $h$  is the scale height of the disk given by  $h = c/\Omega_K$ ; and  $\Omega_K = (GM/r^3)^{1/2}$  is Keplerian angular velocity. For the simplicity,  $\sigma \propto r^{-\alpha}$ ,  $c \propto r^{-\beta}$ , respectively.  $c$  are assumed as In the minimum-mass solar nebula model proposed by Hayashi *et al.* (1985), for example, the surface density and the sound velocity is given by

$$\begin{cases} \sigma = 1.7 \times 10^3 \left( \frac{r}{1\text{AU}} \right)^{-3/2} [g/cm^3], \\ c = 9.9 \times 10^4 \left( \frac{r}{1\text{AU}} \right)^{-1/4} [cm/s]. \end{cases} \quad (3)$$

The angular velocity  $\Omega_{\text{gas}}$  is obtained from the  $r$ -component of Eq. (1) as (Tanaka *et al.* 2002)

$$\Omega_{\text{gas}} = \Omega_K \left[ 1 + (\alpha + 2\beta) \frac{h^2}{r^2} + \beta \frac{z^2}{r^2} \right]^{1/2}, \quad (4)$$

In the derivation of Eq. (4), we neglected the term of  $O(z^4/r^4)$  and the higher. This approximation is valid even for the investigation of the gas drag effect on highly inclined orbits because we do not have to consider the gas drag (and the nebula gas) at a high altitude such as  $z \sim h$ .

We consider the planetesimals of which size is larger than 1 km. For such large body, the gas drag force is described by Newton's law. Thus, the gas drag force per unit mass can be written as (Adachi *et al.* 1976)

$$F_i = A\rho \mid \mathbf{u} \mid u_i, \quad (5)$$

with  $A = C_D \pi d^2 / 2m$ , where  $d$  is the radius of particle, and  $C_D$  is the non-dimensional drag coefficient of which value is  $0.5 \leq C_D \leq 1.5$ . The relative velocity between the body and the gas,  $\mathbf{u}$ , is given by  $\mathbf{u} = \mathbf{v} - \mathbf{v}_{\text{gas}}$ , where  $\mathbf{v}$  is the velocity of the body. Hence, using Eqs. (10), (4), and (5), we can evaluate the drag force on the body by the nebula disk. The ratio of this gas drag force and gravity force of central star is  $4.4 \times 10^{-4}$ . The gas drag force is much smaller than the gravity force, and the time scale of change in orbital elements caused by gas drag force is much longer than Kepler time. Adachi *et al.* (1976) estimated this time scale for body with mass  $m$  and density  $\rho_{\text{mat}}$  as

$$\tau_0 = \frac{1}{A\rho v_K} = \frac{365.}{C_D} \left( \frac{m}{10^{21}\text{g}} \right)^{1/3} \left( \frac{a}{1\text{AU}} \right)^{-1-\alpha} \left( \frac{\rho_{\text{mat}}}{3\text{gcm}^{-3}} \right)^{2/3} \left( \frac{\rho_0}{10^9\text{gcm}^{-3}} \right)^{-1} T_K, \quad (6)$$

where  $T_K$  is Kepler time of the body. In outer region,  $\tau_0$  is much larger.

### 3.2. GENERAL EXPRESSIONS FOR THE CHANGE IN $a$ , $e$ AND $i$

We examine the time variations of semi-major axis  $a$ , eccentricity  $e$ , and inclination  $i$  for a body due to the gas drag force. Such time variations are described by using Gauss's equations:

$$\begin{cases} \frac{da}{dt} = \frac{2}{na} \left( F_R \frac{ae}{\eta} \sin f + F_\phi \frac{a^2 \eta}{R} \right), \\ \frac{de}{dt} = \frac{\eta}{na} [F_R \sin f + F_\phi (\cos f + \cos E)], \\ \frac{di}{dt} = \frac{1}{na\eta} F_\zeta \frac{R}{a} \cos(f + \omega), \end{cases} \quad (7)$$

where  $f$ ,  $n$ , and  $\omega$  are the true anomaly, the mean motion, and the argument of pericenter, respectively and  $R$ -,  $\phi$ - and  $\zeta$ -directions are the cylindrical coordinate system on the basis of the orbital plane of the particle. We express the gas velocity in this coordinate,

$$\begin{cases} v_{\text{gas},R} = 0, \\ v_{\text{gas},\phi} = r\Omega_{\text{gas}}(r, z)\frac{\cos i}{w}, \\ v_{\text{gas},\zeta} = -r\Omega_{\text{gas}}(r, z)\frac{\cos(f + \omega)\sin i}{w}, \end{cases} \quad (8)$$

where  $w = [1 - \sin^2(f + \omega)\sin^2 i]^{1/2} = r/R$ .

Using Eq. (5) and gas and particle velocity, we can rewrite Eqs. (7),

$$\begin{cases} \frac{da}{dt} = -A\rho u \frac{2a}{\eta^2} [1 + 2e \cos f + e^2 - (1 + e \cos f)^{3/2} \kappa \cos i], \\ \frac{de}{dt} = -A\rho u \left[ 2 \cos f + 2e - \frac{2 \cos f + e + e \cos^2 f}{(1 + e \cos f)^{1/2}} \kappa \cos i \right], \\ \frac{di}{dt} = -A\rho u \frac{\cos^2(f + \omega)}{\eta^2(1 + e \cos f)^{1/2}} \kappa \sin i, \end{cases} \quad (9)$$

where  $\rho$ ,  $\kappa$ , and  $u$  are given by

$$\rho(r, z) = \rho_0(a) \left( \frac{\eta^2 w}{1 + e \cos f} \right)^{-\alpha} \exp \left( -\frac{a^2 \eta^4 \sin^2(\omega + f) \sin^2 i}{2h^2(1 + e \cos f)^2} \right), \quad (10)$$

$$\kappa = \frac{\Omega_{\text{gas}}}{\Omega_K w^{3/2}}, \quad (11)$$

$$\begin{aligned} u &= |\mathbf{u}| \\ &= \frac{v_K(a)}{\eta} [1 + 2e \cos f + e^2 - 2(1 + e \cos f)^{3/2} \kappa \cos i + (1 + e \cos f) \kappa^2 w^2]^{1/2}. \end{aligned} \quad (12)$$

Since the gas drag force being much smaller than gravity force of the central star for the body, as before, we can assume that  $a$ ,  $e$ ,  $i$  and  $\omega$  is constant in one orbital period. Thus, in the R. H. Ss. of Eqs. (9) - (12), only the true anomaly  $f$  is dependent on  $t$ . We consider the changes in  $a$ ,  $e$  and  $i$  averaged over one orbital period. The averaged changes are defined as

$$\left\langle \frac{da}{dt} \right\rangle = \frac{1}{T_K} \int_0^{T_K} \frac{da}{dt} dt = \frac{1}{2\pi} \int_0^{2\pi} \frac{da}{dt} \frac{\eta^3}{(1 + e \cos f)^2} df. \quad (13)$$

We also treat with  $e$  and  $i$ , using same averaging.

Adachi *et al.* (1976) derive the averaged changes in  $a$ ,  $e$  and  $i$ , in the case that  $e$  and  $i$  are much smaller than 1. they consider that  $\Omega_K = \Omega_{\text{gas}}(1 - \eta_d/2)^{1/2}$  and that  $\eta_d$  is depend



on  $a$ . They derive the change in  $a$ ,  $e$ , and  $i$  in the case that  $e$ ,  $i$ , and  $\eta_d$  is much smaller than the unity. Inaba *et al.* (2000) modified their results in the form

$$\begin{cases} \frac{\tau_0}{a} \left\langle \frac{da}{dt} \right\rangle = -2[(0.97e)^2 + (0.64i)^2 + \eta_d^2]^{1/2} \eta_d, \\ \frac{\tau_0}{e} \left\langle \frac{de}{dt} \right\rangle = -[(0.77e)^2 + (0.64i)^2 + (1.5\eta_d)^2]^{1/2}, \\ \frac{\tau_0}{i} \left\langle \frac{di}{dt} \right\rangle = -\frac{1}{2}[(0.77e)^2 + (0.85i)^2 + \eta_d^2]^{1/2}, \end{cases} \quad (14)$$

where

$$\tau_0 = \frac{365}{C_D} \left( \frac{m}{10^{21} \text{g}} \right)^{1/3} \left( \frac{a}{1 \text{AU}} \right)^{11/4} \left( \frac{\rho_{\text{mat}}}{3 \text{gcm}^{-3}} \right)^{2/3} T_K.$$

Equations (14) are valid for small  $e$  and  $i$ . If  $i$  is large enough, the body goes at a high altitude where the gas density is extremely low. They do not consider such the effects. We consider them, and attack the problem. As the result, we derive the changes, in the case that  $e$  is large or  $i$  is large.

### 3.3. CASE OF LARGE ECCENTRICITY AND SMALL INCLINATION

Here we consider the case that  $e$  is almost equal to the unity and  $i$  is much smaller than the unity. Expanding Eqs. (9) with respect to  $(1 - e^2)$ , and keeping only the lowest order terms of  $(1 - e^2)$ , we can rewrite Eqs. (9) as

$$\begin{cases} \frac{da}{dt} = -2 \frac{a}{\tau_0} \eta^{-2\alpha-3} \Phi(f), \\ \frac{de}{dt} = -\frac{1}{\tau_0} \eta^{-2\alpha-1} \Phi(f), \\ \frac{di}{dt} = -\frac{i}{\tau_0} \eta^{-2\alpha-1} \Psi(f), \end{cases} \quad (15)$$

where

$$\begin{cases} \Phi(f) = (1 + \cos f)^{\alpha+3/2} (2 - \sqrt{1 + \cos f}) \sqrt{3 - 2\sqrt{1 + \cos f}}, \\ \Psi(f) = (1 + \cos f)^\alpha \sqrt{3 - 2\sqrt{1 + \cos f}} \cos^2(f + \omega). \end{cases} \quad (16)$$

In Eqs. (15), we neglected the terms of  $O(i^2)$  and the higher. The changes in  $a$ ,  $e$ , and  $i$  caused by the large  $e$  are much larger than the velocity difference between the nebula gas and the plane body, we can consider that  $\mathbf{v}_{\text{gas}}$  is equal to  $\mathbf{v}_K$ . To take the orbital average on

Eqs. (15), we only have to integrate  $\Phi(f)$  and  $\Psi(f)$  over the orbital period. That is,

$$\left\{ \begin{array}{l} \left\langle \frac{da}{dt} \right\rangle = \frac{a}{\pi\tau_0} (1 - e^2)^{-\alpha} \bar{\Psi}, \\ \left\langle \frac{de}{dt} \right\rangle = \frac{1}{2\pi\tau_0} (1 - e^2)^{-\alpha+1} \bar{\Psi}, \\ \left\langle \frac{di}{dt} \right\rangle = \frac{i}{2\pi\tau_0} (1 - e^2)^{-\alpha+1} (\bar{\Phi}_1 + \bar{\Phi}_2 \sin^2 \omega). \end{array} \right. \quad (17)$$

Through numerical integration of Eq. (16),  $\bar{\Psi}$  and  $\bar{\Phi}$  are obtained as functions of  $\alpha$ . In the minimum-mass solar nebular model,  $\alpha$  is  $11/4$  and, then,  $\bar{\Psi} = 4.94$  and  $\bar{\Phi} = 1.75 + 0.98 \sin^2 \omega$ .

Figures 8 show the change in  $a$ ,  $e$ , and  $i$ , as the function of  $e/(1 - e^2)$ , in the case of  $i = 0.01$ . The full circles, the dash line, and the short dash line show the result of the numerical orbital integration, Eqs. (17), and Eqs. (14) (Adachi *et al.* 1976), respectively. when  $e/(1 - e^2)$  is small, that is,  $e \ll 1$ , Eqs. 14 (Adachi *et al.* 1976) are consistent with the numerical result. In the case that  $e$  is large, that is,  $1 - e^2 \ll 1$ : large  $1/(1 - e^2)$ , Eqs. 14 (Adachi *et al.* 1976) are consistent with the numerical result. We can derive the changes in  $a$ ,  $e$ , and  $i$ , in the case that  $e$  is large.

### 3.4. CASE OF LARGE INCLINATION

Next, we consider highly inclined such that  $\sin i$  is much larger than  $h/a$ . Bodies with such a large inclination penetrate the nebula disk twice (near the ascending and the descending nodes) in an orbital period. Then we only have to examine the gas drag effect at the penetrations. First, We consider the penetration near the ascending node ( $f \simeq \omega$ ). Since, in highly inclined orbit, the duration of a penetration is much shorter than the orbital period, the only density  $\rho$  changes during a penetration in the R. H. Ss. of Eqs. (9). Thus we put  $f = -\omega$  in others factors of Eqs. (9) for penetrations near the ascending node and expand  $\rho$  with  $f + \omega$ . As a result, we have

$$\left\{ \begin{array}{l} \frac{da}{dt} = -\frac{2a}{\tau_0\eta^3} G(\omega) \exp \left( -\frac{a^2\eta^4(f + \omega)^2 \sin^2 i}{2h^2(1 + e \cos f)^2} \right), \\ \frac{de}{dt} = -\frac{1}{\tau_0\eta} H(\omega) \exp \left( -\frac{a^2\eta^4(f + \omega)^2 \sin^2 i}{2h^2(1 + e \cos f)^2} \right), \\ \frac{di}{dt} = -\frac{\sin i}{\tau_0\eta^3} I(\omega) \exp \left( -\frac{a^2\eta^4(f + \omega)^2 \sin^2 i}{2h^2(1 + e \cos f)^2} \right), \end{array} \right. \quad (18)$$

where

$$\begin{cases} G(\omega) = \tilde{r}^{-\alpha} \tilde{u} [1 + 2e \cos \omega + e^2 - (1 + e \cos \omega)^{3/2} \cos i], \\ H(\omega) = \tilde{r}^{-\alpha} \tilde{u} \left[ 2(e + \cos \omega) - \left( \cos \omega + \frac{\cos \omega + e}{1 + e \cos \omega} \right) \cos i \sqrt{1 + e \cos \omega} \right], \\ I(\omega) = \tilde{r}^{-\alpha+1} \tilde{u} \sqrt{1 + e \cos \omega}, \end{cases} \quad (19)$$

and

$$\begin{cases} \tilde{r} = \frac{\eta^2}{1 + e \cos \omega}, \\ \tilde{u} = \sqrt{2 + 3e \cos \omega + e^2 - 2(1 + e \cos \omega)^{3/2} \cos i}. \end{cases} \quad (20)$$

In the above we also assumed that  $v_{\text{gas}} = v_K$  since the relative velocity is large in this case, too. Since only the exponential functions includes  $f$  in Eqs. (18), we can take the time integration them. If we put  $\omega = -\omega - \pi$  in (18), we have the changes during penetration near the descending node. Considering the changes during two penetration, we can derive the averaged change as:

$$\begin{cases} \left\langle \frac{da}{dt} \right\rangle = -\frac{a}{\tau_0} \frac{\sqrt{2}h_0}{\sqrt{\pi}a\eta^4 \sin i} [\tilde{r}(\omega)^{\gamma+1} G(\omega) + \tilde{r}(\omega + \pi)^{\gamma+1} G(\omega + \pi)], \\ \left\langle \frac{de}{dt} \right\rangle = -\frac{1}{\tau_0} \frac{h_0}{\sqrt{2\pi}a\eta^2 \sin i} [\tilde{r}(\omega)^{\gamma+1} H(\omega) + \tilde{r}(\omega + \pi)^{\gamma+1} H(\omega + \pi)], \\ \left\langle \frac{di}{dt} \right\rangle = -\frac{1}{\tau_0} \frac{h_0}{\sqrt{2\pi}a\eta^4} [\tilde{r}(\omega)^{\gamma+1} I(\omega) + \tilde{r}(\omega + \pi)^{\gamma+1} I(\omega + \pi)], \end{cases} \quad (21)$$

When  $e \ll 1$ , we can rewrite Eq. (25) as

$$\begin{cases} \left\langle \frac{da}{dt} \right\rangle = -\frac{2\sqrt{2}h_0(1 - \cos i)^{3/2}}{\sqrt{\pi}\tau_0 \sin i}, \\ \left\langle \frac{de}{dt} \right\rangle = -\frac{\sqrt{2}h_0(1 - \cos i)^{1/2}}{\sqrt{\pi}\tau_0 a \sin i} \left[ (2 - \cos i) + 2 \left( \alpha - \gamma - \frac{1}{4} \right) (1 - \cos i) \cos^2 \omega \right] e, \\ \left\langle \frac{di}{dt} \right\rangle = -\frac{\sqrt{2}h_0(1 - \cos i)^{1/2}}{\sqrt{\pi}\tau_0 a}. \end{cases} \quad (22)$$

We also compare the numerical result with the analytical equations, in this case. Figures 9 are the changes in  $a$ ,  $e$ , and  $i$ , as a function of the initial  $i$ , in the case of the initial  $e = 0.01$  and  $h_0(a)/a = 4.7 \times 10^{-2}$ . When  $i$  is large, Eqs.(25) are consistent with the numerical results. Figures 10 are the same as Figs.9 except for the initial  $e = 0.9$ . This figures show Eqs.(25) are also consistent for the large  $i$  in the case that the initial  $e$  is large.

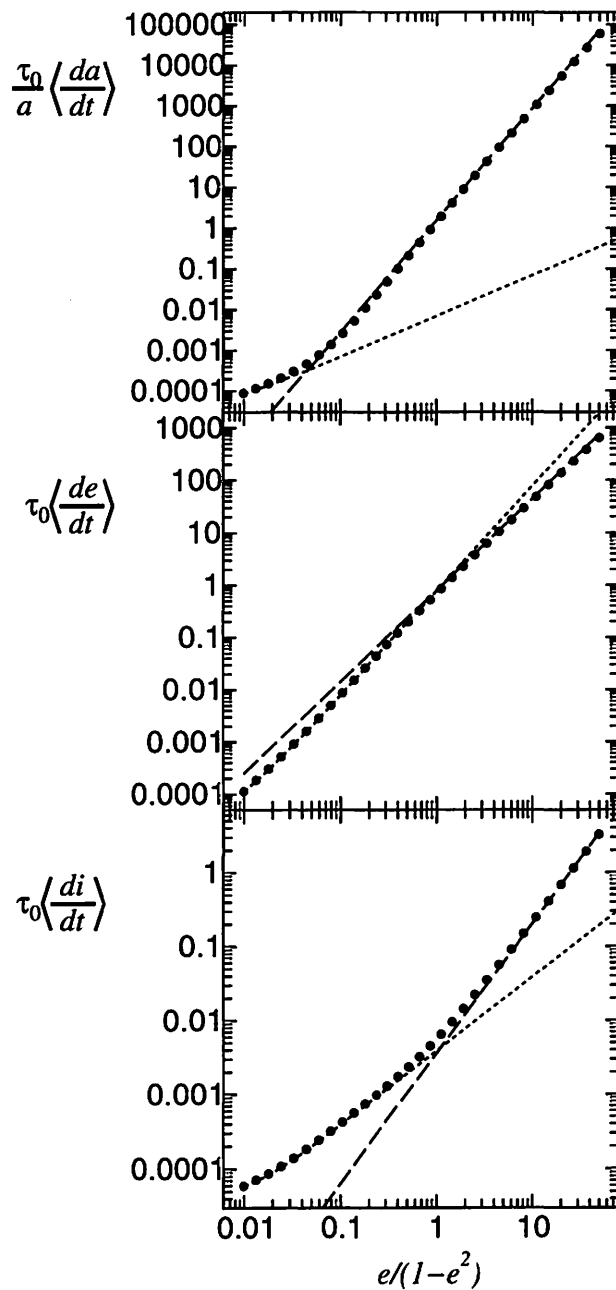


Fig. 8.— The changes in  $a$ ,  $e$ , and  $i$ , as a function of the initial  $e/(1 - e^2)$ , in the case with  $i = 0.01$ . The full circle plots, the dash lines, and the short dash lines show the numerical result, Eqs. (17), and Eqs. (14), respectively.

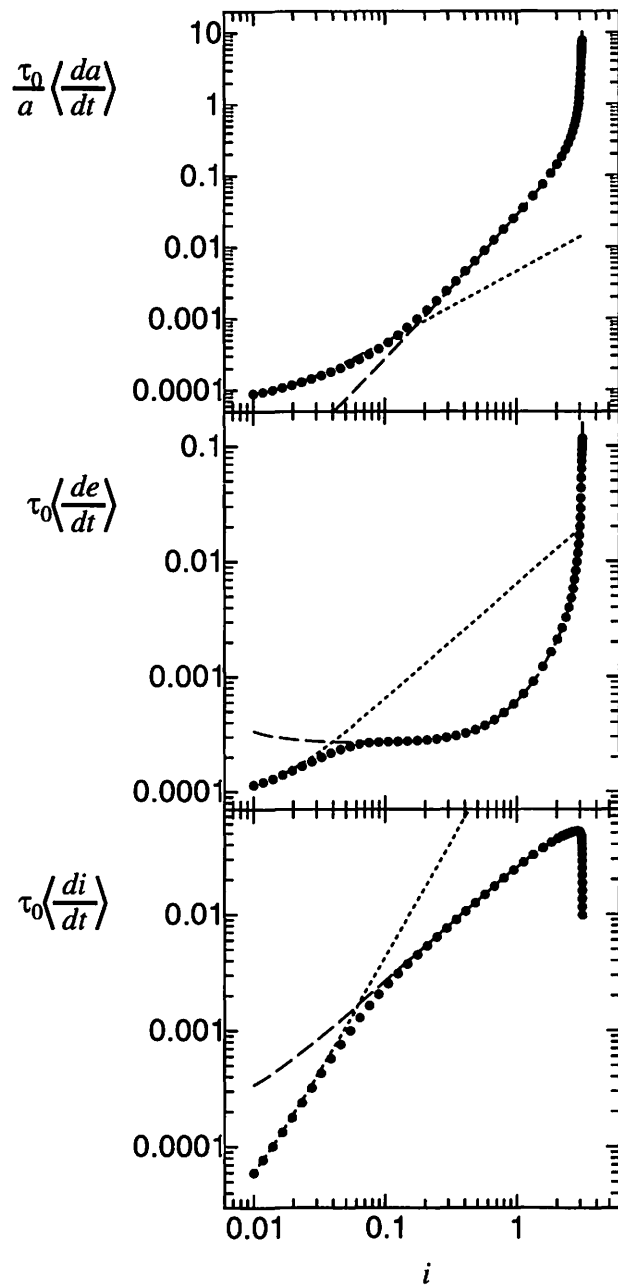


Fig. 9.— The changes in  $a$ ,  $e$ , and  $i$ , as a function of the initial  $i$ , in the case with initial  $e = 0.01$ . The plots and the lines mean the same as Figs. 8.

### 3.5. NEW EQUATION

We derive the analytical equations for the variation of orbit in two cases, and we connect these equations to derive new equations. To derive the equation of the changes in any  $e$  and  $i$ , we modify Eqs. (14), (17), and (25) by the following ways.

$$\begin{cases} \dot{a}_1 = -2\frac{a}{\tau_0} \left[ \left( 0.97 \frac{e}{1-e^2} \right)^2 + (0.64i)^2 + \eta^2 \right]^{1/2} \eta, \\ \dot{e}_1 = -\frac{e}{\tau_0} \left[ \left( 0.77 \frac{e}{1-e^2} \right)^2 + (0.64i)^2 + (1.5\eta)^2 \right]^{1/2}, \\ \dot{i}_1 = -\frac{i}{2\tau_0} \left[ \left( 0.77 \frac{e}{1-e^2} \right)^2 + (0.85i)^2 + \eta^2 \right]^{1/2}, \end{cases} \quad (23)$$

$$\begin{cases} \dot{a}_2 = \frac{a}{\pi\tau_0} \left( \frac{e}{1-e^2} \right)^\alpha \bar{\Psi}, \\ \dot{e}_2 = \frac{1}{2\pi\tau_0} \left( \frac{e}{1-e^2} \right)^{\alpha-1} \bar{\Psi}, \\ \dot{i}_2 = \frac{i}{2\pi\tau_0} \left( \frac{e}{1-e^2} \right)^{\alpha-1} (\bar{\Phi}_1 + \bar{\Phi}_2 \sin^2 \omega). \end{cases} \quad (24)$$

$$\begin{cases} \dot{a}_3 = -\frac{a}{\tau_0} \frac{\sqrt{2}h_0}{\sqrt{\pi}a\eta^4 \sin i} [\tilde{r}(\omega)^{\gamma+1}G(\omega) + \tilde{r}(\omega + \pi)^{\gamma+1}G(\omega + \pi)], \\ \dot{e}_3 = -\frac{1}{\tau_0} \frac{h_0}{\sqrt{2\pi}a\eta^2 \sin i} [\tilde{r}(\omega)^{\gamma+1}H(\omega) + \tilde{r}(\omega + \pi)^{\gamma+1}H(\omega + \pi)], \\ \dot{i}_3 = -\frac{1}{\tau_0} \frac{h_0}{\sqrt{2\pi}a\eta^4} [\tilde{r}(\omega)^{\gamma+1}I(\omega) + \tilde{r}(\omega + \pi)^{\gamma+1}I(\omega + \pi)], \end{cases} \quad (25)$$

Using these expression, we can write in the case with  $i \ll a/h$ ,

$$\left\langle \frac{da}{dt} \right\rangle = \sqrt{(\dot{a}_1)^2 + (\dot{a}_2)^2} \quad (26)$$

$$\left\langle \frac{de}{dt} \right\rangle = \sqrt{(\dot{e}_1)^2 + (\dot{e}_2)^2} \quad (27)$$

$$\left\langle \frac{di}{dt} \right\rangle = \sqrt{(\dot{i}_1)^2 + (\dot{i}_2)^2} \quad (28)$$

In the case with any  $e$ , The orbital change for small and large limit of  $i$  can be expressed by Eqs. (25) and (28). We can write the variation in the case with  $e \ll 1$ ,

$$\left\langle \frac{da}{dt} \right\rangle = \sqrt{(\dot{a}_1)^2 + (\dot{a}_3)^2} \quad (29)$$

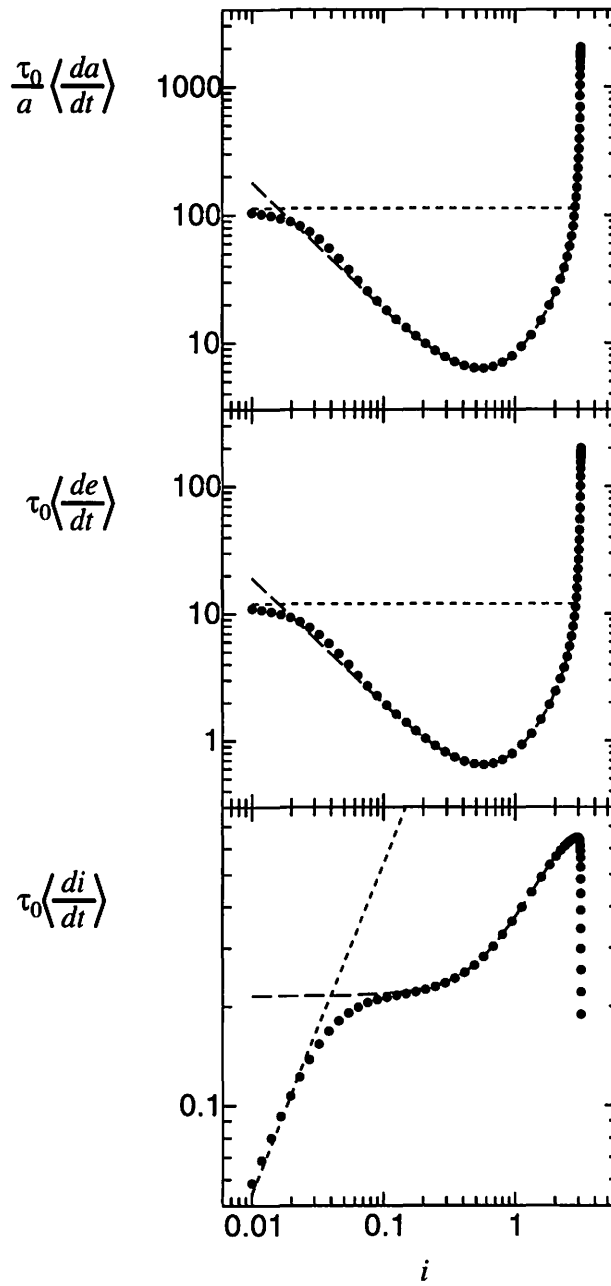


Fig. 10.— The same as Fig. 9 except for initial  $e$  is 0.9. The plots are the same as Figs. 8. The dash lines and short dash lines show Eqs. (25) and Eqs. (17), respectively.

$$\left\langle \frac{de}{dt} \right\rangle = \frac{1}{\sqrt{\frac{1}{(\dot{e}_1)^2} + \frac{1}{(\dot{e}_3)^2}}}, \quad (30)$$

$$\left\langle \frac{di}{dt} \right\rangle = \frac{1}{\sqrt{\frac{1}{(\dot{i}_1)^2} + \frac{1}{(\dot{i}_3)^2}}}, \quad (31)$$

and in the case with  $e \sim 1$ ,

$$\left\langle \frac{da}{dt} \right\rangle = \frac{1}{\sqrt{\frac{1}{(\dot{a}_2)^2} + \frac{1}{(\dot{a}_3)^2}}}, \quad (32)$$

$$\left\langle \frac{de}{dt} \right\rangle = \frac{1}{\sqrt{\frac{1}{(\dot{e}_2)^2} + \frac{1}{(\dot{e}_3)^2}}}, \quad (33)$$

$$\left\langle \frac{di}{dt} \right\rangle = \frac{1}{\sqrt{\frac{1}{(\dot{i}_2)^2} + \frac{1}{(\dot{i}_3)^2}}}, \quad (34)$$

We confirm the validity of the new equations numerically. Figures 8, 9, and 10 show that Equations are compared with the numerical results.

## 4. AFTER THE STELLAR ENCOUNTER

### 4.1. Neptune scattering

If the perihelion distance  $[= a(1 - e)]$  of the planetesimal is smaller than 30 AU that is the semi-major axis of present Neptune, the planetesimal may be ejected by Neptune. The stellar encounter pump up  $e$  of the planetesimals, some planetesimals with  $a > 40$  AU are ejected.

Figures 11 and 12 show the planetesimals' orbit after the stellar encounter with  $D = 120$  AU which is removed one's perihelion distance is smaller than 30 AU. If  $i_*$  is small (in the case of  $30^\circ$ ,  $60^\circ$ ),  $e$  is highly pumped-up. Many bodies are ejected. On the other hand, high inclination ( $i_*$ ) encounter result in small  $i$  and they are not ejected. After ejecting the planetesimals with perihelion distance  $< 30$  AU, We plot  $i$  of the planetesimals with  $a$  between 40 AU to 50 AU as a function with  $e$  of themselves. In the case with  $i_* = 60^\circ$  and  $\omega_* = 90^\circ$ , It similar to EKBO distribution (there are the planetesimals whose  $i$  are higher than their  $e$ ).



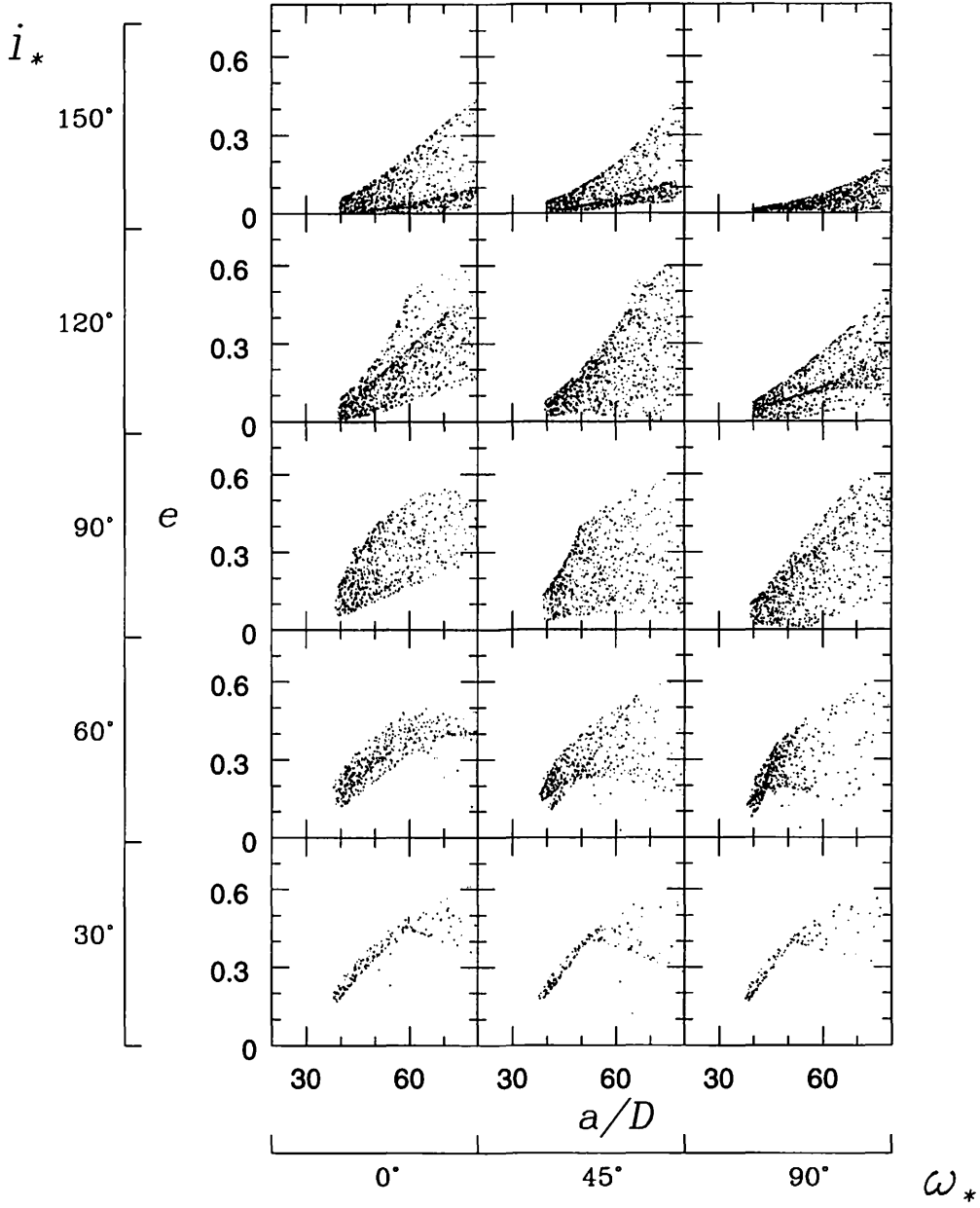


Fig. 11.—  $i$  of the planetsimals with  $a$  between 40 AU and 50 AU and perihelion ( $a(1 - e)$ )  $> 30$  AU as a function of  $e$ , after a stellar encounter with  $D = 120$  AU. We select 1000 bodies with initial  $a$  between 40 AU to 80 AU.

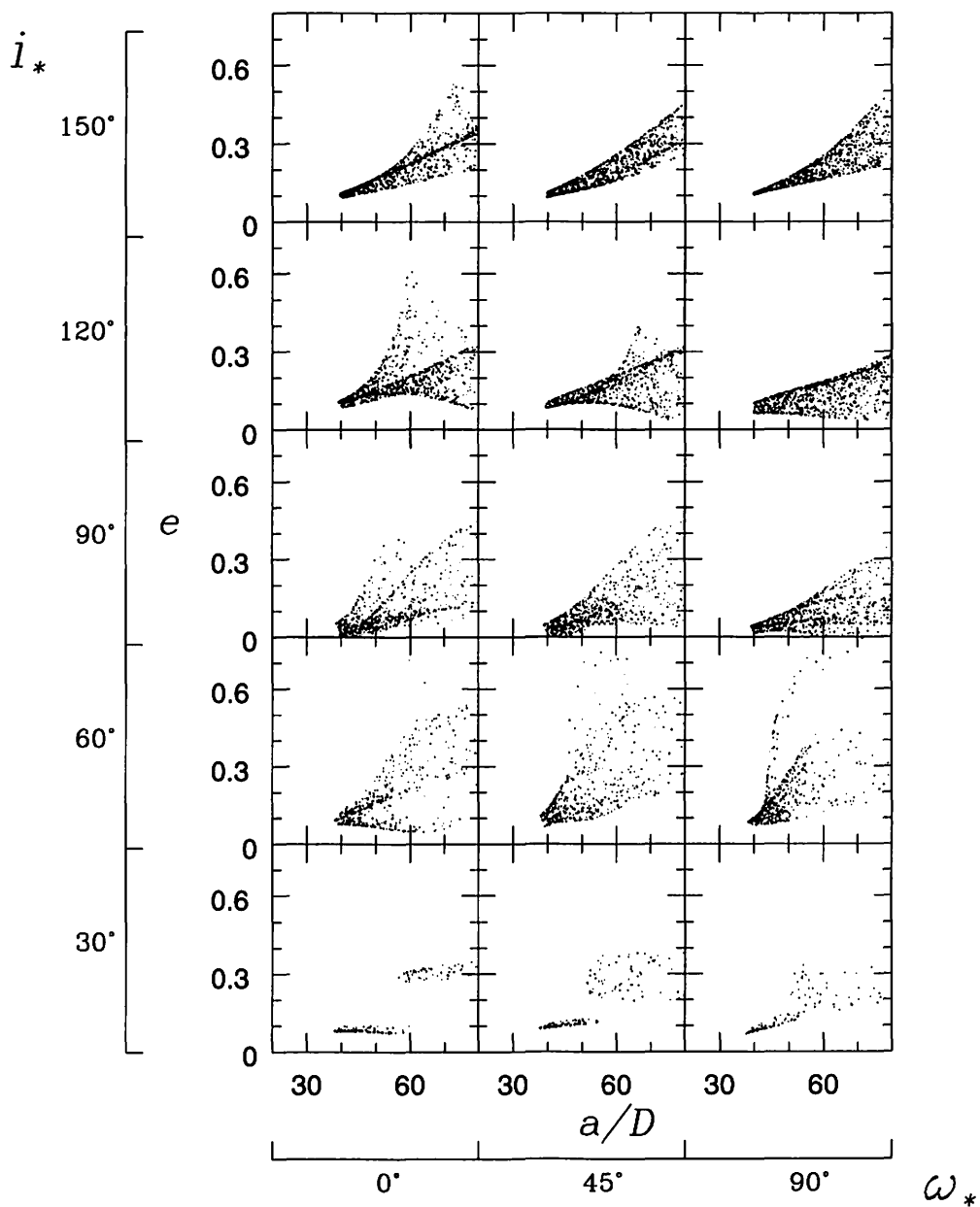


Fig. 12.—  $i$  of the planetsimals with  $a$  between 40 AU and 50 AU as a function of  $e$ , after a stellar encounter with  $D = 120$  AU. We select 1000 bodies with initial  $a$  between 40 AU to 80 AU.

## 4.2. gas drag effect

We calculate the change in  $a$ ,  $e$ , and  $i$  due to the gas drag, using Eqs. 31, 34, 28, and 25 and the data table of the variation for no limit equation. The variation of  $a$  due to gas drag is larger than that of  $e$  and  $i$ , in the case with high  $e$ , because the effect of gas drag tend to changes  $a$  of the body with high  $e$  to its perihelion distance. The gas drag affect on the planetesimals with near perihelion distance. We investigate the effect of gas drag, assuming nebula gas density to decay as  $\exp(t/t_{\text{decay}})$ . Figures 13 and 14 show  $e$  and  $i$  of the planetesimals under the gas drag setting parameters as  $t_{\text{decay}} = 5 \times 10^7$  years and  $\tau = 5 \times 10^6 (m/10^{22}\text{g})^{1/3}$ , as a function of thief  $a$ . The gas drag effect is not only that the planetesimals with high  $e$  and any  $i$  fall into inner region, but also that some planetesimal's  $e$  are dumped and  $i$  remain high. As the result, the bimodal distribution of EKBO can be explained by the stellar encounter with  $D \sim 120$  AU,  $i_* \sim 60^\circ$ , and  $\omega_* \sim 90^\circ$  (see Fig. 15).

## 5. CONCLUSION

We investigate the effects of a stellar encounter on a planetesimal disk and the evidence of the stellar encounter on the distribution of Edgeworth-Kuiper belt objects (EKBO). The classical EKBO are so exited that there are some perturbations in past. The distribution of inclination ( $i$ ) of the classical EKBO are statically two peak in 0.05 and 0.3: The group with large  $i$  have the feature of  $i \gg e$ , and the other  $i \sim e$ .

We considered that a disk of massless particles (planetesimals) orbiting a primary star encounters a passing single star. Encounter parameters are pericenter distance of the encounter ( $D$ ), the argument of perihelion ( $\omega_*$ ), eccentricity ( $e_*$ ) and inclination ( $i_*$ ) of the orbit of the passing star, and the mass ratio ( $M_*$ ) of the passing star's mass to the primary one. We show that a stellar encounter can pump up  $e$  and  $i$  of EKBO. However, in the case of most successive sets of encounter parameters, most of  $e$  of EKBO are as large as  $i$ .

We consider that a stellar encounter occur in early time of planet formation. We also investigate the effect of the gas drag and Neptune scattering, after the stellar encounter. These affect on the bodies with small perihelion distance well, resulting in the bodies falling into inner region or scattering. As the perihelion distance is  $= a(1 - e)$ , the only bodies with high  $e$  are ejected from Edgeworth-Kuiper belt region, and the ones with high  $i$  remain if they have small  $e$ . As the result, We can show that the stellar encounter with some sets of parameters can explain EKBO's bimodal distribution:  $i \gg e$  and  $i \sim e$ .

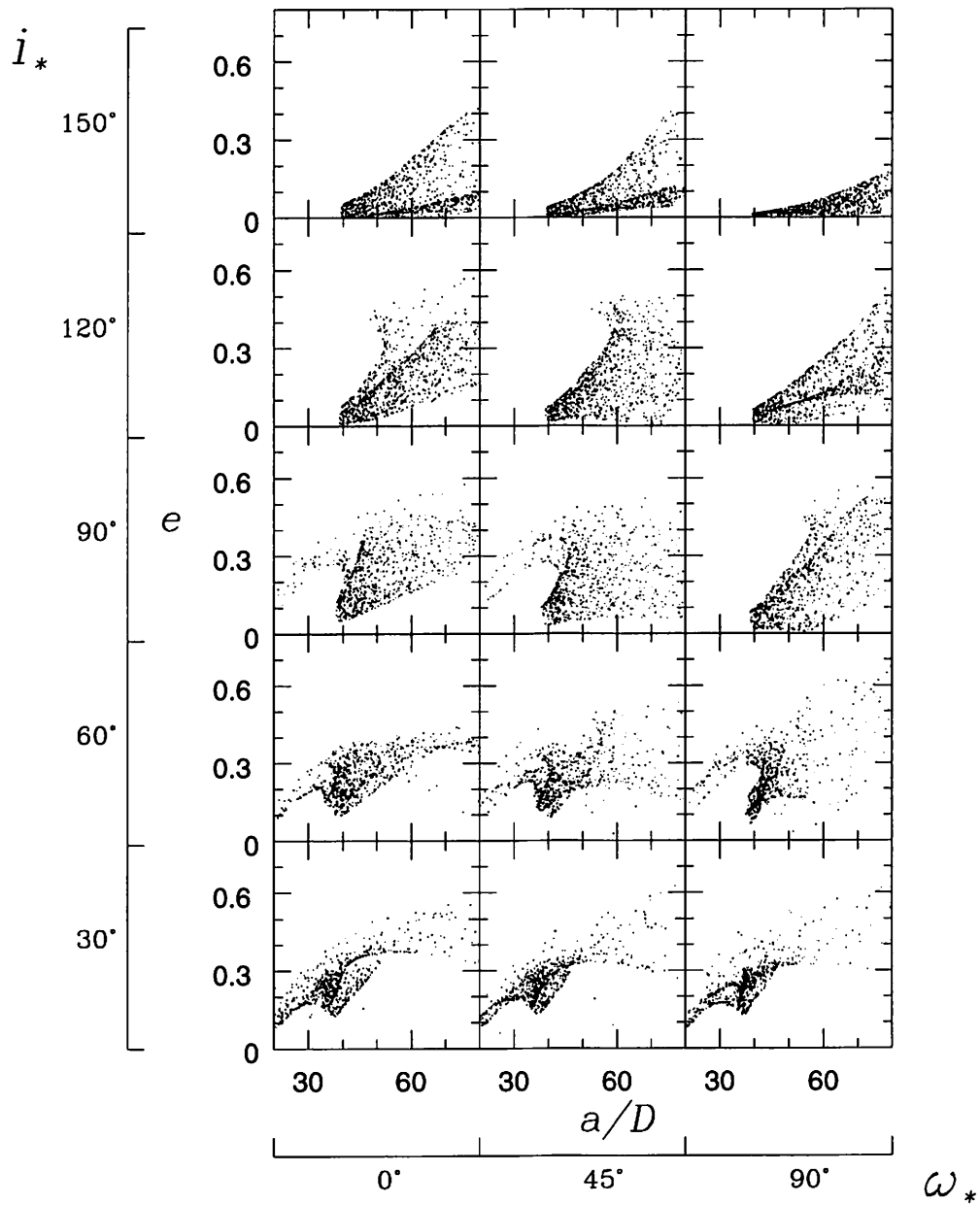


Fig. 13.— The effect of gas drag.  $i$  of the planetsimals with  $a$  between 40 AU and 50 AU AU as a function of  $e$ , after a stellar encounter with  $D = 120$  AU. We select 1000 bodies with initial  $a$  between 40 AU to 80 AU.

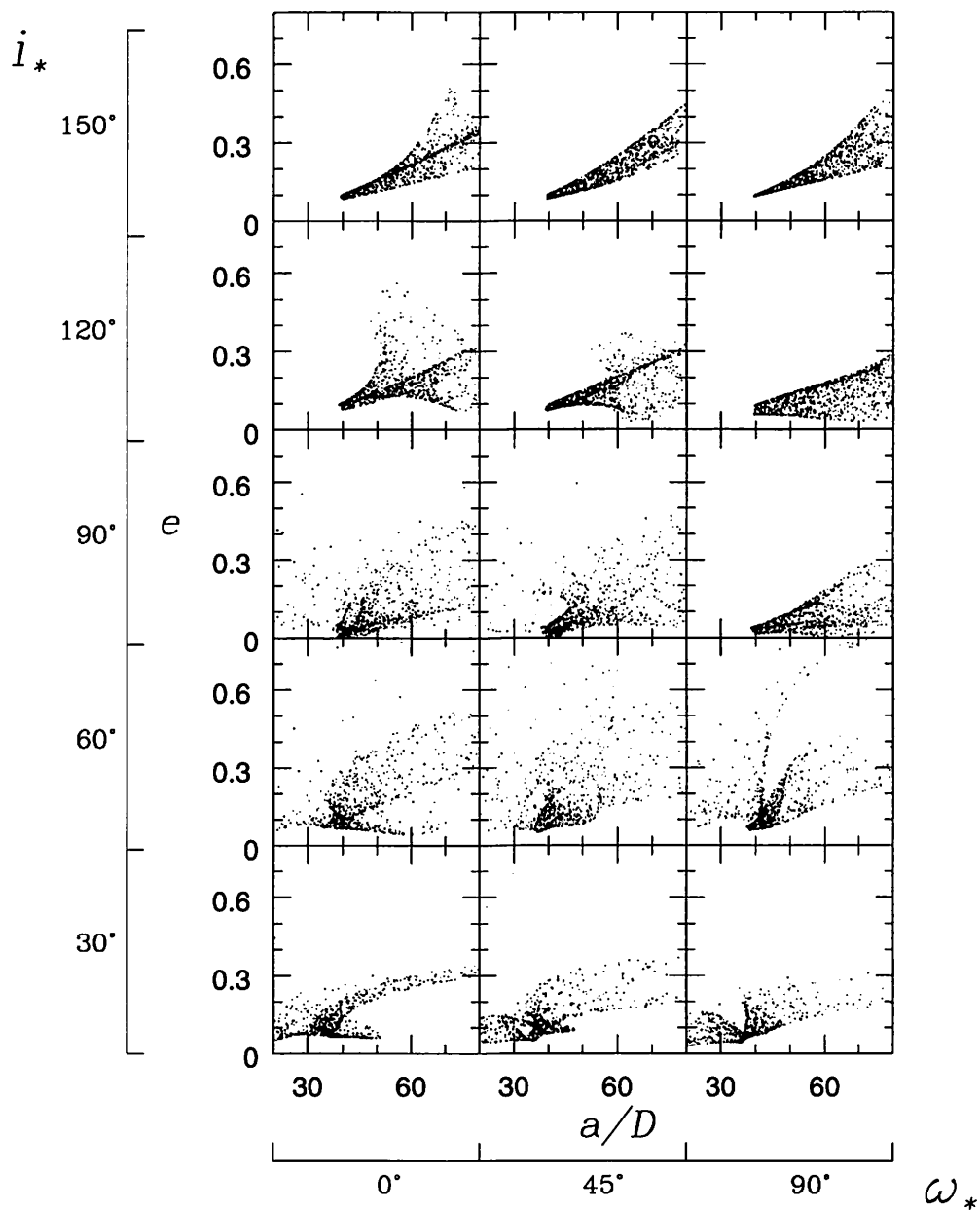


Fig. 14.— The effect of gas drag.  $i$  of the planetsimals with  $a$  between 40 AU and 50 AU as a function of  $a$ , after a stellar encounter with  $D = 120$  AU. We select 1000 bodies with initial  $a$  between 40 AU to 80 AU.

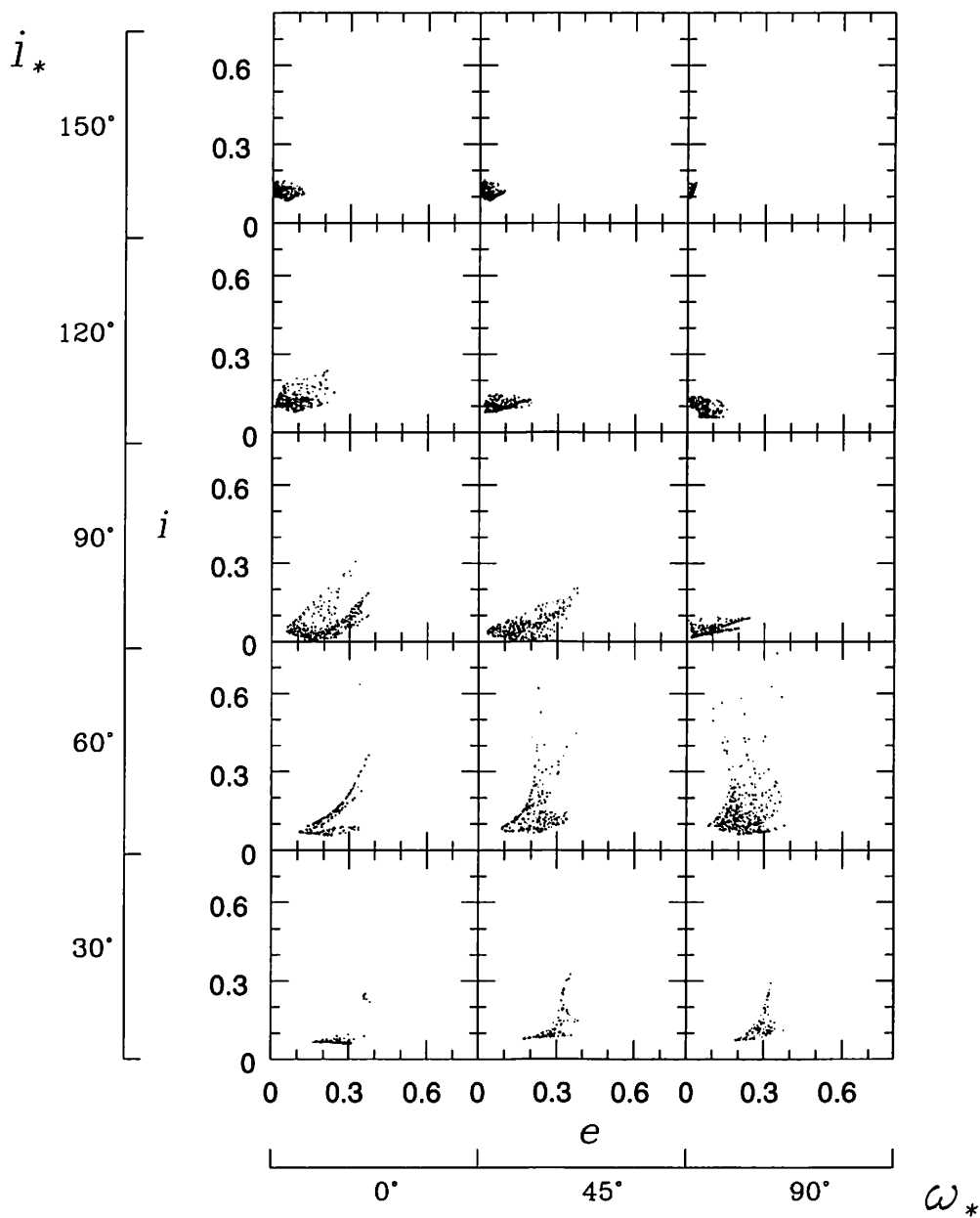


Fig. 15.— The effect of gas drag.  $i$  of the planetsimals with  $a$  between 40 AU and 50 AU as a function of  $e$ , after a stellar encounter with  $D = 120$  AU. We select 1000 bodies with initial  $a$  between 40 AU to 80 AU.

## REFERENCES

- Adachi, I., C. Hayashi, and K. Nakazawa, 1976. The gas drag effect on the elliptical motion of a solid body in the primordial solar nebula. *Prog. Theor. Phys.* **56** 1756-1771.
- Brown, M. 2001. The Inclination Distribution of the Kuiper Belt. *Astrophys. J.* **121** 2804-2814.
- Hayashi, C., K. Nakazawa, and Y. Nakagawa 1985. Formation of the solar system. In *Proto-stars and Planets II* (D. C. Black and M. S. Matthews, Eds.), pp. 1100-1153. Univ. of Arizona Press, Tucson.
- Ida, S., J. Larwood, and A. Burkert 2000. Evidence for Early Stellar Encounters in the Orbital Distribution of Edgeworth-Kuiper Belt Objects. *Astrophys. J.* **528** 351-356
- Inaba, S., H. Tanaka, K. Nakazawa, G. Wetherill, and E. Kokubo 2001. High-Accuracy Statistical Simulation of Planetary Accretion: II. Comparison with N-Body Simulation. *Icarus* **149** 235-250
- Kobayashi, H. and S. Ida 2001. The Effects of a Stellar Encounter on a Planetesimal disk. *Icarus* **153** 416-429
- Malhotra, R. 1995. The Origin of Pluto's Orbit: Implications for the Solar System beyond Neptune. *Astron. J.* **110** 420-429.
- Ngasawa, M. and S. Ida 2000. Sweeping Secular Resonances in the Kuiper Belt Caused by Depletion of the Solar Nebula. *Astron. J.* **120** 3311-3322
- Petit, J., A. Morbidelli, and G. Valsecchi 1999. Large Scattered Planetesimals and the Excitation of the Small Body Belts. *Icarus* **141** 367-387
- Tanaka, H., T. Takeuchi, and W. Ward Three-Dimensional Interaction between a Planet and an Isothermal Gaseous Disk. I. Corotation and Lindblad Torques and Planet Migration. *Astrophys. J.* **565** 1257-1274

# Formation of Low-Mass Multiple Satellites.

Takaaki Takeda

*Theoretical Astrophysics Division, Astrophysics Division, National Astronomical  
Observatory of Japan. E-mail: takedatk@th.nao.ac.jp*

## 1. Introduction

The most favored hypothesis of the origin of the Moon is giant impact hypothesis. It is considered that a Mars-sized protoplanet had collided with the proto-Earth and the Moon is formed from a debris disk, which is splashed by the impact. Using  $N$ -body method, it was shown that a single Moon would be formed from a massive debris disk, which is initially confined within the Roche limit (Ida et al. 1997). The Moon accretion process is as follows. As random velocity of disk particles damps, they begin to form aggregates. Within the Roche limit, tidal effect shears them apart, and spiral pattern develops in the disk. The pattern enhances radial spreading of the disk. Beyond the Roche limit, tips of spiral arms form lunar seeds, and they collide each other and grow. A seed nearest to the Roche limit grows preferentially by the disk particles diffused out. Eventually, other satellite seeds collide to the largest one or are scattered to the Roche limit, and a single Moon remains. The Moon scatters the disk particles to the Earth.

In this work, we performed  $N$ -body simulations of the evolution of debris disks of various mass. We found that a single satellite would be formed only from such a massive disk as the protolunar disk, while multiple satellites would be formed from a less massive disk. In a less massive disk, the satellite forms a gap between the disk and itself. Satellite mass is regulated by the gap formation condition. We found that as the initial disk mass decreases, the satellite mass decreases more rapidly. If the initial disk is less massive, disk-satellite interaction pushes the formed satellite outward rather than the satellite scatters the disk to the central planet. In this case, next satellite is formed near the Roche limit again. We performed  $N$ -body simulations of two-satellite formation cases and semi-analytically derived the condition for the multiple satellite formation. We extrapolate our result to a much less massive disk, and found that multiple satellite system, which is similar to the satellites of outer planets of solar system, may be formed from such a disk.

## 2. Numerical Method

We numerically integrate the orbits of disk particles using Hermite integrator of second-order. Being near the Roche limit, physical radius of particles are comparable to the particles' gravitational radius. Since particles do not approach closely, second-order integrator is enough.

Outside the Roche limit, particles form aggregates. We performed two sets of simulations, adopting different models to express aggregates. In set A, we adopted a rubble

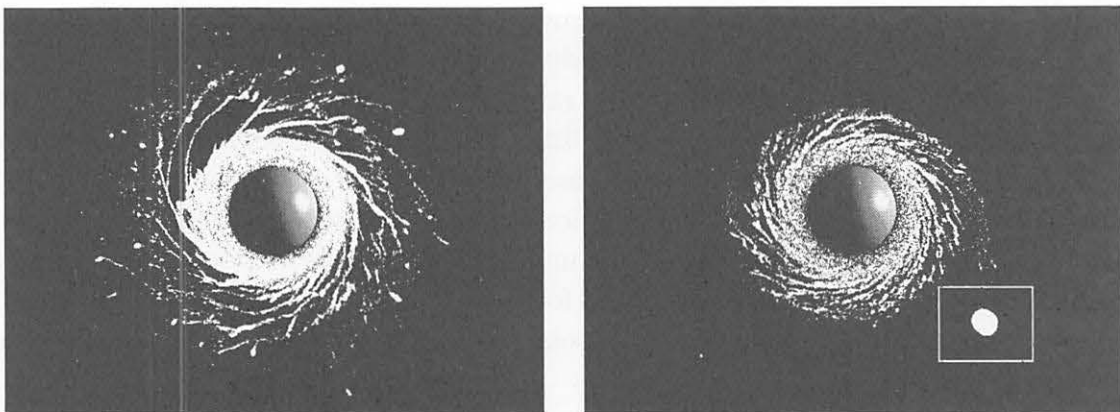


pile model to express the aggregates. When we detect an overlapping of particles, we change the velocity of particles according to restitution coefficient, and simply pushed apart the particles up to the sum of particles' radius. We had found that the simulation result does not depend on restitution coefficient much. We adopted normal restitution coefficient as 0.1. In set A, we performed simulations with various initial disk mass, from 0.01 to  $0.046 M_c$ , where  $M_c$  is mass of central planet, and follows the evolution of the disk until a satellite forms a gap and its growth stops.

After a satellite is formed, it slowly migrates outward. However, rubble pile model takes larger computer power, and it is not good for following longer time evolution. Thus, we perform another set of simulations including artificial accretion process. In set B, we initially put a satellite seed outside the Roche limit, which is smaller than the expected final satellite mass. When we detect a collision between disk particle and satellite seed, we artificially merge them. In set B, a satellite is expressed by one large particle. In this way, we follow orbital evolution of the first formed satellite for longer time, until the first satellite is pushed away enough and the next satellite is formed.

### 3. Results

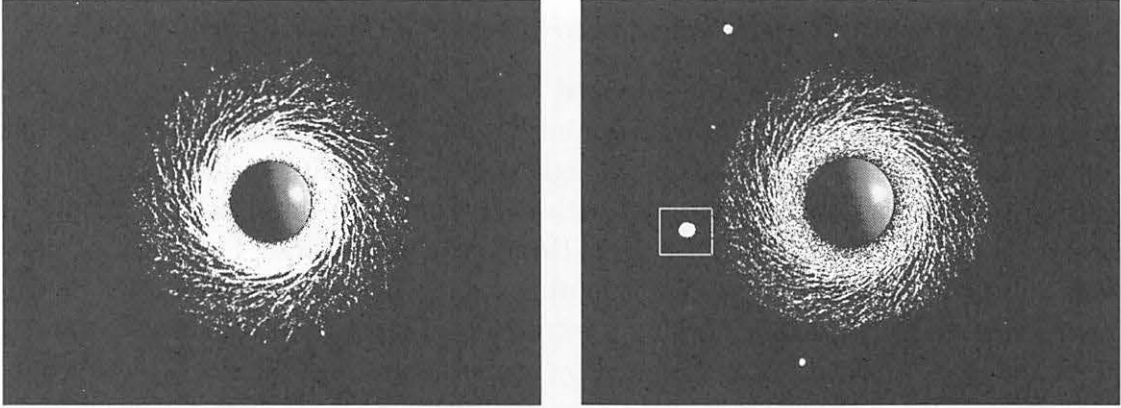
We show typical evolution of the disk. Figures 1 are face-on snapshots of a run of set A with initial disk mass 0.046. Left panel shows snapshot at  $t = 8$ , where the unit of time is Kepler time at Roche limit radius ( $\sim 2.9$  planet radius). Spiral arms develop in the disk, and several satellite seeds are formed beyond the Roche limit. Right panel shows snapshot at  $t = 36$ . Satellite seeds collide each other and form a single large satellite eventually.



Figures 1

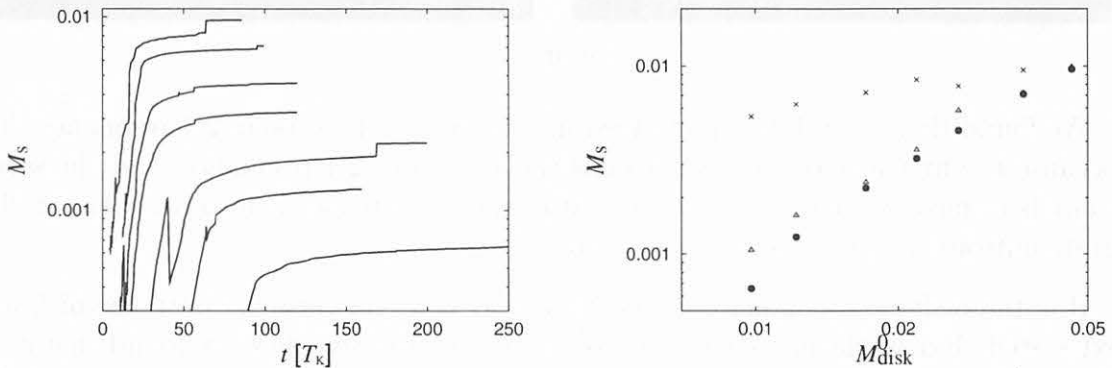
We show the typical evolution of less massive disk. The left panel of Figures 2 is a snapshot at  $t = 20$  of a run with initial disk mass  $0.017 M_c$ . In a case that an initial disk is less massive, scale of spiral structure is small, and time scale of viscous spreading of the disk is much longer. Right panel shows a snapshot at  $t = 106$ . The mass flow from the disk to the satellite is nearly stopped by this time. The largest satellite seed is formed just outside Roche limit. Though several satellite seeds still remain, the second largest aggregate has only 20% mass of the largest one. They would collide to the largest one or

be scattered to the central planet eventually.



Figures 2

We show the growth rate of the largest satellite seed in Figures 3 (left panel). Initially, the satellite grows rapidly. When a gap is formed between the disk and the satellite, mass supply from the disk stops and the growth rate of the satellite becomes much smaller. Slow growth of satellite at this stage mostly comes from particles or small satellite seeds scattered beyond the Roche limit at the early stage of disk evolution. We continued the simulation  $2 T_{\text{stop}}$  at least, where  $t_{\text{stop}}$  is the time at which rapid growth of the satellite stops. By that time, a clear gap is formed between the satellite and the disk and mass flow from the disk to the satellite almost stops. In the right panel, we show the mass of satellite and disk at the end of simulations of set A. Horizontal axis  $M_{\text{disk}}^0$  is initial disk mass. Filled circles are mass of the largest aggregate, and triangles are all mass outside the gap, including the mass of the satellite. The particles outside the gap would collide to the satellite or be scattered eventually. Final satellite mass would be between a filled circle and a triangle. Sometimes, the second largest satellite seed with a considerable mass remains in a horse shoe orbit of the largest one. Crossed points represent a disk mass inside the gap at the same time. As the initial disk mass decreases, satellite mass decreases more rapidly, and mass ratio between the remaining disk and the satellite increases.



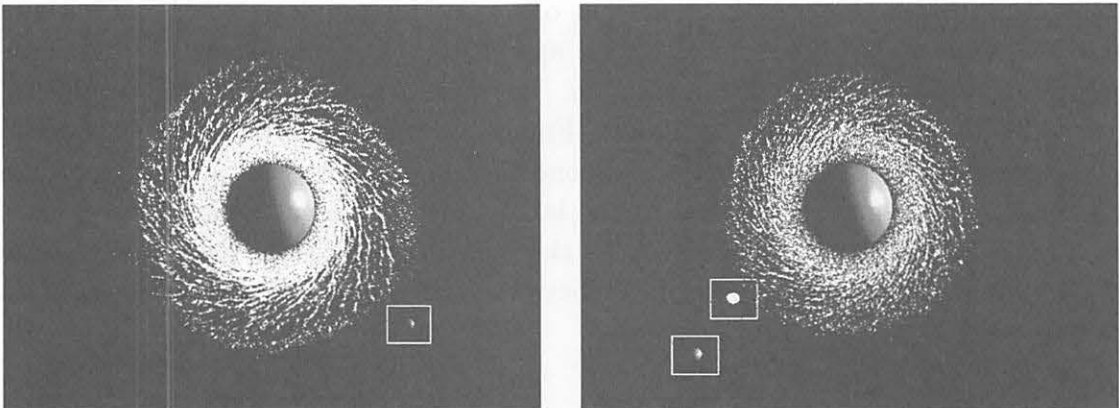
Figures 3

This tendency comes from that the angular momentum transfer rate within a debris disk is proportional to  $\Sigma^3$ , where  $\Sigma$  is surface density (Takeda and Ida 2001, daisaka et al.

2001), while the disk-satellite interaction is proportional to surface density and satellite mass (e.g. Goldreich and Tremaine 1979). Thus, the satellite mass required to form a gap is  $\propto \Sigma^2$ , so that disk-satellite mass ratio increases with decreasing disk mass.

Since angular momentum is transferred from the disk to the satellite, the satellite migrates outward, while the disk mass diminishes as the disk particles fall to the planet in compensation. If disk mass remains enough, a second satellite formation would occur. Since disk-satellite mass ratio increases with decreasing disk mass, second satellite formation would occur in a less massive disk. However, it takes larger cpu time to simulate all the way of the satellite migration adopting rubble pile model. Thus, we introduce another model (set B). In this model, we introduce merger of the satellite seed and disk particles, so that a satellite is represented by one large particle, and its orbital evolution is much easier. We should put the seed near the Roche limit. However, if we put it too close to the Roche radius, it enters the Roche limit. So we put it at about  $1.2R_R$  initially. We have confirmed that the mass of the satellite formed in runs of set B are consist with that of set A, as long as the initial location of satellite seed is well chosen.

We see the snapshots of a run of set B with initial disk mass 0.017. Initially, the satellite seed grows as the disk radially spreads. As the satellite mass increases, it begins to form a gap between itself and the disk. The left side panel of Figures 4 shows a snapshot at  $t = 50$ . After that, the satellite migrates outward by the disk-satellite interaction. When the satellite migrates enough, a second satellite begins to form outside the Roche limit.



Figures 4

We found that a satellite pair formed in this way tends to be in 2:1 resonance. This is because the first formed satellite widens the gap up to 2:1 resonance, and the second satellite is formed at the edge of the gap. After the satellite pair is formed, both satellites migrate outward keeping the resonant state.

If initial disk mass is much smaller, this process would continue until the disk mass is exhausted. Extrapolating our results to a much less massive disk, we found that a disk with  $\sim 0.003M_c$  may form 4 to 5 satellites of mass with 0.001 to 0.0001 $M_c$ . This satellite system resembles satellite systems of outer planets of the solar system. As a future work, we will investigate the evolution of further less massive disks. Also, the stability of the multiple satellite system formed in this way must be investigated, including the effect of tide from the central planet.

# Orbital Stability of a Protoplanet System in the Nebular Gas: Dependence on the masses of Protoplanets

K. IWASAKI, H. EMORI, H. TANAKA, AND K. NAKAZAWA

*Department of Earth and Planetary Sciences, Faculty of Science  
Tokyo Institute of Technology, Meguro-ku, Tokyo 152-8551, Japan*

## Abstract

We investigated the orbital instability of a protoplanet system, in which five protoplanets whose mass is  $M_p$  are distributed with an equal separation distance,  $\Delta\tilde{a}_0$ . In this paper, the cases where  $M_p$  is  $10^{-9}M_\odot$  and where  $M_p$  is  $10^{-5}M_\odot$  were studied and, combining our results with the works on the cases where  $M_p$  is  $10^{-7}M_\odot$ , which were already done in the past, we confirmed that the following properties of an orbital instability are valid, regardless of  $M_p$ .

- (1) The logarithm of the orbital instability time,  $T_{\text{inst}}$ , of a system without a gas disk increases in proportion to  $\Delta\tilde{a}_0$ .
- (2) In the presence of the nebular gas, the instability time of a system becomes extremely large compared with  $T_{\text{inst}}$  by the effect of the drag force due to the nebular gas and the system substantially doesn't experience an orbital instability, when  $\Delta\tilde{a}_0$  is larger than a critical separation distance,  $(\Delta\tilde{a}_0)_{\text{crit}}$ .
- (3) The value of  $(\Delta\tilde{a}_0)_{\text{crit}}$  becomes large with a decrease in the surface density of the nebular gas.

Furthermore, we obtained a semi-analytical expression for  $(\Delta\tilde{a}_0)_{\text{crit}}$ , partly using our simulation results. Finally, we applied our results to the orbital stability of a protoplanet system in the Jovian planet region ( $a = 5\text{AU}$ ,  $M_p = 10^{-5}M_\odot$ ). In the presence of the minimum mass solar nebula,  $(\Delta\tilde{a}_0)_{\text{crit}}$  is estimated to be 4.3, which is smaller than a typical separation distance of a realistic protoplanet system ( $\simeq 10$ ). Thus, in the presence of the minimum mass solar nebula, the orbital instability of a protoplanet system never occurs in the Jovian planet region.

## 1 Introduction

In the standard scenario of planetary formation, terrestrial planets and solid cores of Jovian planets are formed through accretion of planetesimals whose initial sizes are 10 – 100km. (Safronov 1972, Greenberg et al. 1978, Hayashi et al. 1985). The process of planetary formation from planetesimals is divided into plural stages. In the beginning of planetesimal growth, planetesimals grow uniformly via mutual collisions (orderly growth). As planetesimals become massive, mutual gravity between planetesimals becomes effective and have an important role on the growth mode of planetesimals. Coupling effect of dynamical friction caused by mutual gravity (i.e., energy equipartition) and gravitational focusing (i.e., increase in the collisional cross section) leads to the preferential growth of relatively massive planetesimals. This growth mode is called “runaway growth” (Wetherill and Stewart 1989, Kokubo and Ida 1996).

Next, a small number of massive planetesimals (i.e., protoplanets) formed in the runaway growth stage begin to pump up the random velocities (i.e., orbital eccentricities and inclinations) of surrounding planetesimals, and, as a result, reduce the speed of their own growth. This feedback effect regulates the masses of protoplanets to be nearly equal with each other. This growth mode is

called “oligarchic growth” (Kokubo and Ida 1998). The final mass of a protoplanet,  $M_p$ , which would be formed through oligarchic growth, are dependent on the radial distance from the sun. In the terrestrial planet region ( $a < 2.7\text{AU}$ ),

$$M_p \simeq 0.1 \left( \frac{\sigma_d}{\sigma_d^H} \right)^{\frac{3}{2}} \left( \frac{a}{1\text{AU}} \right)^3 M_\oplus, \quad (1)$$

and, in the Jovian planet region ( $a > 2.7\text{AU}$ ),

$$M_p \simeq 2.8 \left( \frac{\sigma_d}{\sigma_d^H} \right)^{\frac{3}{2}} \left( \frac{a}{5\text{AU}} \right)^3 M_\oplus, \quad (2)$$

where  $a$  is the semi-major axis of a protoplanet and  $\sigma_d$  and  $\sigma_d^H$  are the initial surface density of solid materials and that in the minimum mass solar nebula, respectively (Kokubo and Ida 2000). In the above,  $M_\oplus$  represents the Earth mass. Here, note that the mass of a protoplanet in the Jovian planet region attains to about 10 times as large as that in the terrestrial planet region. These protoplanets are displaced with nearly equal separation distance ( $\simeq 10$ ), when the separation distances are scaled by the Hill radius. The Hill radius,  $r_H$ , is defined as follows:

$$r_H = \left( \frac{2M_p}{3M_\odot} \right)^{\frac{1}{3}} \left( \frac{a_{10} + a_{20}}{2} \right), \quad (3)$$

where  $M_{\odot}$  is the solar mass and  $a_{10}$  and  $a_{20}$  are the semi-major axes of two adjacent protoplanets, respectively. The random velocities of protoplanets (i.e., eccentricities and inclinations) are very small just after the formation, since the protoplanets are formed, suffering from dynamical friction with the surrounding planetesimals. In other words, the protoplanets become isolated.

In the terrestrial planet region, after the oligarchic growth, the isolated protoplanets increase their orbital eccentricities owing to mutual gravity and experience orbital crossings (i.e., orbital instability) (Chambers et al. 1996, Yoshinaga et al. 1999, Ito and Tanikawa 1999, 2001). Orbital crossings cause protoplanets to collide with each other and start to grow to the present planets again. Recent works on the formation process of terrestrial planets from protoplanets imply that, in order to form terrestrial planets with small eccentricities such as present Earth and Venus, the nebular gas is needed even after the formation of planets (Chambers and Wetherill 1998, Kominami and Ida 2001). However, in the presence of the nebular disk, a protoplanet system is expected to be prevented from undergoing an orbital instability, since the nebular gas has an effect of suppressing eccentricities of protoplanets through gravitational and hydrodynamical interaction (Adachi et al. 1976, Ward 1988, Artymowicz 1993). Thus, it is important to investigate the orbital stability of a protoplanet system in the nebular disk in detail.

Orbital behaviors of a protoplanet system without the gas disk were already studied by Chambers et al. (1996). They considered a system composed of 5 to 20 protoplanets with mass of  $10^{-7}M_{\odot}$  (about one third of Martian mass). Initial orbits of protoplanets are distributed with equal orbital separation in circular and coplanar orbits. Through long term orbital calculations, they discovered that the orbital instability time of a protoplanet system becomes large exponentially with an increase in the separation distance of protoplanets. Furthermore, Iwasaki et al. (2001, 2002) investigated the orbital instability of a protoplanet system set up by Chambers et al., including the drag forces caused by the hydrodynamical interaction (i.e., gas-drag force) and the tidal interaction (i.e., gravitational interaction) with the nebular gas. According to their results, regardless of the kind of the interactions, the onset of an orbital instability is suppressed by the drag force, when the separation distance is larger than a critical separation distance. They also showed that a critical separation distance increases as the surface density of the gas disk decreases, i.e., the gas disk dissipates. Especially, Iwasaki et al. (2002) elucidated that the surface density of the gas disk must decrease to about 0.1% of that of the minimum mass solar nebula for the onset of an orbital instability of a protoplanet system with a typical separation distance ( $\simeq 10$  Hill radius), when the gravitational interaction with the gas disk is taken into account.

On the other hand, in the Jovian planet region, protoplanets begin to capture the surrounding nebular gas gravitationally and become gaseous giant planets such as Jupiter and Saturn (Mizuno 1980, Bodenheimer and

Pollack 1986, Ikoma et al. 2000). The formation time of these gaseous planets depends strongly on a solid core mass (i.e., a protoplanet's mass). Especially, Ikoma et al. (2000) pointed out that, in order to form massive gaseous envelopes of the present Jovian planets ( $\simeq 100$  to  $1000 M_{\oplus}$ ) around protoplanets within the nebular life time ( $\simeq 10^7$ – $10^8$  year), a protoplanet's mass which would be formed through oligarchic growth must exceed a critical mass ( $\simeq 5$  to  $10 M_{\oplus}$ ). However, from equation (2), we find that a typical mass of a protoplanet formed in the Jovian planet region when  $\sigma_d = \sigma_d^H$  is smaller than the above-mentioned critical mass. This means that a protoplanet formed in the minimum mass solar nebula can not become the present gaseous giant planets alone.

There are two possible ways to resolve this difficulty in the formation of Jovian planets. One is to enhance initial surface density of solid materials,  $\sigma_d$ , from that of the minimum mass solar nebula,  $\sigma_d^H$ . The other is to consider the collision and accretion between protoplanets. Here, we must again discuss about the orbital stability of a protoplanet system in the nebula gas, if we think about the latter possibility (i.e., collisional process between protoplanets). However, the above-mentioned works on the orbital stability of a protoplanet system in the terrestrial planet region (Chambers et al. 1996, Iwasaki et al. 2001, 2002) cannot be applied straightly to this case. This is because their works are limited to the case where the masses of protoplanets are  $10^{-7}M_{\odot}$  which is about  $0.03M_{\oplus}$ .

In this study, we investigated the orbital stability of a protoplanet system in the nebula gas through orbital calculations, changing the masses of protoplanets extensively. The tidal (gravitational) interaction between a protoplanet and the nebular gas is taken into account as a drag force proportional to the random velocity of a protoplanet (Iwasaki et al. 2002). A protoplanet system is set up in the same way in Chambers et al. (1996), i.e., protoplanets with same mass are distributed with equal orbital separation and their initial orbital eccentricities and inclinations are set to be zero.

## 2 Model of Orbital Calculations

We investigated the orbital behaviors of  $n$  protoplanets (in this study,  $n = 5$ ) revolving around the central star in the nebular gas through numerical calculations. The forces acting on a protoplanet in our numerical simulations are the gravitational force from the central star, the mutual gravity between protoplanets, and, the tidal (gravitational) interaction with the nebular gas. The tidal interaction between a protoplanet and the nebular gas is included as a drag force proportional to the random velocity,  $u$ , of a protoplanet. Here, the random velocity,  $u$ , represents a velocity of a protoplanet in coordinates rotating with a Keplerian circular velocity in the semi-major axis of the protoplanet,  $v_K$ , i.e.,

$$u = v_p - v_K, \quad (4)$$

where  $v_p$  is a velocity of the protoplanet in an inertial space. Thus, the drag force caused by the gas disk,  $F_{\text{grav}}$

is expressed by

$$F_{\text{grav}} = -\frac{1}{\tau_D} u, \quad (5)$$

where  $\tau_D$  is a numerical constant.

The value of  $\tau_D$  is obtained by calculating the tidal torque, which is acting on the density waves excited in the nebular gas (Ward 1988, Artymowicz 1993), and its explicit form is (Artymowicz 1993)

$$\tau_D = 8 \times 10^9 \left( \frac{M_p}{1 \times 10^{-7} M_\odot} \right)^{-1} \left( \frac{c_s}{v_K} \right)^4 \times \left( \frac{\sigma_g}{1.7 \times 10^3 \text{ g cm}^{-2}} \right)^{-1} \left( \frac{a}{1 \text{ AU}} \right)^{-\frac{1}{2}} T_K, \quad (6)$$

where  $\sigma_g$  and  $c_s$  represent a gas surface density at 1AU and a sound velocity at the semi-major axis considered, respectively. Furthermore,  $T_K$  denotes a Keplerian period. In the above, we assume that the gas surface density decrease as  $a^{-\frac{3}{2}}$  with a increase in  $a$ . Here, we must mind that  $\tau_D$  depends not only on the physical quantities of the gas disk (i.e.,  $c_s$  and  $\sigma_g$ ) but also on the mass of a protoplanet,  $M_p$ . Generally, the eccentricity of a body suffering from the drag force obeying equation (5) is depressed. The time variation of the eccentricity,  $e$ , of the body in that case is given by (Adachi et al. 1976)

$$\frac{1}{e} \frac{de}{dt} = -\frac{1}{\tau_D}. \quad (7)$$

From the above equation, we understand that  $\tau_D$  is a characteristic damping time of the eccentricity by the drag force.

In order to study the orbital evolution of five protoplanets under the drag force obeying equation (5), we performed orbital calculations for various initial conditions. However, the initial setting of a protoplanet system is the same as that adopted in Chambers et al. (1996) and Iwasaki et al. (2001,2002), i.e.,

- (a) The masses of five protoplanets are equal in each simulation.
- (b) The five protoplanets are placed with equal separation distance scaled by the Hill radius (see equation (3)),  $\Delta \tilde{a}_0$ , from the same semi-major axis, i.e., 1AU.
- (c) The orbital eccentricities and inclinations are set to be zero.
- (d) The position angles of protoplanets are selected at random, under the constraint that the angle between two adjacent protoplanets is set to be larger than 20 degree.
- (e) The value of  $\tau_D$  (i.e., the intensity of the drag force) is not changed with time through each simulation.

We investigated the cases where  $M_p$  is  $10^{-9} M_\odot$  and where  $M_p$  is  $10^{-5} M_\odot$ , in order to see how the property of an orbital stability changes with respect to the masses of protoplanets (the cases where  $M_p$  is  $10^{-7} M_\odot$  were already studied by Chambers et al. (1996) and Iwasaki et

al. (2001,2002)). For each  $M_p$ , the following four cases with different  $\tau_D$  are considered:

$$\tau_D = 1.0 \times 10^3, 3.0 \times 10^3, 9.0 \times 10^3, \text{ and } 2.7 \times 10^4 T_K. \quad (8)$$

We also calculated the cases with different separation distance,  $\Delta \tilde{a}_0$  for each  $M_p$  and  $\tau_D$ . Thus, the separation distance,  $\Delta \tilde{a}_0$ , is changed by a step of 0.2 from 3.6 to 8.8, i.e.,

$$\Delta \tilde{a}_0 = 3.6, 3.8, 4.0, \dots, 8.6, \text{ and } 8.8. \quad (9)$$

Furthermore, for each separation distance,  $\Delta \tilde{a}_0$ , ten cases with different position angles are calculated.

Integration scheme used in this study is a 4th-order  $P(EC)^n$  Hermite scheme (Makino and Aarseth 1992, Kokubo et al. 1998, Kokubo and Makino 1998).

## 3 Results

### 3.1 Orbital Instability of a Gas-free System

Before considering the effect of the drag force caused by a gas disk, we must know the property of the orbital instability of a protoplanet system without a gas disk. As mentioned in the previous sections, the orbital stability of a protoplanet system in a gas-free condition where  $M_p$  is  $10^{-7} M_\odot$  was already studied by Chambers et al. (1996) and their results were re-confirmed by Iwasaki et al. (2001). Through long-term orbital calculations, they showed that a protoplanet system always experiences an orbital instability within  $1 \times 10^7$  year, as long as the normalized separation distance,  $\Delta \tilde{a}_0$  ranges from 3.6 to 8.8. Here, the onset of an orbital instability means the first approach of any two protoplanets within one Hill radius. The logarithm of the time of the onset of an orbital instability,  $T_{\text{inst}}$ , is expressed as a function of  $\Delta \tilde{a}_0$ , i.e., (Chambers et al. 1996)

$$\log_{10}(T_{\text{inst}}/T_K) = b \Delta \tilde{a}_0 + c, \quad (10)$$

where  $b$  and  $c$  are numerical constants. The values of  $b$  and  $c$  are almost constant independently of the number of planets,  $n$ , when  $n$  is equal to or larger than 5 and, for  $n = 5$ , (Iwasaki et al. 2001)

$$b = 0.777 \pm 0.11, \quad (11)$$

and

$$c = -0.154 \pm 0.064. \quad (12)$$

In the present study, we calculated the cases where  $M_p$  is  $10^{-9} M_\odot$  and where  $M_p$  is  $10^{-5} M_\odot$ . Figure 1 shows the relation between the logarithm of  $T_{\text{inst}}$  and  $\Delta \tilde{a}_0$ , when  $M_p$  is  $10^{-5} M_\odot$  and  $n$  is 5. In this figure, ten cases with different initial phase in the position angles are plotted for each  $\Delta \tilde{a}_0$ . The logarithm of  $T_{\text{inst}}$  becomes large with a increase in  $\Delta \tilde{a}_0$  in the same way in the cases where  $M_p$  is  $10^{-7} M_\odot$ . For the cases where  $M_p$  is  $10^{-9} M_\odot$ , such exponential increase of the orbital instability time against  $\Delta \tilde{a}_0$  is also observed. However, the



degree of increasing is different depending on the masses of protoplanets. In figure 1, the least-squares fit given by equation (10) for the present case (the solid line) are also shown, together with the same fitting line for the case where  $M_p$  is  $10^{-7}M_\odot$  (the dotted line). We can easily see that the inclination (i.e., the value of  $b$  in equation (10)) of the fitting line for the present cases is larger than that for the cases where  $M_p$  is  $10^{-7}M_\odot$ . Furthermore, table 1 shows the values of  $b$  and  $c$  in equation (10) obtained by the least-squares fitting of computational runs for the cases where  $M_p$  is  $10^{-9}M_\odot$ , where  $M_p$  is  $10^{-7}M_\odot$ , and where  $M_p$  is  $10^{-5}M_\odot$ . Clearly, the value of  $b$  increases with  $M_p$ . This was already pointed out in Chambers et al. (1996) for a system composed of three protoplanets.

Table 1: The values of  $b$  and  $c$  in equation (10) for a system composed of five protoplanets.

$M_p$	$b$	$c$
$10^{-9}$	$0.639 \pm 0.012$	$0.713 \pm 0.070$
$10^{-7}$	$0.777 \pm 0.011$	$-0.154 \pm 0.064$
$10^{-5}$	$1.06 \pm 0.019$	$-1.68 \pm 0.097$

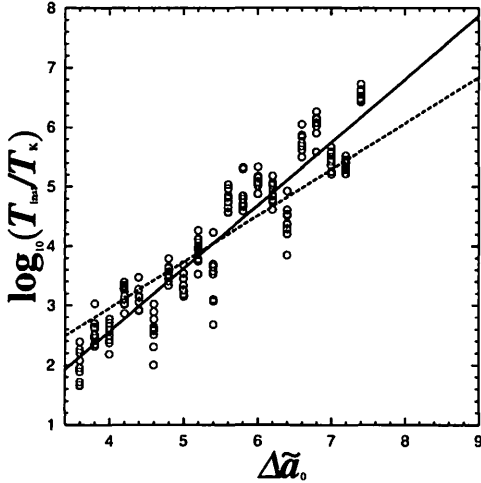


Figure 1: Orbital instability time of a protoplanet system ( $n = 5$ ) without a gas disk,  $T_{\text{inst}}$ , against  $\Delta\tilde{a}_0$ . The mass of a protoplanet is  $10^{-5}M_\odot$ . Ten cases with different initial angle positions are plotted for each  $\Delta\tilde{a}_0$ . The solid line is obtained by a least-squares fitting of all of the plots (see equation (10)). The dotted line shows also the same fitting line for the cases where  $M_p$  is  $10^{-7}M_\odot$  (see equations (11),(12))

### 3.2 Orbital Stability in the Gas Nebula

In this subsection, we investigate the orbital stability of a protoplanet system considered in the previous subsection under the drag force obeying equation (5). The orbital instability time under the influence of the

drag force,  $T_{\text{inst}}^{\text{df}}$ , is expected to become larger than the instability time in a gas-free condition,  $T_{\text{inst}}$ , since the drag force has an effect of suppressing the eccentricities of protoplanets (see equation (7)). Thus, we stopped the orbital calculation, regarding the computational run as a case where the onset of an orbital instability is prevented by the drag force, when the evolutionary time exceeds the cut-off time, which is defined by (Iwasaki et al. 2001, 2002)

$$T_{\text{stop}} = 200 \times T_{\text{inst}}. \quad (13)$$

where  $T_{\text{inst}}$  is the instability time of a system without a gas disk and given by equation (10).

In figure 2, the logarithm of  $T_{\text{inst}}^{\text{df}}$  are plotted against the separation distance,  $\Delta\tilde{a}_0$  for the cases where  $M_p$  is  $10^{-5}M_\odot$  and  $\tau_D$  is  $9.0 \times 10^3 T_K$ . Figure 2 shows that the instability time under the drag force,  $T_{\text{inst}}^{\text{df}}$ , becomes larger than and separates from the instability time of a system without a gas disk,  $T_{\text{inst}}$  (solid line), in other words, a system becomes more stabilized, as  $\Delta\tilde{a}_0$  increases. Thus, the abscissa of figure 2, i.e.,  $\Delta\tilde{a}_0$  can be divided into the three zones, according to the degree of stabilization, i.e., the unstable zone ( $\Delta\tilde{a}_0 < 5.2$ ), the transition zone ( $5.2 \leq \Delta\tilde{a}_0 < 5.6$ ), and the stable zone ( $5.6 \leq \Delta\tilde{a}_0$ ) (Iwasaki et al. 2001, 2002). In the unstable zone, all the ten cases for each  $\Delta\tilde{a}_0$  undergo orbital instabilities within the time nearly equal to  $T_{\text{inst}}$ . In the transition zone, only some cases of the ten cases experience orbital instabilities. In the stable zone, all ten cases for each  $\Delta\tilde{a}_0$  reach the cut-off time,  $T_{\text{stop}}$ , without the onset of an orbital instability. Therefore, the transition zone represents a partition which divides the range of  $\Delta\tilde{a}_0$  into the unstable zone, where a protoplanet system always undergoes an orbital instability like a system without a gas disk does, and the stable zone, where a protoplanet system never experiences an orbital instability.

The position of a transition zone changes, depending on the value of  $\tau_D$ . Table 1 shows the smaller boundary,  $(\Delta\tilde{a}_0)_1$ , and the larger boundary,  $(\Delta\tilde{a}_0)_2$ , of a transition zone for each  $\tau_D$ . The values of  $(\Delta\tilde{a}_0)_1$  and  $(\Delta\tilde{a}_0)_2$  (i.e., the position of a transition zone) move toward large  $\Delta\tilde{a}_0$ , in other words, the unstable zone becomes wide, as  $\tau_D$  increases. This is because an increase in  $\tau_D$  reduce the efficiency of an eccentricity damping by the drag force (see equation (7)) and, as a result, a system becomes more unstable.

The above-mentioned properties of an orbital instability of a system with a gas disk are almost the same for the cases where  $M_p$  is  $10^{-9}M_\odot$ . In table 2, the values of  $(\Delta\tilde{a}_0)_1$  and  $(\Delta\tilde{a}_0)_2$  for the cases where  $M_p$  is  $10^{-9}M_\odot$  are also shown. For these cases, the position of a transition zone becomes large with an increase in  $\tau_D$ .

Figure 3 shows the transition zones for the cases where  $M_p$  is  $10^{-9}M_\odot$ , where  $M_p$  is  $10^{-7}M_\odot$ , and, where  $M_p$  is  $10^{-5}M_\odot$ , against  $\tilde{\tau}_D$ . Here,  $\tilde{\tau}_D$  is defined by

$$\tilde{\tau}_D = \left( \frac{M_p}{1 \times 10^{-7}M_\odot} \right) \tau_D, \quad (14)$$

where  $\tau_D$  is given by equation (6). Thus, the abscissa of figure 3, i.e.,  $\tilde{\tau}_D$ , is independent of the mass of a

Table 2: Boundaries of a transition zone,  $(\Delta\tilde{a}_0)_1$  and  $(\Delta\tilde{a}_0)_2$ .

$\tau_D/10^3 T_K$		1.0	3.0	9.0	27
$M_p = 10^{-9} M_\odot$	$(\Delta\tilde{a}_0)_1$	3.8	4.2	4.8	5.8
	$(\Delta\tilde{a}_0)_2$	5.2	5.2	6.0	6.8
$M_p = 10^{-7} M_\odot$	$(\Delta\tilde{a}_0)_1$	4.6	5.0	5.4	6.2
	$(\Delta\tilde{a}_0)_2$	5.0	5.4	5.8	6.6
$M_p = 10^{-5} M_\odot$	$(\Delta\tilde{a}_0)_1$	4.2	4.8	5.2	5.6
	$(\Delta\tilde{a}_0)_2$	5.2	5.6	5.6	6.6

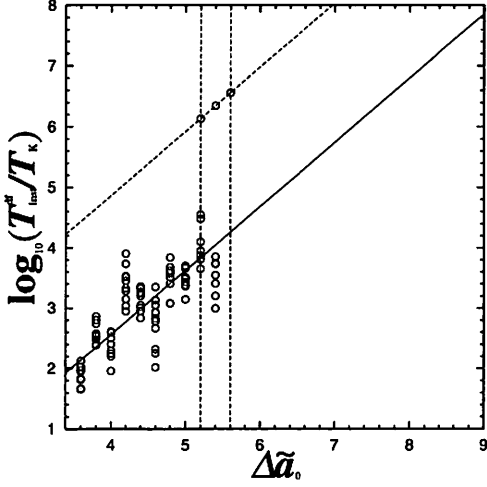


Figure 2: Instability time of a protoplanet system with a gas disk,  $T_{\text{inst}}^{\text{df}}$ , as a function of  $\Delta\tilde{a}_0$ , for the cases where  $M_p$  is  $10^{-5} M_\odot$  and  $\tau_D$  is  $9 \times 10^3 T_K$ . For each  $\Delta\tilde{a}_0$ , ten cases with different initial phase in position angles are plotted like figure 1. Solid line shows the instability time of a system without a gas disk,  $T_{\text{inst}}$ , (see equation 10) and the upper dotted line denotes the cut-off time,  $T_{\text{stop}} (= 200 \times T_{\text{inst}})$ . Two vertical dotted lines represent the boundaries of a transition zone, i.e.,  $(\Delta\tilde{a}_0)_1$  (left), and,  $(\Delta\tilde{a}_0)_2$  (right).

protoplanet,  $M_p$ , and, depends purely on the physical quantities of the gas disk, i.e.,  $c_s$  and  $\sigma_g$ . In figure 3, we also plotted the “critical separation distances”,  $(\Delta\tilde{a}_0)_{\text{crit}}$ , which were introduced in Iwasaki et al. (2001, 2002), in order to point to the position of a transition zone definitely and given by

$$(\Delta\tilde{a}_0)_{\text{crit}} = \frac{(\Delta\tilde{a}_0)_1 + (\Delta\tilde{a}_0)_2}{2}. \quad (15)$$

The lines in figure 3 represent the semi-analytical expressions for a critical separation distance,  $(\Delta\tilde{a}_0)_{\text{crit}}$ , which are derived in the same way in Iwasaki et al. (2002). These expressions are commonly obtained by meeting the requirement that the damping time of a eccentricity due to the drag force caused by the gas disk,  $\tau_D$ , must be equal to the stirring time of an eccentricity by the mutual gravity of protoplanets,  $\tilde{C}T_{\text{inst}}$  ( $\tilde{C}$  is a numerical constant), when  $\Delta\tilde{a}_0 = (\Delta\tilde{a}_0)_{\text{crit}}$ , i.e., (Iwasaki et al. 2002)

$$\tau_D = \tilde{C}T_{\text{inst}} \text{ at } \Delta\tilde{a}_0 = (\Delta\tilde{a}_0)_{\text{crit}}, \quad (16)$$

where  $T_{\text{inst}}$  is given by equation (10). Solving equation (16) for  $(\Delta\tilde{a}_0)_{\text{crit}}$ , we obtain

$$(\Delta\tilde{a}_0)_{\text{crit}} = \frac{1}{b} \left[ \log_{10} \left( \frac{\tilde{\tau}_D}{T_K} \right) + \log_{10} \left( \frac{M_p}{1 \times 10^{-7} M_\odot} \right) \right] - \frac{1}{b} (\log_{10} \tilde{C} + c). \quad (17)$$

Here, in the above equation, we set the value of  $\tilde{C}$  to be 0.3, following Iwasaki et al. (2002). In figure 3, the lines given by equation (17) are in good agreement with the values of  $(\Delta\tilde{a}_0)_{\text{crit}}$  obtained by numerical calculations, regardless of the mass of a protoplanet,  $M_p$ . This implies that equation (17) can be applied to a system composed of protoplanets with an arbitrary (but equal) mass.

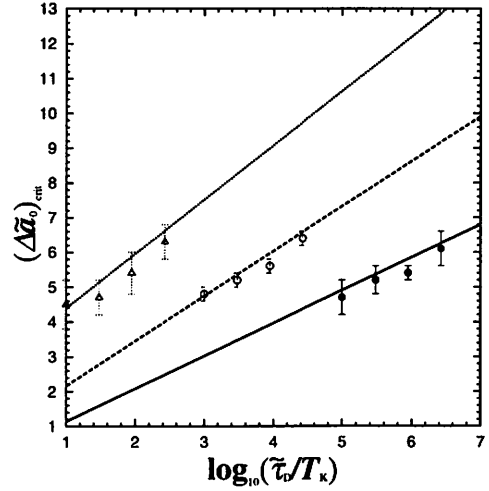


Figure 3: Critical separation distance,  $(\Delta\tilde{a}_0)_{\text{crit}}$ , as a function of  $\tilde{\tau}_D$  for the cases, where  $M_p$  is  $10^{-9} M_\odot$  (triangle), where  $M_p$  is  $10^{-7} M_\odot$  (open circle), and, where  $M_p$  is  $10^{-5} M_\odot$  (filled circle). Vertical error bars represent the width of a transition zone. Three lines, i.e., the dotted line (for the case of  $M_p = 10^{-9} M_\odot$ ), the broken line ( $10^{-7} M_\odot$ ), and, the solid line ( $10^{-5} M_\odot$ ), show the semi-analytical expressions for  $(\Delta\tilde{a}_0)_{\text{crit}}$  given by equation (17).

## 4 Summary and Conclusion

In our present study, we investigated the orbital stability of a system composed of five protoplanets with equal mass, for the cases where  $M_p$  is  $10^{-9} M_\odot$  and where  $M_p$  is  $10^{-5} M_\odot$ . First, we studied the orbital stability of a protoplanet system without a gas disk through orbital calculations. Next, the orbital stability of a system with a gas disk was examined by including the effect of the drag force caused by a gas disk in orbital calculations. Our main results are summarized as follows:

- Orbital stability of a protoplanet system without a gas disk.
- (1) A protoplanet system inevitably undergoes an orbital instability owing to mutual gravity.



- (2) Regardless of the mass of a protoplanet,  $M_p$ , the orbital instability time,  $T_{\text{inst}}$ , increases exponentially with an increase in the separation distance,  $\Delta\tilde{a}_0$ , i.e.,

$$\log_{10}(T_{\text{inst}}/T_K) = b\Delta\tilde{a}_0 + c. \quad (18)$$

However, the values of  $b$  and  $c$  depend on the mass of a protoplanet. Especially, the value of  $b$  increases as  $M_p$  becomes large.

- Orbital stability of a protoplanet system in the nebular gas.

- (1) The drag force due to the nebular gas has an effect of preventing a protoplanet system from undergoing an orbital instability.
- (2) Such a effect of the drag force becomes remarkable, when  $\Delta\tilde{a}_0$  is larger than a critical separation distance,  $(\Delta\tilde{a}_0)_{\text{crit}}$ , and the instability time,  $T_{\text{inst}}^{\text{df}}$ , becomes at least 200 times as large as the instability time of a system without a gas disk,  $T_{\text{inst}}$ .
- (3) The value of a critical separation distance,  $(\Delta\tilde{a}_0)_{\text{crit}}$ , depends on  $\tau_D$  (i.e., the intensity of the drag force) and is estimated by

$$(\Delta\tilde{a}_0)_{\text{crit}} = \frac{1}{b} \left[ \log_{10} \left( \frac{\tau_D}{T_K} \right) - (\log_{10} \tilde{C} + c) \right]. \quad (19)$$

Now, we apply the above results to the orbital instability of a protoplanet system in the Jovian planet region ( $a \simeq 5$  AU). As mentioned in the first section, a protoplanet system must experience an orbital instability in the presence of a gas disk, if we consider the accretion process of protoplanets via mutual collisions, as a possible way to overcome the difficulty in the formation of Jovian planets that the mass of a protoplanet is too small to capture a gaseous envelope. However, the above results show that, in the presence of a gas disk, a protoplanet system never experiences an orbital instability if the separation distance,  $\Delta\tilde{a}_0$ , is larger than a critical separation distance,  $(\Delta\tilde{a}_0)_{\text{crit}}$ . The value of  $(\Delta\tilde{a}_0)_{\text{crit}}$  is obtained by substituting the value of  $\tau_D$  in the Jovian planet region into equation (19). Substituting  $M_p = 10^{-5}M_\odot \simeq 3M_\oplus$ ,  $c_s = 0.05v_K$ , and  $\sigma_g = 1.7 \times 10^3 \text{g/cm}^3$ , (which are the values in the minimum mass solar nebula model) into equation (6), we obtain

$$\tau_D = 2.3 \times 10^2 T_K. \quad (20)$$

From equations (19) and (20),  $(\Delta\tilde{a}_0)_{\text{crit}}$  in this case is given by

$$(\Delta\tilde{a}_0)_{\text{crit}} = 4.3. \quad (21)$$

On the other hand, a typical separation distance,  $\Delta\tilde{a}_0$ , which would be formed through oligarchic growth is about 5 to 10 (Kokubo and Ida 2000), i.e., is larger than  $(\Delta\tilde{a}_0)_{\text{crit}}$  in the above equation. Thus, in the presence of the minimum mass solar nebula, a protoplanet system never experiences an orbital instability. Naturally, the value of  $(\Delta\tilde{a}_0)_{\text{crit}}$  can become larger than the

above value, if the nebular gas dissipates, i.e., the surface density of the nebular gas decreases. However, in order to form a gaseous giant planet such as the present Jupiter whose mass is about  $1 \times 10^{-3}M_\odot$ , there must remain a nebular gas whose surface density is as large as that of the minimum mass solar nebula around a protoplanet, if we assume that the present Jupiter captured the nebular gas in its feeding zone (i.e., a ring region with a width of about 6 Hill radius). Therefore, the scenario that protoplanets grew to a critical mass owing to mutual collisions is implausible. It is natural to think that a surface density of the solid mass in the Jovian planet region was two or three times as large as that of the minimum mass solar nebula.

## ACKNOWLEDGMENTS

The authors express their sincere gratitude to S. Ida for continuous encouragement and valuable advice. Numerical computation was helped by SX-5/16 at the Computer Center of Tokyo Institute of Technology.

## References

- Adachi, I., Hayashi, C., & Nakazawa, K. 1976, *Prog. Theor. Phys.* 56, 1756
- Artymowicz, P. 1993, *ApJ* 419, 166
- Bodenheimer, P., & Pollack, J. B. 1986, *Icarus* 67, 391
- Chambers, J. E., & Wetherill, G. W. 1998, *Icarus* 136, 304
- Chambers, J. E., Wetherill, G. W., & Boss, A. P. 1996, *Icarus* 119, 261
- Hayashi, C. 1981, *Prog. Theor. Phys. Suppl.* 70, 163
- Hayashi, C., Nakazawa, K., & Nakagawa, Y. 1985, in *Protostars and Planets II*, ed D. C. Black, M. S. Matthews (Tucson, University of Arizona Press), 1100
- Ikoma, M., Nakazawa, K., & Emori, H. 2000, *ApJ* 537, 1013
- Ito, T., & Tanikawa, K. 1999, *Icarus* 139, 336
- Ito, T., & Tanikawa, K. 2001, *PASJ* 53, 143
- Iwasaki, K., Tanaka, H., Nakazawa, K., & Emori, H. 2001, *PASJ* 53, 321
- Iwasaki, K., Emori, H., Nakazawa, K., & Tanaka, H. 2002, *PASJ* 54, in press
- Kokubo, E., & Ida, S. 1996, *Icarus* 123, 180
- Kokubo, E., & Ida, S. 1998, *Icarus* 131, 171
- Kokubo, E., & Ida, S. 2000, *Icarus* 114, 247
- Kokubo, E., & Makino, J. 1998, in *Proceedings of the 30th symposium on celestial mechanics*, ed T. Fukushima, T. Ito, T. Fuse, H. Umehara, p248
- Kokubo, E., Yoshinaga, K., & Makino, J. 1998, *MNRAS* 297, 1067
- Kominami, J., & Ida, S. 2001, *Icarus*, in press
- Mizuno, H. 1980, *Prog. Theor. Phys.* 64, 544
- Ward, W. R. 1988, *Icarus* 73, 330
- Wetherill, G. W., & Stewart, G. R. 1989, *Icarus* 77, 330
- Yoshinaga, K., Kokubo, E., & Makino, J. 1999, *Icarus* 139, 328

# Dynamical Stability of Planetary System of GJ876

Hiroshi Kinoshita and Hiroshi Nakai  
National Astronomical Observatory

## Abstract

The main-sequence star GJ 876 was found to have at least two planets from the precise Doppler measurements made at the Lick and Keck observatories. If the two planets are moving in the same plane, which lies in the line of sight, the planetary system of GJ 876 is stabilized by the 2:1 mean motion resonance. As the inclination of the orbital plane to the line of sight decreases from 90 degrees, the planetary mass increases and the mutual perturbation between two planets becomes large. However the planetary system of GJ 876 continues to be stabilized by the corotation of pericenters of the two planets together with the mean motion resonance. This stabilization mechanism is explained by the secular perturbation theory.

## 1 Introduction

So far 77 extrasolar planets have been discovered since 1995 and among them 7 multiple planetary systems ( $\nu$  Andromedae, GJ876, HD168443, HD82943, HD74156, 47 Uma) are confirmed.  $\nu$  Andromedae has three planets, the orbital period of the most inner planet is only 4.6 days and its mass is  $0.68 M_J$ . The orbital periods of two outer planets are 241.3 days and 1299 days and their masses are  $1.94 M_J$  and  $4.02 M_J$ , respectively (California & Carnegie Planet Search, 2002). Since the perturbation to the most inner planet from the outer planets are weak, the motion of the most inner planet is stable. On the other hand the mutual interaction between outer two planets might become strong and unstable because of a close approach due to their large eccentricities. The orbital motion of outer two planets, however, is stabilized by the alignment of the pericenters of two planets, which makes to avoid a close approach between them (Nakai and Kinoshita 2000, 2001, Rivera and Lissauer 2000, and Kinoshita and Nakai 2002). The planetary systems of GJ876 and HD82943 are stabilized by the 2:1 mean motion resonance. As for other four planetary systems the mutual distance is large and the mutual perturbation is weak and their orbital motions are stable.

Marcy et al. (2001) discovered from precise Doppler observations during six years from the Lick and Keck Observatories at least two planets orbiting GJ 876. With the assumption

that the two planets do not disturb each other and their orbits are Keplerian, Marcy et al. determined the following dynamical parameters: masses of  $M \sin i = 0.56$  and  $1.89 M_J$ , orbital periods of  $P=30.1$  and  $61.0$  days, semimajor axes of  $a = 0.13$  and  $0.21$  AU, and eccentricities of  $e=0.27$  and  $0.10$ , respectively. Marcy et al. (2001) performed numerical simulations with minimum mass ( $i = 90^\circ$ ) and twice their minimum masses ( $\sin i = 1/2$ ) under the assumption that the orbital elements in the above are the osculating elements and found that the planetary system of GJ 876 is stable by the 2:1 mean motion resonance.

We assumed the minimum planet mass ( $\sin i = 1$ ) and carried out the numerical simulations for this planetary system with the initial conditions, which are different from those of Marcy et al. (2001) and investigated the orbital stability. We put the planet 2 (the outer planet) on its pericenter at the initial epoch and searched a stable configuration with changing the initial mean anomaly  $l_1$  of planet 1 (the inner planet) from  $0$  to  $360$  degrees. We found that when  $l_1$  is in the range of  $-60 < l_1 < 60$ , the planetary system is in the state of 2:1 mean motion resonance (the critical argument  $\sigma = \lambda_1 - 2\lambda_2 + \varpi_1$ ), which avoids a close approach between two planets and makes the planetary system stable. At the initial epoch the pericenter of planet 2 is very close to that of planet 1. However, in these stable planetary systems the pericenters of two planets move independently and not corotating as in the case of planetary system of  $\nu$  Andromedae.

Since the determination of the planetary mass has ambiguity of  $\sin i$ , which cannot be determined from the Doppler observation under the assumption of two Keplerian orbital fitting. If  $i$  is small, then the planetary mass becomes large and the mutual perturbation grows and the system might be unstable. By changing  $\sin i$  we investigate the stability of the planetary system in the 2:1 mean motion resonance and found that the alignment of pericenter stabilizes the planetary systems and explained this stability mechanism by semi analytical secular perturbation.

## 2 Mass Dependence of the Stability of the Planetary System of GJ 876

The orbital parameters, which Marcy et al. (2001) determined from the Doppler observations, are shown in Table 1. These parameters were determined under the assumption that there is no mutual perturbation between two planets. As for the mass determination from Doppler observation the combination  $M \sin i$ , where  $i$  is the inclination of the orbital plane to the line of sight, is determined and  $i$  and  $M$  cannot be separately determined. We define a mass factor  $m_f$  as

$$m_f = \frac{1}{\sin i}, \quad (1)$$

Table 1: Orbital Parameters (Marcy et al. 2001)

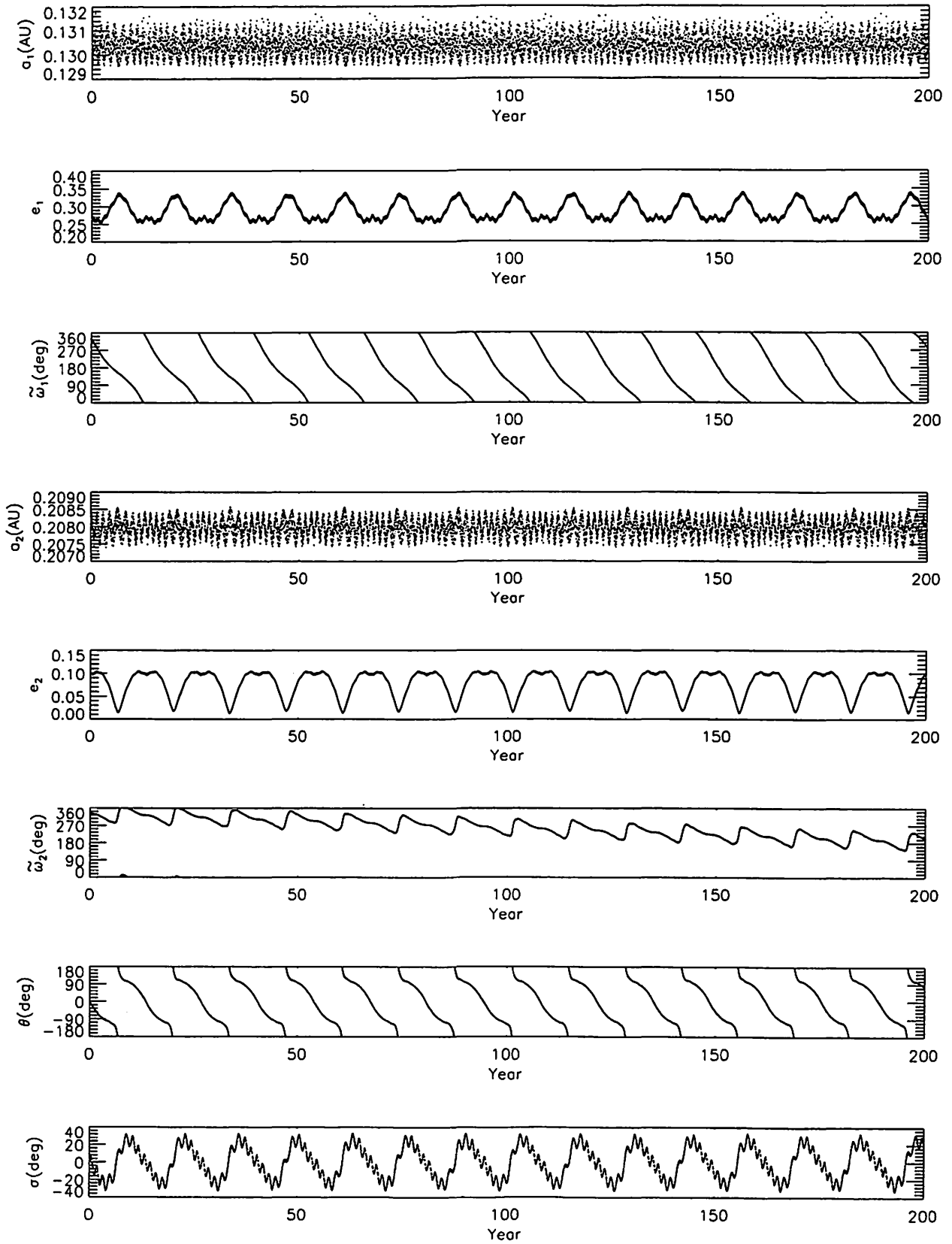
Parameters	Inner	Outer
Orbital Period P(day)	30.12	61.02
Eccentricity ( $e$ )	0.27	0.10
$\varpi$ (deg)	330	333
Periastron Time(JD)	2450031.4	2450106.2
$M \sin i(M_J)$	0.56	1.89
$a$ (AU)	0.130	0.208

where we assume the two planets move in the same plane. Since even in the case of  $m_f = 1$  and not in the state of 2:1 mean motion resonance the planetary system is unstable (Kinoshita and Naka 2001, 2002b), we choose the critical argument  $\sigma = \lambda_1 - 2\lambda_2 + \varpi_1 = 0$  at the epoch. ( $\lambda_1 = 42.46, \lambda_2 = 186.23, \varpi_1 = 330.00, t_{epoch} = 2451362.7426$ , which lies in the period of the observations at Lick and Keck). We use an extrapolation method and a symmetric multiple step as a numerical integrator.

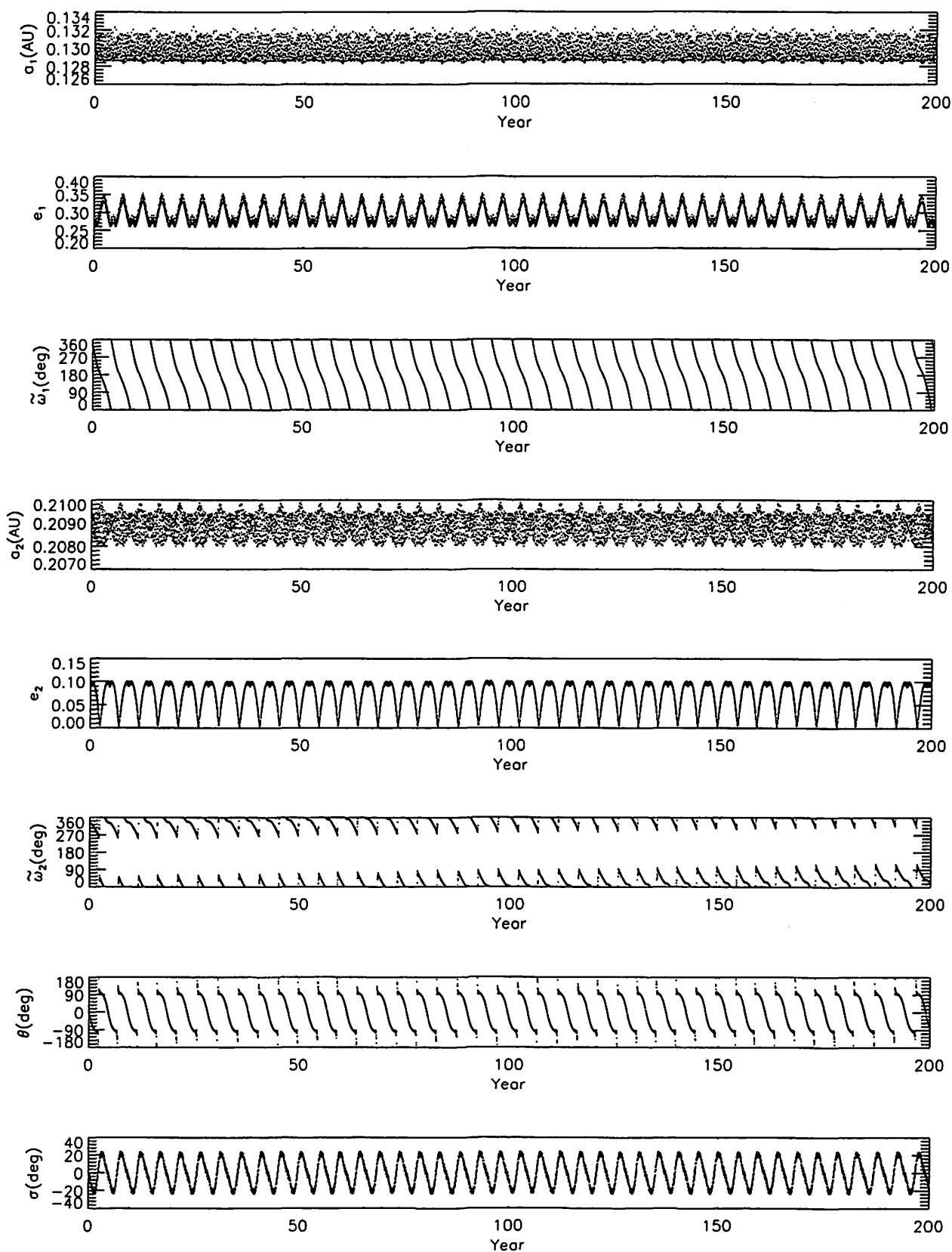
At first we integrated the planetary orbits for 200 years with increasing the mass factor  $m_f$  from 1 to 20 ( $i = 2.9^\circ$ ) by one and show their orbital elements as for  $m_f = 1, 3, 5, 10$ , and 11 in Figure 1. The upper three panels show the semimajor axis  $a_1$ , the eccentricity  $e_1$ , and the longitude of pericenter  $\varpi_1$ , and then the next three panels show the corresponding quantities for the outer planet, and the next two panels show the difference of the pericenter's longitudes ( $\theta = \varpi_1 - \varpi_2$ ) and the critical argument  $\sigma = \lambda_1 - 2\lambda_2 + \varpi_1$ . Even though the planetary system with  $m_f > 10$  is in the 2:1 mean motion resonance at the epoch, the 2:1 mean motion resonance state is soon disrupted and the eccentricities of the planets become large and then the close approach takes place and the planetary system becomes unstable. In order to know the measure of the close approach we calculate the mutual distance of two planets with the relative Hill radius. Here we choose the following three types of Hill radius:

$$\begin{aligned}
 R_{H1} &= \left( \frac{M_1 + M_2}{3M_c} \right)^{\frac{1}{3}} \left( \frac{a_1 + a_2}{2} \right), \\
 R_{H2} &= \left( \frac{M_1 + M_2}{3M_c} \right)^{\frac{1}{3}} \left( \frac{Q_1 + q_2}{2} \right), \\
 R_{H3} &= \left( \frac{M_1 + M_2}{3M_c} \right)^{\frac{1}{3}} \left( \frac{r_1 + r_2}{2} \right),
 \end{aligned} \tag{2}$$

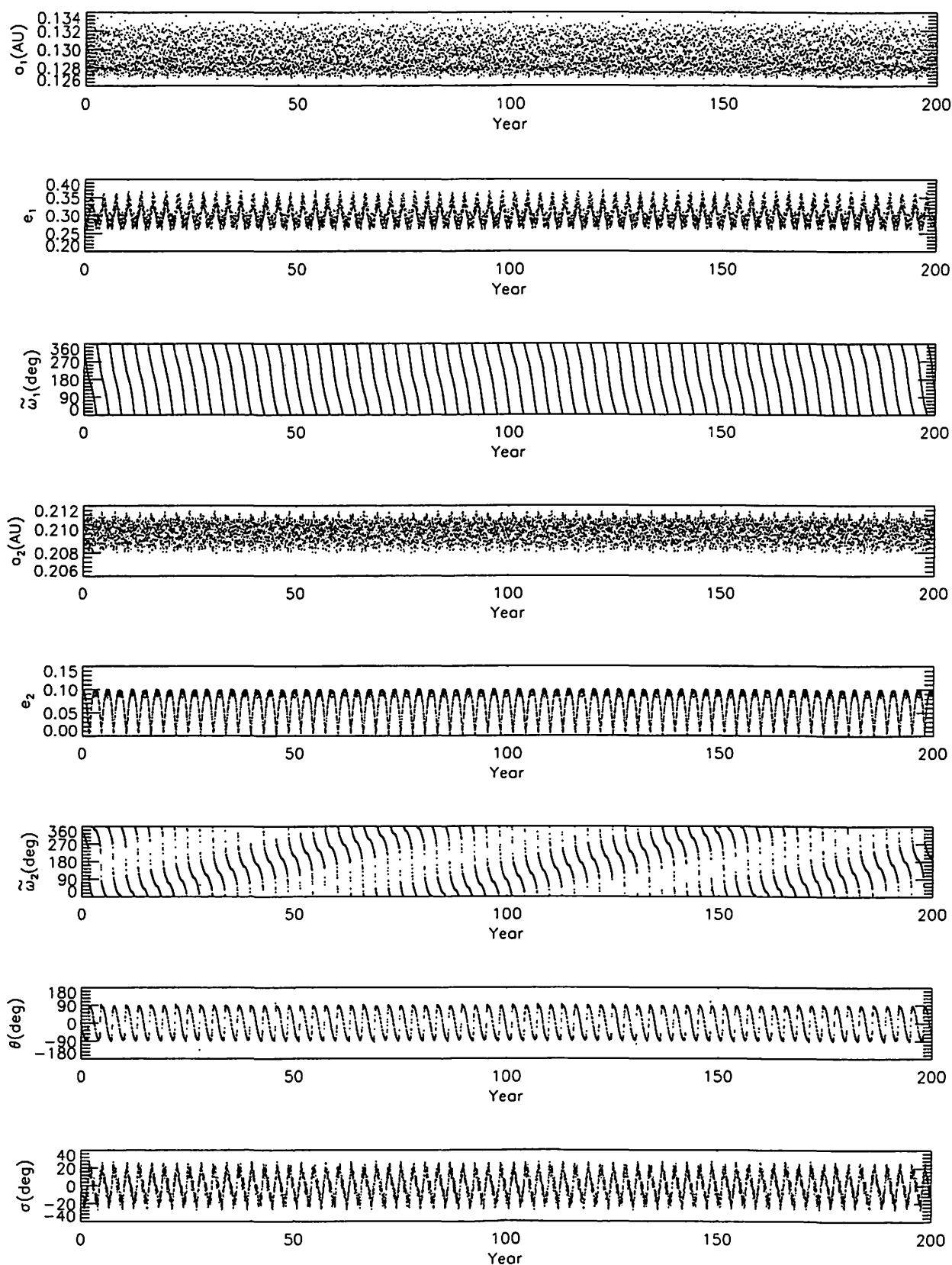
where  $Q_1 = a_1(1 + e_1)$  (the apocentric distance of the inner planet)  $q_2 = a_2(1 - e_2)$  is the pericentric distance of the outer planet). We normalized the mutual distance by these three types of relative Hill radius and found there is no qualitative difference among three normalization. Therefore, we only show the moralized by  $R_{H3}$  as for  $m_f = 1, 5, 10$ , and 11 in



**Figure 1-1.**— The orbital elements of two planets for 200 years ( $m_f = 1$ ): The upper six panels represent the semi-major axis, the eccentricity, and the longitude of the pericenter, respectively. The last two panels show the angles  $\theta = \varpi_1 - \varpi_2$  and the critical argument  $\sigma = \lambda_1 - 2\lambda_2 + \varpi_1$ .



**Figure 1-2.**—The orbital elements of two planets ( $m_f = 3$ ). The explanations for the vertical axes are same as for Figure 1-1.



**Figure 1-3.**—The orbital elements of two planets ( $m_f = 5$ ). The explanations for the vertical axes are same as for Figure 1-1.

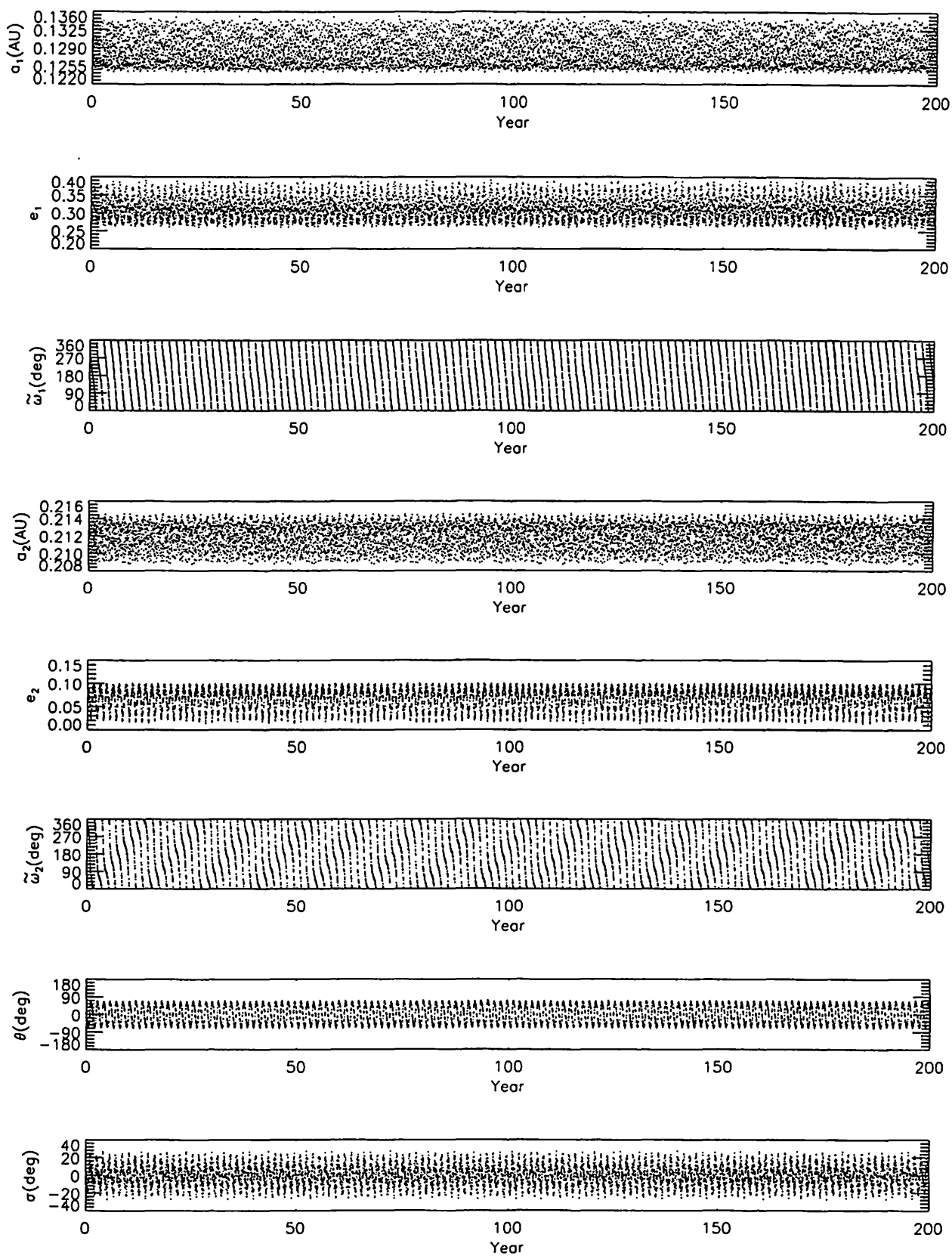
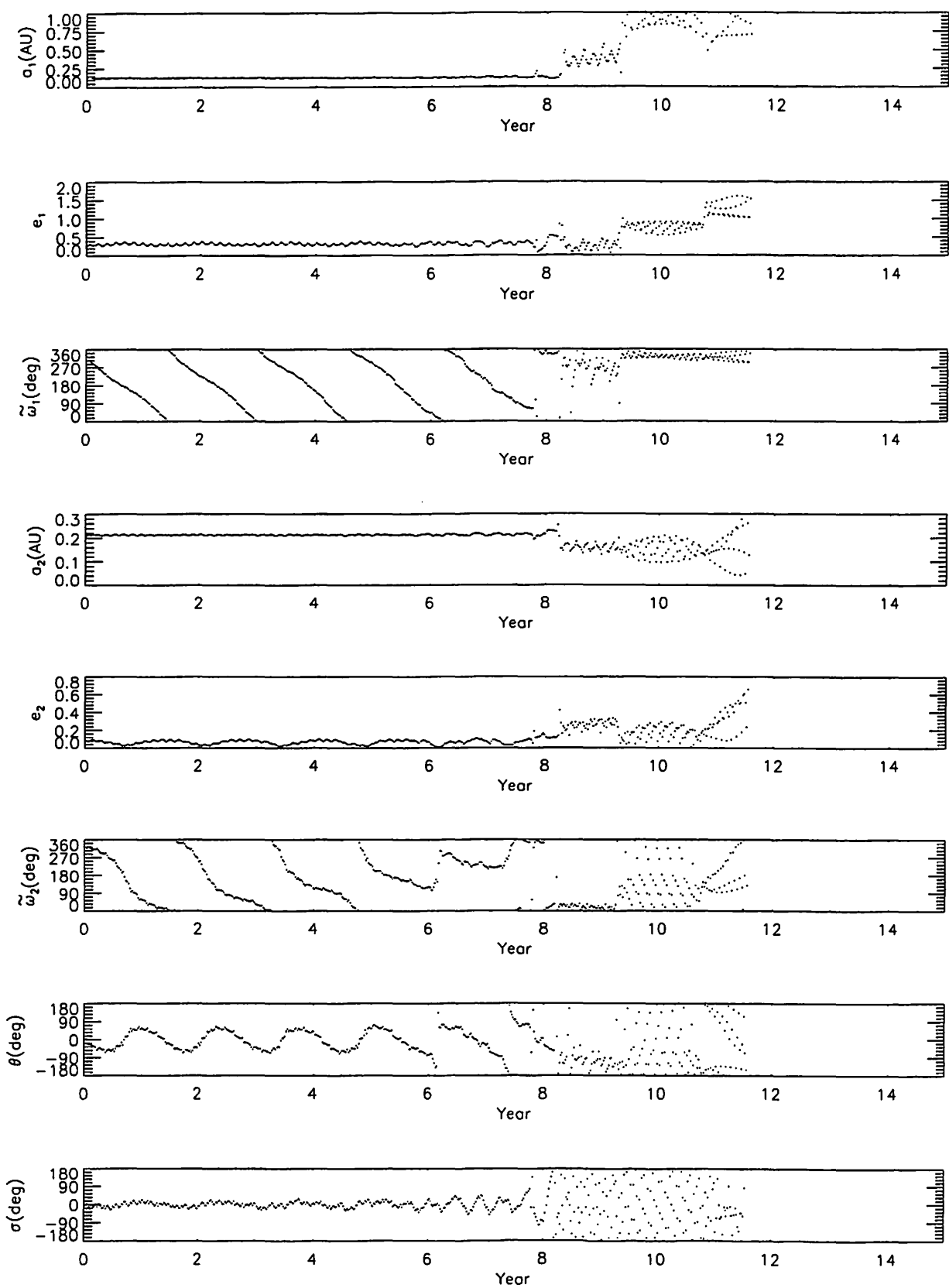


Figure 1-4.—The orbital elements of two planets ( $m_f = 10$ ).The explanations for the vertical axes are same as for Figure 1-1.





**Figure 1-5.**—The orbital elements of two planets ( $m_f = 11$ ). The explanations for the vertical axes are same as for Figure 1-1.

Figure 2. We define

$$\beta = \Delta/R_{H3}, \quad (3)$$

where  $\Delta$  is the minimum distance between the two planets. When the planetary system is stable, the parameter  $\beta$  becomes small from about 4 to 2. When the planetary system is unstable, the parameter  $\beta$  becomes smaller than 1, which means that the planets enter the Hill sphere and the orbits change largely and a close approach takes place and then the planetary system becomes unstable.

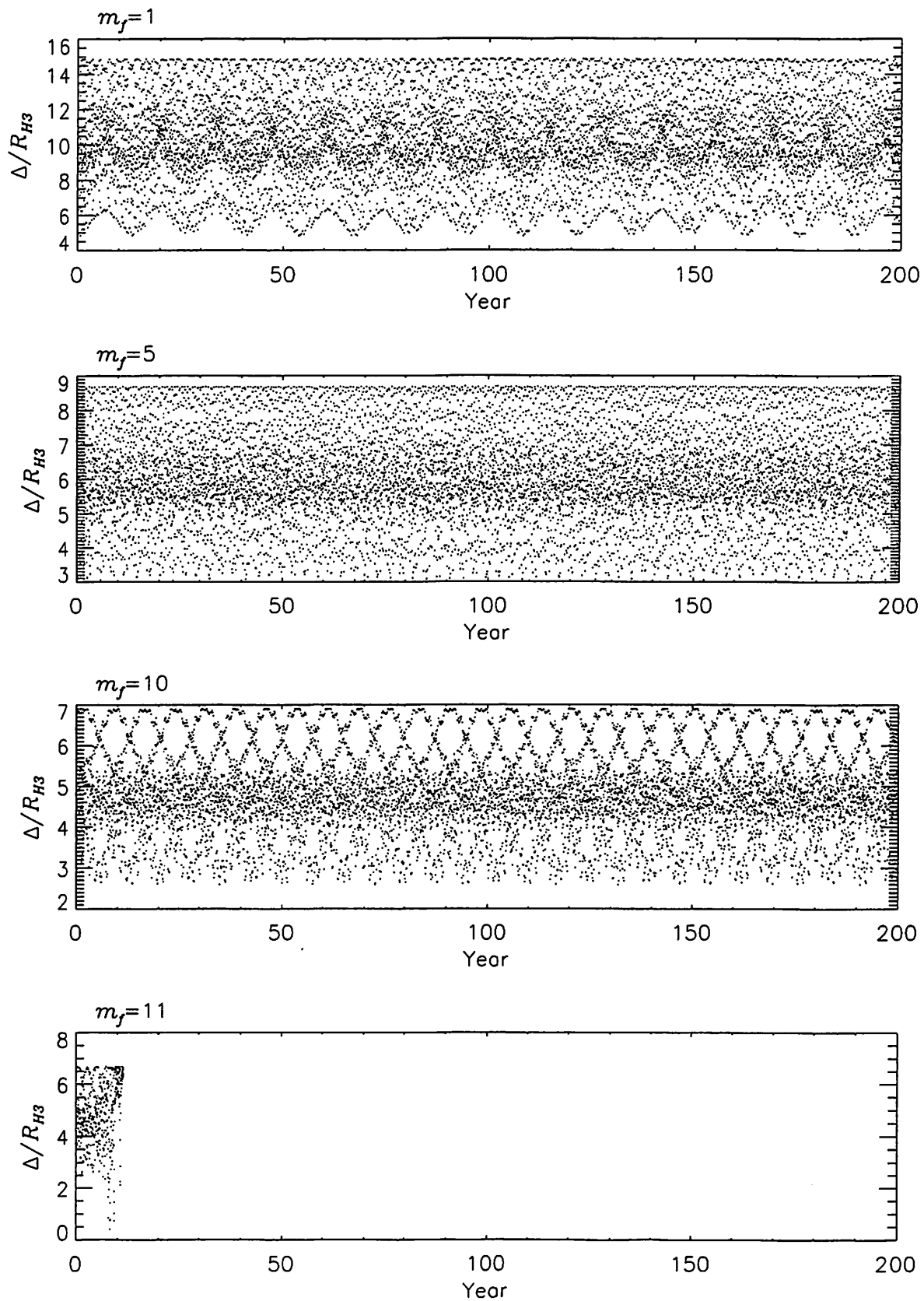
The notable features of the orbital elements for the stable planetary systems are 1) From around  $m_f = 2$ , the difference of pericenter longitudes  $\theta$  circulates for some time and then librates. The circulation and libration of  $\theta$  take place periodically. As  $m_f$  increases, the time of the libration of  $\theta$  becomes longer and for  $m_f = 10$  the angle  $\theta$  only librates.

2) Both the longitudes of pericenters are retrograde and the periods of the circulation of  $\varpi_1$  and  $\varpi_2$  becomes shorter as the mass factor  $m_f$  increases.

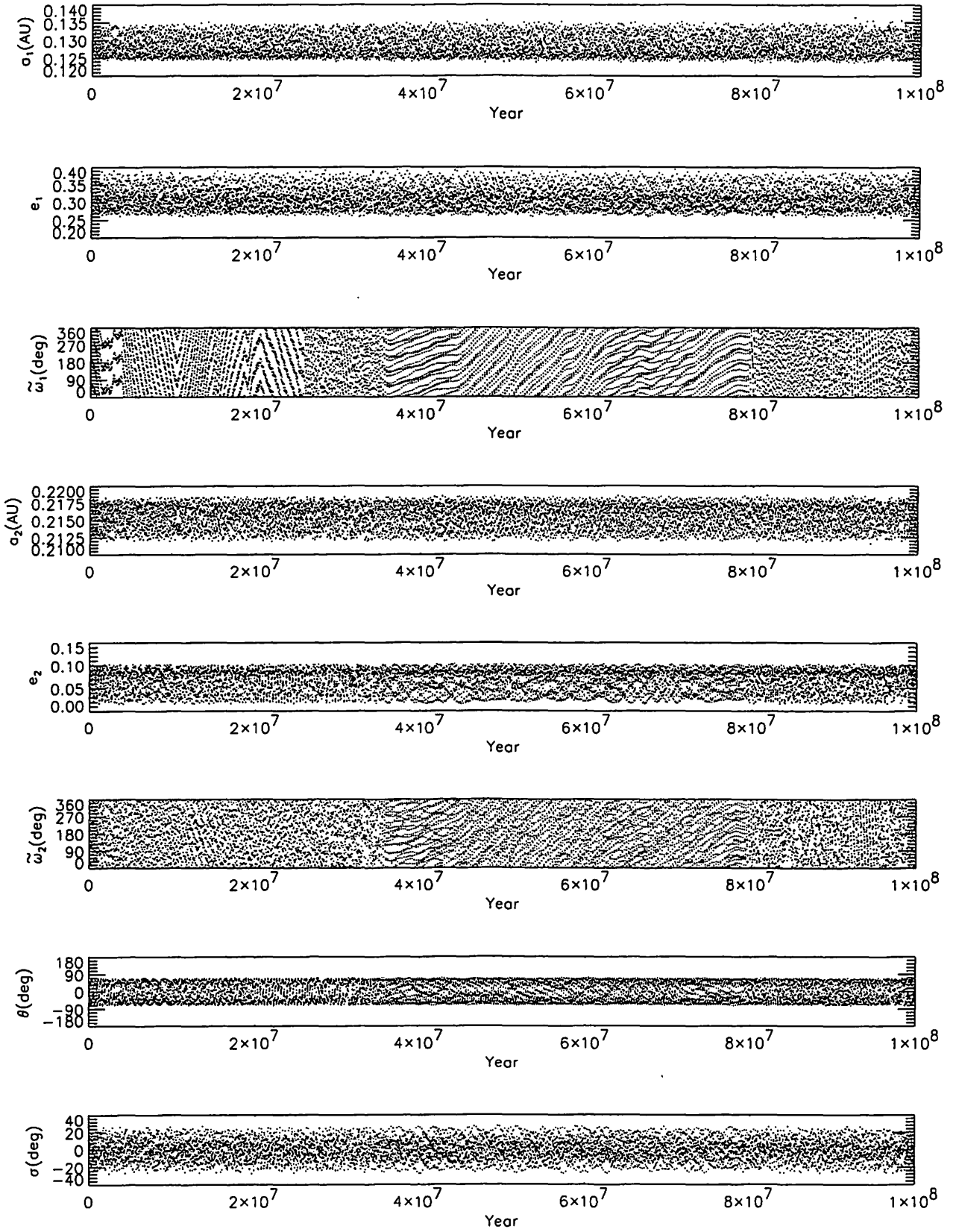
Then we integrated the planetary system, which is stable for 200 year, for  $10^8$  years and found that the planetary systems ( $1 \leq m_f \leq 10$ ) are stable and the 2:1 mean motion resonance and the corotation of the pericenters is conserved. Figure 3 show the orbital elements for the first five million years of  $10^8$  years integration as for  $m_f = 10$ . Figure 4 shows the instability time with respect to the mass factor  $m_f$ . Since the orbital elements as for  $1 \leq m_f \leq 10$  in Figure 3 show no indication of irregularity, the planetary system could plausibly be stable over  $10^8$  years.

### 3 Role of the Alignment of Pericenters

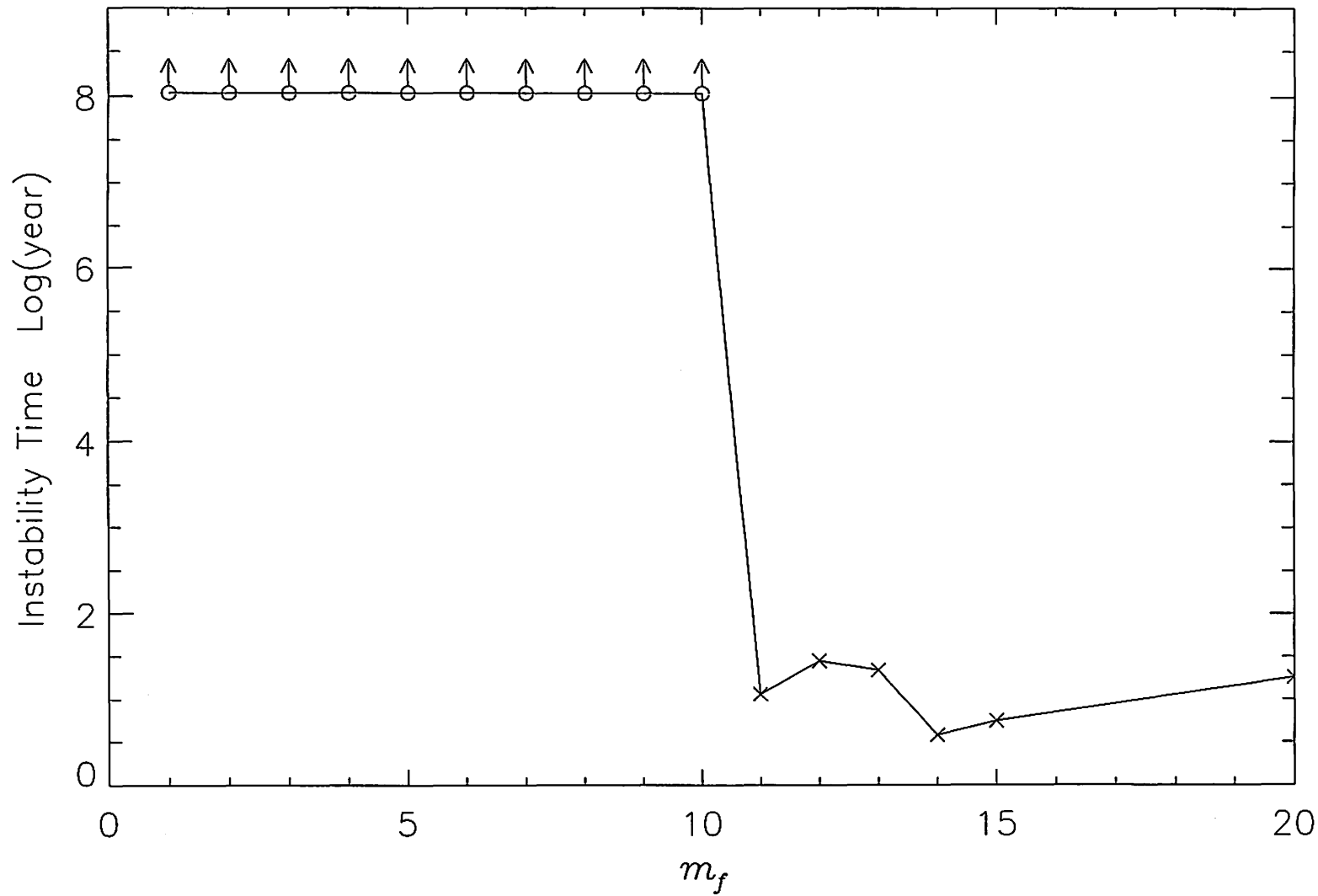
As mentioned in the previous section, around from  $m_f = 4$  the angle  $\theta$  almost librates, which means that the the pericenters of two planets corotate. The initial state in the computation in Section 2 the pericenters are very close. We, therefore, choose  $\varpi_2^* = \varpi - 180^\circ$  as the initial condition and integrated the orbits for 200 years with  $m_f = 1, 2, 3, 4$ . Figure 5 shows the evolution of the orbital elements for the case of  $m_f = 1, 3, 4$ . The angle  $\theta$  almost librates for the case of  $m_f = 2, 3$ . However for the planetary system with  $m_f = 4$  the pericenters moves totally move independently and around  $t = 30$  years the 2:1 mean motion resonance is disrupted and the planetary system becomes unstable. Figure 6 shows the mutual distance normalized by the mutual Hill radius  $R_{H3}$ . Since the pericenters move independently, at some time the apocenter of the inner planet and the pericenter of the outer planet are in the same direction, when  $\beta$  becomes smaller than 1, which takes place around  $t = 30$  years and the 2:1 mean motion resonance is broken.



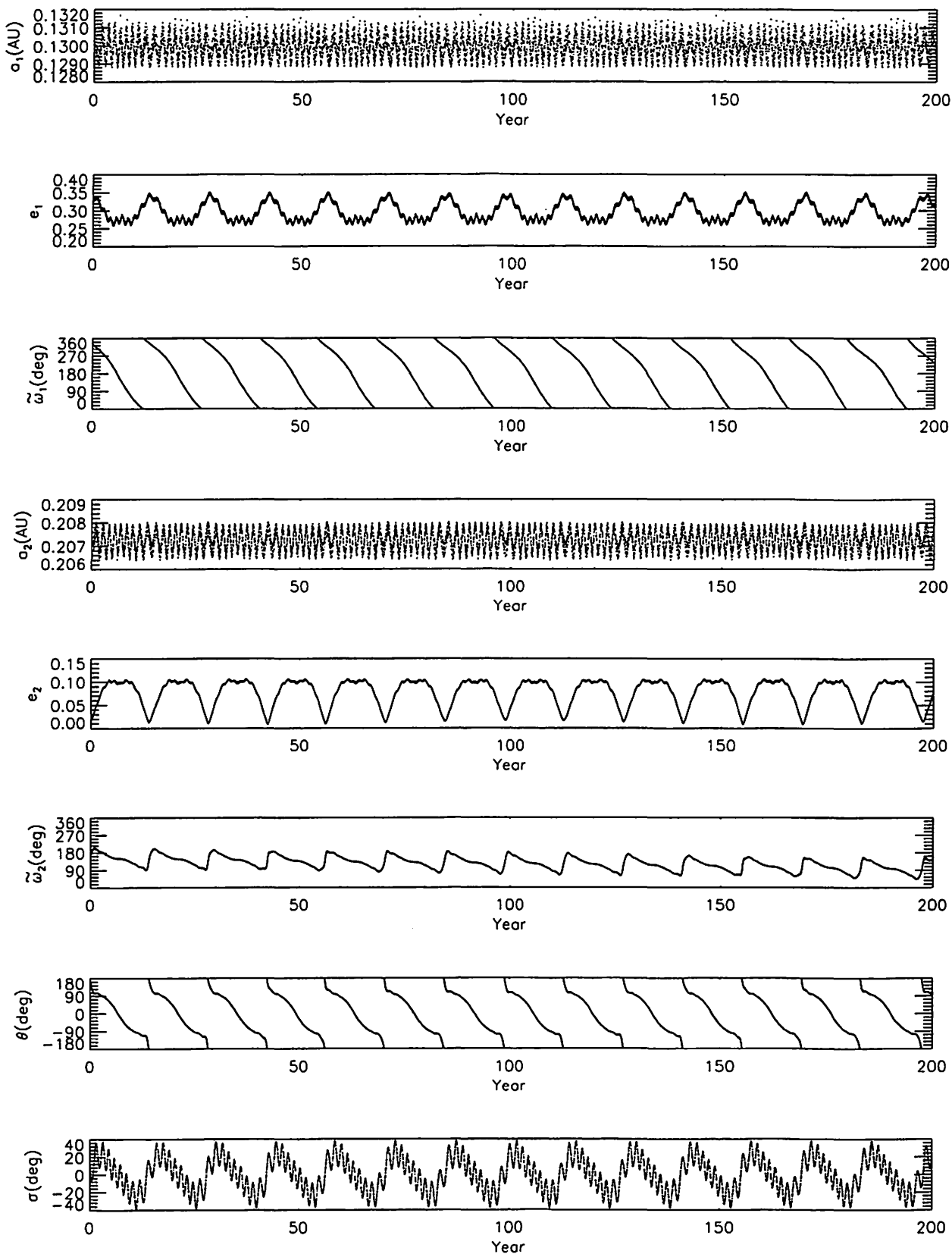
**Figure 2.**— (Mutual Distance)/(Relative Hill Radius( $R_{H3}$ )) for  $m_f=1,5,10$ , and 11.



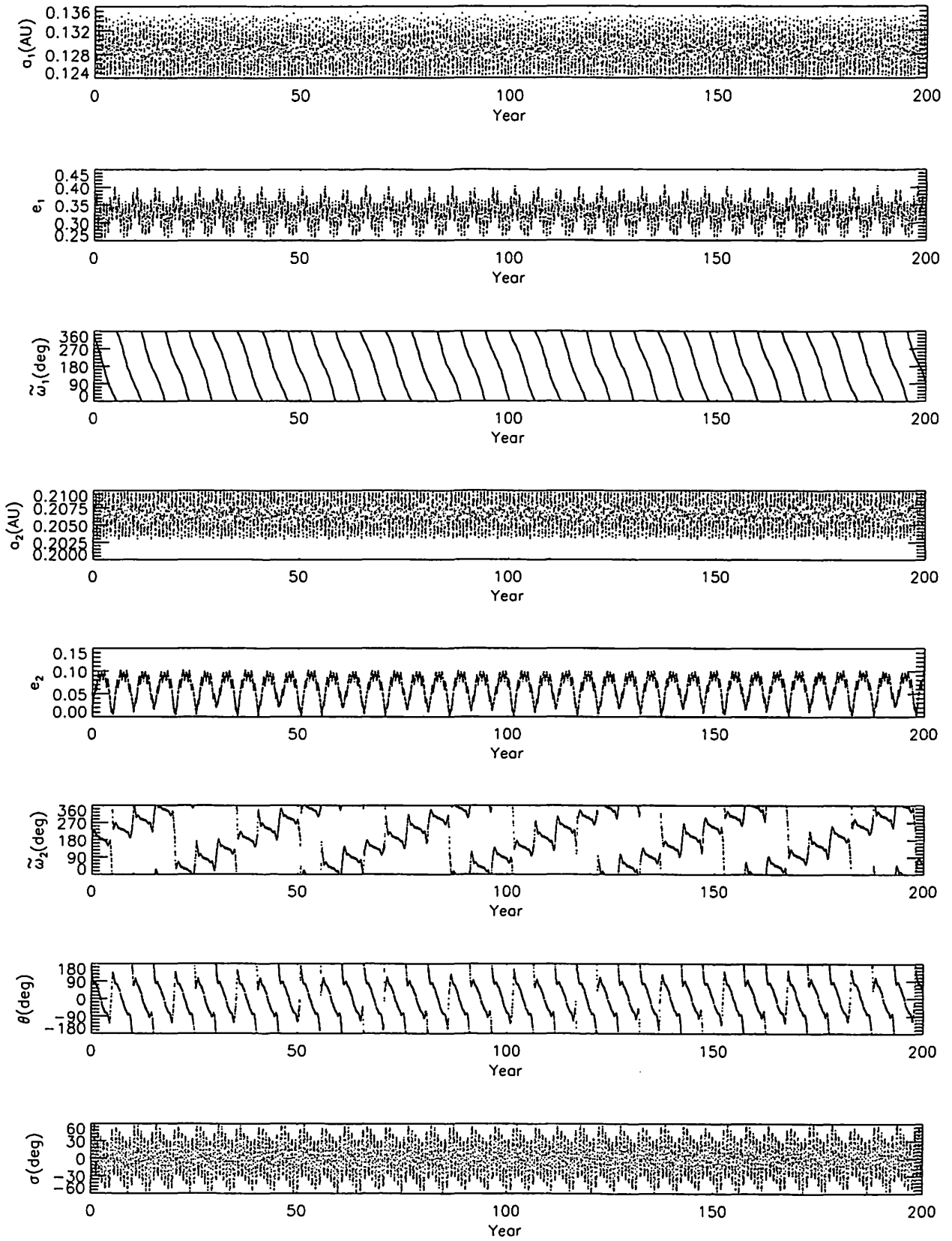
**Figure 3.**—The orbital Elements of two planets for 100 million years integration ( $m_f = 10$ ). The explanations for the vertical axes are same as for Figure 1-1.



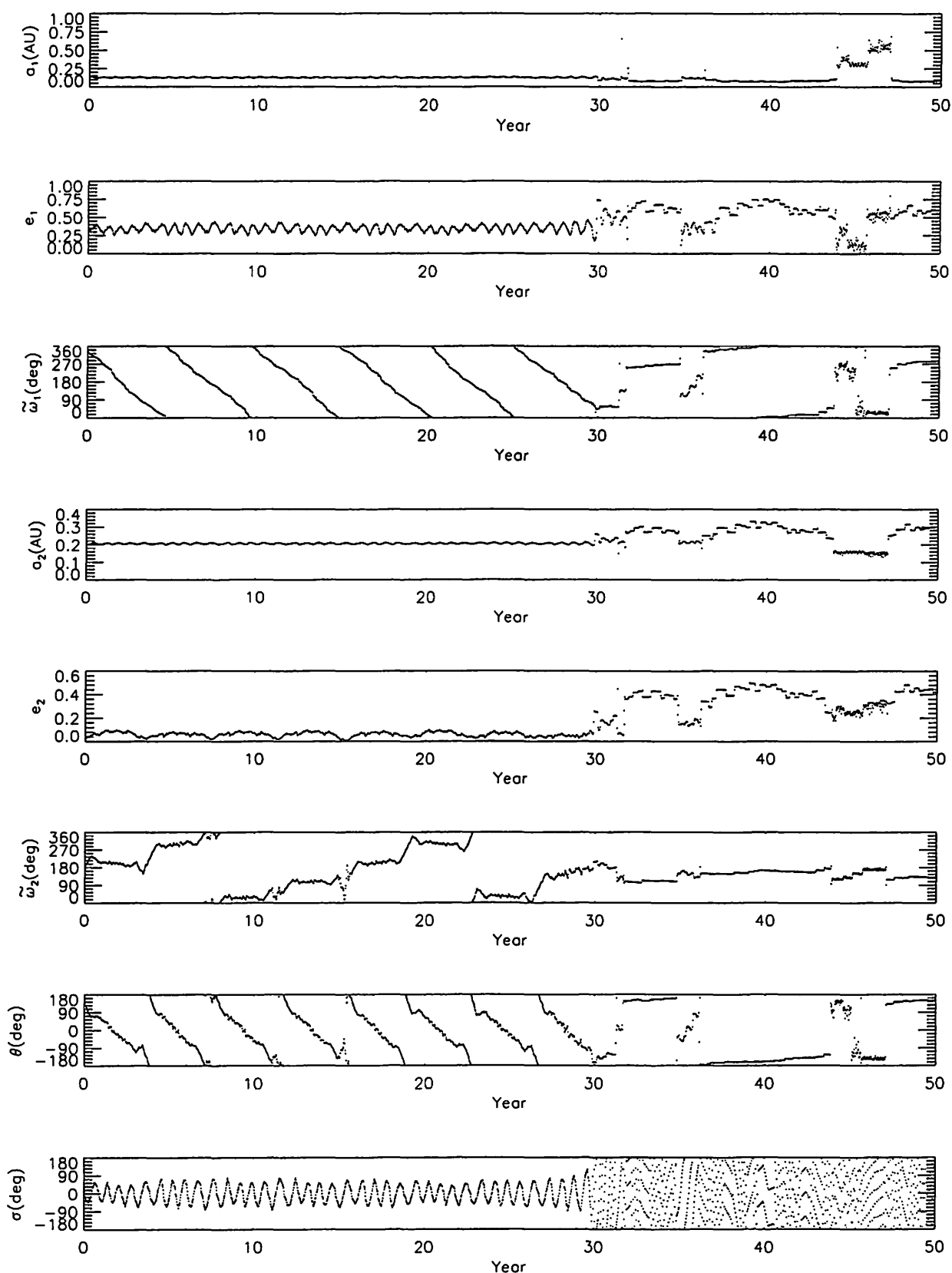
**Figure 4.**—The instability time for various  $m_f$ . The open circle means the system is stable till  $10^8$  years and the arrow indicates that the system might plausibly be stable over  $10^8$  years. The cross indicates that the eccentricity of one planet exceeds one.



**Figure 5-1.**—The orbital elements of two planets ( $m_f = 1$ ) with the initial condition with  $\varpi_2^* = \varpi_1 - 180^\circ$  with keep other elements as same as in the case of Figure 1. The explanations for the vertical axes are same as for Figure 1-1.



**Figure 5-2.**—The orbital elements of two planets ( $m_f = 3$ ) with the initial condition with  $\varpi_2^* = \varpi_1 - 180^\circ$  with keep other elements as same as in the case of Figure 1. The explanations for the vertical axes are same as for Figure 1-1.



**Figure 5-3.**—The orbital elements of two planets ( $m_f = 4$ ) with the initial condition with  $\omega_2^* = \omega_1 - 180^\circ$  with keep other elements as same as in the case of Figure 1. The explanations for the vertical axes are same as for Figure 1-1.



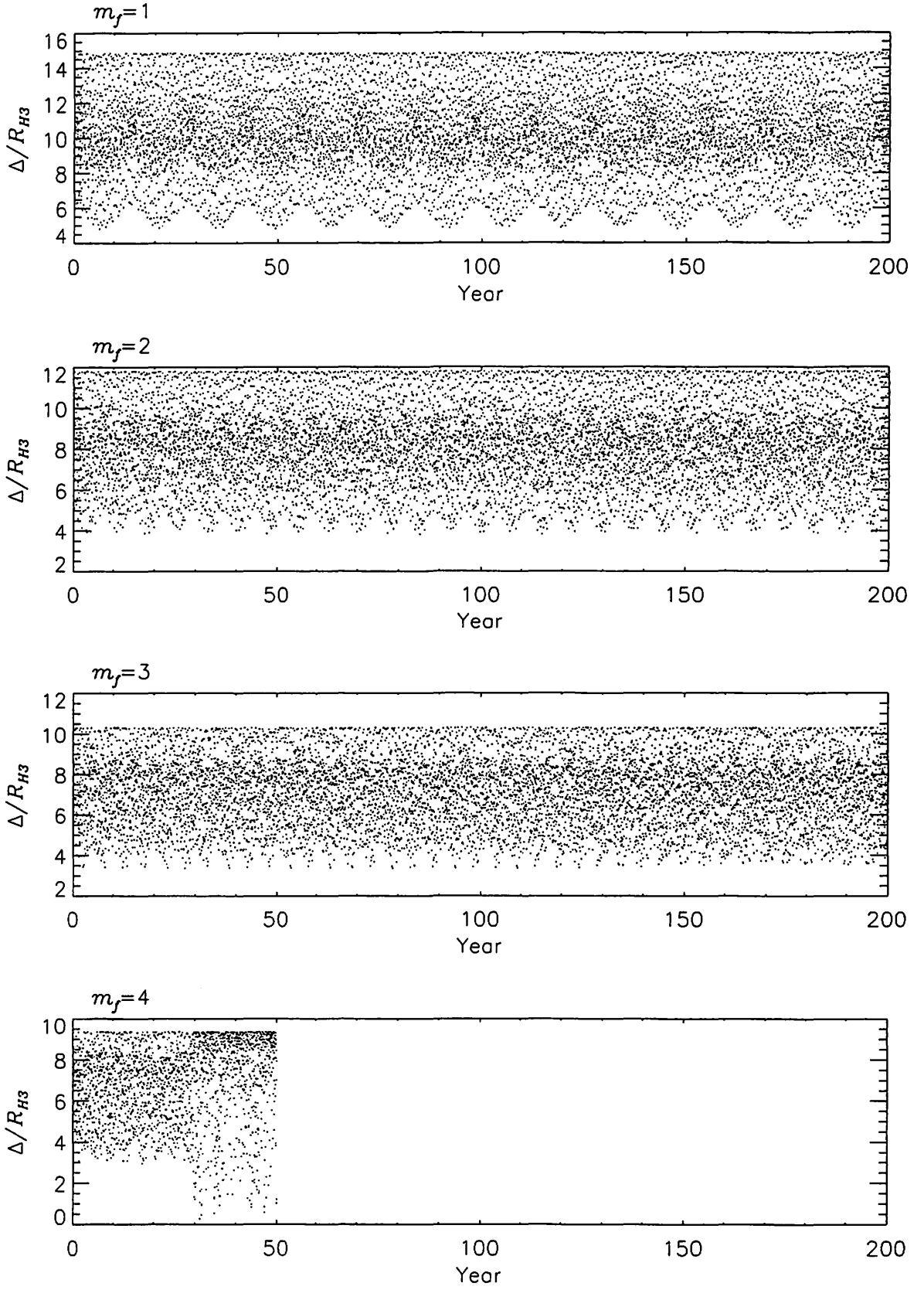


Figure 6.— (Mutual Distance)/(Relative Hill Radius( $R_{H3}$ )) for  $m_f = 1, 2, 3$ , and 4.

Table 2: Orbital Parameters (Laughlin and Chambers,2001a)

Parameter	Inner	Outer
Period(day)	30.13	61.58
Eccentricity $e$	0.226	0.025
$\varpi$ (deg)	156	70
Mean Anomaly(deg)	277	31
$a$ (AU)	0.1297	0.2092

## 4 The Stability with Use of the Orbital Elements Determined by Laughlin and Chambers

As mentioned in Setion 2, Marcy et al. (2001) determined the orbital parameters (Table 1) assuming that both planets are Keplerian and the mutual perturbations are not taken into account. With use of the Doppler Observations, which are given in the paper by Marcy et al. (2001) Laughlin and Chambers (2001a) determined the orbital parameters (Table 2) taking account into the mutual perturbations.

In determination of the parameters of Table 2 Laughlin and Chambers assumed that both planets are coplanar and  $\sin i = 1$  is one for both planets. The orbital parameters are osculating elements, which includes the perturbation. They also made orbital determination including  $\sin i$  as fitting parameters.

With use of the orbital parameters of Table 2 we carried out a similar numerical simulation as in Section 2 by changing the mass factor  $m_f$ . Figure 8 show the orbital elements as for  $m_f = 1, 4$ , and 5. We found until  $m_f = 4$  the 2:1 mean motion resonance state is conserved and the corotation time of the pericenters becomes longer as  $m_f$  increases. From around  $m_f = 5$  the 2:1 mean motion resonance is disrupted and the planetary system becomes unstable. Figure 8 shows the instability time. In making Figure 8 we integrated the orbits for  $10^8$  years, which indicates the planetary system ( $1 \leq m_f \leq 4$ ) might plausibly stable over  $10^8$  years.

The origin of the difference of the instability times of Figure 4 and 8 originates from the fact the initial difference of pericenters of planets  $\theta$  by Marcy et al. (2001) is only 3 degrees but that by Laughlin and Chambers (2001a) is as big as 86 degrees. Laughlin and Chambers (2001b) gave another sets of the orbital parameters, in which  $\sin i$  is determined as one of the fitting parameters. In these new orbital parameters, the pericenters are very

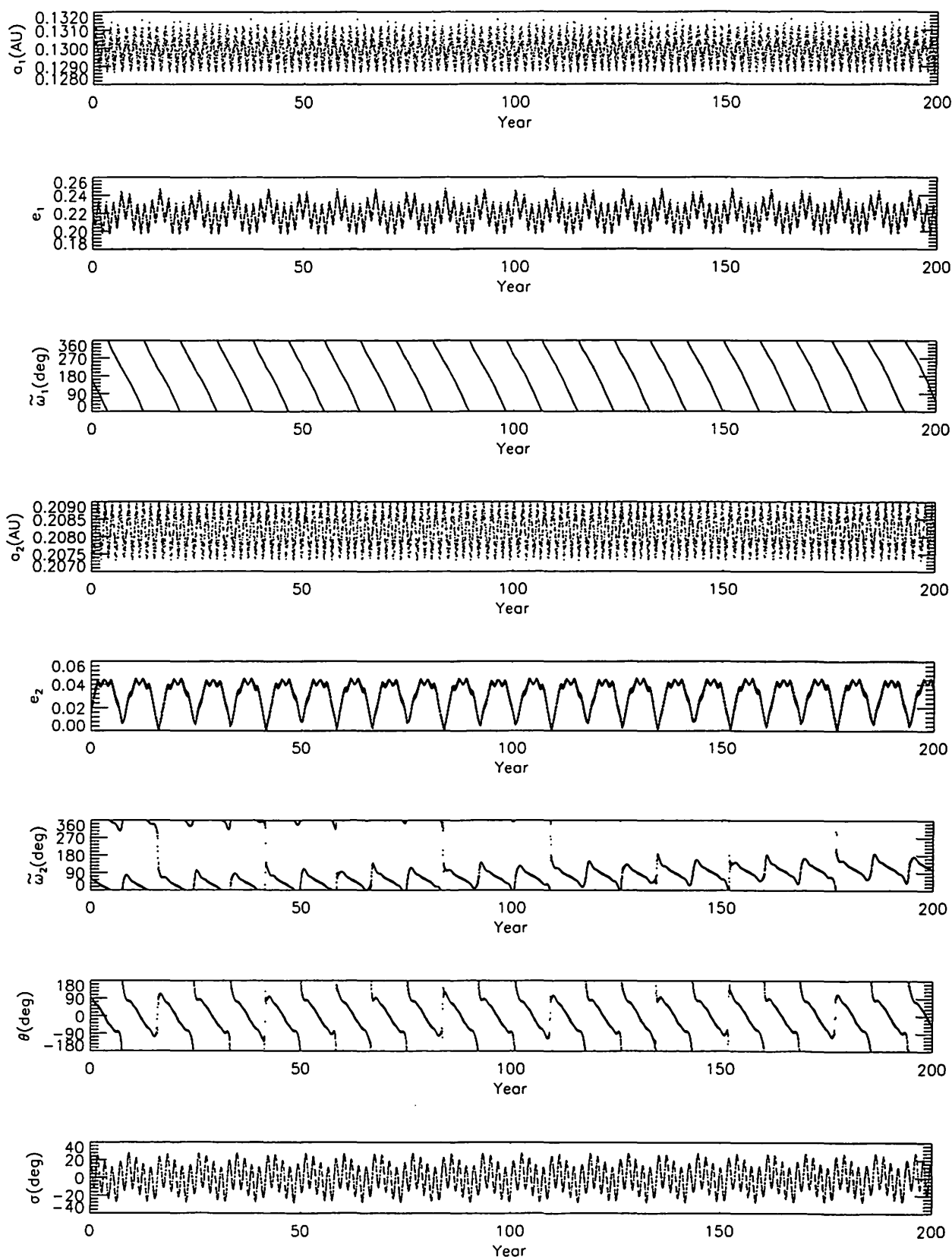
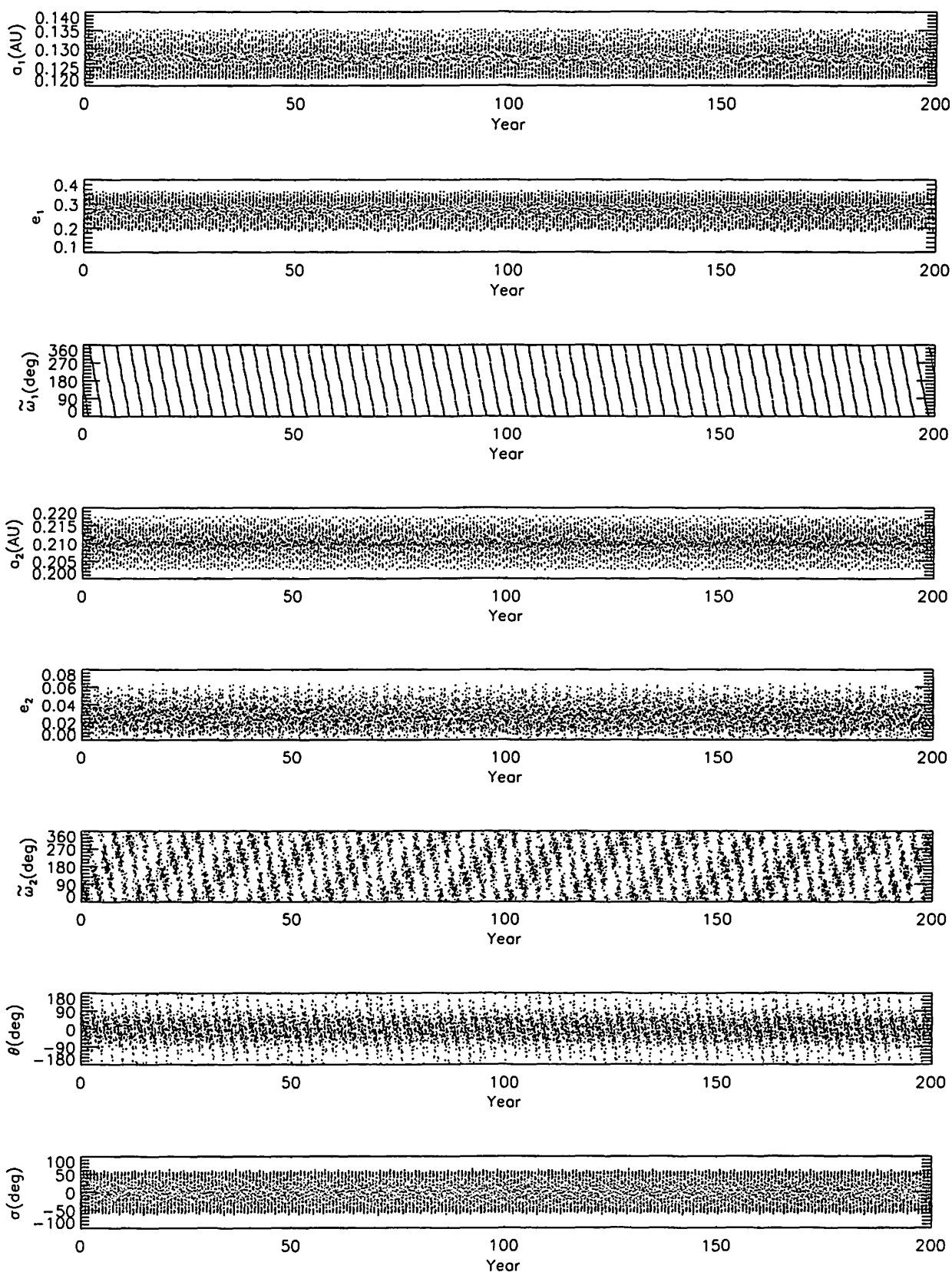
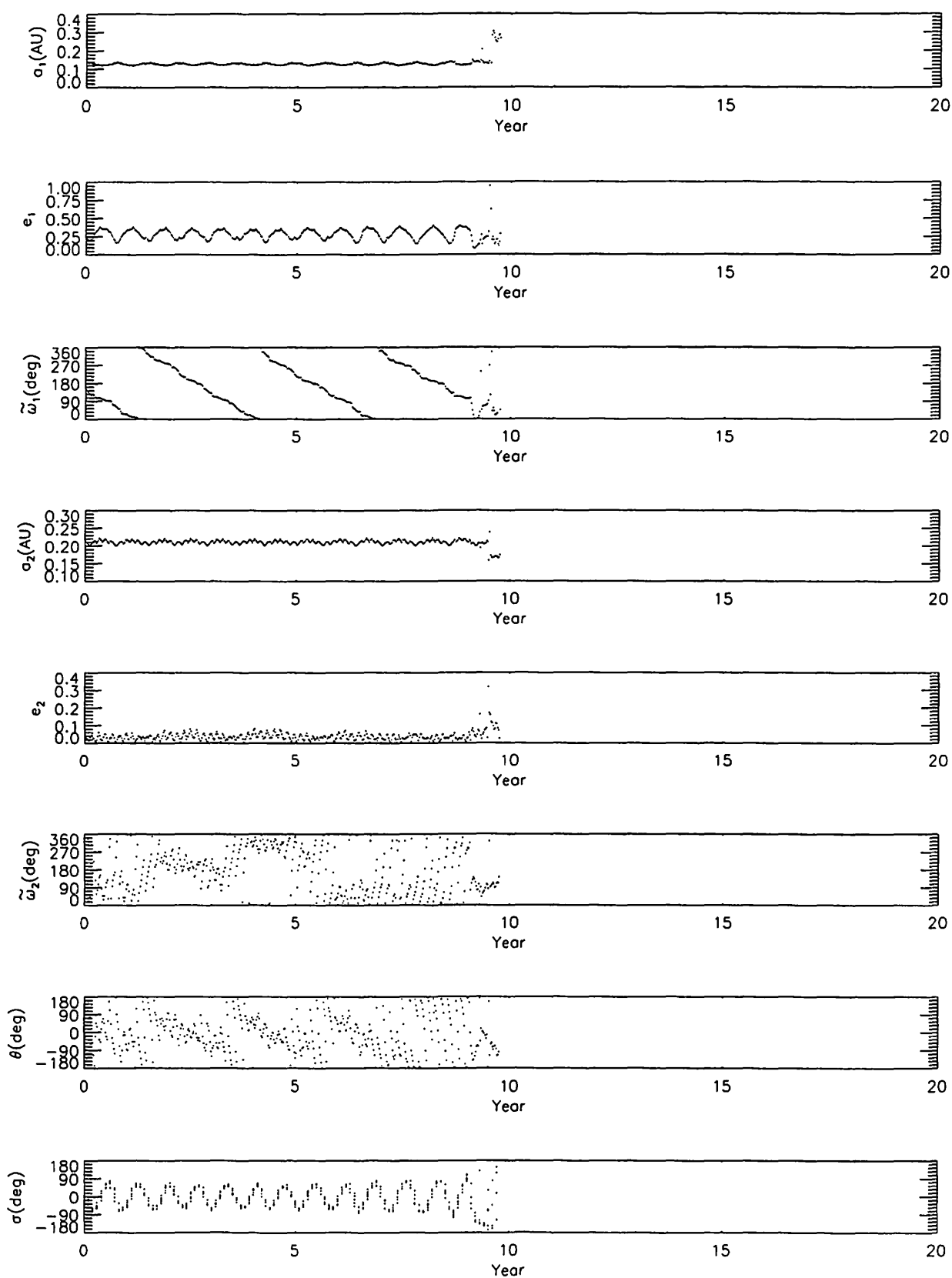


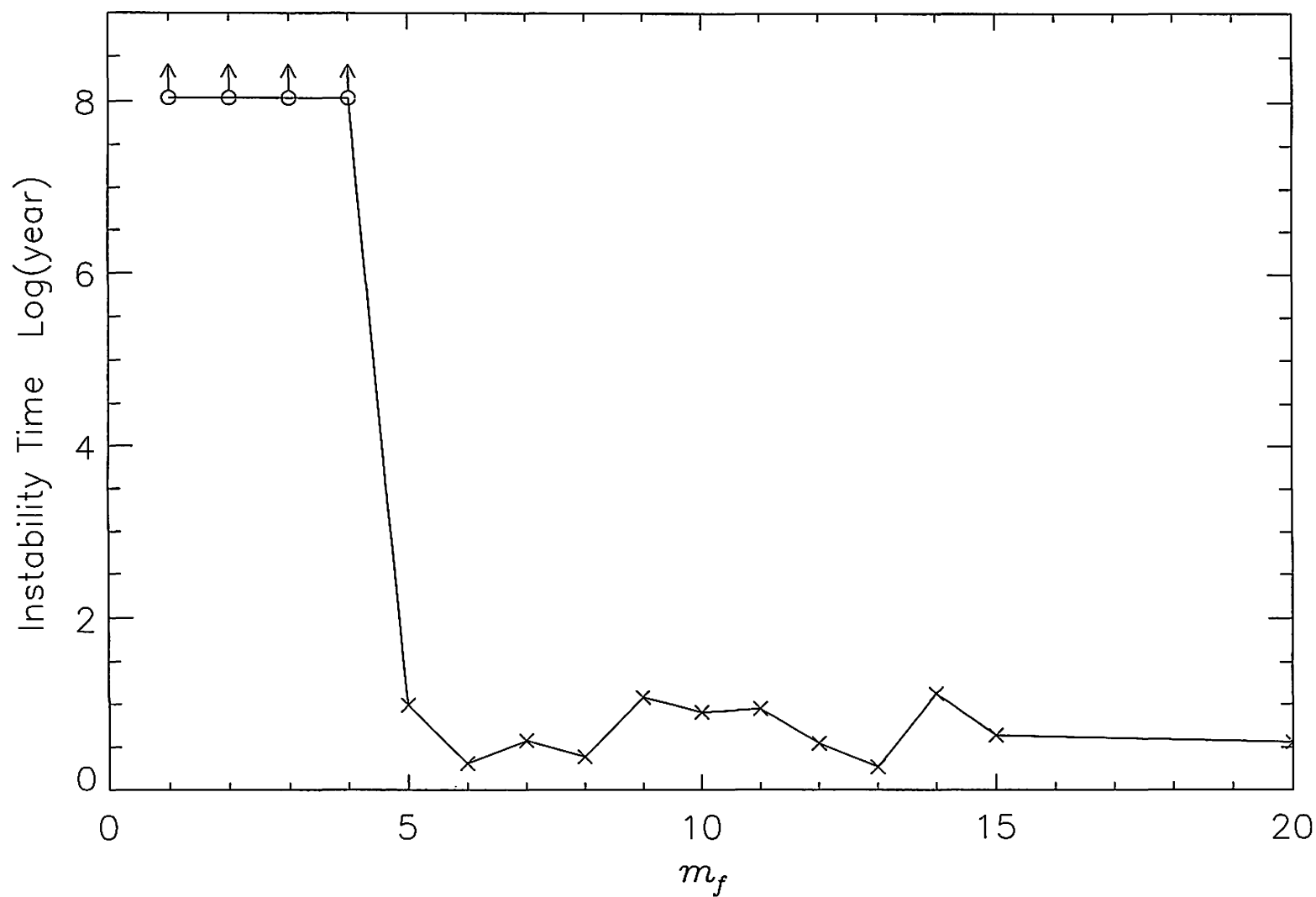
Figure 7-1.—Orbital elements with Laughlin and Chambers's initial values ( $m_f = 1$ ). The explanations for the vertical axes are same as for Figure 1-1.



**Figure 7-2.**—Orbital elements with Laughlin and Chambers's initial values ( $m_f = 4$ ). The explanations for the vertical axes are same as for Figure 1-1.



**Figure 7-3.**—Orbital elements with Laughlin and Chambers's initial values ( $m_f = 5$ ). The explanations for the vertical axes are same as for Figure 1-1.



**Figure 8.**—The instability time for various  $m_f$  for Laugh and Chambers's orbital elements. The open circle means the system is stable till  $10^8$  years and the arrow indicates that the system might plausibly be stable over  $10^8$  years. The cross indicates that the eccentricity of one planet exceeds one.

close by about 10 degrees.

## 5 Secular Perturbation

The Hamiltonian for the coplanar case is

$$F = F(a_1, a_2, e_1, e_2, \varpi_1, \varpi_2, \lambda_1, \lambda_2). \quad (4)$$

Since we discuss the case of 2:1 mean motion resonance case, after short periodic terms the Hamiltonian takes the following form:

$$F^* = F^*(a_1, a_2, e_1, e_2, \sigma, \theta), \quad (5)$$

where  $\sigma = \lambda_1 - 2\lambda_2 + \varpi_1$  ( the critical argument) and  $\theta = \varpi_1 - \varpi_2$ . When both the eccentricities of both planets are small, the Hamiltonian after the elimination of short-periodic terms takes the following form:

$$\begin{aligned} F^* = & +\frac{1}{8}(e_1^2 + e_2^2)(2\alpha D_\alpha + \alpha^2 D_\alpha^2)A_0 \\ & +\frac{1}{4}e_1e_2(2 - 2\alpha D_\alpha - \alpha^2 D_\alpha^2)A_1 \cos \theta \\ & -\frac{1}{2}e_1(4 + \alpha D_\alpha)A_2 \cos \sigma \\ & +\frac{1}{2}e_2(3 + \alpha D_\alpha)A_1 \cos(\sigma - \theta), \end{aligned} \quad (6)$$

where  $\alpha = a_1/a_2$ , and  $A_0 = b_{1/2}^{(0)}$ ,  $A_1 = b_{1/2}^{(1)}$ ,  $A_2 = b_{1/2}^{(2)}$ , which are Laplace coefficients, and  $D_\alpha = \partial/\partial\alpha$ . By using the recurrence formulae for the Laplace coefficients (Brouwer and Clemence, 1961)),  $F^*$  takes the following form:

$$F^* = B_1(e_1^2 + e_2^2) + B_2e_1e_2 \cos \theta + B_3e_1 \cos \sigma + B_4e_2 \cos(\sigma - \theta), \quad (7)$$

where

$$\begin{aligned} B_1 &= \frac{1}{8}\alpha b_{3/2}^{(1)} \\ B_2 &= -\frac{1}{4}\alpha b_{3/2}^{(2)} \\ B_3 &= -\frac{1}{2}(6b_{1/2}^{(2)} + \alpha(b_{3/2}^{(3)} - \alpha b_{3/2}^{(2)})) \\ B_4 &= \frac{1}{2}(4b_{1/2}^{(1)} + \alpha(b_{3/2}^{(2)} - \alpha b_{3/2}^{(1)})) \end{aligned} \quad (8)$$

The degree of freedom of the new Hamiltonian is reduced from four to two. However this Hamiltonian is not integrable. Here we assume the critical argument  $\sigma = 0$ , from which we have  $a_1$  and  $a_2$  are constant. Then the degree of the freedom of the new Hamiltonian is reduced one:

$$F^* = F^*(e_1, e_2, \theta), \quad (9)$$

with conservation of the angular momentum:

$$m_1\sqrt{a_1(1 - e_1^2)} + m_2\sqrt{a_2(1 - e_2^2)} = \text{const} = H. \quad (10)$$

the Hamiltonian  $F^*$  takes the following form:

$$F^* = B_1(e_1^2 + e_2^2) + B_2e_1e_2 \cos \theta + B_3e_1 + B_4e_2 \cos \theta, \quad (11)$$

The last two terms of equation (11) originates from the 2:1 mean motion resonance. The equation of motion from the Hamiltonian (11) is not integrable because of the existence of the two terms. With use of (10) we eliminate  $e$  from  $F^*$  (11) and we have

$$F^* = F^*(e_1, \theta, H) \quad (12)$$

and we can draw the level curves of the Hamiltonian and know the global behavior of  $e_1$  and  $\theta$ . As we see from Figure 1, the eccentricity becomes large and the expanded Hamiltonian (11) is not appropriate for drawing the level curves. Therefore we numerically averaged the original Hamiltonian (4) under the condition of the critical argument  $\sigma = 0$  and the angular momentum conservation (11) and then get numerically the averaged Hamiltonian (12) with the parameter of the angular momentum  $H$ . The technical of this numerical averaging is described in the paper by Kinoshita and Nakai(1985), in which the case of  $\sigma \neq 0$  is also discussed. We draw the contour map of the Hamiltonian (12) taking  $\theta$  as the horizontal axis and  $e_1$  or  $e_2$  as the vertical axis with the parameter  $H$ , which is determined from the initial conditions. Figure 9 shows the contour maps for  $m_f=1,5$ , and 10 and the numerical solutions are plotted on these maps. Since the solution includes the short-periodic terms, the averaged solution is shown in Figure 10. These figures show the good agreement between the numerical solutions and the semi-analytical secular solutions.

## 6 Summary and Discussions

The planetary system of GJ 876 are stabilized by two mechanisms 1) 2:1 mean motion resonance and 2) the alignment of the pericenters. Among two stabilizing mechanism, the necessity of the alignment of pericenters depends mainly on the planetary mass. As the planetary mass increases till  $m_f = 10$ , the duration of the corotation of the pericenters becomes longer for the stabilization. For the planetary system with  $m_f > 10$  the alignment of the pericenter is broken and then the 2:1 mean motion resonance is disrupted and the system becomes unstable.

With use of the Doppler measurements of the main-sequence star GJ 876 ( Marcy et al. 2001) Rivera and Lissauer (2001) determined the orbital parameters for the two planets that accounts for the mutual perturbation between the planets. Their orbital determination method is a Levenberg-Marquardt minimization algorithm, which is one of the nonlinear parameter fitting method and was also used by Laughlin and Chambers ( 2001a and 2001b). They seeked the best fit solutions from the larger parameter space



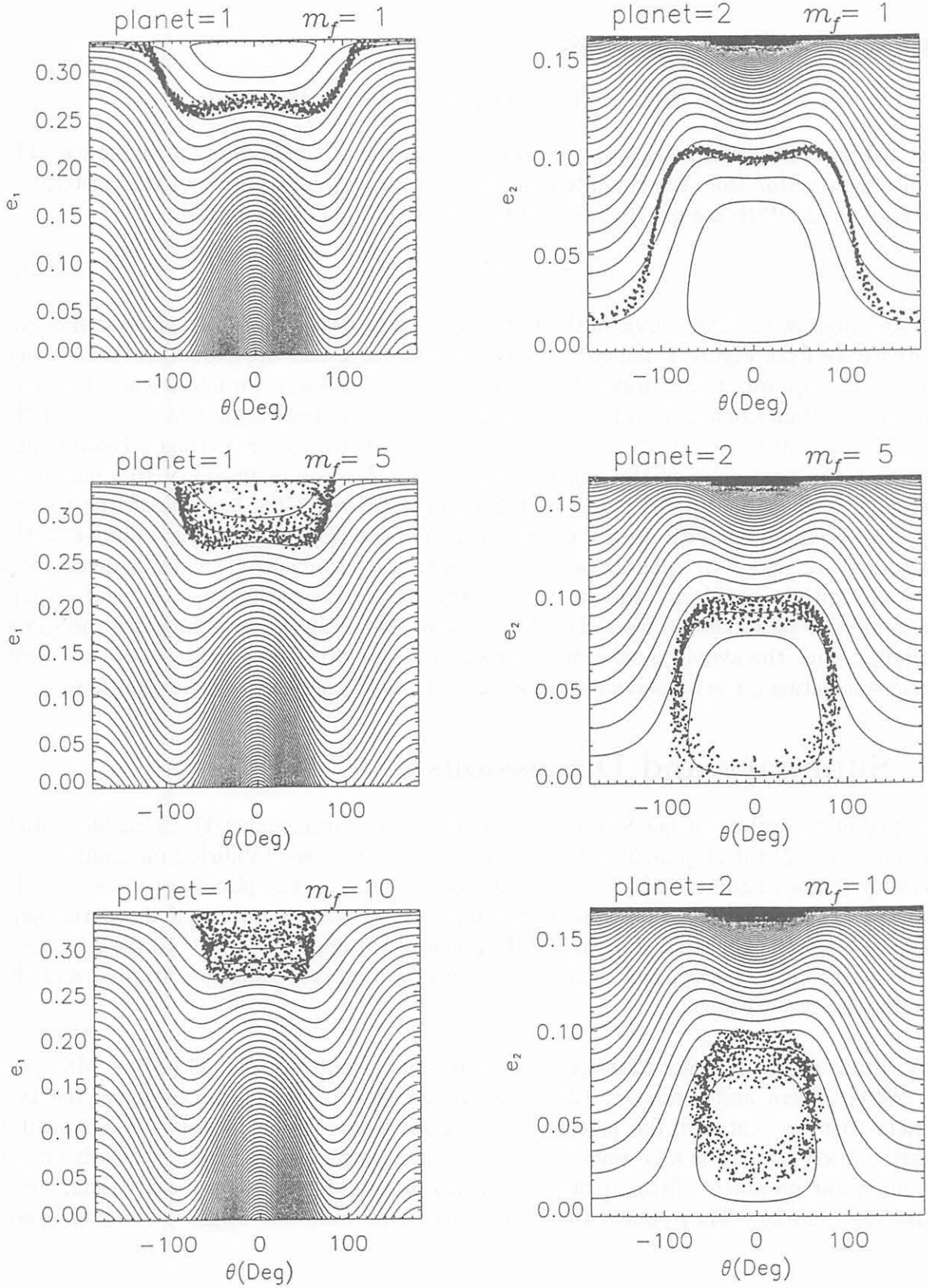


Figure 9.—Equi-Hamiltonian curves and osculating orbital elements for  $m_f=1,5$ , and 10.

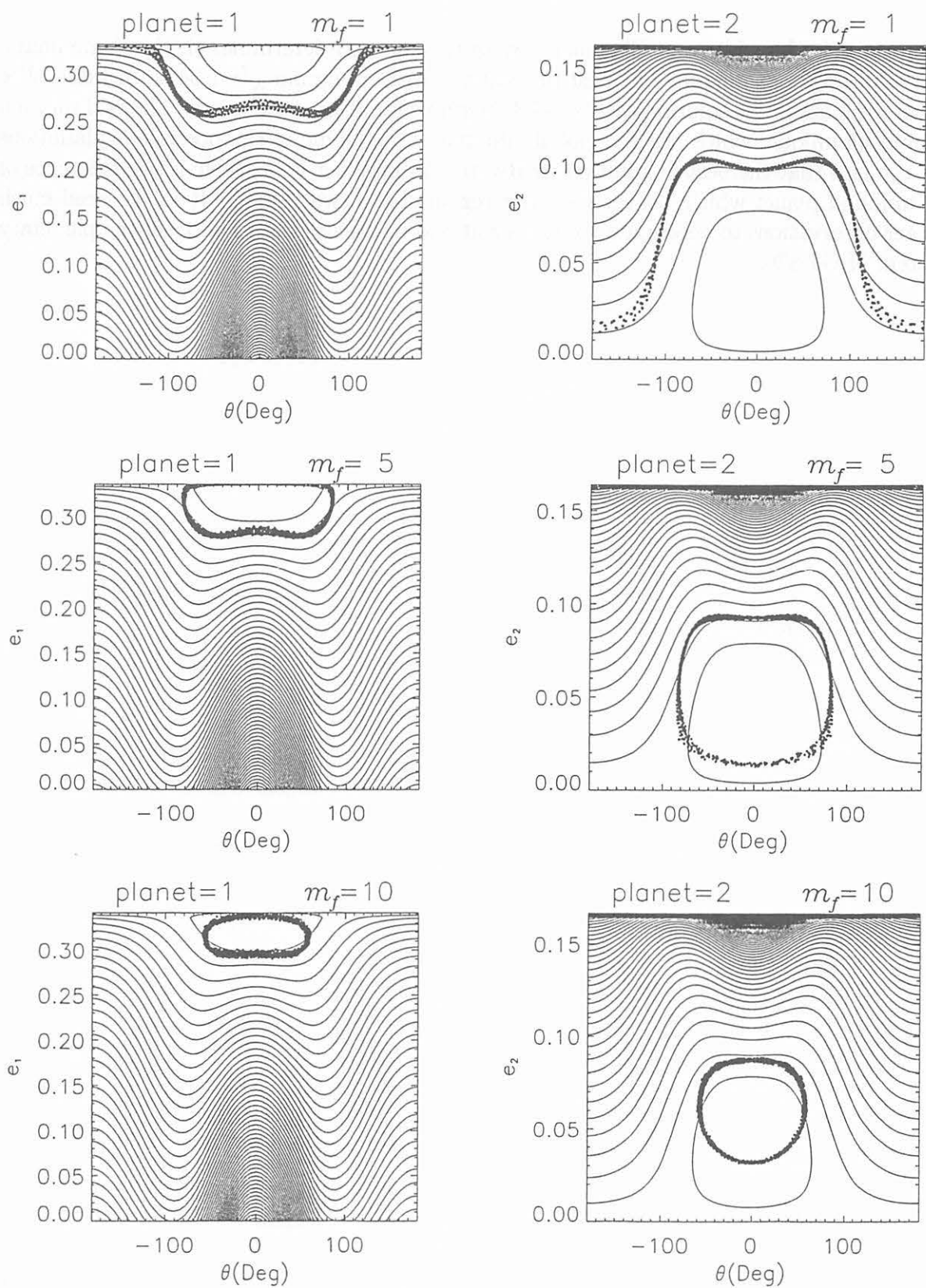


Figure 10.—Equi-Hamiltonian curves and mean orbital elements for  $m_f = 1, 5$ , and 10.

than those by Laughlin and Chambers. Since this method determines the local minimum, the solution is not unique. The best fit solutions with minimum  $\chi^2=1.34$  is unstable. The other local minimum solution with  $\chi^2=1.43$  is stable for at least  $10^8$  years. The real solution should be unique, which has the global minimum of  $\chi^2$ . One of the possible explanations for the fact that the best fit solution by Rivera and Lissauer is unstable is the existence of an unfound planet which locates the outer region of the two planets. If so, we need much longer observations to determine the true orbit and to discuss the stability of the planetary system of GJ 876.

## references

- Brouer, D. and Clemence, G.M. 1961, *Methods of Celestial Mechanics*, Academic Press.
- California & Carnegie Planet Search, 2002, <http://exoplanets.org/almanacframe.html>
- Kinoshita, H. & Nakai, H. 1985, *Celest. Mech.*, **36**, 391-407.
- Kinoshita, H. & Nakai, H. 2001, *PASJ*, **53**, L25-L26.
- Kinoshita, H. & Nakai, H. 2002a, in *Proc. 202th IAU Symp.*, in press.
- Kinoshita, H. & Nakai, H. 2002b, in *Proc. 33rd Symp. Celest. Mech.*, 166-179.
- Laughlin, G. & Chambers, J. E. 2001a, preprint.
- Laughlin, G. & Chambers, J. E. 2001b, *ApJ*, **551**, L109-L113.
- Marcy, G. W., Butler, R. P., Fischer, D., Vogt, S. S., Lissauer, J. J., & Rivera, E. J. 2001, *ApJ*, **556**, 296-301.
- Nakai, H. & Kinoshita, H. 2000, in *Proc. 32th Symp. Celest. Mech.*, 206-215.
- Nakai, H. & Kinoshita, H. 2002, in *Proc. 33rd Symp. Celest. Mech.*, 180-188.
- Rivera, E. J. & Lissauer, J. J. 2000, *ApJ*, **530**, 454-463.
- Rivera, E. J. & Lissauer, J. J. 2001, *ApJ*, **558**, 392-412.

# 「Schwarzschild 解」の意味するもの

## A comprehension for the solution due to Schwarzschild

井上 猛 (京都産業大学)

T. Inoue

Kyoto Sangyo University

**Abstract.** It is widely accepted that the general theory of relativity established by Einstein (1915) solved totally the problem of the excess advance which exists in the longitude of the perihelion of Mercury. This excess advance was found by Le Verrier (1859) when he tried to adjust his theory of the motion of Mercury to the observations of the passage of Mercury in front of the disk of the Sun. The same excess advance was equally necessary in the theory of Newcomb (1895). That is, Newcomb confirmed the Le Verrier's result.

The present author pointed out that the theory of Le Verrier for the motion of Mercury contains tiny errors and that because of this he was not able to well predict the passage phenomena of Mercury (Inoue, 1992). We revealed, in the case of Newcomb, that his theory also contains the same kind of errors (Inoue, 2000).

If one corrects the errors, one can obtain perfect theories of the motion of Mercury. This means that there never exist discrepancies between the theory of the motion and the passage observations in nineteen century.

Then, what did Einstein and Schwarzschild solve? There was no problem to be solved ! In order to legitimately understand, we propose the following postulate *"The Schwarzschild's solution gives exactly the motion of the two-body problem in the Newtonian mechanics"*.

Let us describe the system of Schwarzschild by  $(R, \Phi)$  and that of Newton by  $(r, \phi)$ . The equations of the motion for Mercury take the form as follows in the system of Schwarzschild :

$$\frac{d^2 R}{dt^2} - R \left( \frac{d\Phi}{dt} \right)^2 + \frac{\mu}{R^2} = \frac{2m\mu}{R^3} + \frac{3m}{R^2} \left( \frac{dR}{dt} \right)^2 - 2m \left( \frac{d\Phi}{dt} \right)^2$$

$$\frac{1}{R} \frac{d}{dt} \left( R^2 \frac{d\Phi}{dt} \right) = \frac{2m}{R} \frac{dR}{dt} \frac{d\Phi}{dt}$$

For the problem of two bodies, the equations are given by the following forms :

$$\frac{d^2 r}{dt^2} - r \left( \frac{d\phi}{dt} \right)^2 + \frac{\mu}{r^2} = 0$$

$$\frac{1}{r} \frac{d}{dt} \left( r^2 \frac{d\phi}{dt} \right) = 0$$

With the aid of functions  $\rho$  and  $\sigma$ , we are able to combine the two systems  $(R, \Phi)$  and  $(r, \phi)$  as follows :

$$R = r + m\rho = r \left[ 1 + \frac{m}{r} \left\{ \frac{2}{e^2} (1 - \xi) - \frac{3u}{\eta} e \sin f \right\} \right]$$

$$\Phi = \phi + m\sigma = \phi + \frac{m}{p} \left\{ 3 \left( f - \frac{u}{\eta} \xi^2 \right) + \frac{2+2\xi+e^2}{e^2} e \sin f \right\}$$

The detailed explanations of the quantities are given in the text (in Japanese).

The relations satisfy 'the Schwarzschild condition' : when the quantity  $r$  tends to infinity, the difference between  $R$  and  $r$  tends to zero. It is the same for the difference between  $\Phi$  and  $\phi$ .

In the Newtonian mechanics, the elliptic orbit does not move and there is no advance in the longitude of the perihelion  $\varpi$ . Therefore, one revolution after the longitude  $\phi$  takes simply the value :  $\phi = \varpi + 2\pi$ .

In the Schwarzschild system, the circumstances are different. At the moment of start, the quantity  $\Phi$  obtains the value  $\varpi$ . But, so as to just return to the 'perihelion', one should give the following value for the quantity  $\Phi$  :

$$\Phi = \left\{ \varpi + \frac{6\pi m}{p} \right\} + 2\pi .$$

久しく心に留めて居て、人にも知って欲しいと願って居るものがある。それは、私かに尊崇して止まない平山清次先生がお書きになったものである。此処に添付するのは、甚だ御無礼な事かも知れないが「積年の懐いの吐露」として、御寛恕の程を請う次第である。

昭和五年(1930)十月、東京天文臺内に本部を置く日本天文學會が、日本天文學會編輯の下に、本格的な研究誌としての『日本天文學會要報』なるものを世に送り出さむとした。その「第一號」に、理事長の平山清次先生が、「發刊の辭」を述べて居られるのである。此処に全文を引用するのであるが、残念な事に“文字上の制約”から、本文の通りに打つ事が出来ない憾みがある。それ以外は、先生のお書きになったものの通りである。

## 發 刊 の 辭

日本天文學會理事長

理學博士      平   山   清   次

現代の日本の學者が歐文で學術上の論文を書くのは、外國人がそれを読む事を豫期する結果、便利を考へてさうするのであつて、徳川時代の儒者が自己の學識を表したい為めに殊更難解の字句を並列し文章を飾つた事と混同してはならない。英佛獨等の外國語によるよりも望ましいのは 에스ぺ란토 の如き世界語を用ひる事であるが歐米諸國の學者が眞に此事に目醒めぬ限り實行は出來かねる。

國旗は 國家の符徴に外ならないが 言語は事實上民族を結束する締縄であり其表徴である。それであるから民族が其個性を維持せんとする志望を失はない限り言語は尊重されなくてはならない。海外より弘く有益な知識を吸収する為めに外國文を読む事を怠ってはならぬが、然し其為めに自國語を蔑視したりしてはならない。

それであるから學術上の論文は、如何に世界的のものであつても、それを自國語で書くのは當然であつて、更にそれを外國文に綴つて弘く海外に發表するのが最も適當な處置である。事情によりそれを自國語だけに止めるのは自由であるが、外國文で發表して自國語で發表しないのは大なる誤と言はなければならぬ。

のみならず外國文と自國文と内容は同一でも著者の意想が自國文の方に適切に表はれるのは當然である。従つて其方に原作としての價値の多い事は無論である。

日本天文學會の役員の主張は實に此點に在る。其結果、定會の議決を経て實現したのは此要報である。事は小さいながらも思慮は十分に深く實行の意志は強固である。筆者は理事長として茲に役員を代表し此主意を會員諸氏並に關係ある識者に傳へる事を大なる光榮とするものである。

日本の天文學の前途は有望である。それが吾々の此新らしい計畫によって順當なる發達を遂げ、深遠なる宇宙と靈妙なる自然とが一層明確に示さるゝ事を得れば誠に幸である。

1. Le Verrierの研究(1859)に端を発する「水星近日点黄経 $\omega$ に於ける余剰の永年変化の問題」は、Einsteinの一般相対性理論(1915)に依って、完全解決を見た事になって居る。就中、重力場方程式に対するSchwarzschildの美麗な解(1916)は、Eddingtonの著書(1924)等を通じて広く知られて居る処である。「当該問題解決法」は、“現代の常識”である。

斯うした折りに「Le Verrierの研究には誤りがあった。その故に《余剰の永年変化》の出現を見る事になった」との結論を得た(Inoue, 1992)。その後の調査で「Newcombの研究(1895)にも、Le Verrier同様の誤りが在る。その為に一世紀当り43秒角の不等性が生ずる事になって居る」のが明らかになった(井上、2000)。

我々の主張をその儘に受け容れるならば、『件の問題そのものが存在しなかった』事になる。これは“現代の常識”に反する。然し、『我々の結果』を否定しなければならない理由を、我々は見つける事が出来ないのである。そこで『Newton力学に於ける二体問題が表わす処のものは、Einstein理論に於ける Schwarzschild解が与えるものである』と捉える事にする。斯うする事は可能である。此の事を以下で見て行く。

実は、十二年ほど前の事になるが、次の様な問題を「理論天文学」の後期試験に課して居る。

◎平面運動であるとして、水星近日点黄経に於ける余剰前進の問題に関連させて以下を問う。

(甲) Newton力学に於ける二体問題は次の形に書ける：

$$\frac{d^2 r}{dt^2} - r \left( \frac{d\phi}{dt} \right)^2 = - \frac{\mu}{r^2} \quad , \quad \mu \equiv G (m_{\odot} + m^*) \quad ;$$

$$\frac{d}{dt} \left( r^2 \frac{d\phi}{dt} \right) = 0 \quad : \quad r^2 \frac{d\phi}{dt} = h \quad .$$

(乙) Schwarzschildの謂わば一体問題は次の様に与えられる：

$$\frac{d^2 R}{d\tau^2} - R \left( \frac{d\Phi}{d\tau} \right)^2 = - \frac{m c^2}{R^2} - 3 m \left( \frac{d\Phi}{d\tau} \right)^2 \quad ,$$

$$\frac{d}{d\tau} \left( R^2 \frac{d\Phi}{d\tau} \right) = 0 \quad : \quad R^2 \frac{d\Phi}{d\tau} = H \quad .$$

問1. 上で、Newtonの場合とSchwarzschildの場合とで文字の書き分けを行なった。その必要性の有無を、理由を付けて述べよ。

問2. 此の問題は、その総てが Le Verrierに依る、水星の太陽面通過に対する条件方程式から出発して居る。此の事に付いて知れる処を述べてみよ。

問3. Schwarzschildの場合には、SiriusやCapella等の連星系に於ける運動をどの様に理解するのが良いか？ 考えを述べよ。

問4. 此の問題は、一般相対論との関連に於て、どの様に扱われるのが望ましいと考えるか？ 考える処を記してみよ。

$$\frac{dt}{d\tau} = \left( 1 - \frac{2m}{R} \right)^{-1/2} \quad , \quad \mu : m c^2 \quad . \quad (90Z14V \text{ のうへ たけし})$$

当時は「Le Verrierの研究に誤りがある」等と云った認識は、皆無であった。然し『文字書き分けの必要性』には考えが及んで居た様である。これが有ったればこそ「Le Verrier理論に誤りがあった」のに気付く事も出来れば「余剰前進問題の解決」に取り組んで行く事も出来たのであった。

上で、 $r$ 及び $\phi$ は動径および経度を表わす。 $R$ 及び $\Phi$ は、それらに対応する「座標」である。Newton力学で言う絶対時間  $t$  は、Schwarzschildの系では「座標時」に相当する。

量 $\tau$ は「固有時」を表わす。量 $c$ は「光の速さ」を表わし、量 $m$ は、長さの単位を有する“微小な定数”を表わす。量 $G$ は万有引力定数を表わし、量 $m_{\odot}$ は太陽の質量を、また量 $m^*$ は水星(惑星)の質量を表わす。

此の時に、量 $\mu$ を量 $mc^2$ に等しいと置くか否かに付いては、大いに「議論の余地」の有る処と考へて居る。先に引用した Eddingtonの場合には量 $mc^2$ は、 $Gm_{\odot}$ に等しいと捉へるのが妥当と考えられる(同著:81頁~83頁)。ところが別の著者の場合には、此の量 $mc^2$ は「量 $\mu$ 即ち量 $G(m_{\odot}+m^*)$ に等しいと理解すべきである」と考えられるのである(Brumberg, 1991:1頁, 2頁, 5頁, 76頁, 82頁)。

水星の場合には、質量 $m^*$ が、太陽の質量 $m_{\odot}$ の $1.7 \times 10^{-7}$ 倍でしかないのであるから「量 $m^*$ の存在を無視する」と云う見方が有る様である。然し、それは的を射たものとは言ひ難い。何となれば、我々が問題として居る「余剰の前進量 $\delta \varpi_{(s)}$ 」の大きさは「太陽の重力下で水星が公転運動する量」の $8 \times 10^{-8}$ 倍でしか無いからである。詰り、 $1.7 \times 10^{-7}$ の大きさの量を無視するのならば、 $8 \times 10^{-8}$ の大きさの物は更に積極的に無視しなければならない事になるからである。

此の問題は、真剣に取り組まなければならないものなのであるが、多くを考へる事なく此処では、「単純に、量 $\mu$ は量 $mc^2$ に等しい」と置いて、以下の議論に移る事にする。

2. 先ずは、Einstein理論に於ける Schwarzschild解が与える「惑星運動の基本方程式」の導出を試みる。これを、行なうのに Eddingtonの著書を参考にした。同著の86頁及び87頁から一部を抜き書きしてみる。此処で、 $ds$ は $cd\tau$ の事である。

Differentiating with respect to  $\phi$ , and removing the factor  $\frac{d}{d\phi} \frac{1}{r}$ ,

$$\frac{d^2}{d\phi^2} \frac{1}{r} + \frac{1}{r} = \frac{m}{h^2} + \frac{3m}{r^2} \quad (39.61),$$

with

$$r^2 \frac{d\phi}{ds} = h \quad (39.62).$$

Compare these with the equations of a Newtonian orbit

$$\frac{d^2}{d\phi^2} \frac{1}{r} + \frac{1}{r} = \frac{m}{h^2} \quad (39.71),$$

with

$$r^2 \frac{d\phi}{dt} = h \quad (39.72).$$

In (39.61) the ratio of  $\frac{3m}{r^2}$  to  $\frac{m}{h^2}$  is  $\frac{3h^2}{r^2}$ , or by (39.62)

$$3 \left( r \frac{d\phi}{ds} \right)^2.$$

For ordinary speeds this is an extremely small quantity — practically three times the square of the transverse velocity in terms of the velocity of light. For example, this ratio for the earth is .0000 0003. In practical cases the extra (以下 87頁) term in (39.61) will represent an almost inappreciable correction to the Newtonian orbit (37.71).

前節で指摘した様に、只今の Eddingtonの謂が一般には「何の疑問も抱かれる事なく」受け入れられて居るのである。此の方程式群が導かれる基となったものは、当然の事乍ら「Schwarzschildの計量」及び「測地線の方程式を解いて得られる積分」である。そこで再び同著から少しく引用して置く事にしよう。但し、「平面問題」に限定してのものとする。

$$r = 1 - \frac{2m}{r} \quad (38.7),$$

$$ds^2 = -r^{-1}dr^2 - r^2d\phi^2 + rc^2dt^2 \quad (38.8).$$



以上は、85頁に所載のものである。此処で Eddingtonは、光の速さ  $c$  を「速さの単位」に選んだので、(38・8)式に記した  $c$  は、彼の表式には現われては居ない。続いて、86頁から積分表式を引用するが、此処でも我々流に書き換えたものを記す事にする。

$$r^2 \frac{d\phi}{ds} = h \quad (39\cdot41) ,$$

$$\frac{dt}{ds} = \frac{1}{c r} \quad (39\cdot42) .$$

「基本となる方程式群は 既にEddingtonが与えて居るではないか」と、改めての導出を訝る向きもあるであろうが、我々は『文字の書き分けの必要性』を主張して居るのであるから、此処は「いちから」見て置きたいのである。

此の時に問題となるのが「Newton力学に於ける絶対時間  $t$ 」と「Einstein理論に於ける座標時  $t$ 」の間の関係である。此処では、「同一の文字“ $t$ ”」を「考えも無しに用いて居る」かの印象を与えるかも知れない。然し、『これらを同一視しても良い』と云うのが既に広く行なわれて居るので、此処はそれに依拠した迄である。

先ずは(38・7)式(38・8)式を、我々の文字記号で書いて置く。

$$(1) \quad r = 1 - \frac{2m}{R} .$$

$$(2) \quad ds^2 = c^2 d\tau^2 = r c^2 dt^2 - r^{-1} dR^2 - R^2 d\Phi^2 .$$

同様にして、(39・41)式(39・42)式も書き換えて置く。

$$(3) \quad R^2 \frac{d\Phi}{d\tau} = H \quad ; \quad R^2 \frac{d\Phi}{dt} r^{-1} = H .$$

$$(4) \quad \frac{dt}{d\tau} = \frac{1}{r} = r^{-1} .$$

以下に計算をして行く。

$$c^2 = r c^2 \frac{dt^2}{d\tau^2} - r^{-1} \frac{dR^2}{d\tau^2} - R^2 \frac{d\Phi^2}{d\tau^2}$$

$$c^2 = r^{-1} c^2 - r^{-3} \frac{dR^2}{dt^2} - \frac{H^2}{R^2}$$

此の表式を、独立変数  $t$  で微分する。然る後に各項を、因子  $-2r^{-3} \frac{dR}{dt}$  で整除すれば次の表式に到達するであろう。

$$\frac{d^2 R}{dt^2} - \frac{H^2}{R^3} r^3 + \frac{m c^2}{R^2} r - \frac{3m}{R^2} \left( \frac{dR}{dt} \right)^2 r^{-1} = 0$$

運動方程式の体裁を整える為に、(3)式を用いて積分定数の  $H$  を消去した形に導く。

$$(5) \quad \frac{d^2 R}{dt^2} - R \left( \frac{d\Phi}{dt} \right)^2 r + \frac{m c^2}{R^2} r - \frac{3m}{R^2} \left( \frac{dR}{dt} \right)^2 r^{-1} = 0 .$$

以下では、先程「問題とした」量  $m c^2$  を単純に量  $\mu$  に等しいと置く事にする。量  $m$  は充分に「小さい」として、これの自乗以上の量は総て無視する事とする。斯くして、次の Einstein理論に於ける「動径」に関する運動方程式が得られる。

$$(6) \quad \frac{d^2 R}{dt^2} - R \left( \frac{d\Phi}{dt} \right)^2 + \frac{\mu}{R^2} = \frac{2m\mu}{R^3} + \frac{3m}{R^2} \left( \frac{dR}{dt} \right)^2 - 2m \left( \frac{d\Phi}{dt} \right)^2 .$$

これに対応する「偏角(経度)」に関する運動方程式は、(3)式を次の形に書いた後に時間  $t$  で両辺を微分して導く事が出来る。

$$R^2 \frac{d\Phi}{dt} = r H = H \left( 1 - \frac{2m}{R} \right)$$

$$(7) \quad \frac{1}{R} \frac{d}{dt} \left( R^2 \frac{d\Phi}{dt} \right) = \frac{2m}{R} \frac{dR}{dt} \frac{d\Phi}{dt} .$$

上記の(6)式および(7)式の表わすものが「Einstein理論に於ける平面惑星運動に対する基本方程式」である。勿論、此処で対象として居るのは、太陽と一個の惑星のみである。これらの表式が、只今の立場で惑星運動を論ずる際の基本方程式である事を確認するには以下の文献を参照すれば充分であろう。

Brumberg : *Essential Relativistic Celestial Mechanics* p. 82 (3.1.49)式

平山清次 : 天体力学 一般摂動論 p. 71 (1.3)式

此の「二体問題」に対する「Newton力学に於ける基本方程式」は、「文字書き分け」の主張に基づいて、次の形に書かれる事になる。

$$(8) \quad \frac{d^2 r}{dt^2} - r \left( \frac{d\phi}{dt} \right)^2 + \frac{\mu}{r^2} = 0 ,$$

$$(9) \quad \frac{1}{r} \frac{d}{dt} \left( r^2 \frac{d\phi}{dt} \right) = 0 .$$

初等的な事柄を長々と述べて来たが、我々の言わむとする処を此処で強調して置こう。

「Einstein理論に於ける基本方程式」(6)式および(7)式が表わす惑星運動は「Newton力学に於ける基本方程式」(8)式および(9)式が与えるものに完全に一致する。

3. 問題の、「Einstein理論の系を記述する量 (R, Φ)」と「Newton力学の系を記述する量 (r, φ)」とを次の関係で結び付ける。

$$(10) \quad R = r + m \rho ,$$

$$(11) \quad \Phi = \phi + m \sigma .$$

此処に、量 ρ および量 σ は「以下の条件」を満たすべく決定されなければならない未知量である。

(6)式および(7)式を記述する量 R および量 Φ を、(10)式および(11)式の関係に依って量 r および量 φ に置き換えるなら、直ちに(8)式および(9)式が得られる

未知量の ρ および σ を求めるのに、方程式(6)式および(7)式を積分した形を用いる。そこで、先ずは(7)式から考える事にしよう。これは、次の形に書き換える事が出来る。

$$\left( R^2 \frac{d\Phi}{dt} \right)^{-1} \frac{d}{dt} \left( R^2 \frac{d\Phi}{dt} \right) = \frac{2m}{R^2} \frac{dR}{dt}$$

従って、容易に積分する事が出来る。結果は以下の如しである。

$$(12) \quad R^2 \frac{d\Phi}{dt} = H \cdot \exp \left( -2 \frac{m}{R} \right) , \quad (H : \text{積分定数}) .$$

これを、(6)式に代入し、量 m の自乗以上の微小量を見捨てる事により次の表式が得られる。

$$\frac{d^2 R}{dt^2} - \frac{H^2}{R^3} + \frac{\mu}{R^2} = \frac{2m\mu}{R^3} + \frac{3m}{R^2} \left( \frac{dR}{dt} \right)^2 - \frac{6mH^2}{R^4} .$$

因子に微小量 $m$ を有する項では、「Newton力学」に於ける「二体問題の解表式」を、量 $R$ および量 $\Phi$ で書き表わしたものをを用いても構わないであろう。そこで、エネルギー積分の次の形のものを取り入れる。

$$\frac{1}{2} \left\{ \left( \frac{dR}{dt} \right)^2 + \left( R \frac{d\Phi}{dt} \right)^2 \right\} - \frac{\mu}{R} = E , \quad (E : \text{積分定数}) .$$

斯くして、最上段の表式は、容易に積分可能な次の形へと導かれる事になった。

$$(13) \quad \frac{d^2 R}{dt^2} - \frac{H^2}{R^3} + \frac{\mu}{R^2} = \frac{8m\mu}{R^3} + \frac{6mE}{R^2} - \frac{9mH^2}{R^4} .$$

積分の結果は以下の通りである。此処で、量 $C$ は積分定数である。

$$(14) \quad \left( \frac{dR}{dt} \right)^2 + \frac{H^2}{R^2} - \frac{2\mu}{R} = - \frac{8m\mu}{R^2} - \frac{12mE}{R} + \frac{6mH^2}{R^3} + C .$$

これに、先に設定した関係式(10)式(11)式を代入して行く訳である。この時、次の表式を用いるのは当然の事である。

$$(15) \quad R = r \left( 1 + \frac{m\rho}{r} \right) ,$$

$$(16) \quad \frac{dR}{dt} = \frac{dr}{dt} + m \frac{d\rho}{dt} .$$

以下の計算の流れを容易にする目的で、「Newton力学」に於ける種々の関係式を一覧の形で記して置く事しよう。

$$(17) \quad \frac{1}{2} \left\{ \left( \frac{dr}{dt} \right)^2 + \left( r \frac{d\phi}{dt} \right)^2 \right\} - \frac{\mu}{r} = E , \quad [\text{エネルギー積分}] .$$

$$(18) \quad r^2 \frac{d\phi}{dt} = h ; \quad (h : \text{積分定数}) \quad [\text{角運動量積分}] .$$

$$(19) \quad E = - \frac{\mu}{2a} ; \quad h = \sqrt{\mu p} , \quad p \equiv a(1 - e^2) .$$

此処で、量 $a$ および量 $e$ は楕円軌道の長半径および離心率である。

準備が整ったので(15)式(16)式を(14)式に代入する。この時に定数の間に以下の等式の成立を要請する。

$$(20) \quad C = 2E , \quad H = h .$$

その結果、共通因子の $m$ で整除の後に、未知量 $\rho$ に対する微分方程式として次の形のものが得られる。

$$(21) \quad \frac{dr}{dt} \cdot \frac{d\rho}{dt} - \frac{\mu p}{r^3} \rho + \frac{\mu}{r^2} \rho = - \frac{4\mu}{r^2} - \frac{6E}{r} + \frac{3\mu p}{r^3} .$$

再び、楕円軌道に対して成立する関係式を書き出して置こう。

$$(22) \quad \frac{dr}{dt} = \frac{na}{\eta} e \sin f , \quad \frac{d\phi}{dt} = \frac{df}{dt} = \frac{h}{r^2} = \frac{n\xi^2}{\eta^3} .$$

此処で、量 $f$ は真近点離角である。更に、次の略記号を用いた。

$$(23) \quad n \equiv \sqrt{\left\{ \frac{\mu}{a^3} \right\}} , \quad \eta \equiv \sqrt{1 - e^2} ; \quad \xi \equiv \frac{p}{r} = 1 + e \cos f .$$

初等的な関係式等を並べて来たが、結局は未知量 $\rho$ に対する方程式は次の形に導かれる。

$$(24) \quad \frac{d\rho}{df} \cdot e \sin f + (1 - \xi) \rho = \frac{3\eta^2}{\xi} - 4 + 3\xi .$$

此の方程式を満たす量 $\rho$ は、「人工系の方法(定数変化法)」に依っても解く事が出来て次の形を取る。結果の正当性は、これを(24)式に代入してみれば容易に確かめられる処である。此処で、離心近点離角 $u$ に登場して貰う必要があった。

$$(25) \quad \rho = \frac{2}{e^2} (1 - \xi) - \frac{3u}{\eta} e \sin f \quad .$$

次には、只今の「動径」方向の關係を用いて、「経度」方向の差 $\sigma$ を求めて行く。

$$\begin{aligned} R^2 \frac{d\Phi}{dt} &= (r + m\rho)^2 \cdot \left( \frac{d\phi}{dt} + m \frac{d\sigma}{dt} \right) = \\ &= r^2 \frac{d\phi}{dt} + 2rm\rho \frac{d\phi}{dt} + r^2 m \frac{d\sigma}{dt} + \dots = \\ &= h \cdot \exp \left( -\frac{2m}{R} \right) = h - \frac{2hm}{R} + \dots = \\ &= r^2 \frac{d\phi}{dt} - \frac{2np}{\eta^3} m\xi + \dots \end{aligned}$$

因子 $m$ の一次の項を等置する事に依って、未知量 $\sigma$ に対する方程式を導く事が出来る。

$$\begin{aligned} r^2 \frac{d\sigma}{dt} + 2 \frac{\rho}{r} r^2 \frac{d\phi}{dt} &= - \frac{2np}{\eta^3} \xi \quad , \\ r^2 \frac{d\phi}{dt} &= na^2 \eta \quad ; \quad dt = \frac{\eta^3}{n\xi^2} df \quad . \end{aligned}$$

此処でも、真近点離角 $f$ に関する微分方程式に書いて置く。

$$(26) \quad p \frac{d\sigma}{df} = - 2\xi (1 + \rho) \quad .$$

これも容易に解く事が出来て、次の形の解が得られる。

$$(27) \quad p\sigma = 3 \left( f - \frac{u}{\eta} \xi^2 \right) + \frac{1}{e^2} (2 + 2\xi + e^2) e \sin f \quad .$$

以上で、目的は総て達成された。即ち「Einstein理論の系」と「Newton力学の系」とを結び付けるべく設定した関係式(10)式(11)式が、滞り無く求められたのである。これらを書いてみれば以下の如しである。

$$(28) \quad R = r + m\rho = r \left[ 1 + \frac{m}{r} \left\{ \frac{2}{e^2} (1 - \xi) - \frac{3u}{\eta} e \sin f \right\} \right] \quad ,$$

$$(29) \quad \Phi = \phi + m\sigma = \phi + \frac{m}{p} \left\{ 3 \left( f - \frac{u}{\eta} \xi^2 \right) + \frac{2 + 2\xi + e^2}{e^2} e \sin f \right\} \quad .$$

これで見れば明らかな様に、 $r \rightarrow \infty$ および $p \rightarrow \infty$ に際して、 $R \rightarrow r$ および $\Phi \rightarrow \phi$ が実現し所謂「Schwarzschild条件」を満たして居るのが判る。

4. 量 $R$ が「動径」を表わすものならば、「近日点」に於ては、時間に依る「微分」商は「零」に等しくならなければならない。これを確かめる目的で、(22)式に与えた関係式を改めて書き出して置く。

$$\frac{dr}{dt} = \frac{na}{\eta} e \sin f \quad , \quad \frac{df}{dt} = \frac{n\xi^2}{\eta^3} \quad (22)$$

更に、離心近点離角 $u$ に関する「微分」商も書いて置く。

$$(30) \quad \frac{du}{dt} = \frac{n\xi}{\eta^2} \quad .$$

$$\begin{aligned} (31) \quad \frac{dR}{dt} &= \frac{na}{\eta} e \sin f + m \left\{ \frac{2}{e^2} \cdot \frac{n\xi^2}{\eta^3} e \sin f + \right. \\ &\quad \left. - \frac{3}{\eta} \cdot \frac{n\xi}{\eta^2} e \sin f - \frac{3u}{\eta} \cdot \frac{n\xi^2}{\eta^3} e \cos f \right\} = 0 \quad . \end{aligned}$$

此処で、量  $f$  や量  $u$  は「Newton力学」に於ける「二体問題の解」を記述する量なのであるから、 $f=0$  では  $u=0$  である。此の時、水星（惑星）は「近日点」に在って、動径  $r$  は「近日点距離」を与える。上式から明らかな様に、量  $R$  も、此の時点で「近日点距離」を与える事になって居る。

「Newton力学」に於ける「二体問題」では、一公転の後には  $f=2\pi$  となり  $u$  も  $2\pi$  に等しくなる。当然の事ながら、動径  $r$  は、再び「近日点距離」を与える事になって居る。然し、量  $R$  の方は、上記(31)式を成立させるには、量  $f$  は  $2\pi$  とは少しばかり異なる値を取らなければならない事になって居る。そこで、量  $m$  の大きさの微小量  $\delta_i$  の存在を仮定して、(31)式の成立を図る。量  $m$  の一次の大きさ限定して考えるのは当然の事である。

$$f \equiv 2\pi + \delta_i, \quad ,$$

$$\frac{n a e}{\eta} \delta_i + m \left\{ O(\delta_i) - O(\delta_i) - \frac{6\pi}{\eta} \cdot e \cdot \frac{n(1+e)^2}{\eta^3} \right\} = 0 \quad ;$$

$$(32) \quad \delta_i = \frac{3m}{p} \cdot \frac{(1+e)^2}{\eta} \times 2\pi \quad .$$

同様の事を、「経度」の方に付いても見て行く。「Newton力学」に於ける「二体問題」では、経度  $\phi$  は、近日点黄経  $\varpi$  を用いる時は次の形に与えられる。

$$(33) \quad \phi = \varpi + f \quad .$$

従って、 $f=0$  に於ては  $\phi=\varpi$  である。量  $\Phi$  の方とは言えば、(29)式から明らかな様に此の時は、 $\Phi=\varpi$  である。

「Newton力学」では、初め「近日点」に在った水星（惑星）は、量  $f$  が  $f=2\pi$  となった時に一公転を終え、経度  $\phi$  は計算上  $\phi=\varpi+2\pi$  となる。端的に言えば、 $\phi=\varpi$  の儘なのである。従って、「近日点」に《前進》も無ければ《後退》も無い。

同じ事が、量  $\Phi$  に付いても言えるか？ 言えない !!! 「Einstein理論」の枠組みでは「動径」 $R$  が、再度「近日点距離」を与える時には、量  $f$  は  $f=2\pi+\delta_i$  なる値を取るのであった。それ故に、「経度」 $\Phi$  に対しては、(29)式を次の様に扱わなければならない事になって居る。

$$\Phi = \varpi + (2\pi + \delta_i) + \frac{m}{p} \left[ 3 \left\{ 2\pi - \frac{2\pi}{\eta} (1+e)^2 \right\} + O(\delta_i) \right] =$$

$$= \varpi + 2\pi + 2\pi \times \frac{3m}{p} \cdot \frac{(1+e)^2}{\eta} + \frac{m}{p} \times 3 \left\{ 2\pi - \frac{2\pi}{\eta} (1+e)^2 \right\} \quad ,$$

$$(34) \quad \Phi = \varpi + 2\pi + \frac{6\pi m}{p} \quad .$$

此処でも、量  $2\pi$  の存在は無視して考えても構わない。然し、等式  $\Phi=\varpi$  は、明らかに成立し得ず、量  $m$  の大きさの「前進」を有する事になって居る。此の「前進量」を与える表式は、将に「Einsteinの一般相対性理論」が、「水星近日点黄経  $\varpi$  に於ける余剰の永年変化の問題を解いた」とする「表式」そのものである。

以上で、我々が、本小論で『論じたいと考えた処のもの』は総て述べ終った事になる。

小論を閉じるに当り、フランス航空宇宙研究所のマルシャル博士に敬意と謝意を表して置きたい。博士には、当該問題に付いて、折々に、議論に乗って頂き懇切な助言を受けて来た。そのお蔭で「より深く」此の問題を考える事も出来たと云う次第である。

Nous sommes très heureux de dire à haute voix que Docteur Christian MARCHAL, Office national d'études et de recherches aérospatiales (ONERA) France, nous a donné beaucoup de suggestions utiles grâce auxquelles nous avons pu achever et perfectionner les études actuelles dont nous remercions de tout notre cœur.

本小論で扱った Schwarzschildの系に関しては、此処とは異なる観点から論じた事がある（井上, 1989）。

#### 参考文献

1. Brumberg, V. A. : 1991, *Essential Relativistic Celestial Mechanics*  
Adam Hilger
2. Eddington, A. S. : 1924, *The Mathematical Theory of Relativity*  
Cambridge University Press
3. Einstein, A. : 1915, *Preuss. Akad. Wiss. Berlin, Sitzber.* 47
4. 平山清次 : 1930, 天體力學 : 一般摂動論 岩波書店
5. 井上 猛 : 1989, 第23回天体力学研究会集録, p. 156
6. Inoue, T. : 1992, *Proceedings of the Twenty-Fifth Symposium on Celestial Mechanics*, p. 205
7. 井上 猛 : 2000, 第32回天体力学N体力学研究会集録, p. 147
8. Le Verrier, U. J. : 1859, *Annales de l'Observatoire Impérial de Paris*, V
9. Marchal, C. : 2000, Private communications
10. Newcomb, S. : 1895, *Astronomical Papers of the American Ephemeris*, VI
11. Schwarzschild, K. : 1916, *Sitzber. Deut. Akad. Wiss. Berlin, Kl. Math. -Phys. Tech.*

# Subsystems in a stable planetary system I. A classification

Kiyotaka TANIKAWA and Takashi ITO<sup>1</sup>

National Astronomical Observatory of Japan  
Mitaka 181-8588, Japan

## Abstract

Our planetary system is dynamically stable for the lifetime of the solar system according to the long-term numerical integrations of planetary orbits (Ito & Tanikawa, 2002). we discuss various forms of subsystems of a stable planetary system which may contribute to maintain the system stability. It is well known that resonances play an important role in such a long time scale. We stress that contrary to the restricted problem such as the stability of asteroids and comets, multi-planet subsystems may have variety of mechanisms for keeping stability.

## §1. Introduction

The solar system is dynamically stable for at least five billion years in the past and four billion years in the future (Ito & Tanikawa, 2002; hereafter IT2002). The orbital elements of nine planets are almost constant within a small range of variability for 9 billion years. Only the eccentricity and/or inclination of Mercury might have a slight touch of secular change. We expect that if the solar system of nine planets would become unstable, then Mercury would be the first that makes a close encounter with Venus or Earth and possibly dives into the Sun or is scattered away from the solar system.

For the moment, we are safe to say that the solar system is stable. We are apt to ask: Why is the solar system stable? This is a dangerous question because the question requires a 'true' reason of stability. We instead ask: How is the solar system stable? This is a feasible problem. We can ask more specifically. 'What kind of stabilization mechanisms are working? 'What kind of realization and what kind of characteristics do these interactions have?' 'Have these interactions continued to exist from the formation era of the system? 'Are these interactions common to extrasolar systems?

In the present paper and the one follows this, we address these questions and answer some of them or at least show the direction of study to be taken. Five-billion-year integrations of our solar system provided us information not only on the stability of the system but also on mechanisms acting on various subgroups of the system. These are the clues for unveiling the mechanisms of keeping stability. The purpose of the present paper is to classify the groupings of planets, to examine the stabilizing mechanisms of our solar system and justify the validity of the classification. As is well-known one of the stabilizing mechanisms is the grouping of individual bodies. The whole system becomes stable upon making subgroups if otherwise unstable. One famous and typical example is the hierarchical structure of gravitating systems. We will see in this report other types of grouping in the solar system.

---

<sup>1</sup> email: tanikawa.ky@nao.ac.jp and tito@tabby.mtk.nao.ac.jp

We start with classifying and counting subsystems. We do not claim that we count up all subsystems. A number of subsystems may have been overlooked.

- (1) (mean motion) Resonant pair
  - 1.1. Binary planets (Sister planets)
  - 1.2. Other 1:1 resonant pair
  - 1.3. Other (mean motion) resonant pairs
  - 1.4. Resonant triples
- (2) Close neighbors (Cousin planets)
- (3) Planetary groups
- (4) Independent planetary subsystems

In the following section, we describe individual subsystems and examine their stability where it is possible. We specifically pay attention to the last three subgroups of the above list.

## §2. More about subsystems

In this section, we consider subsystems in the Solar System together with possible subsystems which can be existent in other solar systems. If the latter subsystems do not exist in other solar systems, then this gives strong constraints on the formation process or mechanism of planetary systems. In some cases, we alter the existent subsystems and compare the stability of the whole solar system with and without this alteration. In other cases, we add perturbations from outside the Solar System. This stability analysis is only possible because we are sitting outside the solar system and in front of computers.

### 2.1. Resonant pair planets

If plural planets are in mean motion resonance(s), then we call these *resonant multiples*. If the number of planets involved is two, these are called a resonant pair. As is well-known, there are an infinite number of resonances because there are an infinity of rational numbers. It is also known that higher resonances have smaller resonance regions. Thus, the number of resonant multiples are not so many though, up to the observational limit, every pair or triple, etc is in resonance.

#### 2.1.1. Binary planets

A pair of planets in 1:1 mean motion resonance is called a (1:1 mean motion) resonant pair or simply *binary planets*. The member bodies have strongest interaction and connection. Examples are the Earth–Moon system and Pluto–Charon system. There may be an objection that these are not the planet–planet pairs. We do not know the reason why in our solar system there are no planet–planet binaries. We do not know why Jupiter and Saturn do not have satellites of mass comparable with that of the Earth or Venus. We may find these combinations in extrasolar planetary systems. The formation process of planetary systems somehow prohibited these combinations in our solar system.

Essentially, there are three kinds of binary planets.

#### 1) Pluto–Charon system

In this system, orbits of both planets are sometimes convex and other times concave to the Sun. The semi-major axis is  $17 R_{\text{pluto}}$ , the eccentricity is equal to zero, and the



period is 6.387 days. The masses are  $M_{\text{charon}} = 0.08M_{\text{pluto}}$ . The acceleration of Charon due to the Sun is

$$a_{\text{charon}} = GM_{\odot}/r_{(\odot-\text{charon})}^2 = 3.79 \times 10^{-6} \text{m/s}^2, \quad (1)$$

whereas the acceleration due to Pluto is

$$a_{\text{charon}} = GM_{\text{pluto}}/r_{(\text{pluto}-\text{charon})}^2 = 2.63 \times 10^{-3} \text{m/s}^2. \quad (2)$$

The acceleration from Pluto is far larger than that from the sun. So, the orbit of Charon is always concave to Pluto.

## 2) Earth–Moon system

The orbital semi-major axis is  $60R_{\oplus}$ , the eccentricity is 0.0549, and the orbital period is 27.3 days. The masses satisfy  $M_{\text{moon}} = 0.020M_{\oplus}$ . The acceleration of the Moon due to the Sun is

$$a_{\text{moon}} = GM_{\odot}/r_{(\odot-\text{moon})}^2 = 5.93 \times 10^{-3} \text{m/s}^2, \quad (3)$$

whereas the acceleration due to the Earth is

$$a_{\text{moon}} = GM_{\oplus}/r_{(\oplus-\text{moon})}^2 = 2.70 \times 10^{-3} \text{m/s}^2. \quad (4)$$

The orbit of the Moon is always concave to the sun. The interaction between components is weaker than the Pluto–Charon system.

## 3) Retrograde binary planets

Mikkola & Innanen (1997) coined the term 'quasi satellite' for the object orbiting the mother planet with unusually large semimajor axis. Wiegert et al. (2000) surveyed the stability of these orbits around outer four giant planets and found that low inclination orbits survive for the age of the solar system. This orbit is fairly large compared with a prograde satellite orbit. So a tentative system will have a relative distance greater than the Earth–Moon system. We classify these objects as a *retrograde binary*.

### 2.1.2. Other 1:1 mean motion resonant pair

#### 1) Tadpol-type pair

Typical examples are Trojan asteroids. If, instead of an asteroid, there is an object of planetary mass, the system can be called a planetary pair. There can be a pair that one of the members explores more wide area. But the position is always in front of the main member or in its back side. It seems that the stability of this pair is not well analysed.

#### 2) Horseshoe pair

Theoretically, there can be horseshoe type pairs. However, there is a problem of long-term stability. We do not so far find any pair of comparable masses in our solar system. Small objects are found in this position. It is interesting to see the stability including perturbative forces from other planets.

### 2.1.3. Neptune–Pluto system

This is a very special system maintaining its long-term stability using plural resonant mechanisms (Kinoshita & Nakai, 1995, 1996; IT2002). In the shortest time scale, this is a 3:2 mean motion pair. It roughly resumes the original configuration every five hundred

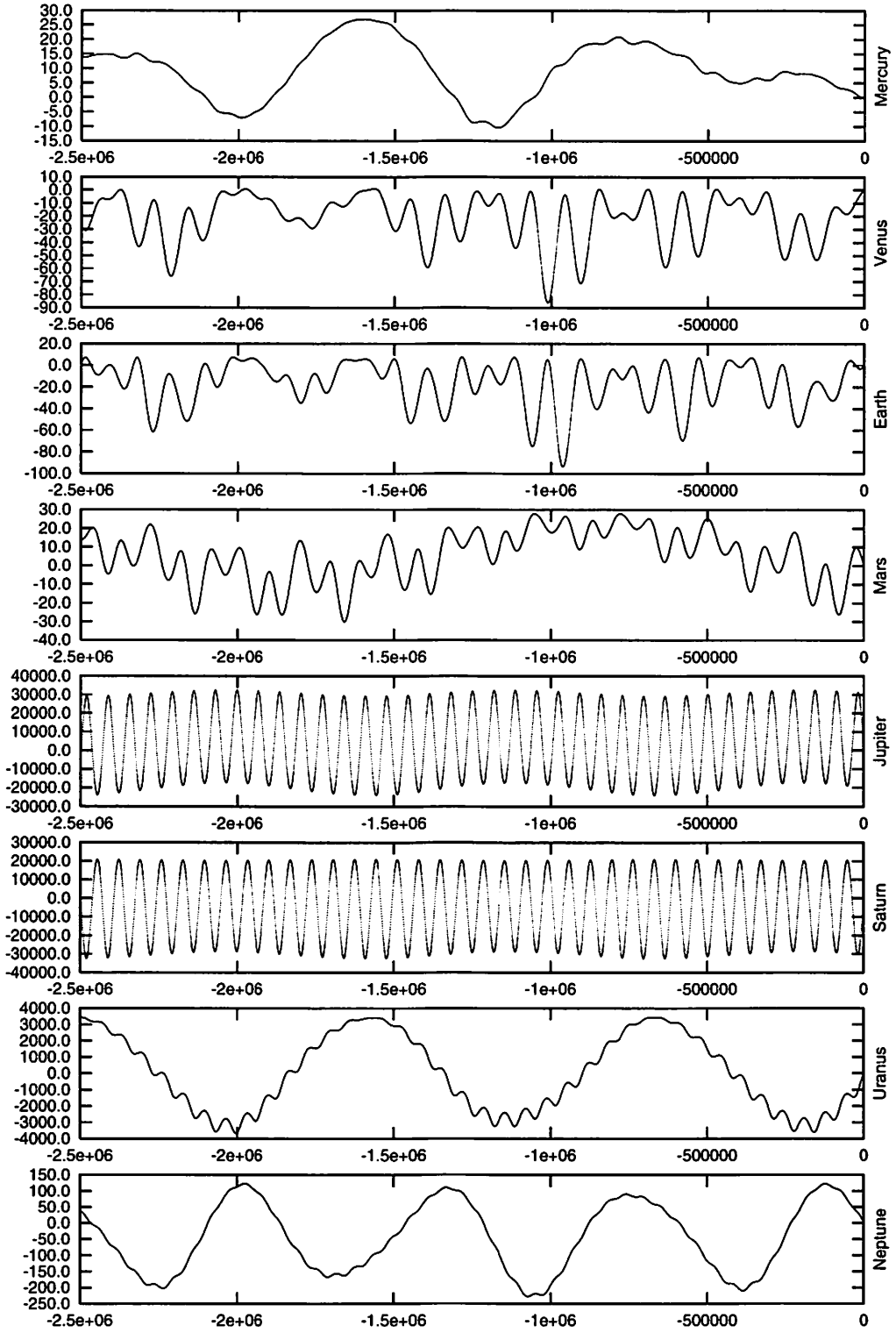


Figure 1: Variation of the angular momentum of eight planets for 2.5 million years from Brower & van Woerkom's (1950) theory of secular perturbation. From the top, Mercury, Venus, Earth, Mars, Jupiter, Saturn, Uranus, and Neptune. The ordinate represents the total angular momentum. The unit is  $10^{-12} M_{\odot} \text{AU}^2 \text{day}^{-1}$ .

years. In a more microscopic view, the critical argument  $\theta_1 = 3\lambda_P - 2\lambda_N - \varpi_P$  librates around  $180^\circ$ . The period is  $2 \times 10^4$  year. Thus, the system returns to the original configuration every  $2 \times 10^4$  years.

Pluto's argument of perihelion  $\omega_P = \theta_2 = \varpi_P - \Omega_P$  librates around  $180^\circ$ . The period is  $3.8 \times 10^6$  years. The system resumes its original configuration, which takes into account the relative position of the ascending node, every  $3.8 \times 10^6$  years. This is a very long periodicity. The longitude of Pluto's node referred to the longitude of Neptune's node,  $\theta_3 = \Omega_P - \Omega_N$ , circulates. The period of circulation is equal to the period of  $\theta_2$  libration. When longitudes of ascending nodes of Neptune and Pluto coincide ( $\theta_3 = 0$ ), Pluto's inclination becomes maximum, its eccentricity minimum, and argument of perihelion  $90^\circ$ . When  $\theta_3 = 90^\circ$ , Pluto's inclination becomes minimum, its eccentricity maximum, and argument of perihelion  $90^\circ$  again. This was confirmed by Milani et al. (1989).

There is a longer periodicity. An argument  $\theta_4 = \varpi_P - \varpi_N + 3(\Omega_P - \Omega_N)$  librates around  $180^\circ$ . The period is  $5.7 \times 10^8$  years. IT2002 showed that  $\theta_4$  varies between librations and circulation in  $O(10^{10})$ -year timescale. This is one of the longest timescales ever known.

## 2.2 Close neighbors

In our solar system, the Earth-Venus system, if these two can be called a system at all, occupies a special position as a subsystem. The Earth and Venus are planets of similar character. Nonetheless, these two planets have not explicitly been regarded as a dynamical pair. According to the long-term integrations of our planetary system (IT2002), these planets have interesting behaviors. These are not in any low order mean motion resonance. In the secular perturbation theory (Brouwer & van Woerkom, 1950), their orbital angular momenta have negative correlation in a short time-scale ( $\sim$  million years), i.e., if one planet obtains the angular momentum, the other loses the angular momentum, and *vice versa* (the second and third panels of Fig. 1). This is also observed in our numerical integrations. In a longer time-scale (billion years), the orbital angular momenta seem to have a positive correlation (Laskar, 1994; IT2002), i.e., if, for example, one planet obtains the angular momentum, the other also obtains the angular momentum. Figure 2 (the second and third panels) shows one of the numerical results taken from IT2002. This means that in a longer time-scale, these two planets behave synchronously against the perturbation from outside. They are repulsive each other in a shorter time scale but move together against outer forces in a longer time scale. These may be alternatively called cousin planets.

So far, our statement is of a qualitative nature. In order to quantitatively confirm the above statement and to see the difference from other conceivable pairs like Jupiter-Saturn and Uranus-Neptune, we take correlations of orbital elements using the results of long-term integrations of IT2002.

A simplest method of taking correlation would be, after removing the trend from a time series, that is, after subtracting an average value, to assign + (resp. -) to the data at  $t + \Delta t$  if the value of a parameter at  $t + \Delta t$  is larger than or equal to (resp. less than) that at  $t$ . Pick up two time series and take the data according to the elapse of time. If at  $t$  both data have +, then we add 1. If at  $t$  one has + and the other has -, then we add -1. If at  $t$  both data have -, then we add 1. We sum up 1 and -1 made from two time series and divide by the number of data. This quantity  $\rho$  is the easiest index for the strength of correlation. If  $\rho$  is positive and large, then a positive correlation. If  $\rho$  is negative and large in absolute value, then a negative correlation. If the absolute value of

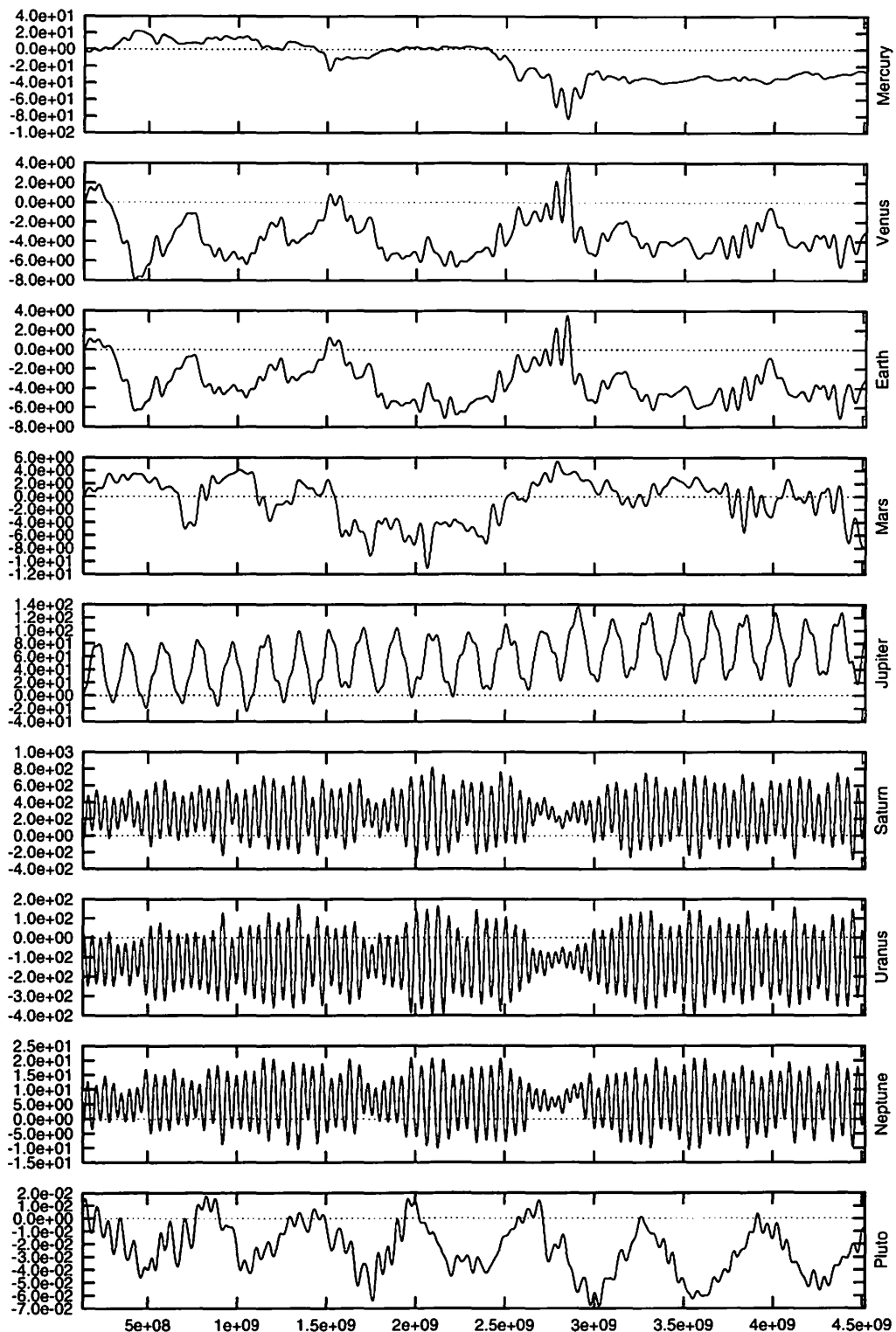


Figure 2: Variation of the angular momentum of nine planets for 4.5 billion years showing qualitatively the correlation of orbital elements (reproduced from IT2002). From the top, Mercury, Venus, Earth, Mars, Jupiter, Saturn, Uranus, Neptune, and Pluto. The ordinate represents the total angular momentum.

$\rho$  is small, then the correlation is weak.

The second easiest method which we have adopted is as follows. As before, we subtract the trend and obtain and average for each time series. From each data  $x_i$ , we get a normalized value  $X_i = (x_i - \bar{x})/|\bar{x}|$ . Then correlation  $\rho$  is defined by

$$\rho = \frac{1}{N} \sum_{i=1}^N X_i Y_i$$

where  $X_i$  and  $Y_i$  are normalized data from two time series.

Table I. Correlation of the orbital energy and angular momentum in Earth–Venus, Jupiter–Saturn and Uranus–Neptune pairs.

N+1	Energy	Ang. Mom.	z-component
Venus–Earth	−9.5563e-01	8.9389e-01	9.0139e-01
Jupiter–Saturn	−4.7369e-01	−1.1749e-02	−9.9988e-03
Uranus–Neptune	9.6988e-02	−9.4638e-01	−9.4363e-01
N+2			
Venus–Earth	−9.1239e-01	7.8515e-01	8.3140e-01
Jupiter–Saturn	−2.0947e-01	2.2497e-03	5.9993e-03
Uranus–Neptune	−3.7995e-02	−9.0689e-01	−9.0414e-01
N+3			
Venus–Earth	−7.5304e-01	9.9014e-01	9.9272e-01
Jupiter–Saturn	−5.7535e-01	5.5392e-01	5.6135e-01
Uranus–Neptune	−2.3225e-01	−1.5655e-01	−1.4569e-01
N-1			
Venus–Earth	−9.0385e-01	9.1585e-01	8.6786e-01
Jupiter–Saturn	−3.1795e-01	−6.6656e-03	−4.6659e-03
Uranus–Neptune	1.2198e-01	−9.0218e-01	−8.9585e-01
N-2			
Venus–Earth	−8.8844e-01	8.5159e-01	8.8302e-01
Jupiter–Saturn	−5.0593e-01	−4.3422e-02	−4.1708e-02
Uranus–Neptune	5.6563e-02	−9.6215e-01	−9.5958e-01
N-3			
Venus–Earth	−8.2674e-01	9.5758e-01	8.5702e-01
Jupiter–Saturn	−8.1674e-01	3.8252e-01	3.8395e-01
Uranus–Neptune	−3.7880e-01	−1.8083e-01	−1.9054e-01
short			
Venus–Earth	−2.0819e-01	−1.7899e-02	−1.5059e-01
Jupiter–Saturn	−6.6097e-02	−4.5458e-01	−4.5498e-01
Uranus–Neptune	−3.3998e-03	−4.2998e-03	−3.3998e-03

We have several time series of long-term integrations of our planetary system. In addition, there are several dynamical and physical quantities for which correlations can be considered. This time, we consider the correlation of the energy, angular momentum and the z-component of the angular momentum. The results are shown in Table I. Let us first explain the meaning of  $N \pm i$  ( $i = 1, 2, 3$ ) and 'Short' in the table.  $N + i$  ( $i = 1, 2, 3$ ) is for the future and  $N - i$  ( $i = 1, 2, 3$ ) is for the the past. Except for  $N - 3$  for which the Moon is neglected, the Earth–Moon barycenter has been adopted in the integrations (see IT2002, Table I). 'Short' is obtained from the initial 10 million years of  $N + 2$ . It is to be noted that data series  $N \pm i$  have passed a low-pass filter, whereas 'short' data have not.

The first feature we point out is that the correlation is strongest in the Earth–Venus system. This confirms our first impression looking at Figs. 1 and 2 that the Earth–Venus may form a pair. The second characteristics seen in the table is that the Earth–Venus

system has reverse correlations in long and short time scales. The contrast between long and short time scales may be more conspicuous if we take a shorter time series for the shorter data. In 10 million times scale, the correlation seems in a transition, i.e., the correlation is rather weak in particular for the angular momentum. The correlation of the energy is negative in both cases, that is, if the orbital energy of the Earth increases, then the energy of Venus decreases in any time scale, and *vice versa*. The correlation of the angular momentum has a different character. In the shorter time scale, it has the same character as that of the energy. However, in a longer time scale, the angular momentum correlation is positive. It means that both planets gain and lose their angular momenta synchronously. Both planets behave as a unit against external perturbations. The third characteristics seen in Table I is rather unexpected. The Uranus–Neptune pair has strong correlations for the angular momentum. Correlations are negative in the long time scale, whereas the correlations are very small in the short time scale. The energy correlation is weak in both time scales. Thus, Uranus and Neptune move independently in a short time scale, whereas in the long time scale, their motions are related. When Uranus gain the angular momentum, then Neptune loses and vice versa. As a pair, the connection is weaker in the Uranus–Neptune pair than in the Earth–Venus pair.

It is interesting that Mercury and Mars seem to have a positive correlation for the angular momentum in the long-term. However this is a consequence of the fact that each planet behaves with negative correlations to the Earth–Venus system. We cannot say that Mars and Mercury constitute a system.

### 2.3 Planetary groups

So far, most of the efforts on the stability problem of our solar system have been concentrated on the motion of small bodies like satellites, asteroids, and rings. Recently, Kuiper-belt objects are added to the list. In other word, most studies treated the stability problem as a restricted problem in the sense that bodies for which the stability of motion is examined are assumed to be massless, and give no responding force to perturbers (see the review papers by Wisdom, 1987; Lissauer, 1999; Lecar et al., 2001).

There have occasionally been numerical studies on the stability of the whole solar system. 350 years of Eckert (1951) was the first. Then, in the chronological order, 120 thousand years of Cohen & Hubbard (1965), 1 million years of Cohen et al. (1973), 5 million years of Kinoshita & Nakai (1984), 100 Myr of Nobili et al. (1989), and so on. The longest record at present (2002) is  $\pm 5$  billion years integration of 9 planets by IT2002. IT2002 have shown that our solar system is stable 5 billion years in the past and 4 billion years in the future.

There have been a small number of systematic studies of the stability of planets as a group in our solar system. One of the first trials has been done by Gladman (1993) and later by Chambers et al. (1996). Ito & Tanikawa (1999 or IT1999) noticed the importance of these ideas to apply to the actual solar system. IT1999 were the first to examine the possibility that Jupiter form prior to the formation of terrestrial planets. The idea was that the existence of Jovian planets, especially the existence of Jupiter may affect the formation process of terrestrial planet group and even now affect the stability of this group.

Ito & Tanikawa (2001 or IT2001) also noticed that the so-called outer planetary group or Jovian planets form a subsystem which are not affected from other groups. Indeed, the motion of the Jovian planetary group may be not altered if there is no terrestrial planet

group. On the other hand, the terrestrial planet group receives the secular perturbation of Jupiter ( $\sim 300$  thousand year periodicity of the motion of Jupiter's perihelion) and make it uninfluential by sharing the effect of perturbation (IT1999).

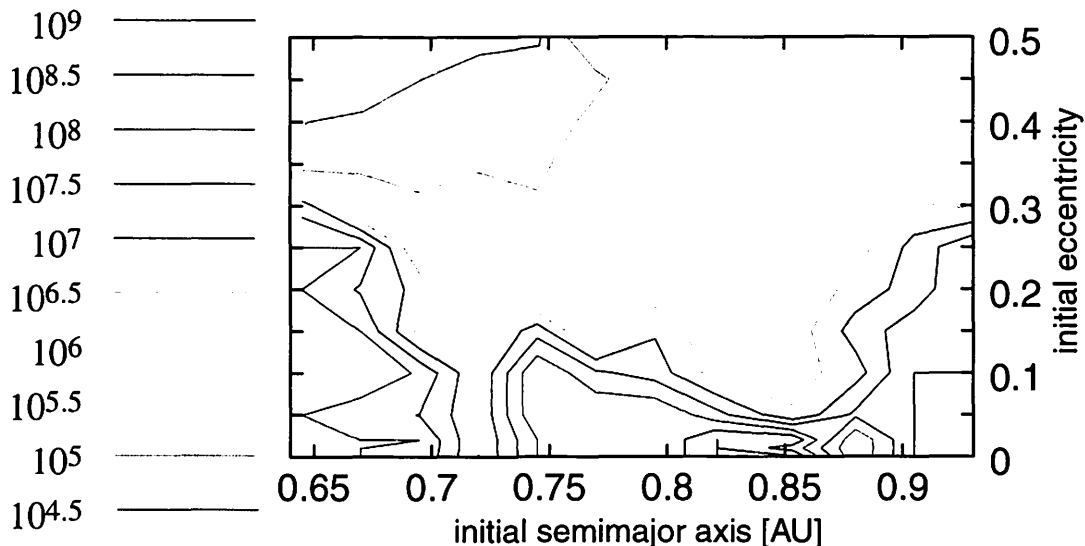


Figure 3: The contours of the same instability time. The results of the numerical integrations with the Earth and Venus merged. The ordinate and abscissa are the semimajor axis and eccentricity of a merger. The inclination is related to  $e$  by  $2I = e$ .

In order to check the meaning and role of this kind of subsystem, one somehow alters the terrestrial system and see what happens. Innanen et al. (1998) examined the dynamical stability of the inner solar system by removing each of terrestrial planets in turn. They found that a drastic phenomenon takes place when the Earth-Moon system is removed. The eccentricity of Venus oscillates with large amplitude. If Mercury is included, then the large oscillation of Mercury's eccentricity takes the role of Venus. Venus is nearly at the position of secular resonance from Jupiter. Innanen et al. (1998) interpreted that the Earth-Moon system suppress the secular resonance. IT1999 gave a slightly different interpretation to this phenomenon. Terrestrial planets share the effect of the secular perturbation from Jupiter. Eccentricities of all the inner planets increase. The mechanism is simple. Due to the perturbation of Jupiter, the eccentricity of Venus tends to increase. The eccentricity of the Earth increases according to the long time scale behavior of correlations in Table I. In other words, the Earth shares the increase of the eccentricity of Venus. The enhancement of the eccentricity is weakened by sharing.

Yet another method of checking stability is to merge two of the planets. To conserve the total mass, total energy and angular momentum is a too restrictive condition, so we consider two cases in which the total mass is conserved and either energy or angular momentum is conserved. We survey the stability of the altered system around the specified position. The position of the merger will be at  $a \simeq 0.855$  AU if the merger has the mass of the Earth and Venus and the orbital energy of the Earth and Venus neglecting the eccentricity. Indeed, we put the merger at various places between  $a = 0.645$  AU and  $a = 0.930$  AU with  $e$  ranging from 0 to 0.5. The number of different sets of initial parameters are more than 150. To reduce the number of integrations, we assumed  $e = 2I$ .

The result of stability analyses is depicted in Fig. 3. Here the abscissa is the semi-major axis and the ordinate is the eccentricity or inclination of the merger. Instability is meant if the orbit of some planet crosses the orbit of another planet. In our case, always Mercury does this. In the figure, the contours of equi-instability time are drawn. The lines with purple color corresponds to 100 my, i.e., the system is stable until 100 my. We need longer integration times to see the final fate of the system. We see two unstable intervals of semi-major axis centered at around  $a = 0.72$  and  $a = 0.88$  with  $e = 0$ . These are the positions of secular resonance from outer planets. Figure 4 shows this. In Fig. 4(top), the secular resonant motion of Mercury with the merger at  $a = 0.73$  is shown. The position of the merger is close to the resonance with Jupiter. Here instead of the increase of the mergers's eccentricity, Mercury's eccentricity increases. The number of points is small because Mercury soon becomes unstable. Figure 4(middle) and (bottom) show the secular resonant motion of Mercury with Jupiter and Uranus when the merger is at  $a = 0.88$ . We stopped the integration at  $t = 10^8$  years in most cases. Between  $a = 0.8$  and  $a = 0.85$  and for small  $e$ , integrations are done until  $t = 10^9$  years. The stable area in Fig.3 becomes smaller if we extend the integration time.

We can conclude that if the Earth and Venus merged in the early solar system, Mercury would have in high probability escaped away. The individual roles of the Earth and Venus as they occupy the present positions in the terrestrial zone contribute the maintenance of the stability of our planetary system. As pointed in IT1999, terrestrial planets keep their stability by sharing and weakening the secular perturbation from Jupiter. The Earth-Venus system plays the distributor of the effects of perturbation.

## 2.4 Independent planetary subsystems

Innanen et al. (1997) carried out interesting numerical simulations. They put a companion star of  $0.02 \sim 0.5M_{\odot}$  at 400AU from the Sun with a circular but inclined orbit from the invariant plane of the solar system with various inclinations. They wanted to see the stability of a planetary system in a binary. They took as a representative case the Jovian planetary system (Jupiter, Saturn, Uranus, and Neptune) and tried to see its behavior under the perturbation of a companion star. It was initially expected that Kozai mechanism will independently drive the inclination variations of planets and hence planets would soon experience close encounters. For a suitable parameter set, however, the Jovian planetary group shows a stability. The Kozai mechanism of individual planets is suppressed and the motions of the ascending node of the planets synchronize. They called this synchronous state a *dynamical rigidity*.

In our context, two experiments of Innanen et al. (1997, 1998) can be used as a tool for checking the strength of connection among planets against internal and external perturbations. The first experiment (Innanen et al., 1997) is a check for the unity of the system against external perturbation. The second experiment (Innanen et al., 1998) can be arranged to test the internal rigidity of the system. We will explain these in the corresponding subsections. In both cases, if there is no rigidity in two subsystems and these two are stable, then we can regard that these two systems are independent.

We give here a preliminary result in a sense that the number of experiments are not enough. We divide our planetary systems into three groups: [1] Mercury, Venus, Earth, Mars, and Jupiter; [2] Saturn; [3] Uranus, Neptune, and Pluto. We will carry out numerical integrations of orbits with and without the second group (Saturn).

### 2.4.1 Rigidity against internal perturbation



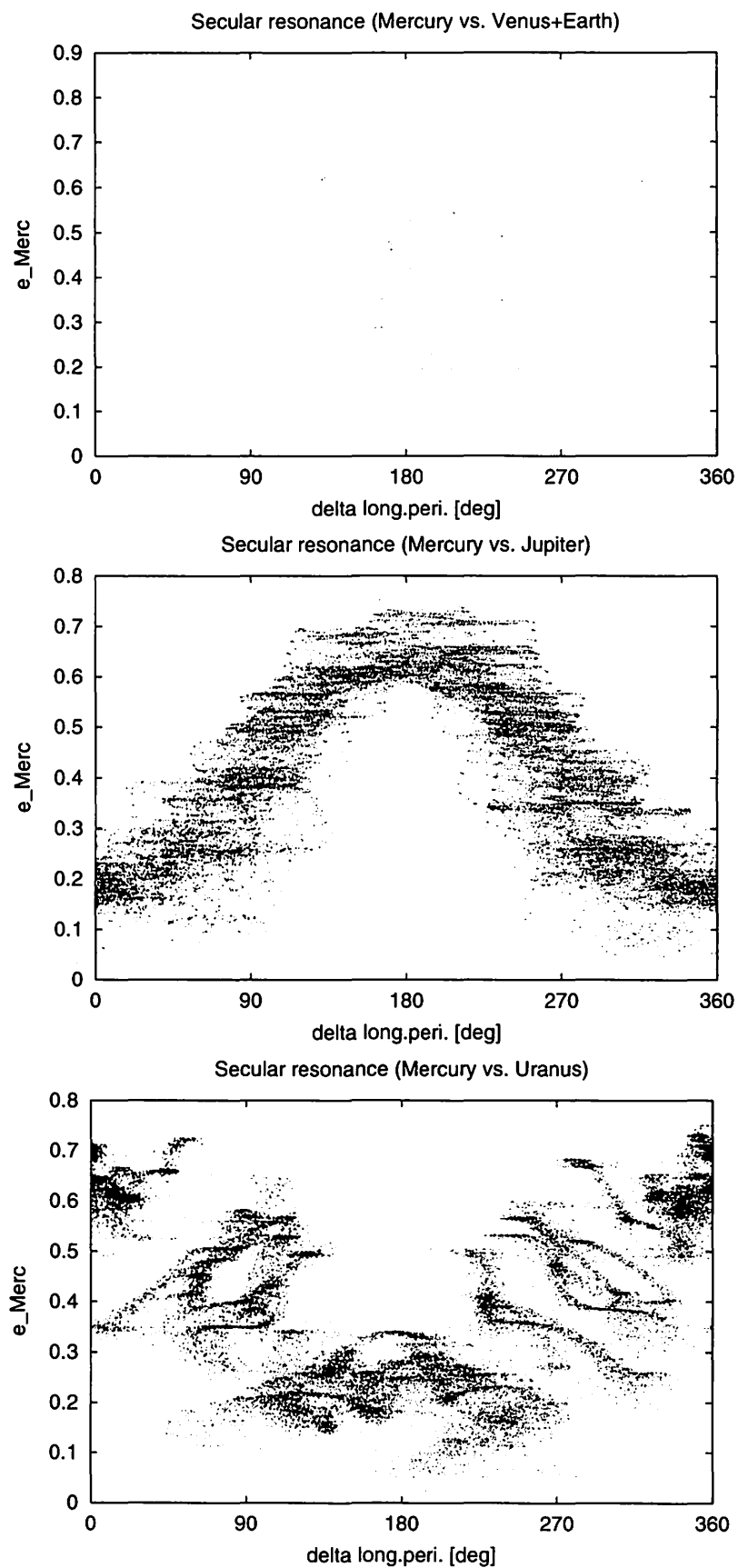


Figure 4: Secular resonance of Mercury (top) with the mergr at  $a \sim 0.73\text{AU}$  and (middle) with Jupiter when the merger is at  $a \sim 0.88\text{AU}$ , (bottom) with Uranus when the merger is at  $a \sim 0.88\text{AU}$ .

We carry out two integrations over 100 million years. In the first calculation, Group [1] and Groups [2] + [3] are inclined by  $20^\circ$ . In the second calculation, Groups [1] and [3] are inclined  $20^\circ$  each other, whereas Group [2] is not included (see Fig. 5(a)).

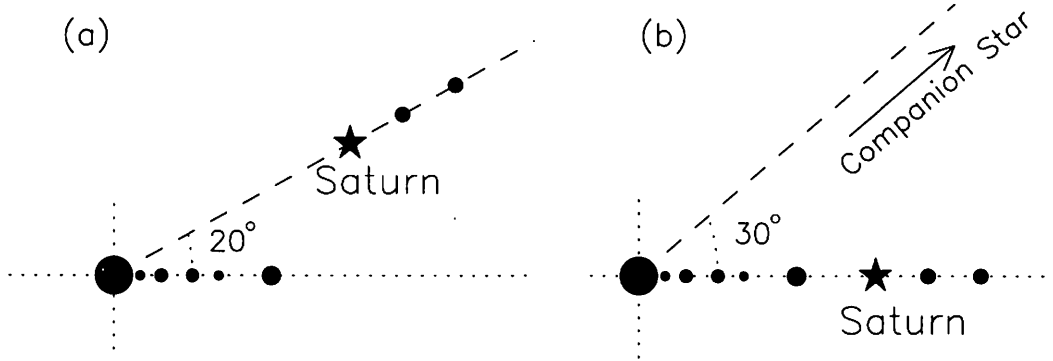


Figure 5: Initial configurations of planets for checking rigidity, (a) rigidity against internal perturbation, (b) rigidity against external perturbation.

Possible results are : (i) The system in the first calculation is unstable, whereas the system in the second calculation is stable; (ii) Both systems are unstable; (iii) Both systems are stable. Result (i) implies that Saturn plays the role of pivot to connect outer and inner planetary systems and the planetary system comprises two independent subsystems without Saturn. Result (ii) implies that the connection among planets are strong enough irrespective of the existence or non-existence of Saturn and the planetary system is unstable with other configurations. Result (iii) implies the parameter used in this investigation is not suitable to check the independence of subsystems.

Numerical results are shown in Fig.6. At around  $t = 4.4 \times 10^6$  years, the system with Saturn becomes unstable in the sense that Mercury's eccentricity becomes as large as 0.8 and more (the upper panel of Fig. 6) and its inclination approaches  $60^\circ$  (the lower panel of Fig. 6). In the system without Saturn, the eccentricity of Mercury oscillates stably between 0.18 and 0.23, and the inclination has similar variations (Fig. 7). Result (i) is attained. Saturn plays the role of a pivot between outer and inner planets. It is to be noted that the terrestrial planets join the rigidity of the whole planetary system. This again confirms the rigidity of terrestrial planets against the perturbation of Jupiter.

#### 2.4.2 Rigidity against external perturbation

We use Innanen et al. (1997)'s numerical experiments as a method of measuring the strength of connection among planets against external perturbations. Rigidity implies the unity of the system. If there is no rigidity in two subsystems and these two are stable, then we can regard that these two systems are independent.

We carry out two integrations: one with Saturn and the one without Saturn. In both cases, a perturbing star of mass  $0.2M_\odot$  moves on the circular orbit of 500AU with inclination  $30^\circ$ . The initial configuration is shown in Fig. 5(b). The integration time is 100 million years. The results are shown in Figs. 8 and 9. Interestingly, both systems are stable and have synchronous motion of nodes  $\Omega$ . Though initially the nodes are distributed random on each orbits of planets, they converge and gather. Inclinations of planets change together (middle panels of both figures). The only noticeable difference is that the oscillation amplitude of Mercury's eccentricity is smaller when Saturn is not

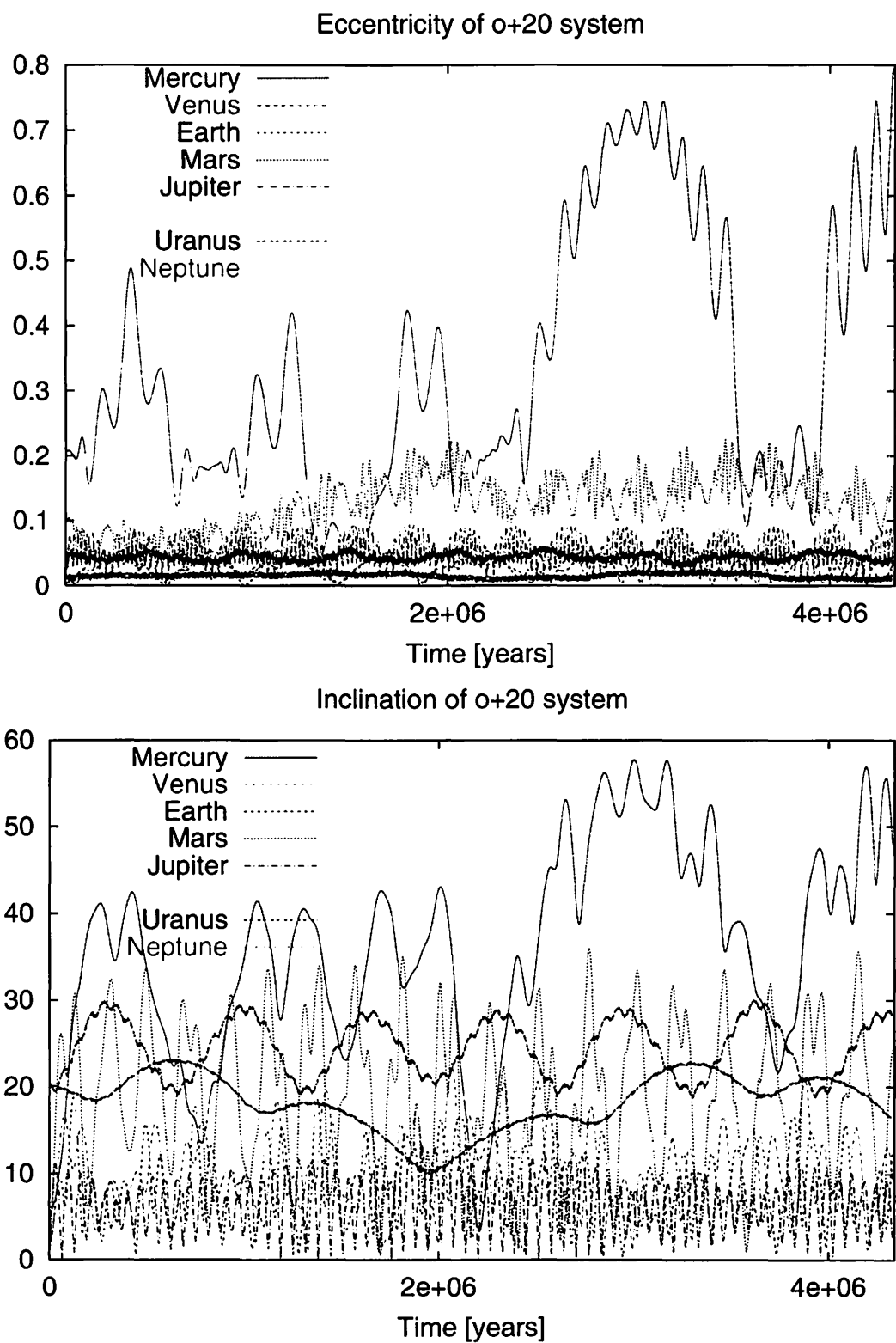


Figure 6: Rigidity against internal perturbations. The planetary system of Fig. 5 (a) is unstable if Saturn exists. Variations of eccentricities (top) and inclinations (bottom) of planets.

existent. Jupiter affects the motion of Uranus and Neptune without the connection of Saturn. This is rather a surprising fact.

The difference of the results of sections 2.4.1 and 2.4.2 indicates that the independence of planetary groups should be examined more carefully. The different checking methods may discriminate subtle differences of the systems otherwise overseen. The existence or the non-existence of Jupiter may be more explicitly related to the independence of outer and inner planetary systems.

### 3. Discussions and Summary

We examined several groupings of planets. Some of the groupings are not real in our solar system. The existence of a particular grouping of planets should have reflected the formation process of planets and planetary systems. Thus for example, there are no binary planets in our solar system. There are binary asteroids and binary Kuiper-belt objects. Pluto and Charon may be interpreted as a binary Kuiper-belt objects. The Earth–Moon system may be conceivable as binary planets. However, the hypothesis of a giant impact which is most successful at present in explaining the formation of the Moon presupposes the existence the Earth prior to the impact. The hypothesis is not compatible with binary planets. We may need to consider different formation processes to obtain a binary of comparable masses. Can a binary of comparable size be produced through a giant impact?

Binary planets are in a sense a paradoxical objects. Suppose a multi-planet planetary system. In general, the system becomes unstable if two of the members make a close approach. However, if the two are close enough and continued to be close enough, then these two constitute a subsystem and the whole system becomes stable once again. The difference is that through many body interactions, gravitational potential energy is released from the binary and is given to the remaining constituents of the system as kinetic energy. If the whole system has negative enough energy so that kinetic energy does not cancel out the total energy, then the system remains stable. So binaries of small masses can be possible to form. Dissipative media may help to form binaries by absorbing the energy and scatter itself away. Planetesimals are one of the candidates. Thus, a planetary system can release energy when it forms. The Oort cloud may be interpreted as an object resulted from the energy release.

As a subsystem, cousin planets sit between binaries (sister planets) and planetary groups. Binaries can be regarded as a single body because the motion of the center of gravity frequently replaces the individual motions of the components. Cousin planets turn out in the present study to be important first because one component suppresses the orbital instability of the other component and secondly because they transmit the perturbation to other members of a larger subgroup to stabilize this larger subgroup. The Earth–Venus is a good example. Uranus–Neptune may be weak cousin planets. Numerical experiments confirm that the  $N$ -body considerations are important. Secular perturbation theory predicts the position of secular resonances of massless particles. In the non-massless perturbed system, secular perturbations may be nullified or the position of secular resonance may be moved off the system.

A planetary group is a collection of loosely connected mutually dynamically dependent planets. A typical example is the terrestrial planets. As Innanen et al. (1998) showed, if one of the important members is void, the terrestrial system becomes unstable against external perturbation. In the present paper, we showed that if two of the members merge

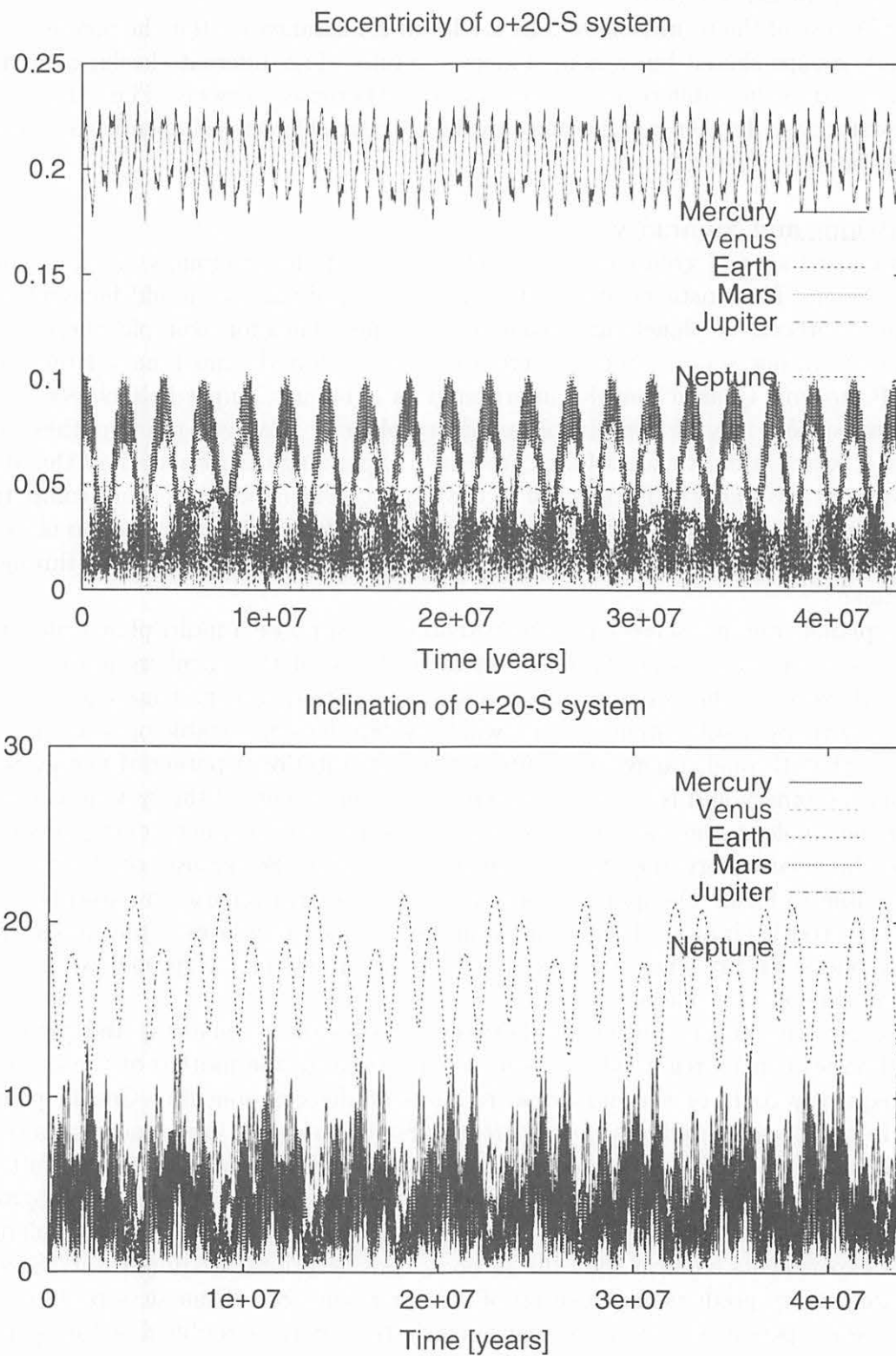


Figure 7: Rigidity against internal perturbations. The planetary system of Fig. 5(a) is stable if Saturn does not exist. Variations of eccentricities (top) and inclinations (bottom) of planets.

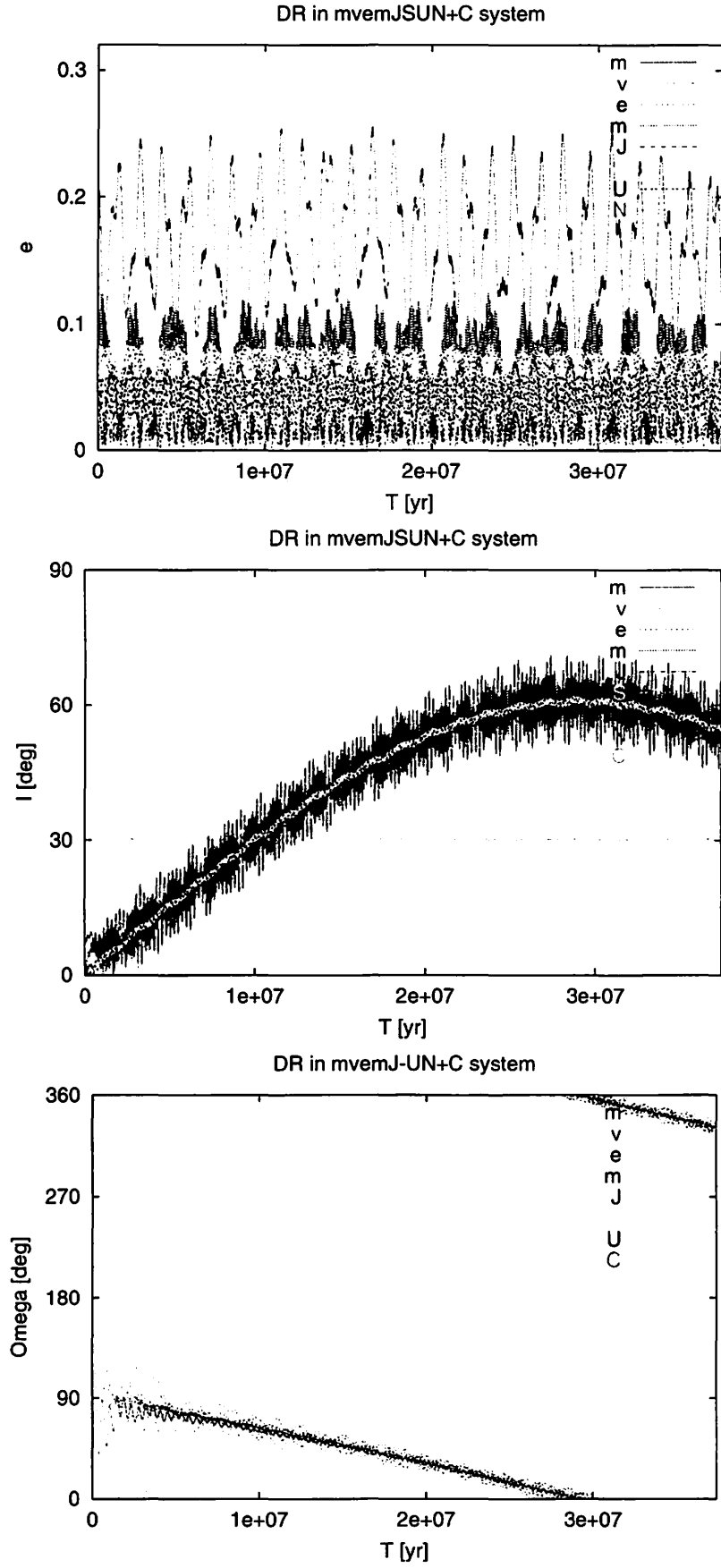


Figure 8: Rigidity of the planetary system with Saturn against external perturbations. Variations of eccentricities (top), inclinations (middle), and nodes (bottom) of planets.

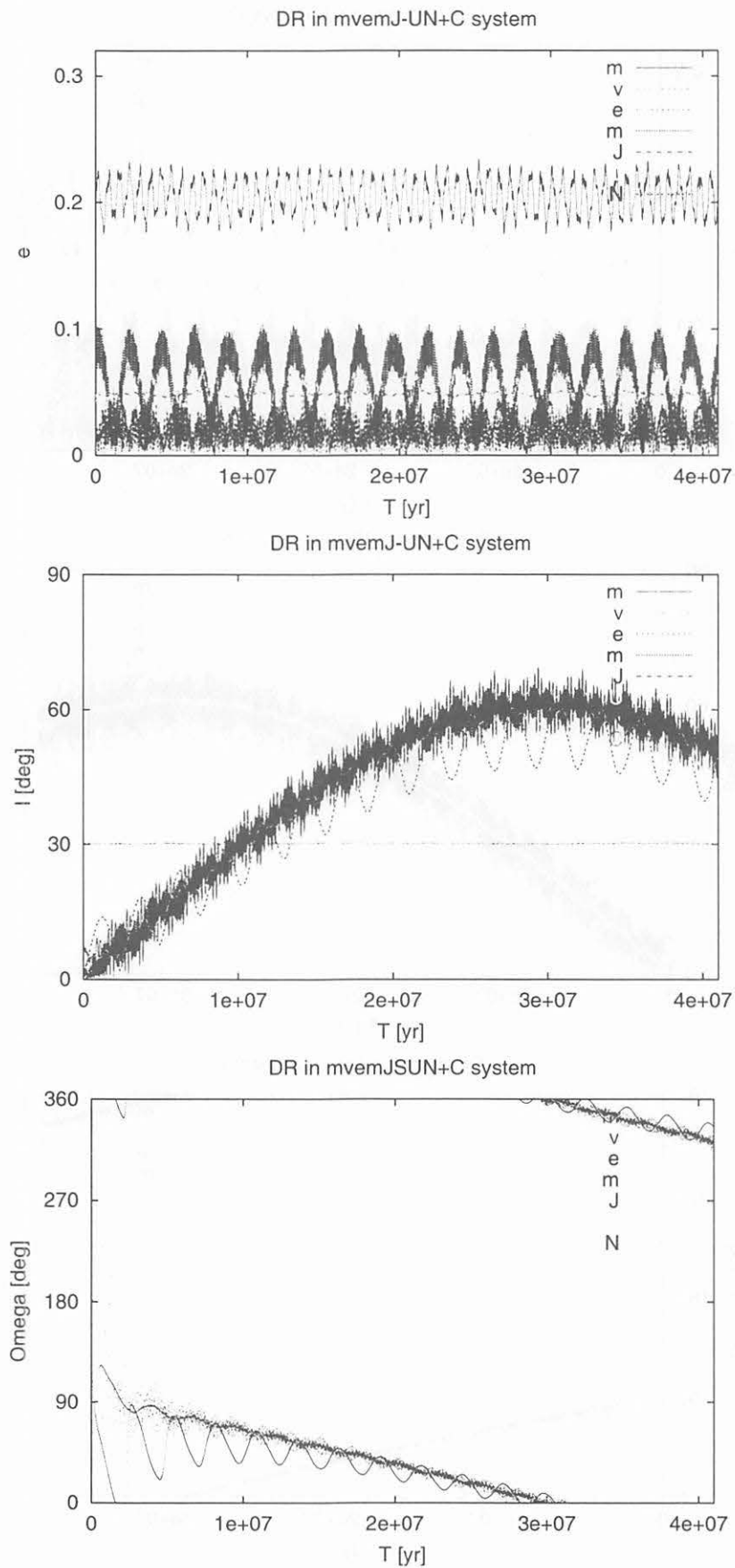


Figure 9: Rigidity of the planetary system without Saturn against external perturbations. Variations of eccentricities (top), inclinations (middle), and nodes (bottom) of planets.

into one, the system becomes unstable also against external perturbation. Even if one of the members is in strong secular resonance with perturbing bodies, all the members share the perturbation to stabilize the whole subgroup. The configuration itself contributes to the stability of the system. This gives a strong constraint to the formation process.

If there are two independent planetary subsystems around a star, We may say that planetary formation processes took place twice. These may be expected in a binary stellar system. One planetary system is around one member of the binary and the other system is around both member stars.

## References

- [1] Brouwer, D. and van Woerkom, A.: The secular variations of the orbital elements of the principal planets, *Astron. Pap. Amer. Ephemeris. Naut. Alm.* **13**, 81–107 (1950).
- [2] Chambers, J.E., Wetherill, G.W., and Boss, A.P.: The stability of multi-planet systems, *Icarus* **119**, 261–268 (1996).
- [3] Cohen, C.J. and Hubbard, E.C.: Libration of the close approaches of Pluto to Neptune *Astron. Journal* **70**, 10 (1965)
- [4] Cohen, C.J., Hubbard, E.C., and Oesterwinter, C.: Planetary elements for 10000000 years, *Celest. Mech.* **7**, 438–448 (1973),
- [5] Eckert, W.J.: Numerical theory of the five outer planets, *Astron. Journal* **56**, 38 (1951).
- [6] Gladman, B.: Dynamics of systems of two close planets, *Icarus* **106**, 247–263 (1993).
- [7] Innanen, K.A., Zheng, J.Q., Mikkola, S., and Valtonen, M.J.: The Kozai mechanism and the stability of planetary orbits in binary star systems, *Astron. J.* **113**, 1915–1919 (1997).
- [8] Innanen, K., Mikkola, S., and Wiegert, P.: The Earth–Moon System and the Dynamical Stability of the Inner Solar System, *Astron. Journal* **116**, 2055–2057 (1998).
- [9] Ito, T. and Tanikawa, K. (IT1999): Stability and instability of the terrestrial protoplanet system and their possible roles in the final stage of planet formation, *Icarus*, **139**, 336–349, 1999.
- [10] Ito, T. and Tanikawa, K. (IT2001): Stability of terrestrial protoplanet systems and alignment of orbital elements, *Publ. Astron. Soc. Japan* **53**, 143–151 (2001).
- [11] Ito, T. and Tanikawa, K. (IT2002): Long-term integrations and stability of planetary orbits in our solar system, *Mon. Not. R. Astron. Soc.* **335**, 2002 (in press).
- [12] Kinoshita, H. and Nakai, H.: Motions of the perihelions of Neptune and Pluto, *Celest. Mech.* **34**, 203–217 (1984).
- [13] Kinoshita, H. and Nakai, H.: The motion of Pluto over the age of the solar system, in *Dynamics, ephemerides and Astrometry in the Solar System*, Kluwer Academic Publishers, Dordrecht, pp.61–70 (1995).



- [14] Kinoshita, H. and Nakai, H.: Long-term behavior of the motion of Pluto over 5.5 billion years, *Earth, Moon, Planets* **72**, 165–173 (1996).
- [15] Laskar, J.: Large-scale chaos in the solar system *Astron. Astrophys.* **287**, L9–L12 (1994).
- [16] Lecar, M.L., Franklin, F.A., Holman, M.J., and Murray, N.W.: Chaos in the solar system, *Ann. Rev. Astron. Astrophys.* **39**, 581–631 (2001).
- [17] Lissauer, J.: Chaotic motion in the solar system, *Reviews of Modern Phys.* **71**, 835–845 (1999).
- [18] Mikkola, S. and Innanen, K.: Orbital Stability of Planetary Quasi-Satellites in *The Dynamical Behaviour of our Planetary System*, Proceedings of the Fourth Alexander von Humboldt Colloquium on Celestial Mechanics, Kulwer Academic Publishers, edited by R. Dvorak and J. Henrard, p.345 (1997).
- [19] Milani, A., Nobili, A.M., Carpino, M.: Dynamics of Pluto, *Icarus* **82**, 200–217 (1989).
- [20] Nobili, A.M., Milani, A., Carpino, M.: Fundamental frequencies and small divisors in the orbits of the outer planets, *Astron. Astrophys.* **210**, 313–336 (1989).
- [21] Wiegert, P., Innanen, K., and Mikkola, S.: The stability of quasi satellites in the outer solar system, *Astron. J* **119**, 1978–1984 (2000).
- [22] Wisdom, J.: Urey prize lecture: Chaotic dynamics in the solar system, *Icarus* **72**, 241–275 (1987).

# Evolution of obliquity of a terrestrial planet due to gravitational perturbation by a giant planet

Keiko Atobe<sup>1</sup>, Takashi Ito<sup>2</sup>, and Shigeru Ida<sup>1</sup>

1. Departement of Earth and Planetary Sciences,  
Tokyo Institute of Technology, Japan

2. Astronomical Data Analysis Computer Center,  
National Astronomical Observatory, Japan

## Abstract

We have studied the change in planetary obliquity near a spin-orbit resonance through numerical calculations and analytical arguments. To clarify basic process of the obliquity evolution, we have considered a simple system that consists of a host star, a hypothetical terrestrial planet, and a hypothetical giant planet. When the precession rate of the spin axis of the terrestrial planet coincides with the frequency of secular variations in its orbital inclination, spin-orbit resonance occurs and the obliquity of the terrestrial planet has large variations. We investigated time evolution of the obliquity near the resonance through numerical calculations of secular precession equations as well as analytical arguments. We derived the resonance width semi-analytically. Using this result, we predict the resonance region as a function of semi-major axis for a given giant planet.

## 1 Introduction

In general, the orientation of the planets' spin axis is not fixed, but changes all the time. Because of their equatorial bulge, planets are subject to torques arising from the gravitational forces of their satellites, host star and other planets. This causes precessional motion of the spin axis. Since the planets' orbits exhibit secular variations induced by gravitational perturbations exerted by other planets, the obliquity of the planets (the angle of the spin axis relative to the orbital plane) generally changes periodically, too. At present, Earth's spin has a precessional period of about 26,000 years, and its obliquity varies by  $\pm 1.3$  degrees around the mean value of 23.3 degree. Such obliquity variations would affect the planet's global climate through insolation change.

Ward (1974) and Ward & Rudy (1991) showed that large  $\sim \pm 10$  degrees variations of the obliquity of Mars are caused by the spin-orbit resonance, employing the secular precession equations (Eq. (1)). Here, spin-orbit resonance means that the precession rate of the spin axis coincides with one of the eigenfrequencies of secular variations in the orbital inclination.

Overlapping of spin-orbit resonances may cause chaotic variations of the obliquity of terrestrial planets (Ward 1992, Laskar & Robutel 1993, Touna & Wisdom 1993). The maximum oscillation amplitude of orbital inclination at a spin-orbit resonance was approximately derived by Ward (1993), through a nonlinear analysis of the secular precession equations. If changes

in planetary orbital inclinations are quasi-periodic, secular perturbation theory (e.g., Brouwer & Clemence 1961) predicts locations of the spin-orbit resonances.

Laskar (1989) showed that orbital evolution of the terrestrial planets has Lyapunov time scale  $\sim 10^7$  years, which may imply that the orbital evolution is chaotic on a timescale  $\sim 10^7$  years (although their orbits are globally stable). The chaotic orbital evolution of planets results in more complicated obliquity changes (Laskar & Robutel 1993).

With Fourier spectrum of time-dependent eigenfrequencies for orbital inclination variations obtained by Laskar (1990), Laskar & Robutel (1993) integrated the secular precession equations with wide ranges of initial obliquity ( $\epsilon_0$ ) and the precession parameter ( $\alpha$ ). They found large chaotic regions in the  $\epsilon_0$ - $\alpha$  plane and suggested that the obliquity of all the terrestrial planets except the Earth in the Solar system could have experienced large and chaotic variations.

Since the procedure to find the chaotic regions by Laskar & Robutel (1993) is rather complicated, it would not be easy to apply their results to more general extrasolar planetary systems. Laskar & Robutel (1993) suggested that arguments based on the spin-orbit resonances can be still used to understand qualitative features of their results. Here we are interested in obliquity variations of rocky planets in habitable zone in extrasolar planetary systems where a gas giant planet(s) has been detected. Large obliquity variations, even if they are not chaotic, may inhibit habitability.

For this purpose, we re-analyze the spin-orbit resonance in more general form. In order to clarify fundamental processes of obliquity evolution, we study a system containing a host star, a hypothetical terrestrial planet and a hypothetical giant planet, in wide parameter ranges. In particular, we investigate the behavior of obliquity near a resonance through analytical arguments and numerical calculations of the secular precession equations. In section 2, we briefly summarize the spin-orbit resonance. We show results of numerical calculations in section 3. In section 4, resonance width is derived semi-analytically. In section 5, we briefly discuss the regions for a hypothetical terrestrial planet where its obliquity variations become large by the spin-orbit resonance.

## 2 Basic Equations and Model

We consider a system containing a host star, a hypothetical terrestrial planet with axisymmetric shape and negligible mass, and a hypothetical giant planet. We assume the two planets initially have circular orbits around a host star. We numerically solve the orbit-averaged Euler equations (secular precession equations) given by Ward (1974):

$$\frac{d\hat{\mathbf{s}}}{dt} = \alpha(\hat{\mathbf{s}} \cdot \hat{\mathbf{n}})(\hat{\mathbf{s}} \times \hat{\mathbf{n}}), \quad (1)$$

where  $\hat{\mathbf{s}}$  is a unit vector in the direction of the spin axis with components

$$s_x = \sin \theta \sin \psi, \quad (2)$$

$$s_y = -\sin \theta \cos \psi, \quad (3)$$

$$s_z = \cos \theta, \quad (4)$$

$\theta$  is the angle between the spin axis and the z-axis, and  $\psi$  is the longitude of the equator of the terrestrial planet in inertial frame. The precessional constant  $\alpha$  is given by

$$\alpha = \frac{3G(C - A)M_1}{2C\omega a^3}. \quad (5)$$

where  $G$  is the gravitational constant,  $M_1$  is the mass of the host star,  $(C - A)/C$ ,  $\omega$  and  $a$  are the dynamical ellipticity, the spin rate, and the semi-major axis of the terrestrial planet, respectively. In our calculation, we assume  $(C - A)/C$  and  $\omega$  are constants.

$\hat{\mathbf{n}}$  is a unit vector normal to the orbital plane,

$$n_x = \sin I \sin \Omega, \quad (6)$$

$$n_y = -\sin I \cos \Omega, \quad (7)$$

$$n_z = \cos I, \quad (8)$$

where  $I$  is the orbital inclination and  $\Omega$  is the longitude of the ascending node of the terrestrial planet. We here adopt the orbital plane of the giant planet as the reference frame. According to secular perturbation theories (e.g., Brouwer & Clemence 1961), the variations in  $I$  and  $\Omega$  are given as

$$I = \text{const.}, \quad (9)$$

$$\Omega = -Bt + \Omega_0, \quad (10)$$

where  $\Omega_0$  is the initial ascending node of the terrestrial planet.  $B$  is

$$B = n \frac{1}{4} \frac{M_2}{M_1} \alpha_2^2 b_{3/2}^{(1)}(\alpha_2), \quad (11)$$

where  $n$  is the mean motion of the terrestrial planet,  $M_2$  is the mass of the giant planet.  $\alpha_2$  is  $a/a_2$ , where  $a_2$  is the semi-major axis of the giant planet.  $b_{3/2}^{(1)}$  is a Laplace coefficient.

The obliquity of the terrestrial planet  $\epsilon$ , the angle between  $\hat{\mathbf{n}}$  and  $\hat{\mathbf{s}}$ , is obtained by

$$\hat{\mathbf{n}} \cdot \hat{\mathbf{s}} = \cos \epsilon. \quad (12)$$

The relationship between the reference plane, orbital plane and equator is schematically shown in Figure. 1.

Substitution of Eqs. (2) to (8) into (1) yields

$$\dot{\theta} \simeq \alpha \cos \epsilon I \sin(\psi - \Omega) + \mathcal{O}(I^2), \quad (13)$$

$$\dot{\psi} \simeq -\alpha \cos \epsilon + \mathcal{O}(I). \quad (14)$$

with the assumption  $I \ll 1$  (Ward 1974; note that definition of  $\psi$  is different). We denote the precession frequency of the spin axis  $\sim -\alpha \cos \epsilon$  by  $p_f$ . When  $I \ll 1$ ,  $\theta \simeq \epsilon$ . Since the sign of  $\dot{\theta}$  changes with frequency  $(\dot{\psi} - \dot{\Omega})$ ,  $\epsilon$  and  $\theta$  usually oscillate with amplitude  $\sim \mathcal{O}(I)$ . However, when  $\dot{\psi} - \dot{\Omega} \simeq -\alpha \cos \epsilon + B \simeq 0$ , oscillation period of  $\epsilon$  becomes very long and  $\epsilon$  has a large amplitude. This is the spin-orbit resonance.

$\hat{\mathbf{n}}$  has dependence as  $\hat{\mathbf{n}} = \hat{\mathbf{n}}(I, Bt)$ . If we scale time by  $\alpha^{-1}$ , Eq. (1) is transformed to

$$\frac{d\hat{\mathbf{s}}}{d\tilde{t}} = [\hat{\mathbf{s}} \cdot \hat{\mathbf{n}}(I, \frac{B}{\alpha}\tilde{t})][\hat{\mathbf{s}} \times \hat{\mathbf{n}}(I, \frac{B}{\alpha}\tilde{t})], \quad (15)$$

where  $\tilde{t} = \alpha t$ . This equation shows that the evolution path of  $\hat{\mathbf{s}}$  is dependent only on the values of  $I$  and  $B/\alpha$ . We will show that the evolution path near a resonance can be written as a form independent of  $I$  and  $B/\alpha$  with further scaling of  $\tilde{t}$  and  $\cos \epsilon$ .

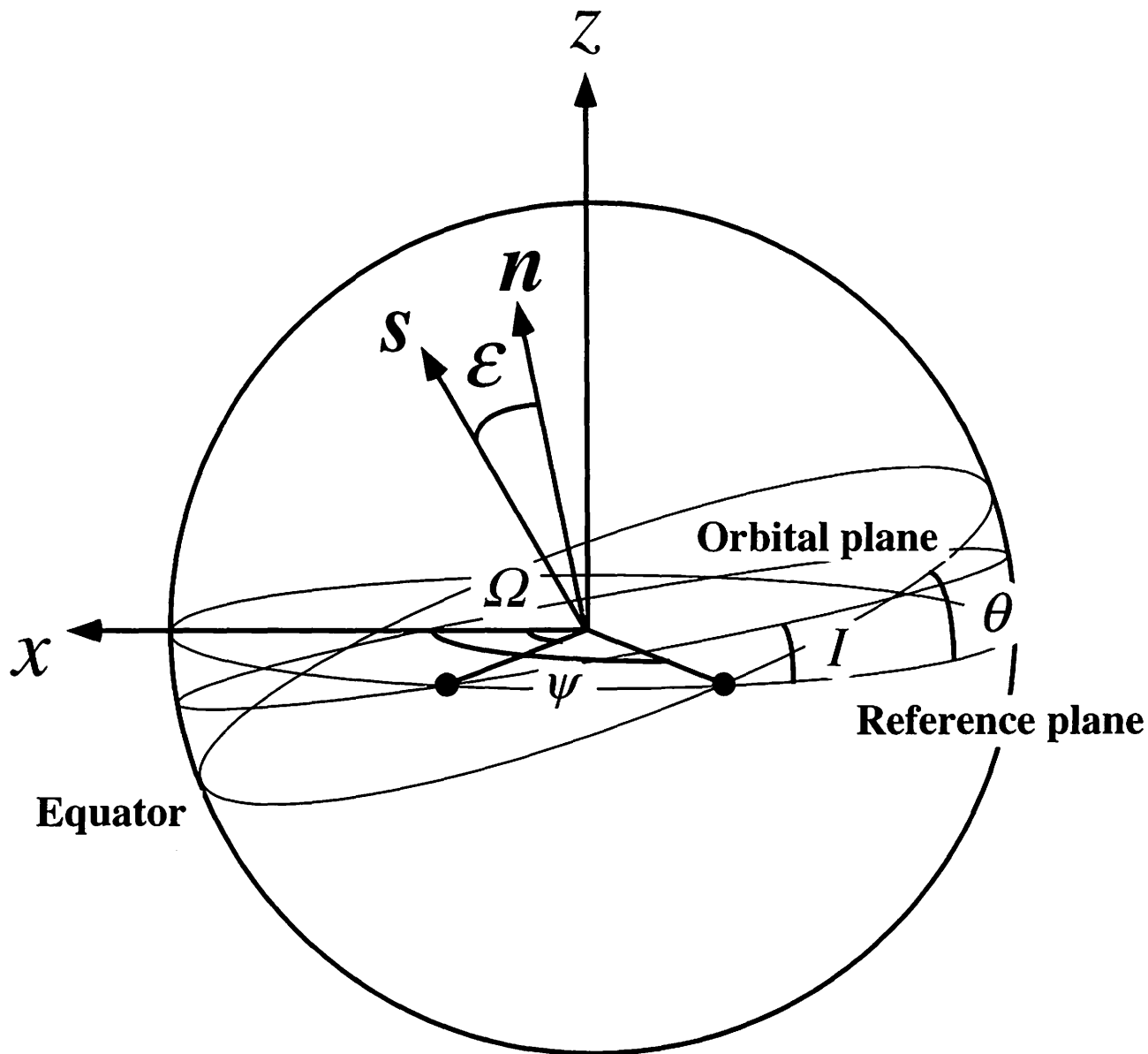


Figure 1: Relationship between the reference plane, orbital plane and equator.  $\hat{s}$  is the unit vector in the direction of the spin axis,  $\hat{n}$  is the unit vector normal to the orbital plane,  $\epsilon$  is the obliquity,  $\theta$  is the angle of the equator relative to the reference plane,  $\psi$  is the longitude of the equator,  $I$  is the orbital inclination, and  $\Omega$  is the longitude of the ascending node.

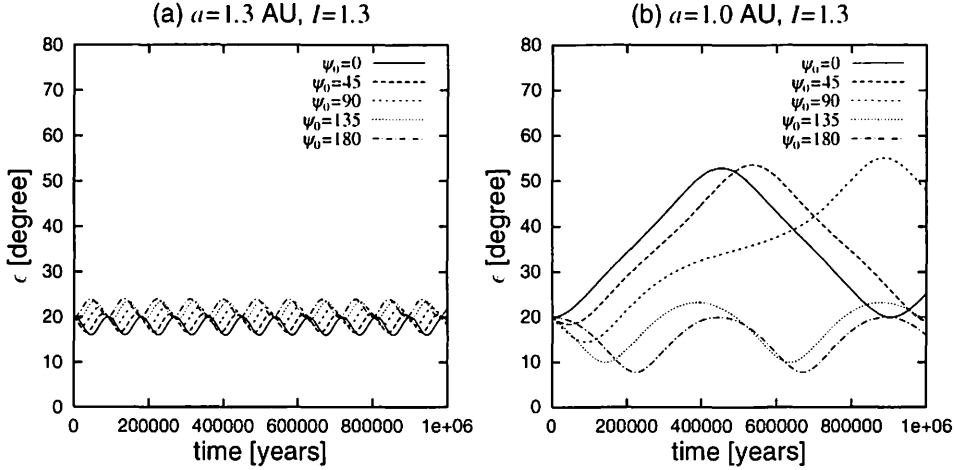


Figure 2: The time evolution of the obliquity  $\epsilon$  of the terrestrial planet.  $\psi_0$  is the initial longitude of the equator of the terrestrial planet,  $\psi$ .

### 3 Numerical Results

We integrate the precession equations (1) over  $10^6$  years, using a fourth order Runge-Kutta scheme. Here we adopt  $(C - A)/A = 0.00335$ ,  $\omega = 7.292 \times 10^{-5}$  rad/year, which are the same as the current Earth's values,  $M_1 = M_\odot$ ,  $a_2 = 5.2$  AU and  $M_2 = 2M_J$  ( $M_J$  is the mass of Jupiter).

Figure 2 shows examples of time evolution of  $\epsilon$ : (a) off-resonance case ( $a = 1.3$  AU and  $I = 1.3^\circ$ ), (b) the case of resonance ( $a = 1.0$  AU and  $I = 1.3^\circ$ ). Initial obliquity  $\epsilon_0$  is  $20^\circ$  in both cases.  $\psi_0 - \Omega_0$ , where  $\psi_0$  is initial precession angle, is  $0^\circ$ ,  $45^\circ$ ,  $90^\circ$ ,  $135^\circ$ , and  $180^\circ$ .  $\epsilon$  oscillates regularly with periods  $\sim 10^6$  years. In the off-resonance case, the variation amplitude is  $\sim \mathcal{O}(I)$ , while that in the resonance case is much larger than  $\mathcal{O}(I)$ .

To investigate the resonance width at  $a = 1.0$  AU, we did similar calculations for different initial obliquity from  $0^\circ$  to  $90^\circ$  with  $1^\circ$  step size. The minimum and maximum values of  $\epsilon$  plotted as a function of  $\epsilon_0$  in Figure 3 (the lower panel).  $p_f/\alpha$  and  $B/\alpha$  are also plotted in the upper panel. The resonance occurs when  $p_f/\alpha \simeq B/\alpha$ . Fig. 3 shows large resonance zone extending from  $15^\circ$  to  $55^\circ$  around exact resonant obliquity  $\epsilon_* \simeq 40^\circ$ , where  $\epsilon_*$  is defined by  $B/\alpha = \cos \epsilon_*$ . Even if it has the same initial obliquity, the amplitude depends on  $\psi_0 - \Omega_0$ .

To explain these features, we investigate evolution of  $x \equiv \psi - \Omega$  and its time derivative  $y \equiv (\dot{\psi} - \dot{\Omega})/\alpha\beta \simeq (-\alpha \cos \epsilon + B)/\alpha\beta$ , where  $\beta = \sqrt{I \cos \epsilon_* \sin \epsilon_*} = \sqrt{I(B/\alpha) \sqrt{1 - (B/\alpha)^2}}$  (Eqs. (10) and (14)). The meaning of the scaling factor  $\beta$  for  $y$  will be clear later. Since  $\alpha$  and  $B$  are fixed, the evolution of  $\epsilon$  is uniquely determined by that of  $y$ . For  $\psi_0 - \Omega_0 = \pi$ , time evolution of  $y$  for the results in Fig. 3 is shown in Figure 4. Trajectories start at  $\psi_0 - \Omega_0 = \pi$  with different  $\epsilon_0$ , that is different initial  $y (= y_0)$ . Trajectories with  $-2 \lesssim y_0 \lesssim 2$  show libration around the center of  $\alpha \cos \epsilon = B$  ( $y = 0$ ) and  $\psi - \Omega = \pi$  ( $x = \pi$ ), while the other trajectories show

circulation. The former cases are “resonance”. The trajectories with libration generally have large periodic variations, that is, large periodic obliquity variations. Note that as mentioned before, Eq. (1) is scaled with  $\alpha$  and hence the contour map of Fig. 4 holds for the cases which have the same  $B/\alpha$ . For the parameters adopted here,  $y = -2$  and  $2$  correspond to  $\epsilon \simeq 15^\circ$  and  $\simeq 55^\circ$ , respectively, which explains the results in the lower panel of Fig. 3. In this case, resonant width in  $y$  is  $\simeq 2$  for that in  $\epsilon \simeq 20^\circ$ .

For other  $I$  and  $B/\alpha (= \cos \epsilon_*)$ , the libration range in  $y$  can change in principle. Resonance width depends on  $\cos \epsilon_*$  and  $I$  as in Figures 5. Since  $\alpha$  and  $B$  are dependent on  $a$ , different  $a$  corresponds to different  $\cos \epsilon_*$ . From these results and the results with other  $\cos \epsilon_*$  and  $I$ , we found the resonance width changes approximately as  $\propto I^{1/2}$  and weakly depends on  $\cos \epsilon_*$ . However, we will show the evolution trajectories on the  $x - y$  plane do not change for other  $I$  and  $B/\alpha$ , at least near a resonance.

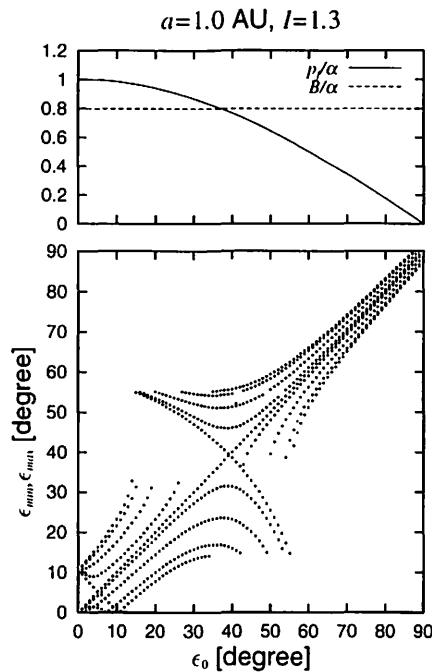


Figure 3: The lower panel is the minimum and maximum values of obliquity  $\epsilon$  over  $10^6$  years calculation as a function of initial obliquity  $\epsilon_0$ .  $p_f$  and  $B$  scaled by  $\alpha$  are plotted in the upper panel.

## 4 Analytical calculation

We derive an analytical solution to Eq. (1) near a resonance. Near a resonance ( $\alpha \cos \epsilon \simeq B$ ), Eqs. (2) – (4), (6) – (8), (13) and (14) lead to

$$\frac{d}{dt}x \simeq \alpha\beta y, \quad (16)$$

$$\frac{d}{dt}y = \frac{d}{dt}(\hat{\mathbf{s}} \cdot \hat{\mathbf{n}}) \simeq -\alpha\beta \sin x. \quad (17)$$

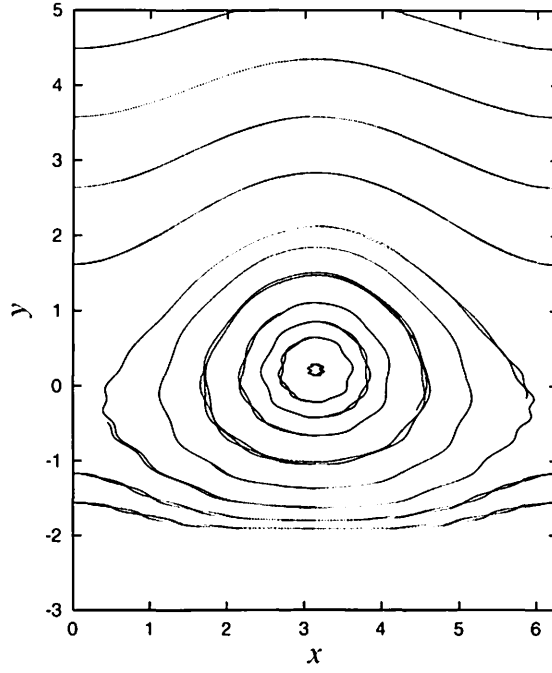


Figure 4: Trajectories on  $x - y$  plane for the results with  $\psi_0 - \Omega_0 = \pi$  in Fig. 3. Different trajectories correspond to runs with different  $y_0(\epsilon_0)$ .

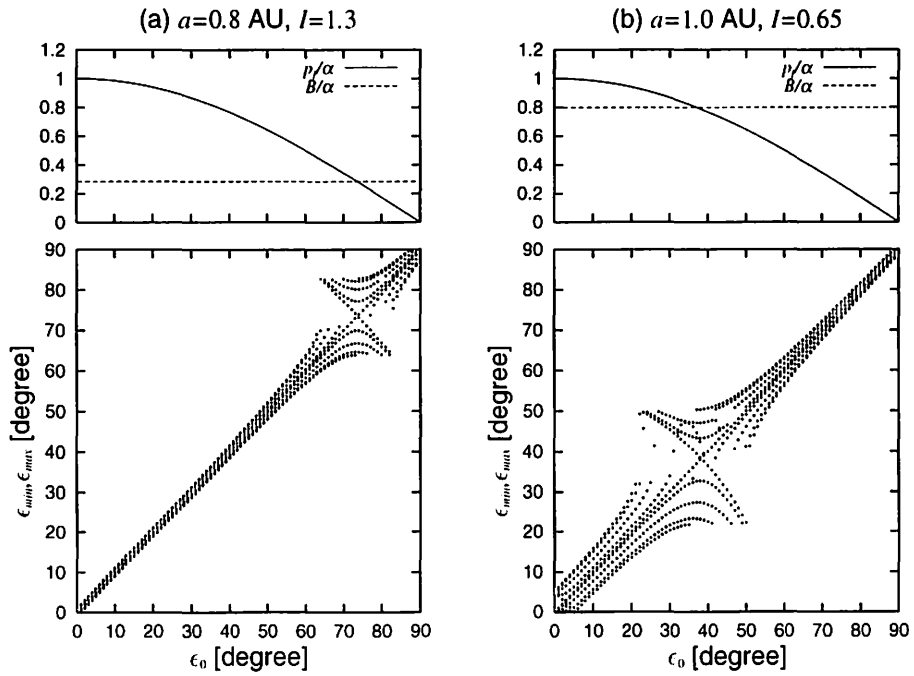


Figure 5: The same figures as Fig.2 expect for  $a$  or  $I$ . (a)  $a = 0.8$  AU and  $I = 1.3^\circ$ . (b)  $a = 1.0$  AU and  $I = 0.65^\circ$ .



In the right hand side of Eqs. (16) and (17), we retain the lowest order terms of  $I$ , assuming  $I \ll 1$ . With scaled variables

$$\tilde{t} = \alpha \beta t, \quad (18)$$

Eqs. (16) and (17) are reduced to

$$\frac{d}{d\tilde{t}}x \simeq y, \quad (19)$$

$$\frac{d}{d\tilde{t}}y \simeq \sin x, \quad (20)$$

which are independent of  $I$  and  $B/\alpha$ . The equilibrium points are  $y = 0$  and  $\sin x = 0$  ( $x = 0, \pi$ ).  $x = 0$  is unstable against small displacements.  $x = \pi$  is stable.

Let  $\delta x$  and  $\delta y$  be small displacements from the equilibrium point,  $x = \pi$  and  $y = 0$ . Eqs. (19) and (20) gives

$$\frac{d}{d\tilde{t}}\delta x \simeq \delta y, \quad (21)$$

$$\frac{d}{d\tilde{t}}\delta y \simeq -\delta x. \quad (22)$$

The solution is

$$\delta x = C \cos(\tilde{t} + \gamma), \quad (23)$$

where  $C$  and  $\gamma$  are constants of integration which are determined by initial conditions  $\psi_0$  and  $\epsilon_0$ . Substituting Eq. (23) into Eq. (21), we have

$$\delta y = C \sin(\tilde{t} + \gamma). \quad (24)$$

Eqs. (23) and (24) represent librating motion centered at  $x = \pi$  and  $y = 0$  with libration period  $2\pi$  in unit of  $\tilde{t}$ . The scaling factor  $\alpha\sqrt{I \cos \epsilon_* \sin \epsilon_*}$  for  $\tilde{t}$  appears to be angular velocity for the libration.

As mentioned in Sec. 3, if  $\alpha$ ,  $B$  and  $I$  are given,  $\epsilon$  uniquely corresponds to some value  $y$ . Therefore Figure 4 explains the evolution of  $\epsilon$ . If a starting point is on closed trajectory,  $x$  and  $y$  librate around the equilibrium point and  $\epsilon$  exhibits large variation. Equations (19) and (20) have an integration as

$$H = \frac{1}{2}y^2 + \cos x, \quad (25)$$

where  $H$  is constant. Different values of  $H$  corresponds to different trajectories. In Fig. 4, we start with  $\cos x = -1$ . Thus, the trajectory with  $y_0$  corresponds to the contour of  $H = y_0^2/2 - 1$ . Trajectories of libration correspond to  $-1 \leq \cos x \leq 1$  at  $y = 0$ , which is equivalent to  $-1 \leq H \leq 1$ . Since  $|y|$  takes a maximum value  $= \sqrt{2(H+1)}$  at  $\cos x = -1$  ( $x = \pi$ ) for each trajectory,

$$|y_{\max}| = 2. \quad (26)$$

In Fig. 6, these analytical estimates are compared with numerical results in Fig. 4, which agree with each other. Numerical results with other  $I$  and  $B/\alpha$  also agree with the analytical estimates.

Therefore we derived the resonance width semi-analytically as

$$|\delta \cos \epsilon|_{\max} \simeq 2\sqrt{I \cos \epsilon_* \sin \epsilon_*}. \quad (27)$$

The dependence of  $|\delta \cos \epsilon|_{\max}$  on  $\epsilon_*$  and  $I$  as well as the magnitude in the above almost completely agrees with the numerical results. Ward et al. (1979) derived  $|\delta \epsilon| \sim \sqrt{I/\tan \epsilon_*}$  at a resonance, through higher order expansion of Eq. (1), although the detailed derivation is not presented. Although his result does not have any dependence on  $\psi_0$  and  $\epsilon_0 - \epsilon_*$  and it includes uncertainty of a numerical factor, it is consistent with Eq. (27) in the limit of  $\delta \epsilon \rightarrow 0$  expect for a numerical factor.

## 5 Conclusion and Discussion

We have investigated the evolution of obliquity through analytical arguments and numerical calculations. We re-analyzed the spin-orbit resonance in a more general form. We considered a system containing a host star, a hypothetical terrestrial planet, and a hypothetical giant planet, and calculated the evolution of obliquity  $\epsilon$  of the terrestrial planet in wide ranges of parameters  $I, B, \alpha$ , and initial conditions of  $\epsilon$  and  $\psi - \Omega$ , where  $I$  is the orbital inclination,  $B$  is the frequency of the orbital variation,  $\alpha$  is the precessional constant,  $\epsilon$  is the obliquity,  $\psi$  is the precession angle, and  $\Omega$  is the longitude of the ascending node.

We found the following results:

1. Evolution of obliquity is described by a contour map on the plane of  $x = \psi - \Omega$  and  $y = (\dot{\psi} - \dot{\Omega})/\alpha \sqrt{I(B/\alpha) \sqrt{1 - (B/\alpha)^2}}$ . Different contours correspond to different initial conditions of  $\epsilon$  and  $\psi - \Omega$ . The contour map does not depend on  $I, B$ , and  $\alpha$ .
2. In the librating region centered at a resonant point,  $\cos \epsilon_* = B/\alpha$  and  $\psi - \Omega = \pi$  ( $x = \pi, y = 0$ ), the obliquity has variations with large amplitudes (Fig. 3 and 4).
3. The width of libration region is  $|\delta y|_{\max} = 2$ , which reads as

$$|\delta \cos \epsilon|_{\max} \simeq 2\sqrt{I \cos \epsilon_* \sin \epsilon_*}. \quad (28)$$

The range of the obliquity variation near a resonance  $\epsilon_*$  is

$$\cos^{-1}(\cos \epsilon_* + |\delta \cos \epsilon|_{\max}) \leq \epsilon \leq \cos^{-1}(\cos \epsilon_* - |\delta \cos \epsilon|_{\max}). \quad (29)$$

Note that the width of resonance region does not explicitly depend on the masses and semi-major axis of the terrestrial planet and the giant planet. In Figure 6, we plot this range as a function of  $a$ , in the case of  $I = 1.3^\circ$  and the other parameters given in section 4. If another giant planet is considered, the large variation regions are expressed by the superposition of spin-orbit resonances due to individual eigenfrequencies of orbital change of the terrestrial planet. If the resonant regions overlap, the obliquity variations could be chaotic.

Many extrasolar giant planets have been found around nearby solar-like stars. Substituting the masses and the semi-major axes of such planets, we can obtain the resonance regions where obliquity variations are large by the spin-orbit resonance if planetary ellipticities and spin rates are given. In a system with a giant planet with relatively large semi major axis, Earth-like planets (small rocky planets) may exist inside the orbits of the giant planet. For life to exist in such a Earth-like planet, the planet may need to have not only H<sub>2</sub>O ocean but also orbit and obliquity with small variation to keep the climate stable. Assuming probable values of ellipticity and the spin rate, we may evaluate the probability of existence of such “habitable” planets in extrasolar systems. We will address this issue in next paper.

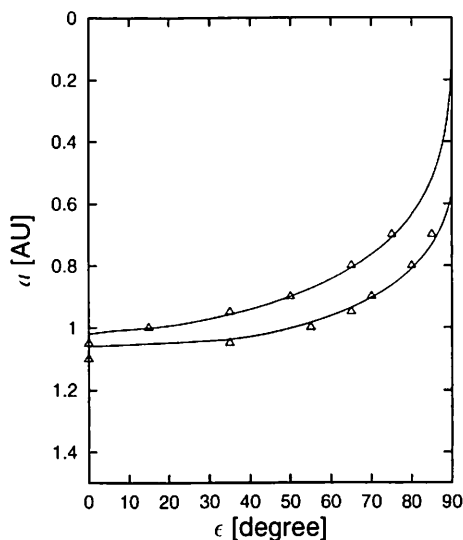


Figure 6: Resonance region in the case of  $I = 1.3^\circ$ . The area between the two solid lines expresses a resonance region. Triangles show the numerical results. Solid lines express the analytical expression given by Eq. (29)

## References

- [1] Laskar, J. 1988. Secular evolution of the Solar System over 10 million years. *Astron. Astrophys.* **198**, 341–362.
- [2] Laskar, J. 1990. The chaotic motion of the Solar System: a numerical estimate of the size of the chaotic zones. *Icarus* **88**, 266–291.
- [3] Laskar, J., and P. Robutel 1993. The chaotic obliquity of the planets. *Nature* **361**, 608–612.
- [4] Touma, J., and J. Wisdom 1993. The chaotic obliquity of Mars. *Science* **259**, 1294–1297.
- [5] Ward, W.R. 1974. Climatic Variations on Mars 1. Astronomical theory of insolation. *J. Geophys. Res.* **79**, 3375–3386.
- [6] Ward, W.R. 1979. Present obliquity oscillations of Mars: fourth-order accuracy in orbital  $e$  and  $I$ . *J. Geophys. Res.* **84**, 237–241.
- [7] Ward, W.R. 1992. Long-term orbital and spin dynamics of Mars. In *Mars*, p.298–320. Univ. Arizona Press, Tucson.
- [8] Ward, W.R., and D.J. Rudy 1991. Resonant obliquity of Mars? *Icarus* **94**, 160–164.

# Motion around Triangular Lagrange Points Perturbed by Other Bodies

Hideyoshi ARAKIDA\*, Toshio FUKUSHIMA†

## Abstract

We constructed an analytical theory of the motion of the test particle being linear with respect to the magnitude of departure from the Lagrange points, and taking into account the direct effects of the other perturbing bodies including the effects of their eccentricities up to the second order, in the planar restricted  $N$ -body problem. We compared our analytical solution with a numerical integration and confirmed that the solution represents the linear part of the true solution so well that the residuals are only due to the non-linear effect of the primary and the secondary system mainly which we ignored. The results will be useful in designing the orbit of near-future space missions to be located in the vicinity of the triangular Lagrange points.

## 1 Introduction

The five Lagrange points from  $L_1$  to  $L_5$  are special solutions of the three body problem. Especially the triangular Lagrange points,  $L_4$  and  $L_5$ , are important in terms of their linear stability because the centrifugal force and the Coriolis force balance. Since then, there have been vast investigations of the Lagrange points. Of course, they have high potentiality for the long-lived space missions and the astronautical applications. Most studies were developed in the framework of the restricted and general three body problems. In fact, the simplest approach to obtain an approximate solution is to linearize the equation of motion around the Lagrange point in the restricted circular and planar three body problem [44, 60]. This analytical solution quite well coincides with the result of numerical integration. Obviously the next step was to include additional physical effects; the non-linear effect by Bhatnagar [3], Deprit [12], Goździewski [32], Hagel [33], and Papadakis [63]; the effect of eccentricity of the primary and secondary bodies by Danby [9], Erdi [18], Içtiaroglou [41], Kinoshita [45], Selaru [70], and Todoran [81]; the effect of high inclination of the orbit of the test particle by Zhang *et al.* [91]; the effect of radiation pressure due to the primary body by Kumar [49], Lukyanov [53, 54], Ragos [67], Simmons [72], and Todoran [80]; the effect of dragging force by Murray [61]; the effect of  $J_2$  and other higher order gravitational field of the primary body by Kondurar [48], Shrivastava *et al.* [71], and Sharma *et al.* [76, 75]; the effect of variability of mass of the primary body by Horedt [37, 39], Horedt *et al.* [38], and Singh *et al.* [73]; the general relativistic effect of the primary body by Maindl [55] and Maindl *et al.* [56]; and the effect of electro-magnetic force of a charged primary body by Dionysiou [16].

On the other hand, the stability of the motion was analytically studied by Celletti [7], Danby [9], Deprit *et al.* [13, 15], Giorgilli [28], Howard [40], Kinoshita [46], Roels [68], Whipple [85], and Zagouras *et al.* [88]. Garfinkel considered the motion of Trojan asteroids, mainly taking notice of the tadpoles and horseshoe orbits in the three body problem ([22] to [27]).

While the long term behavior of the motion around the Lagrange points has been mostly studied by numerical integrations; the population of the long-lived asteroids by Melita *et al.*

---

\*E-Mail : h.arakida@nao.ac.jp

†E-Mail : Toshio.Fukushima@nao.ac.jp

[59], the search for stable orbits of planets by Weibel *et al.* [83], Innanen *et al.* [42], Erdi [18], Zhang *et al.* [89, 90], and Markellos *et al.* [57], and the escape probability from the triangular region by Tsiganis *et al.* [82]. Unfortunately, in the framework of  $N$ -body problem, there are few analytical researches, especially for the perturbations on the third and other perturbing bodies.

As we briefly summarized, the most of the analytical researches on the dynamical behavior around the Lagrange points have been conducted in the framework of the restricted or general three body problems. These treatises rather correspond to solve the free oscillation problem around an equilibrium point in terms of the oscillation dynamics. However, when we consider an actual problems in the  $N$ -body system such as our solar system, there exist not only the primary and the secondary bodies but other perturbing bodies. In the existing works, very few authors discussed the applicability of the results of the restricted and general three body problems to the actual  $N$ -body system. According to our preliminary numerical integrations, the obtained results indicate that for the same initial condition, quite large differences in the motion around  $L_4$  appears between the restricted three body problem and the restricted four body problem, which is the representative of restricted  $N$  body problem. Fig. 1 shows the orbit of a Trojan-like asteroid in the corotational coordinate system of the Sun-Jupiter system. In the experiment, we assumed that both Jupiter and Saturn are moving on the circular orbit. We put Saturn on the positive  $X$  axis at  $t = 0$ . Obviously the solution in the case of Sun-Jupiter-Saturn system is quite different from that of Sun-Jupiter system; about 10 times larger than the latter in the magnitude. Therefore the result obtained from the three body problem does not become the first approximation of the  $N$ -body problem. The differences are no other than the effect

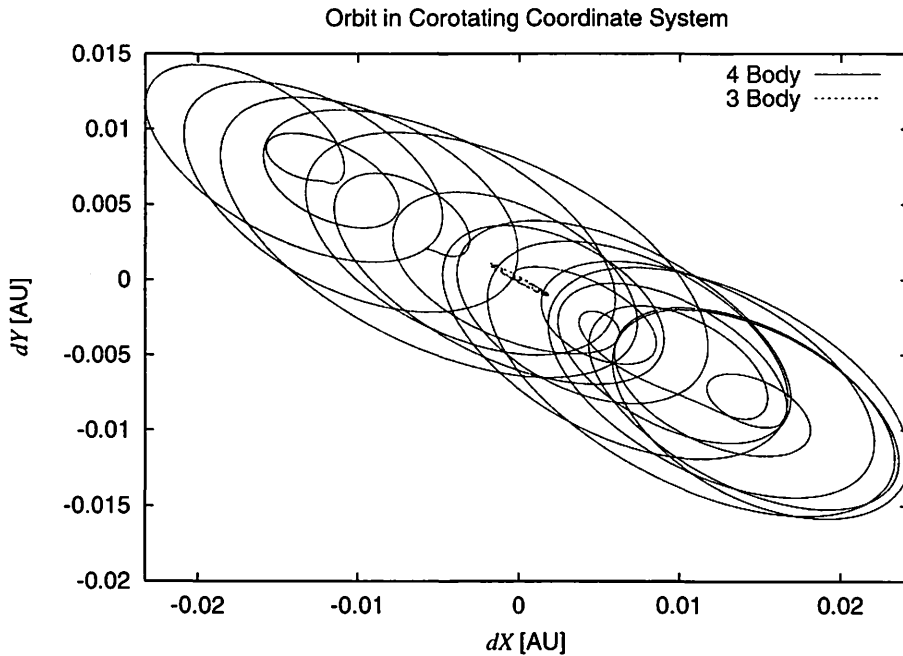


Figure 1: Orbits of Trojan-like asteroid in the corotational coordinate system.

of forced oscillation term due to Saturn. Of course, the dominant gravity force acting on the asteroid is evidently those due to the Sun and Jupiter. However the closer to  $L_4$  the asteroid is, the smaller the influence of the gravity of these two bodies becomes. Therefore, in the vicinity of  $L_4$ , the effect of Saturn plays the key role. In terms of oscillation dynamics, we must first consider the forced oscillation caused by the external bodies like Saturn in this case.

As we mentioned earlier, there are only a few investigations that dealt analytically with the

	Our Solution	Gómez <i>et al.</i> (2001)
Main interests	Sun-Planetary system	Earth-Moon system
Approach	Purely analytical	Semi-analytical
	Expansion of disturbing force	Lie transformation
Non-linear effect	No	Yes
Eccentricity of primary and secondary	No	No
Eccentricity of other bodies	Second order	No
Comparison with numerical integration	Yes	No
Discussion of application limit	Yes	No
Global aspect of orbit	Yes	No

Table 1: Comparison of our theory and that of Gómez *et al.* (2001)

dynamical behaviors around Lagrange points in the framework of  $N$ -body problem (and Gómez *et al.* [29, 30, 31]). In the series of their study, they tried to investigate motion of the collinear and the triangular Lagrange points by both the numerical and the semi-analytical. In their works, the corresponding approach to our study is the bicircular problem in which the motion of test particle moves under the gravitational forces of Earth, Moon, and Sun. They assumed that the orbit of Moon around the Earth is circular, and the Sun moves around the Earth-Moon barycenter in another circular orbit. In that simplified model, they interpreted the procedure to obtain the analytical solution based on the Lie transformation in coordinate. However, they did not solve it completely, and the several figures inserted in their textbook were described with the aid of some numerical processes. Hence their approach is not satisfactory in many points as illustrated in Table 1. Thus as yet, there does exist the purely analytical solution around the Lagrange point. In this work, we will construct a purely analytical theory of the motion of the test particle around the triangular Lagrange points in the framework of  $N$ -body system. We include the effect of direct gravitational force of the third and other perturbing bodies. And we express the solution as an explicit function of time. Then we compare our analytical solution with the results of numerical integration. It is also beneficial to construct an analytical theory, especially for designing the orbits of some space missions. As we mentioned before, some space missions located on the triangular Lagrange are planned. But the present orbital design must carry out the vast of numerical integration in the huge initial condition space. Then the analytical solution is expected that it reduces the considerable time of the numerical integration and restricts the initial condition of numerical integration. The re-evaluation of the orbital region is much easy especially when the spacecraft was put into a wrong orbit. Then it is possible to provide the preliminary constrain for the initial condition of the orbital design.

## 2 Analytical Solution

Let us construct an analytical theory of the motion of a test particle around the triangular Lagrangian point when there exist the perturbations due to extra perturbing bodies. We assume that (1) the orbits of all bodies are coplanar, (2) the orbits of the primary and the secondary around their barycenter are circular, and (3) the extra perturbing bodies are moving around the barycenter of the primary and secondary in non-circular Keplerian orbit. As we showed in Fig. 1, the solution of the three body problem does not become the first approximation of the motion of  $N$ -body problem. Therefore we can no longer regard the effect of the other bodies as the “perturbation” around the Lagrange point as the usual approach of the perturbation theory. Thus we start from the equation of motion in the inertial coordinate system and expand it around the Lagrange point in the coordinate. Then we estimate the magnitude of the expanded

terms and include the dominate terms.

The equation of motion of the test particle is written in the inertial coordinate system as,

$$\begin{aligned}\frac{d^2 \mathbf{r}}{dt^2} &= \mathbf{F}_L(\mathbf{r}, t) + \mathbf{F}_E(\mathbf{r}, t) \\ \mathbf{r} &= \mathbf{r}_0, \quad \dot{\mathbf{r}} = \dot{\mathbf{r}}_0, \quad t = 0,\end{aligned}\tag{1}$$

where

$$\mathbf{F}_L(\mathbf{r}, t) = GM_1 \frac{\mathbf{r}_1 - \mathbf{r}}{|\mathbf{r}_1 - \mathbf{r}|^3} + GM_2 \frac{\mathbf{r}_2 - \mathbf{r}}{|\mathbf{r}_2 - \mathbf{r}|^3}\tag{2}$$

$$\mathbf{F}_E(\mathbf{r}, t) = \sum_{I=3}^N GM_I \frac{\mathbf{r}_I - \mathbf{r}}{|\mathbf{r}_I - \mathbf{r}|^3}.\tag{3}$$

The subscript  $I$  denotes the perturbing bodies,  $G$  is Newton's gravitational constant, and  $M_I$  is the mass of perturbing bodies. In Eq. (1), we separated the perturbing force into two parts. The one is the net effect of the primary and the secondary bodies which vanishes at the Lagrange points in the corotational coordinate frame. The other is the contribution due to the extra perturbing bodies which remain finite at the Lagrange point. Let us introduce a new variable  $\delta \mathbf{r} = \mathbf{r} - \mathbf{r}_L$  where  $\mathbf{r}_L$  denotes the Lagrange point, Then we expand Eq. (1) around  $\mathbf{r}_L$  as,

$$\begin{aligned}\frac{d^2(\mathbf{r}_L + \delta \mathbf{r})}{dt^2} &= \mathbf{F}_L(\mathbf{r}_L + \delta \mathbf{r}, t) + \mathbf{F}_E(\mathbf{r}_L + \delta \mathbf{r}, t) \\ &\approx \mathbf{F}_L(\mathbf{r}_L, t) + \mathbf{F}_E(\mathbf{r}_L, t) + \left( \frac{\partial \mathbf{F}_L(\mathbf{r}, t)}{\partial \mathbf{r}} \right)_{\mathbf{r}_L} \delta \mathbf{r} \\ &\quad + \left( \frac{\partial \mathbf{F}_E(\mathbf{r}, t)}{\partial \mathbf{r}} \right)_{\mathbf{r}_L} \delta \mathbf{r},\end{aligned}\tag{4}$$

Noting that the Lagrange point satisfies the relation,

$$\frac{d^2 \mathbf{r}_L}{dt^2} - \mathbf{F}_L(\mathbf{r}_L, t) = 0,$$

we rewrite the above equation of motion as.

$$\begin{aligned}\frac{d^2 \delta \mathbf{r}}{dt^2} &= \mathbf{F}_E(\mathbf{r}_L, t) + \left( \frac{\partial \mathbf{F}_L(\mathbf{r}, t)}{\partial \mathbf{r}} \right)_{\mathbf{r}_L} \delta \mathbf{r} + \left( \frac{\partial \mathbf{F}_E(\mathbf{r}, t)}{\partial \mathbf{r}} \right)_{\mathbf{r}_L} \delta \mathbf{r} \\ &\approx \mathbf{F}_E(\mathbf{r}_L, t) + \left( \frac{\partial \mathbf{F}_L(\mathbf{r}, t)}{\partial \mathbf{r}} \right)_{\mathbf{r}_L} \delta \mathbf{r}.\end{aligned}\tag{5}$$

Here we estimate the magnitude of each term of Eq. (5) in the Sun-Jupiter-Saturn system. The ratios of  $(\partial \mathbf{F}_E(\mathbf{r}, t)/\partial \mathbf{r})_{\mathbf{r}_L} \delta \mathbf{r}/\mathbf{F}_E(\mathbf{r}_L, t)$  and  $(\partial \mathbf{F}_E(\mathbf{r}, t)/\partial \mathbf{r})_{\mathbf{r}_L} \delta \mathbf{r}/\mathbf{F}_E(\mathbf{r}_L, t)$  are,

$$\begin{aligned}\frac{(\partial \mathbf{F}_E(\mathbf{r}, t)/\partial \mathbf{r})_{\mathbf{r}_L} \delta \mathbf{r}}{\mathbf{F}_E(\mathbf{r}_L, t)} &= 0.014\epsilon \\ \frac{(\partial \mathbf{F}_L(\mathbf{r}, t)/\partial \mathbf{r})_{\mathbf{r}_L} \delta \mathbf{r}}{\mathbf{F}_E(\mathbf{r}_L, t)} &= 4.6 \times 10^{-5}\epsilon,\end{aligned}$$

where

$$\epsilon = \frac{\delta r}{r_L}.$$

Therefore we ignore the third term in the first line in Eq. (5). Next we assume that the solution is split as,

$$\delta \mathbf{r} = \delta \mathbf{r}_L + \delta \mathbf{r}_E,\tag{6}$$

where  $\delta \mathbf{r}_L$  and  $\delta \mathbf{r}_E$  satisfy the following equations, respectively,

$$\frac{d^2 \delta \mathbf{r}_L}{dt^2} = \left( \frac{\partial \mathbf{F}_L}{\partial \mathbf{r}} \right)_{\mathbf{r}_L} \delta \mathbf{r}_L, \quad \delta \mathbf{r}_L = \delta \mathbf{r}_0, \delta \dot{\mathbf{r}}_L = \delta \dot{\mathbf{r}}_0, \text{ at } t = 0, \quad (7)$$

$$\frac{d^2 \delta \mathbf{r}_E}{dt^2} = \mathbf{F}(\mathbf{r}_L, t)_E + \left( \frac{\partial \mathbf{F}(\mathbf{r}, t)_L}{\partial \mathbf{r}} \right)_{\mathbf{r}_L} \delta \mathbf{r}_E, \quad \delta \mathbf{r}_E = 0, \delta \dot{\mathbf{r}}_E = 0, \text{ at } t = 0. \quad (8)$$

In the following sections, we will specifically derive the expression of the solution.

## 2.1 Solution of Free Oscillation

First, we derive the solution of the free oscillation part in the inertial coordinate system. The equation of motion of  $\delta \mathbf{r}_L$  is expressed as

$$\frac{\partial^2 \delta \mathbf{r}_L}{dt^2} = \left( \frac{\partial \mathbf{F}(\mathbf{r}, t)}{\partial \mathbf{r}} \right)_{\mathbf{r}_L} \delta \mathbf{r}_L. \quad (9)$$

Usually the solution is given in the corotational coordinate system. See the derivation given in the literature such as Kinoshita [44] or Murray & Dermott [60]. The equation of motion in the case of the restricted three body problem is expressed in the corotational coordinate system as,

$$\frac{d^2 \delta \bar{X}}{dt^2} - 2n \frac{d\delta \bar{Y}}{dt} = \left( \frac{\partial^2 U}{\partial X^2} \right)_{\mathbf{r}_L} \delta \bar{X}, \quad (10)$$

$$\frac{d^2 \delta \bar{Y}}{dt^2} + 2n \frac{d\delta \bar{X}}{dt} = \left( \frac{\partial^2 U}{\partial Y^2} \right)_{\mathbf{r}_L} \delta \bar{Y}, \quad (11)$$

in which  $U$  is the potential,

$$U = -\frac{1}{2}n^2(X^2 + Y^2) - \frac{GM_1}{|\mathbf{r}_1 - \mathbf{r}|} - \frac{GM_2}{|\mathbf{r}_2 - \mathbf{r}|}. \quad (12)$$

The solution has a form of harmonic oscillator of two modes,

$$\delta \bar{X} = \sum_{\beta=1}^2 C_{\beta} \cos(\omega_{\beta} t + \gamma_{\beta}), \quad \delta \bar{Y} = \sum_{\beta=1}^2 S_{\beta} \sin(\omega_{\beta} t + \gamma_{\beta}). \quad (13)$$

where  $\omega_1, \omega_2$  are the eigenfrequencies expressed as,

$$\omega_1 = \sqrt{\frac{1}{2} \left\{ 1 + \sqrt{1 - 27\nu(1 - \nu)} \right\}} n \quad (14)$$

$$\omega_2 = \sqrt{\frac{1}{2} \left\{ 1 - \sqrt{1 - 27\nu(1 - \nu)} \right\}} n, \quad (15)$$

where

$$\nu = \frac{M_2}{M_1 + M_2}, \quad n = \sqrt{\frac{G(M_1 + M_2)}{a^3}}.$$

In the case of Sun-Jupiter system,  $\omega_1$  and  $\omega_2$  have the period of 11.901 year and 147.42 year, respectively. The amplitude and initial phase depend on the initial condition. We put  $(\delta \bar{x}, \delta \bar{y}, \delta \dot{\bar{x}}, \delta \dot{\bar{y}}) = (\delta \bar{x}_0, \delta \bar{y}_0, \delta \dot{\bar{x}}_0, \delta \dot{\bar{y}}_0)$  at  $t = 0$  as the initial condition and the amplitude and initial phase are given by

$$C_1 = \frac{2n\delta \bar{y}_0 + (\omega_2^2 - a^*)\delta \bar{x}_0}{(\omega_2^2 - \omega_1^2)} \sqrt{\frac{a^{*2}(2n\delta \bar{y}_0 + (\omega_2^2 - a^*)\delta \bar{x}_0)^2}{2n\omega_1\omega_2^2\delta \bar{y}_0 + \omega_1(a^* - \omega_2^2)\delta \bar{x}_0^2} + 1} \quad (16)$$



$$C_2 = \frac{2n\omega_1^2\omega_2\delta\bar{y}_0 + (a^* - \omega_1^2)\omega_2\delta\bar{x}_0}{a^*(\omega_1^2 - \omega_2^2)} \sqrt{\frac{a^{*2}(2n\delta\bar{y}_0 + (\omega_1^2 - a^*)\delta\bar{x}_0)^2}{2n\omega_1^2\omega_2\delta\bar{y}_0 + (a^* - \omega_1^2)\omega_2\delta\bar{x}_0^2} + 1} \quad (17)$$

$$\gamma_1 = \arctan\left(\frac{-(\omega_2^2 - a^*)\omega_1\delta\dot{\bar{x}}_0 + 2n\omega_1\omega_2^2\delta\bar{y}_0}{a^*(2n\delta\dot{\bar{y}}_0 + (\omega_2^2 - a^*)\delta\bar{x}_0)}\right) \quad (18)$$

$$\gamma_2 = \arctan\left(\frac{-(\omega_1^2 - a^*)\omega_2\delta\dot{\bar{x}}_0 + 2n\omega_1^2\omega_2\delta\bar{y}_0}{a^*(2n\delta\dot{\bar{y}}_0 + (\omega_1^2 - a^*)\delta\bar{x}_0)}\right) \quad (19)$$

$$\frac{S_\beta}{C_\beta} = -\frac{\omega_\beta^2 - a^*}{2n\omega_\beta}, \quad \beta = 1, 2, \quad (20)$$

where

$$a^* = -\frac{3}{2}(1 - \sqrt{1 - 3\nu(1 - \nu)})n^2.$$

Thus the solution is expressed in the inertial system as

$$x = X \cos \ell - Y \sin \ell, \quad y = X \sin \ell + Y \cos \ell \quad (21)$$

where  $\ell$  is the longitude of the secondary body and

$$\begin{aligned} X &= a \left\{ \cos\left(\frac{\pi}{3}\right) - \mu + \delta\bar{X} \cos \alpha - \delta\bar{Y} \sin \alpha \right\} \\ Y &= a \left\{ \sin\left(\frac{\pi}{3}\right) + \delta\bar{X} \sin \alpha + \delta\bar{Y} \cos \alpha \right\}. \end{aligned}$$

Here  $\alpha$  is the rotation angle and defined by the relation.

$$\alpha = -\frac{\arctan(-\sqrt{3}(1 - 2\nu))}{2}.$$

Note that the the effect of the initial condition of the motion is completely absorbed by the solution of free oscillation.

## 2.2 Solution of Forced Oscillation

From Eq. (5), the equation of motion of  $\delta\mathbf{r}_E$  becomes

$$\frac{d^2\delta\mathbf{r}_E}{dt^2} = \mathbf{F}_E(\mathbf{r}_L, t) + \left(\frac{\partial\mathbf{F}_L(\mathbf{r}, t)}{\partial\mathbf{r}}\right)_{\mathbf{r}_L} \delta\mathbf{r}_E. \quad (22)$$

As we examined the order of terms in the right-hand side, the first term is the main part of the forced oscillation. In evaluating it, we consider the effect of eccentricity of the extra bodies up to the second order. Note that the eccentricity of Saturn is as small as  $e = 0.0555$ . We expand the position of the perturbing bodies in the coordinate system where  $x$ -axis is the direction of the perihelion, up to the order of  $e_I^2$  as

$$\frac{\eta_I}{a_I} \equiv \sqrt{1 - e_I^2} \sin u_I = \left(1 - \frac{5}{8}e_I^2\right) \sin \ell_I - \frac{1}{8}e_I \sin 2\ell_I + \frac{3}{8}e_I^2 \sin 3\ell_I + \dots \quad (23)$$

$$\frac{\xi_I}{a_I} \equiv \cos u_I = -\frac{1}{2}e_I + \left(1 - \frac{3}{8}e_I^2\right) \cos \ell_I + \frac{1}{2}e_I \cos 2\ell_I + \frac{3}{8}e_I^2 \cos 3\ell_I + \dots \quad (24)$$

By rotating them by the angle  $\varpi$ , the longitude of perihelion, we obtain the expression in the inertial coordinate system as,

$$x_I = \xi_I \cos \varpi - \eta_I \sin \varpi, \quad y_I = \xi_I \sin \varpi + \eta_I \cos \varpi. \quad (25)$$

Since we assume that the orbit of the secondary body is circular, we approximate

$$\begin{aligned}\frac{r}{r_I} &= \left(\frac{a}{a_I}\right) \frac{1}{1 - e_I \cos u_I} \approx \frac{a}{a_I} (1 + e_I \cos u_I + e_I^2 \cos^2 u_I) + \dots \\ &\approx \frac{a}{a_I} \left(1 - \frac{e_I^2}{2} + e_I \cos \ell_I + \frac{3}{8} e_I^2 \cos 2\ell_I\right).\end{aligned}\quad (26)$$

Then we expand the denominator of  $F_E$  by using Legendre polynomials up to the order of  $(r/r_I)^7$  as,

$$\begin{aligned}\frac{1}{|\mathbf{r}_I - \mathbf{r}|^3} &= \frac{1}{r_I^3} \left[ 1 + 3 \cos S_i \frac{r}{r_I} + \frac{3}{2} (5 \cos^2 S_I - 1) \left(\frac{r}{r_I}\right)^2 \right. \\ &\quad + \frac{5}{2} (7 \cos^3 S_I - 3 \cos S_I) \left(\frac{r}{r_I}\right)^3 \\ &\quad + \frac{15}{8} (21 \cos^4 S_I - 14 \cos^2 S_I + 1) \left(\frac{r}{r_I}\right)^4 \\ &\quad + \frac{21}{8} (33 \cos^5 S_I - 30 \cos^3 S_I + 5 \cos S_I) \left(\frac{r}{r_I}\right)^5 \\ &\quad + \frac{7}{17} (429 \cos^6 S_I - 495 \cos^4 S_I + 135 \cos^2 S_I - 5) \left(\frac{r}{r_I}\right)^6 \\ &\quad \left. + \frac{9}{17} (715 \cos^7 S_I - 1001 \cos^5 S_I + 385 \cos^3 S_I - 35 \cos S_I) \left(\frac{r}{r_I}\right)^7 \right],\end{aligned}\quad (27)$$

where  $S_I = \ell_I - \ell = n_I(t - t_{I0}) - n(t - t_0)$ ,  $n_I$  is the mean motion of  $I$ -th perturbing body, and  $t_0$ ,  $t_{I0}$  are the times of perihelion passage of Lagrange point and  $I$ -th perturbing body, respectively. We expanded  $1/|\mathbf{r}_I - \mathbf{r}|^3$  up to the order of  $(r/r_I)^7$  and ignore the third and higher terms with respect to  $e_I$ . Since we assume that the orbit of the secondary body is circular, the position of the triangular Lagrange point ( $L_4$ ) becomes

$$x_L = a \left( \cos \left( \ell + \frac{\pi}{3} \right) - \mu \cos \ell \right), \quad y_L = a \left( \sin \left( \ell + \frac{\pi}{3} \right) - \mu \sin \ell \right) \quad (28)$$

Thus the first term of the right-hand side of Eq. (22) is explicitly expressed as the function of time  $t$  as,

$$\begin{aligned}\begin{pmatrix} F_x \\ F_y \end{pmatrix} &= \frac{1}{a_I^3} \left[ 1 + 3 \cos S_i \left( 1 - \frac{e_I^2}{2} + e_I \cos \ell + \frac{3}{8} e_I^2 \cos 2\ell \right) \frac{a}{a_I} \right. \\ &\quad + \frac{3}{2} (5 \cos^2 S_i - 1) \left( 1 - \frac{2e_I^2}{2} + 2e_I \cos \ell + \frac{6}{8} e_I^2 \cos 2\ell \right) \left(\frac{a}{a_I}\right)^2 \\ &\quad + \frac{5}{2} (7 \cos^3 S_i - 3 \cos S_i) \left(\frac{a}{a_I}\right)^3 \\ &\quad + \frac{15}{8} (21 \cos^4 S - 14 \cos^2 S + 1) \left(\frac{a}{a_I}\right)^4 \\ &\quad + \frac{21}{8} (33 \cos^5 S - 30 \cos^3 S + 5 \cos S) \left(\frac{a}{a_I}\right)^5 \\ &\quad + \frac{7}{17} (429 \cos^6 S - 495 \cos^4 S + 135 \cos^2 S - 5) \left(\frac{a}{a_I}\right)^6 \\ &\quad \left. + \frac{9}{17} (715 \cos^7 S - 1001 \cos^5 S + 385 \cos^3 S - 35 \cos S) \left(\frac{a}{a_I}\right)^7 \right]\end{aligned}$$

$$\times \begin{pmatrix} a_I \left[ \left\{ \left( 1 - \frac{5}{8}e_I^2 \right) \sin \ell_I - \frac{1}{8}e_I \sin 2\ell_I + \frac{3}{8}e_I^2 \sin 3\ell_I \right\} \cos \varpi \right. \\ \left. + \left\{ \frac{1}{2}e_I - \left( 1 - \frac{3}{8}e_I^2 \right) \cos \ell_I - \frac{1}{2}e_I \cos 2\ell_I - \frac{3}{8}e_I^2 \cos 3\ell_I \right\} \sin \varpi \right] \\ - a \left( \cos \left( \ell + \frac{\pi}{3} \right) - \mu \cos \ell \right) \\ a_I \left[ \left\{ \left( 1 - \frac{5}{8}e_I^2 \right) \sin \ell_I - \frac{1}{8}e_I \sin 2\ell_I + \frac{3}{8}e_I^2 \sin 3\ell_I \right\} \sin \varpi \right. \\ \left. - \left\{ \frac{1}{2}e_I - \left( 1 - \frac{3}{8}e_I^2 \right) \cos \ell_I - \frac{1}{2}e_I \cos 2\ell_I - \frac{3}{8}e_I^2 \cos 3\ell_I \right\} \cos \varpi \right] \\ \left. - a \left( \sin \left( \ell + \frac{\pi}{3} \right) - \mu \sin \ell \right) \right]. \quad (29)$$

The partial derivative in Eq. (22) is given by

$$\frac{\partial F_{Lj}(\mathbf{r}, t)}{\partial r_k} = \sum_{I=1}^2 GM_I \left[ -\frac{\delta_{jk}}{|r_{Ik} - r_k|^3} + \frac{3(r_{Ij} - r_j) \otimes (r_{Ik} - r_k)}{|r_{Ik} - r_k|^5} \right] \quad (30)$$

where the subscripts  $j$  and  $k$  represent the coordinate and  $\delta_{jk}$  is the Kronecker's delta. The components of the right-hand side of Eq. (30) become

$$\frac{\partial F_{Lx}}{\partial x} = \frac{1}{a^3} \left[ \frac{\mu}{2} \left\{ 1 + 3 \cos \left( 2 \left( \ell + \frac{\pi}{3} \right) \right) \right\} + 3GM_2 \left\{ \frac{\cos 2\ell}{2} - \cos \left( 2\ell + \frac{\pi}{3} \right) \right\} \right] \quad (31)$$

$$\frac{\partial F_{Ly}}{\partial y} = \frac{1}{a^3} \left[ \frac{\mu}{2} \left\{ 1 - 3 \cos \left( 2 \left( \ell + \frac{\pi}{3} \right) \right) \right\} - 3GM_2 \left\{ \frac{\cos 2\ell}{2} - \cos \left( 2\ell + \frac{\pi}{3} \right) \right\} \right] \quad (32)$$

$$\frac{\partial F_{Lx}}{\partial y} = \frac{3}{a^3} \left[ \frac{\mu}{2} \sin \left( 2 \left( \ell + \frac{\pi}{3} \right) \right) + GM_2 \left\{ \frac{\sin 2\ell}{2} - \sin \left( 2\ell - \frac{\pi}{3} \right) \right\} \right] \quad (33)$$

$$\frac{\partial F_{Ly}}{\partial x} = \frac{\partial F_{Lx}}{\partial y} \quad (34)$$

Now it has become straightforward to solve Eq. (22) since the main term contains only  $t$ . Actually we solve it iteratively. Namely we expand the solution as,

$$\delta \mathbf{r}_E \equiv \sum_{n=0}^{\infty} \delta \mathbf{r}_E^{(n)} = \delta \mathbf{r}_E^{(0)} + \delta \mathbf{r}_E^{(1)} + \delta \mathbf{r}_E^{(2)} + \dots, \quad (35)$$

where  $\delta \mathbf{r}_E^{(n)}$  satisfy the following equations

$$\frac{d^2 \delta \mathbf{r}_E^{(0)}}{dt^2} = \mathbf{F}_E(\mathbf{r}_L, t), \quad (36)$$

$$\frac{d^2 \delta \mathbf{r}_E^{(1)}}{dt^2} = \frac{\partial \mathbf{F}_L(\mathbf{r}_L, t)}{\partial \mathbf{r}} \delta \mathbf{r}_E^{(0)}, \quad (37)$$

$$\frac{d^2 \delta \mathbf{r}_E^{(2)}}{dt^2} = \frac{\partial \mathbf{F}_L(\mathbf{r}_L, t)}{\partial \mathbf{r}} \delta \mathbf{r}_E^{(1)}, \quad (38)$$

...

These equations are directly solved by the double integration as

$$\delta \mathbf{r}_E^{(0)} = \int \left[ \int \mathbf{F}_E(\mathbf{r}_L, t) dt \right] dt, \quad (39)$$

$$\delta \mathbf{r}_E^{(1)} = \int \left[ \int \left( \frac{\partial \mathbf{F}_L(\mathbf{r}, t)}{\partial \mathbf{r}} \right)_{\mathbf{r}_L} \delta \mathbf{r}_E^{(0)} dt \right] dt, \quad (40)$$

...

Here we only consider  $\delta \mathbf{r}_E^{(0)}$  and  $\delta \mathbf{r}_E^{(1)}$  because  $\delta \mathbf{r}_E^{(2)}$  and the higher correspond to the solve the nonlinear terms. Note that the initial condition is given in the form  $\delta \mathbf{r}_E^{(j)} = 0$  at  $t = 0$  since the initial conditions of the whole equation of motion is satisfied by the solution of the free

oscillation part. In the above expression, we ignored the second and higher terms. Thus the final solution is expressed in the series of  $a/a_I$  as,

$$\delta \mathbf{r}_E = \sum_{j=0}^{\infty} \left[ \delta \mathbf{r}_E^{(0j)} + \delta \mathbf{r}_E^{(1j)} \right] \left( \frac{a}{a_I} \right)^j. \quad (41)$$

Actually we expanded the series up to  $(a/a_I)^7$ .

### 3 Numerical Comparisons

From now on, we will show the comparisons between our analytical solution and the numerical integrations. As a test problem of the  $N$ -body system, we consider the motion of test particle in the vicinity of  $L_4$  point of the Sun-Jupiter system under the gravitational forces of the Sun, Jupiter, and Saturn. We assume that the orbits of all the bodies are coplanar. We compare our analytical solution (hereafter noted Analytical) with the numerical solution of the restricted four body problem (hereafter 4 Body), and with the numerical solution of restricted three body problem (hereafter 3 Body). We also assume that the Jupiter's orbit is circular and its semi-major axis is 5.2026 AU. For the Analytical and 4 Body cases, we include the perturbation of Saturn, whose semi-major axis and eccentricity are chosen as 9.5549 AU and 0.0555, respectively. Initially the test particle is located at a point departed by the radius  $\delta r$  from  $L_4$  and the initial position angle  $\Phi$  (see Fig. 2). The numerical solutions were obtained by using the method of variation of parameter based on the KS regularization developed by us [2] based on the results of [1]. And we adopted the Adams method as the numerical integrator.

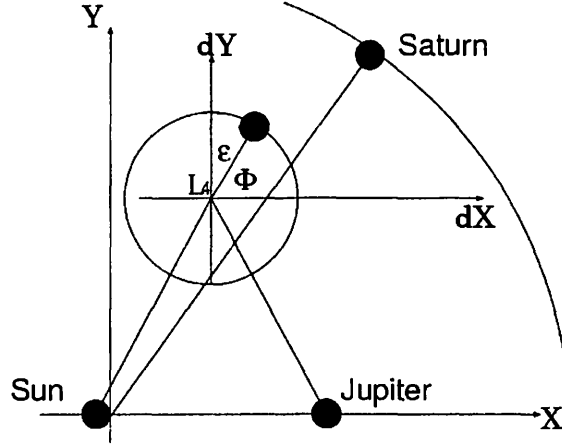


Figure 2: Initial Condition

#### 3.1 Time Evolution of Orbits

First we show the differences in the corotational coordinate system. Fig. 3 shows the time evolution for the 20 revolutions with respect to the Jupiter. This is for the initial condition  $\delta r = 5.2 \times 10^{-5}$  AU, and  $\Phi = 0$  degree. This figure illustrates that the analytical solution agrees well with the numerical integration of 4 Body case. Fig. 4 illustrates the difference between the analytical and numerical (4 Body) solutions for the 20 revolutions of Jupiter's orbit. And Fig. 5 is the same as Fig. 4 but for a longer time span, 1000 revolutions of Jupiter. For the short span,

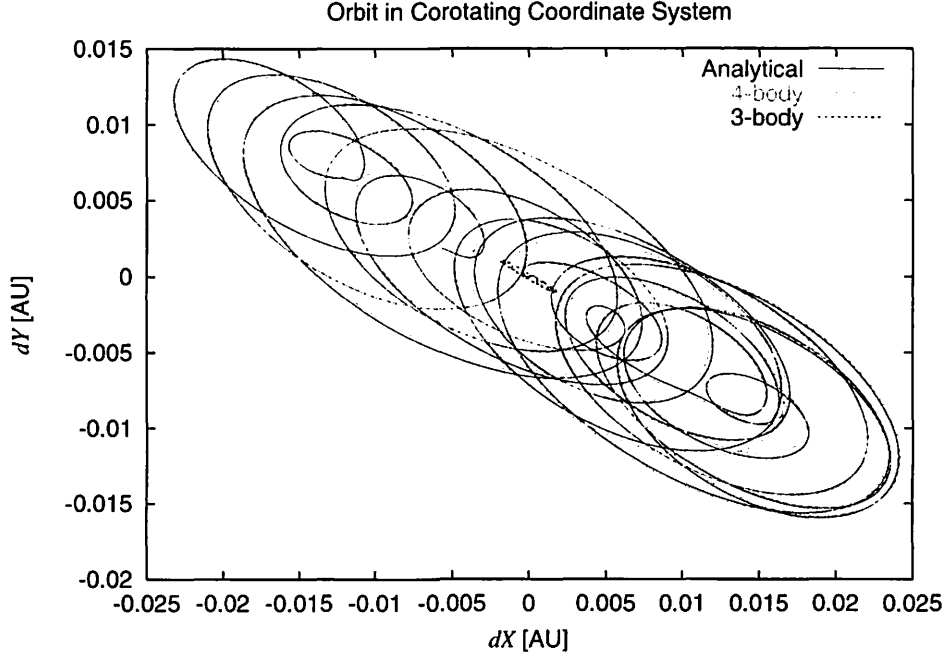


Figure 3: Orbit in the corotational coordinate system :  $\theta = 0$  (deg.)

the maximum error is of the order of  $10^{-4}$ . For the long term, the error is of the form of the mixed secular term and finally becomes of the same order of the size of orbit. Fig. 6 is the same as Fig. 4 but plotted are for a longer period. From these figures, the beating period is almost 3250 orbital periods of Jupiter. The feature of the error in the short term supports an idea that the error is caused by the non-linear effect of the free oscillation term. The observed frequency of the error is around 0.16, which is quite close to the double of the longer eigenfrequency of the free oscillation,  $\omega_2 = 0.0805$ . Fig. 7 plots the frequency analysis of the residual of  $dX$  by Fast Fourier Transform (FFT). Clearly, the effect of  $2\omega_2$  stand out and we can find that the error growth of our analytical solution is mainly occurred by elimination of the nonlinear effect of the tidal force of the primary and the secondary bodies. Also the amplitude of the relative error is of the order of the square of the relative amplitude of the linear solution,  $C_{1,2}$  or  $S_{1,2}$ . Figs. 8 and 12 show the error in radius. Figs. 9 and 13 are the same as Figs. 8 and 12 but for the longitude. Figs. 10 and 11 are the same as Figs. 8 and 9 but plotted are the close-up of Analytical. Unchanged is the error increasing in a mixed secular manner. We note that, when perturbed by other bodies, the behavior of the test particle moving around the triangular Lagrange point hardly depends on the initial condition. Namely the solution of the main part remains the forced oscillation due to the other bodies as we expected.

### 3.2 Comparison with Real System

As the next examination, we compare our solution with the numerical integration where the motion of perturbors is given not by pure Keplerian ones but by the actual planetary ephemeris, DE 405 (hereafter cited by DE). Our analytical solution is constructed based on the simple physical model where the orbit of all the bodies is on the same plane and all the perturbors move on the Keplerian orbit. However in the actual system such as our solar system, the motion of the perturbors is more complicated; the effect of the eccentricity of the secondary, the effect of the inclination of the other perturbors, the influences of other planets as Uranus and Neptune, that say, the orbit of the perturbors is not completely closed. Therefore it is useful to examine

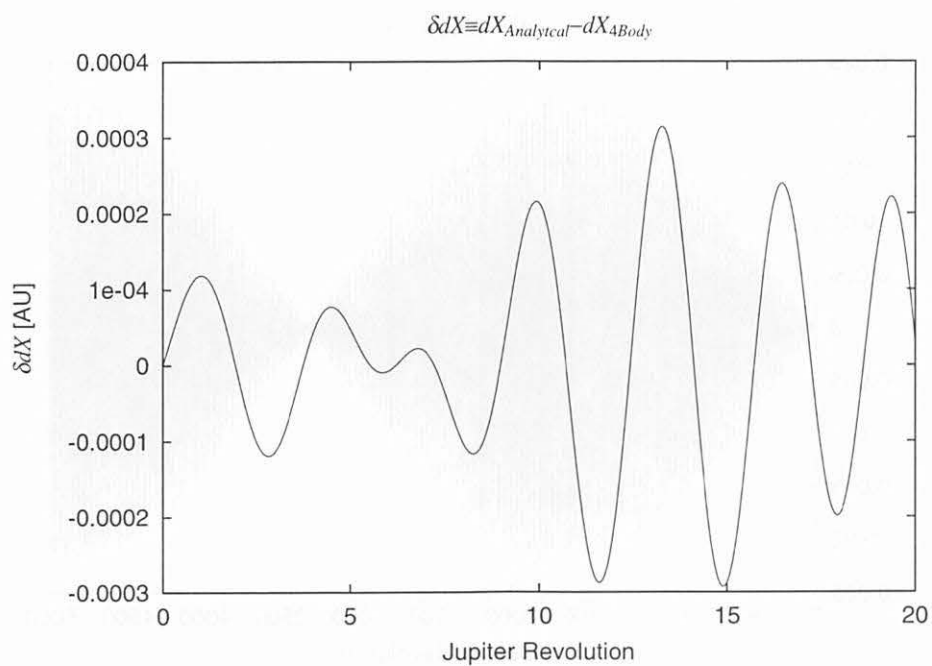


Figure 4: Error of Analytical Solution with respect to Numerical Solution

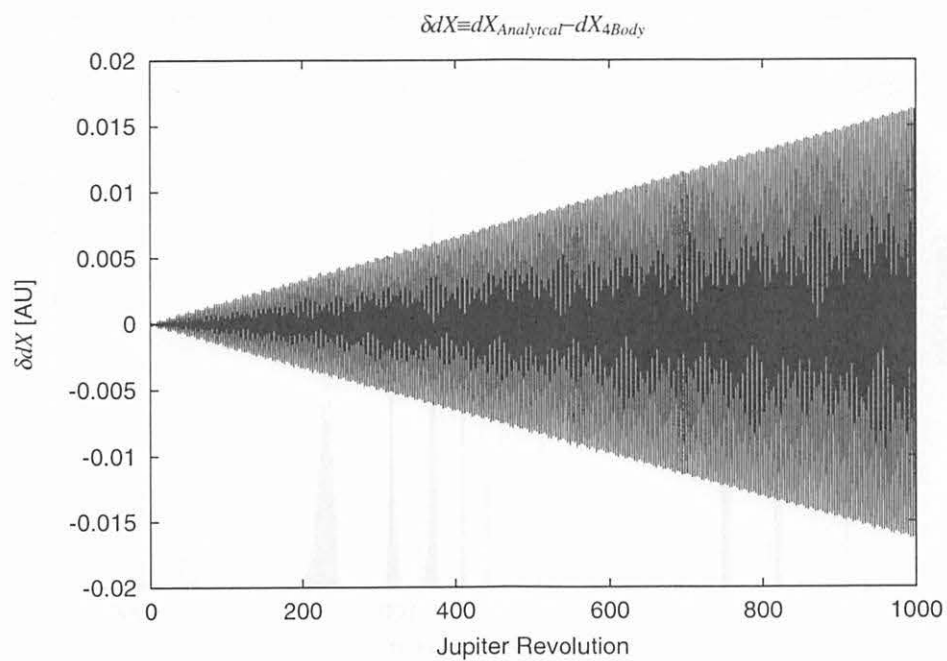


Figure 5: Error of Analytical Solution with respect to Numerical Solution

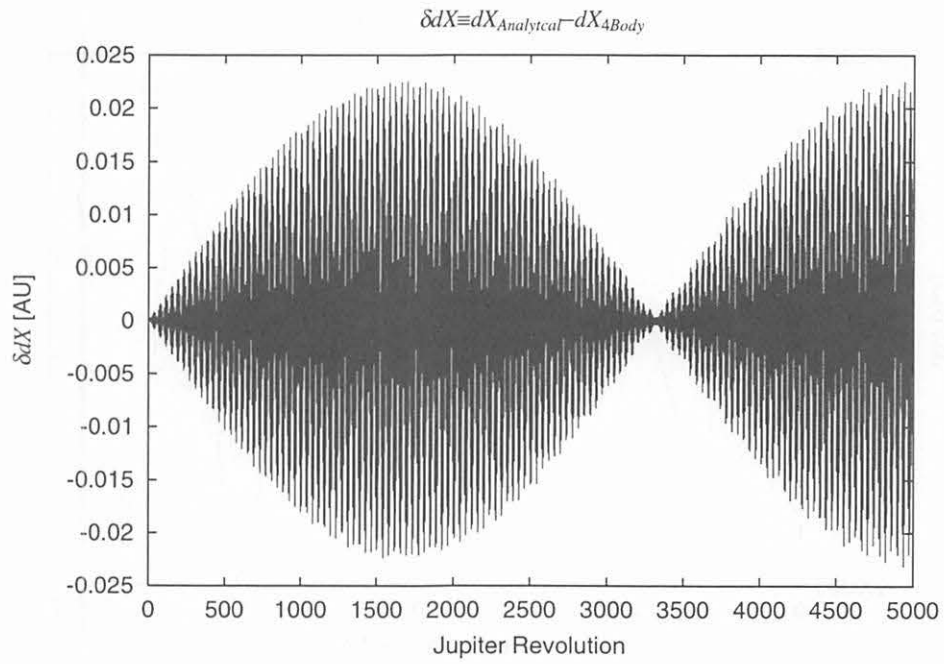


Figure 6: Error of Analytical Solution with respect to Numerical Solution

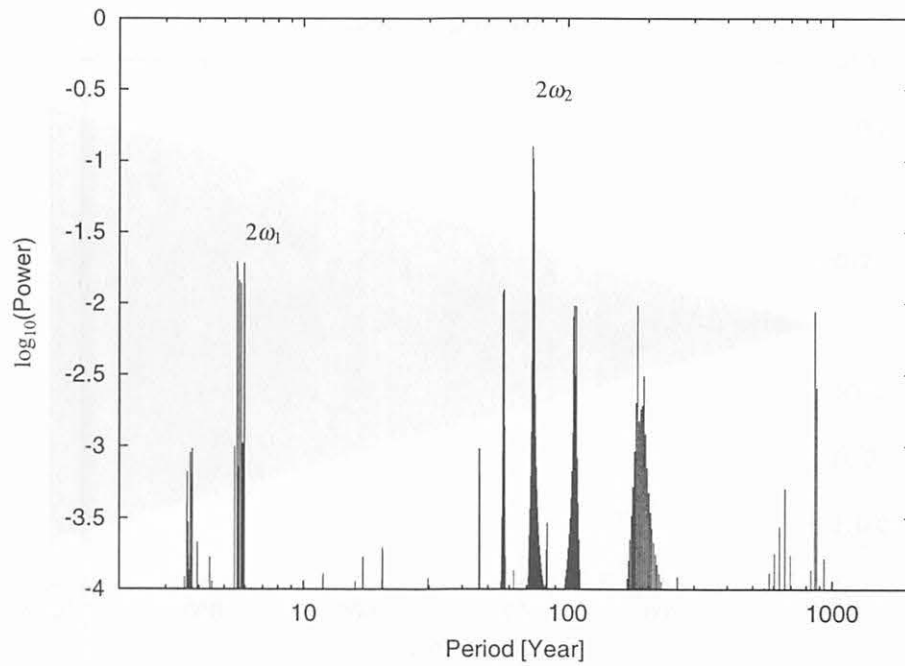


Figure 7: Frequency analysis of residual : Close Up.

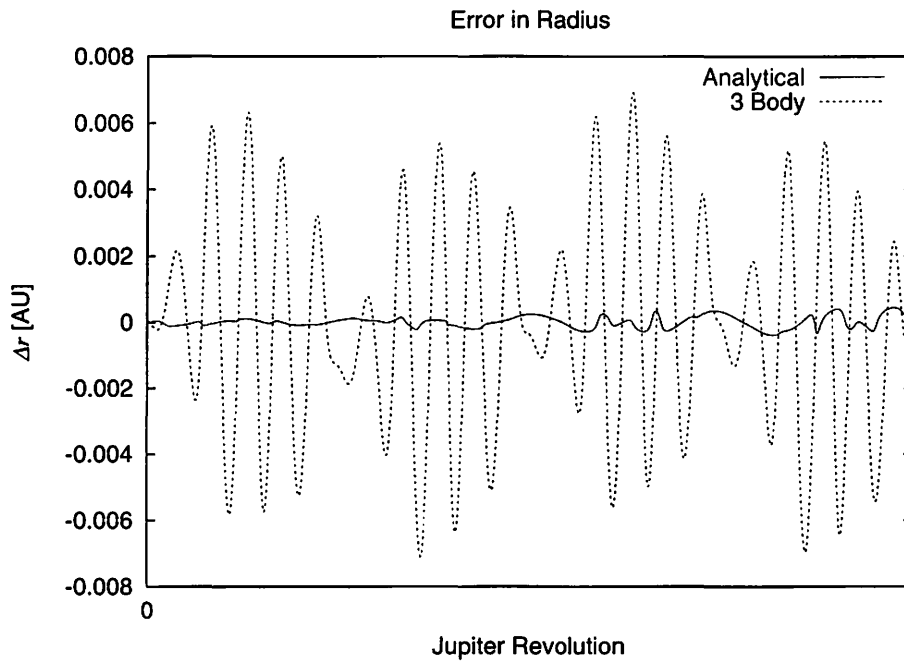


Figure 8: Radial Error of Analytical Solution with respect to Numerical Solution

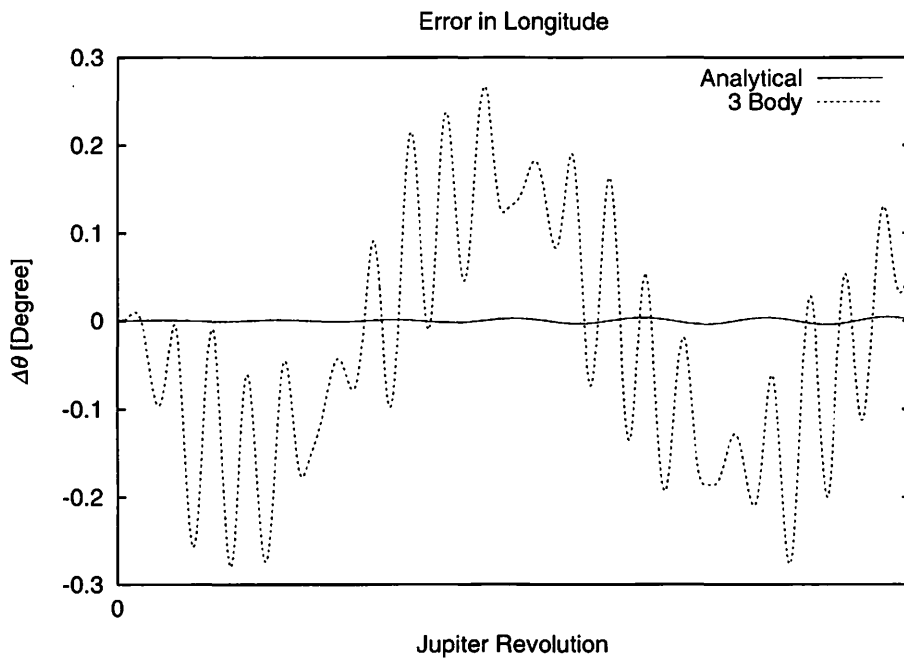


Figure 9: Longitude Error of Analytical Solution with respect to Numerical Solution



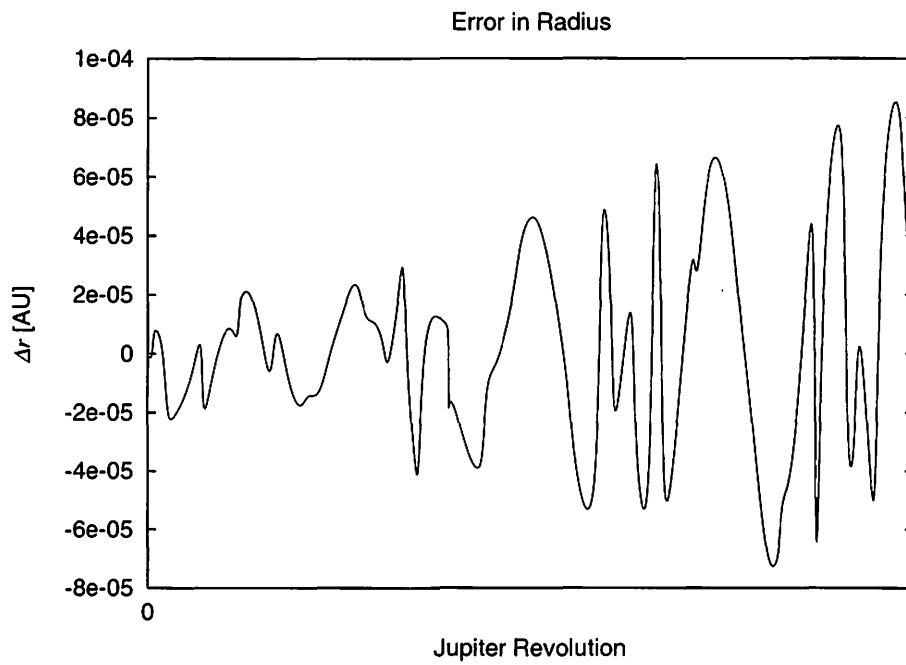


Figure 10: Radial Error of Analytical Solution with respect to Numerical Solution

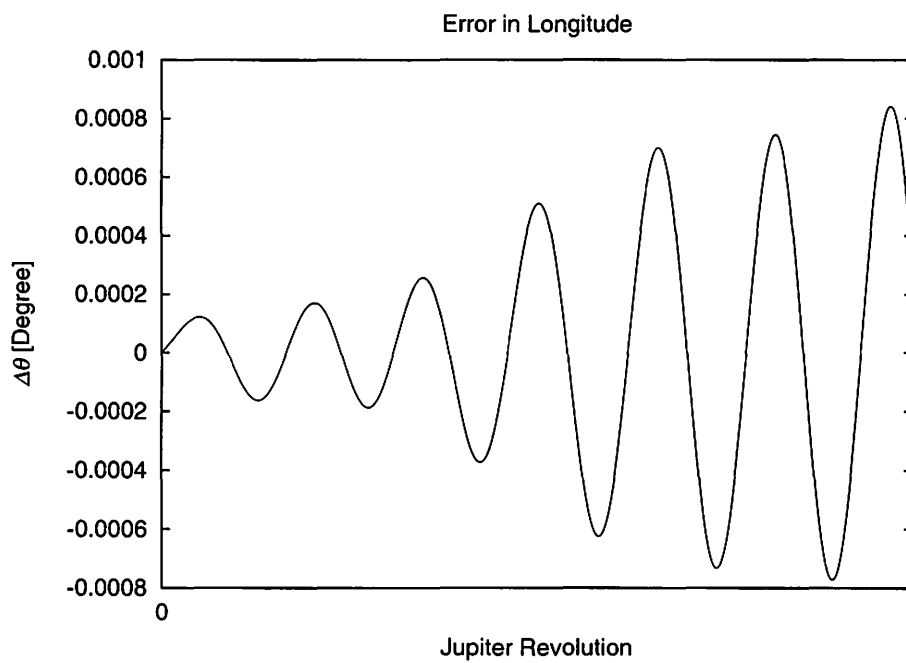


Figure 11: Longitude Error of Analytical Solution with respect to Numerical Solution

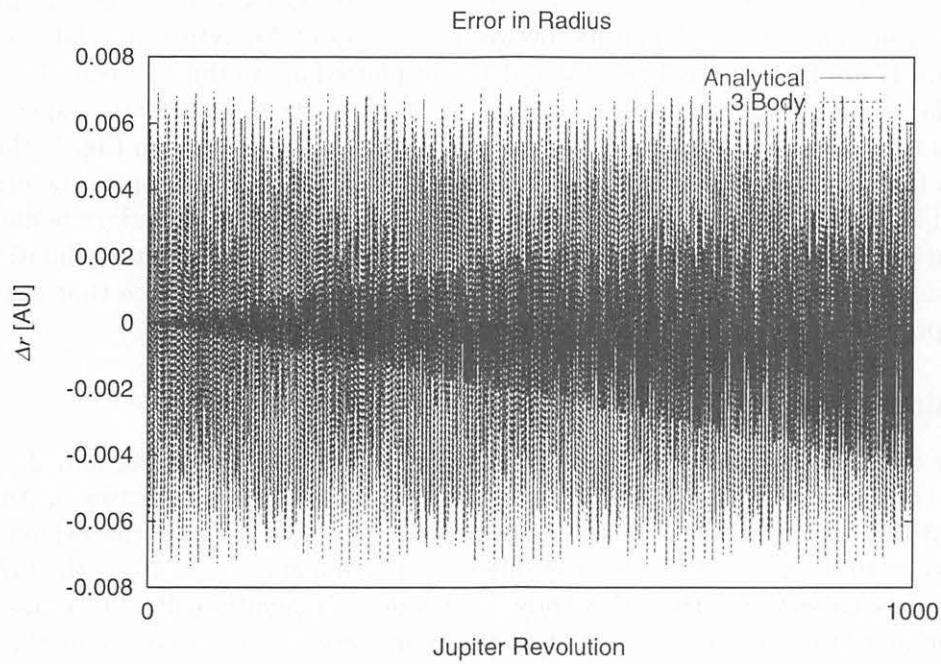


Figure 12: Radial Error of Analytical Solution with respect to Numerical Solution

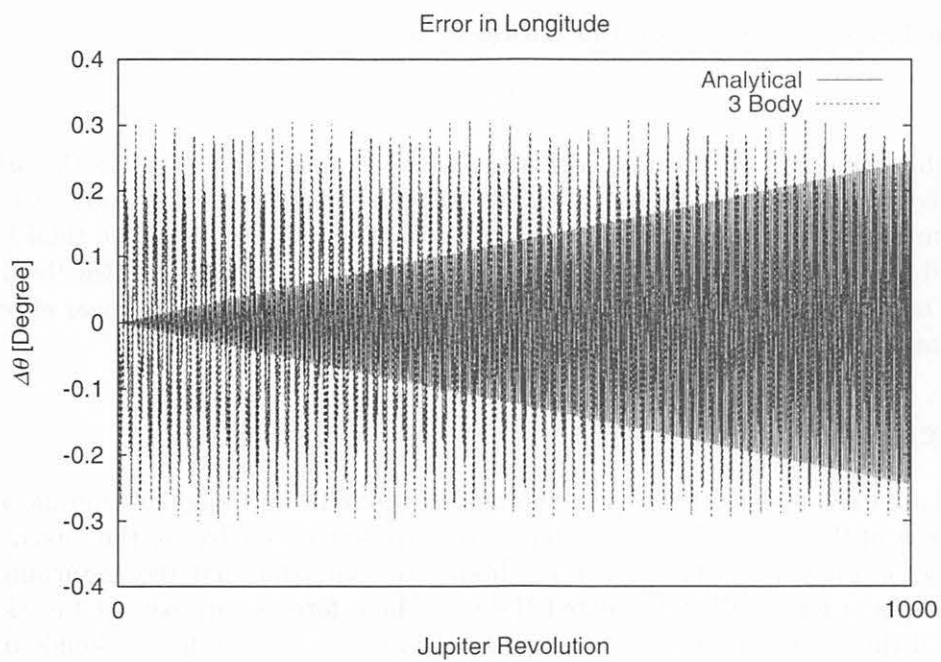


Figure 13: Longitude Error of Analytical Solution with respect to Numerical Solution

how much our solution represent the real motion.

We chose the epoch as JD 2305424.50, and adopted the osculating elements at the epoch in constructing the analytical theory. Since our theory is planar, we regard the Sun-Jupiter plane of DE as our fundamental plane and projected the actual orbit of Saturn on this plane. The initial condition of the test particle was set as  $\delta r = 5.2 \times 10^{-5}$  AU and  $\Phi = 0$ . Figs. 14, and 15 show the deviation in radius and in longitude between the cases of Analytical and DE, respectively. Figs. 16, and 17 are the same as Figs. 14, and 15 but plotted are in the long run. The maximum values of the relative errors in radius and in longitude are 0.0025 and 0.001, respectively. Fig. 18 illustrates the frequency analysis of the residual in  $dX$ . Comparing with Fig. 7, there appear the various frequencies, however, we can realize that the main contribution to the error growth is produced by the lack of the nonlinear effect of the primary and the secondary bodies and then the elimination of the effect of the eccentricity of the secondary body or the inclination of other perturbers is relatively less contribution to the residual. Then, we conclude that our analytical solution represents the quite well the features of true orbit.

### 3.3 Limitation of Application

Finally, we examine the limitation of our analytical solution. Figs.19, 20, and 21 plots of  $\epsilon$ -dependence of the ratio of the width and thickness of orbital region for cases of Analytical, 4 Body, and 3 Body, respectively. For the cases of Analytical and 4 Body, the behaviors are the same up to  $\epsilon = 10^{-4}$ . At  $\epsilon = 10^{-3}$ , a bit of difference occurs, and at  $\epsilon = 10^{-2}$ , the behavior of 4 Body is rather similar to the that of 3 Body. The order of magnitude  $10^{-3}$  is almost the same as the value at which the nonlinear effect of the tidal force of the primary and the secondary bodies and the direct gravitation due to Saturn balance,

$$\mathbf{F}_E(\mathbf{r}_L, t) \sim \frac{1}{2} \frac{\partial^2 F_{Li}(\mathbf{r}_L, t)}{\partial r_j \partial r_k} \delta r_j \delta r_k.$$

For the Sun-Jupiter-Saturn system, this critical value is,

$$\epsilon_c \sim 4.1 \times 10^{-3}.$$

Therefore this value almost coincides with the observed value. Fig. 22 shows the orbits in the corotating coordinate system where  $\epsilon = 10^{-2}$ . Obviously the orbits of 4 Body and of 3 Body have the similar feature. This is because the effect of the non-linearity of the tidal force of the primary and the secondary dominates over the forced oscillation term due to the third and other perturbing bodies. Since our theory is linear, it cannot deal with the non-linear effect and this is the limitation.

## 4 Conclusion

We created a purely analytical theory of the motion around the triangular Lagrangian point in the framework of the planar  $N$ -body system. We considered the effect of the forced oscillation term due to the third and other perturbing bodies by expanding not the disturbing potential but the disturbing force. We represented the disturbing force as an explicit function of time and obtained the correction in position due to the forced oscillation by its double integration. The effect of the eccentricity of the third and other perturbers was taken into account up to the second order. Here we emphasize that our analytical solution is linear theory such that it is easily applicable to calculate the perturbations due to the any number of perturbing bodies though we limit our discussion to the restricted 4 body problem in the main text. For the short period, our solution well coincides with the numerical solution with a relative maximum error less than  $10^{-4}$ . For the long period, the residual between the analytical solution and the

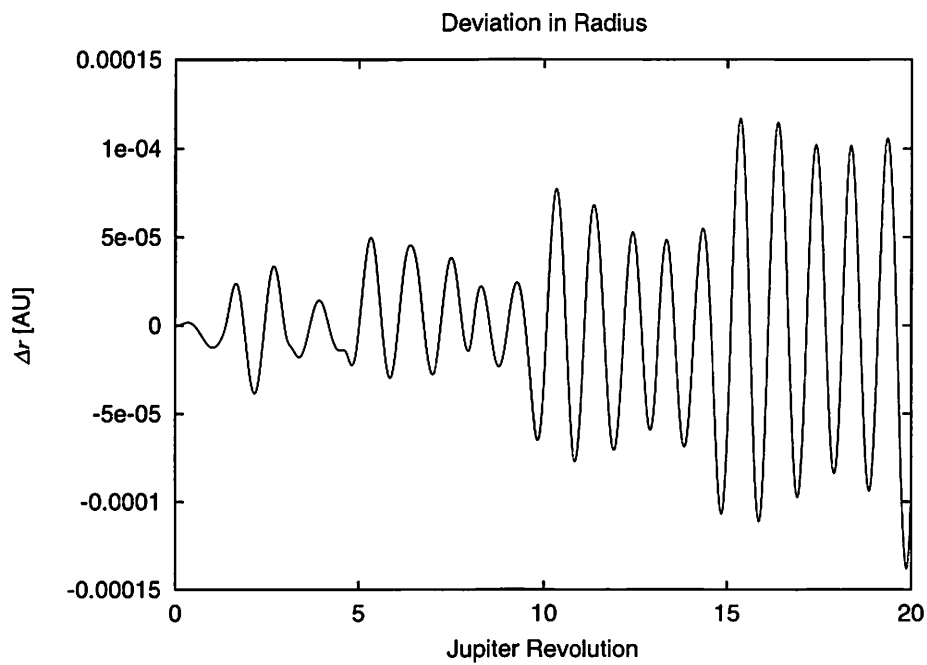


Figure 14: Deviation of Radius between Analytical and DE

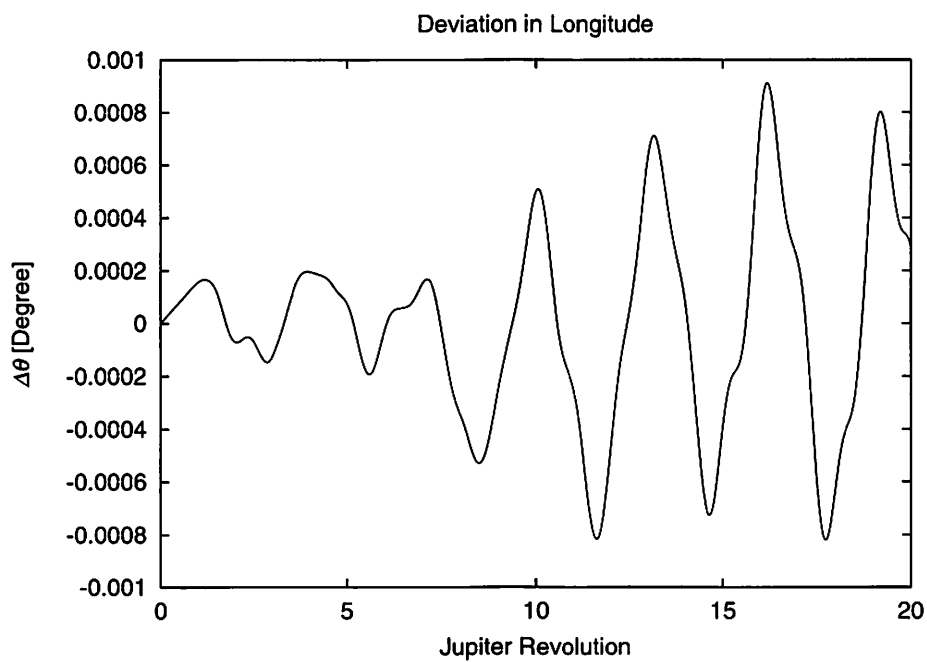


Figure 15: Deviation of Longitude between Analytical and DE

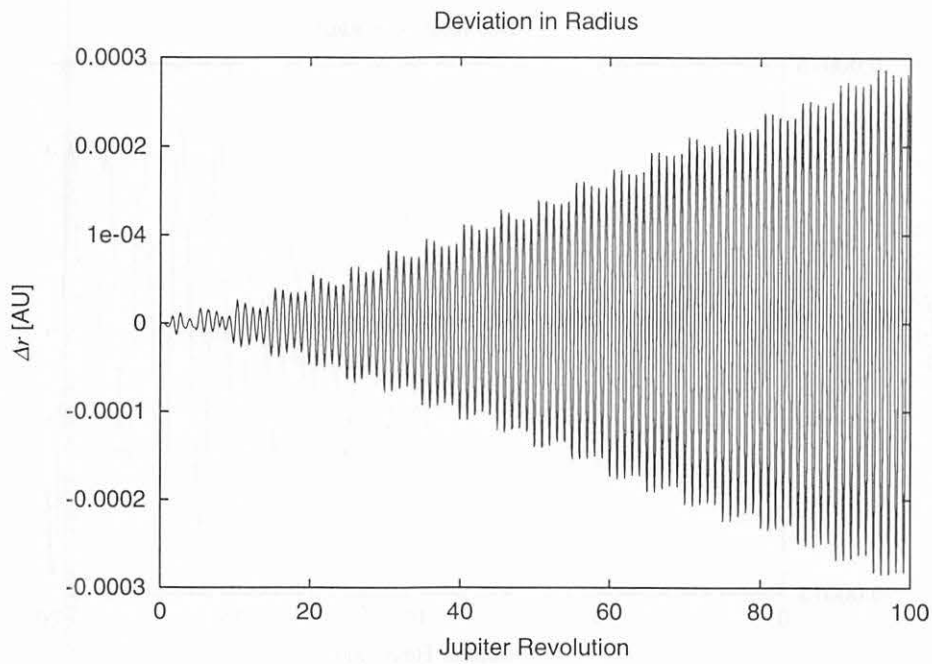


Figure 16: Deviation of Radius between Analytical and DE

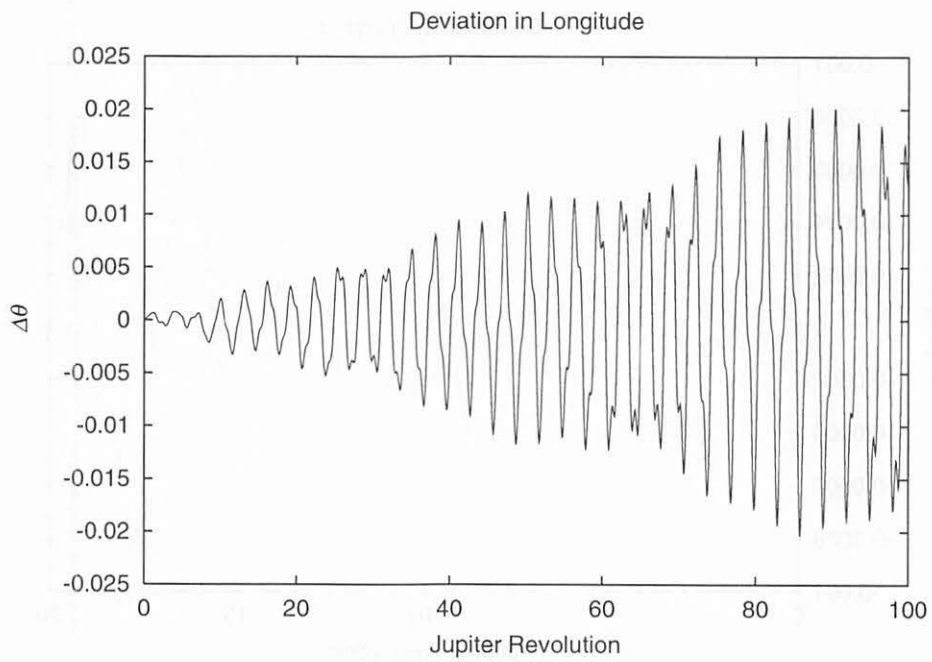


Figure 17: Deviation of Longitude between Analytical and DE

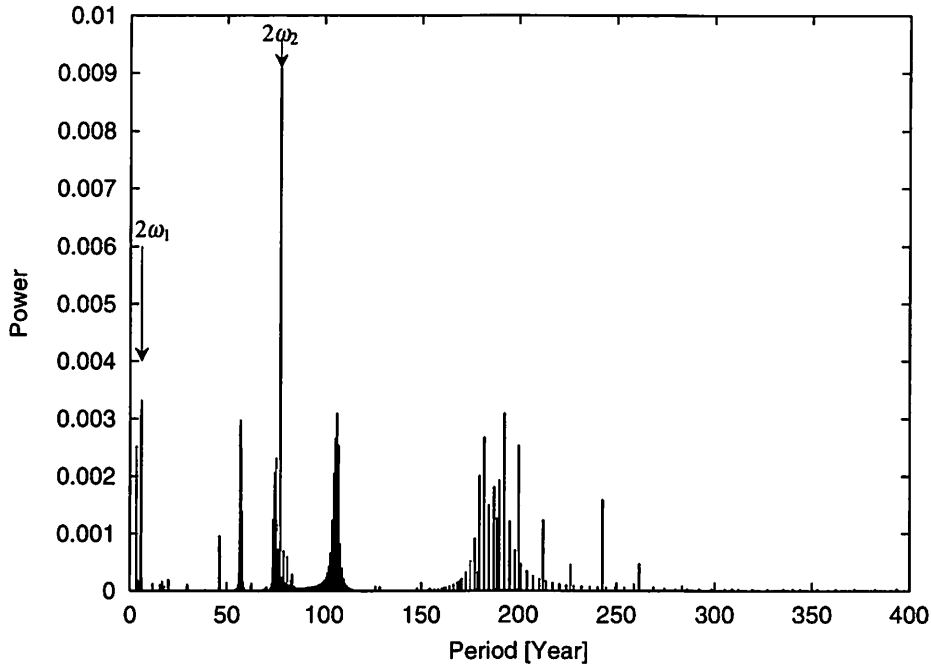


Figure 18: Frequency analysis of residual

numerical solution grows as the mixed secular manner and then beat. By the frequency analysis of the residual, we found that the deviation is mainly caused by the non-linear effect of the tidal force of the primary and the secondary bodies. This fact indicates that the analytical solution we derived is well express the effect of the direct gravitational force due to the other bodies in the linear theory. Then the improvement for the solution must be performed not to the part of the forced oscillation but to that of free oscillation of the primary and the secondary bodies. Further we compared our solution with the numerical one based on the real solar system by using JPL's planetary ephemeris DE405. From this comparison, the main contribution to the residual between the analytical solution and the numerical one based on the DE405 is also mainly caused by the elimination of the nonlinear effect of the tidal force in the part of the free oscillation of the primary and the secondary bodies. Finally we examined the limitation of the application of our analytical solution and realized that when the direct gravitation of the other bodies and the nonlinear effect of the tidal force of the primary and the secondary bodies balance, our solution reach the limit.

In conclusion, for the motion of the test particles around triangular Lagrange points in the restricted  $N$ -body system, the key role is played by the direct gravitational force due to the third and other bodies rather than the non-linear effect of the tidal force of the primary and the secondary bodies. Practically speaking, this range of allowance is quite large. For example, in the case of the Sun-Earth system, it corresponds to 450000 km in the deviation from  $L_4$  by the gravitational influence of Jupiter.

Our solution is especially effective for designing the orbit of some space missions to be put near the triangular Lagrange points such as the gravitational wave detection or the space telescope for observing the near Earth crossing objects. This is because it is expected that the analytical solution limits the initial condition and then reduce considerably the vast of numerical integration and also it makes us estimate easily the orbital region of the spacecrafts without the numerical integration.

As a future work, it is necessary to improve the part of the free oscillation of the primary and

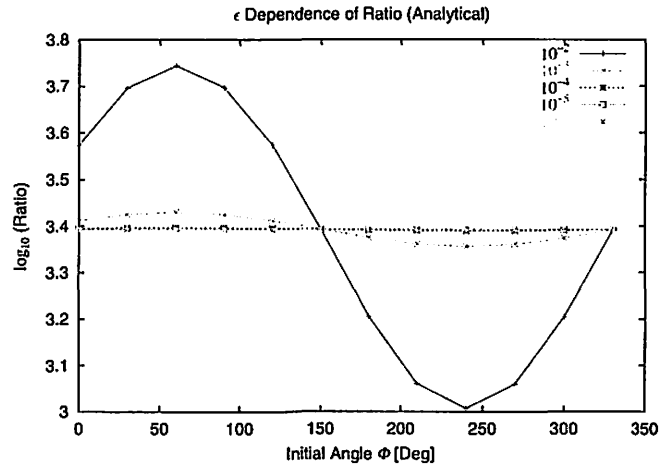


Figure 19:  $\epsilon$  Dependence of Ratio (Analytical)

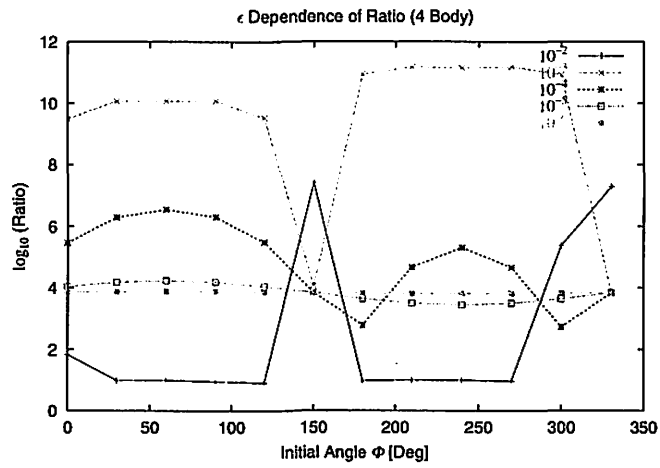


Figure 20:  $\epsilon$  Dependence of Ratio (4 Body)

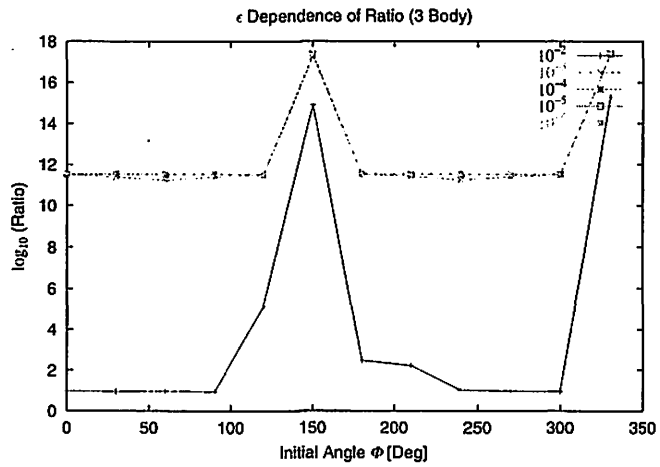


Figure 21:  $\epsilon$  Dependence of Ratio (3 Body)

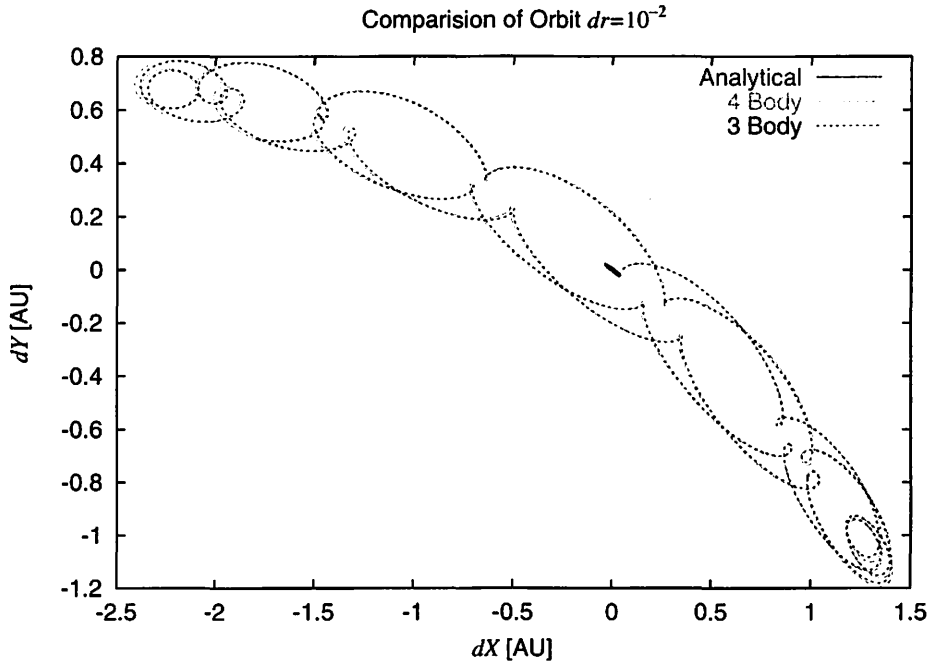


Figure 22: Comparison of Orbit  $dr = 10^{-2}$

the secondary bodies including the nonlinear effect of their tidal force for better agreement with the numerical integration. Since our formalism is restricted to the planar case, the introduction of the inclination of the test particle will be significant in the applying it to the real system such as the motion of the Trojan asteroids.

## References

- [1] Arakida, H., & Fukushima, T., 2000, *AJ*, **120**, 3333
- [2] Arakida, H., & Fukushima, T., 2001, *AJ*, **121**, 1764
- [3] Bhatnagar, K. B., & Gupta, U., 1986, *CeMec*, **39**, 67
- [4] Bhatnagar, K. B., Gupta, U., & Bhardwaj, R., 1994, *CeMDA*, **59**, 345
- [5] Breakwell, J. V., & Brown, J. V., 1979, *CeMec*, **20**, 389
- [6] Brouwer, D., & Clemence, G. M., 1961, "Methods of Celestial Mechanics", Academic Press, New York & London
- [7] Celletti, A., & Giorgilli, A., 1991, *CeMDA*, **50**, 31
- [8] Danby, J. M. A., 1964, *AJ*, **69**, 294
- [9] Danby, J. M. A., 1992, "Fundamentals of Celestial Mechanics" (2nd ed.), Willmann-Bell Inc.
- [10] Danby, J. M. A., 1964, *AJ*, **69**, 165
- [11] Dermott, S. F., & Murray, C. D., 1981, *Icarus*, **48**, 12



- [12] Deprit, A., & Deldvoie, A., 1965, *Icarus*, **4**, 242
- [13] Deprit, A., & Deprit-Bartholome, A., 1967, *AJ*, **72**, 173
- [14] Deprit, A., Henrard, J., & Rom, A. R. M., 1967, *Icarus*, **6**, 381
- [15] Deprit, A., & Palmore, J., 1966, *AJ*, **71** 94
- [16] Dionysiou, D. D., & Stamou, G. G., 1989, *Ap&SS*, **152**, 1
- [17] Dvorak, R., 1992, *CeMDA*, **54**, 195
- [18] Erdi, B., 1978, *CeMec*, **18**, 141
- [19] Erdi, B., 1983, Proceedings of the Seventy-fourth Colloquium “Dynamical trapping and evolution in the solar system” , Gerakini, Greece, August 30-September 2, 1982. Reidel Publishing Co., 165
- [20] Evans, N. W., and Tremaine, S., 1999, *AJ*, **118**, 1888
- [21] Fukushima, T., 1999, *Proc. IAU Coll.* 173, 304, eds: Svoreň, J., Pittich, E.M., and Rickman, H., *Astron. Inst, Slovak. Acad. Sci., Tatranská Lomnica*
- [22] Garfinkel, B., 1976, *CeMec*, **14**, 301
- [23] Garfinkel, B., 1977, *AJ*, **82**, 368
- [24] Garfinkel, B., 1978, *CeMec*, **19**, 259
- [25] Garfinkel, B., 1980, *CeMec*, **22**, 267
- [26] Garfinkel, B., 1983, *CeMec*, **30**, 373
- [27] Garfinkel, B., 1985, *CeMec*, **36**, 19
- [28] Giorgilli, A., & Skokos, C., 1997, *A&A*, **317**, 254
- [29] Gómez, G., Jorba, À., Masdemont, J., & Simó, C., 2001, “Dynamics and Mission Design Near Libration Points, Vol. III Advanced Methods for Collinear Points”, World Scientific, Singapore, New Jersey, London and Hong Kong
- [30] Gómez, G., Jorba, À., Masdemont, J., & Simó, C., 2001, “Dynamics and Mission Design Near Libration Points, Vol. IV Advanced Methods for Triangular Points”, World Scientific, Singapore, New Jersey, London and Hong Kong
- [31] Gómez, G., Llibre, J., Martínez, R., & Simó, C., 2001, “Dynamics and Mission Design Near Libration Points, Vol. II Fundamentals: The Case of Triangular Libration Points”, World Scientific, Singapore, New Jersey, London and Hong Kong
- [32] Goździewski, K., 1998, *CeMDA*, **70**, 41
- [33] Hagel, J., 1995, *CeMDA*, **63**, 205
- [34] Hairer, H., Nørsett, S. P., and Wanner, G., 1987, *Solving Ordinary Differential Equation I*, Springer-Verlag, Berlin
- [35] Henrici, P., 1962, *Discrete Variable Methods in Ordinary Differential Equations*, Wiley & Sons, Inc., New York
- [36] Hiday-Johnston, L. A., & Howell, K. C., 1994, *CeMDA*, **58**, 317

- [37] Horedt, G., 1974, *CeMec*, **10**, 319
- [38] Horedt, G. P., 1984, *CeMec*, **33**, 367
- [39] Horedt, G., & Mioc, V., 1976, *CeMec*, **13**, 421
- [40] Howard, J. E., 1990, *CeMDA*, **48**, 267
- [41] Ichtiaroglou, S., & Voyatzis, G., 1990, *JApA*, **11**, 11
- [42] Innanen, K. A., & Mikkola, S., 1989, *AJ*, **97**, 900
- [43] Jorba, À., 2000, *A&A*, **364**, 327
- [44] Kinoshita, H., 1998, *Celestial Mechanics and Orbital Dynamics*, University of Tokyo Press (In Japanese)
- [45] Kinoshita, H., 1970, *PASJ*, **22**, 373
- [46] Kinoshita, H., 1979, *Proceedings of the Advanced Study "Instabilities in dynamical systems: Applications to celestial mechanics"* Cortina d'Ampezzo, Italy, July 30-August 12, 1978, D. Reidel Publishing Co., 229
- [47] Kinoshita, H., Yoshida, H., and Nakai, H., 1991, *Celest. Mech.*, **50**, 59
- [48] Kondurar, V. T., 1974, *CeMec*, **10**, 327
- [49] Kumar, V., & Choudhry, R. K., 1987, *CeMec*, **40**, 155
- [50] Lambert, J. D., and Watson, I. A., 1976, *J. Inst. Math. Applic.*, **18**, 189
- [51] Lohinger, E., & Dvorak, R., 1993, *A&A*, **280**, 683
- [52] Lukianov, L. G., 1976, *CeMec*, **13**, 455
- [53] Lukyanov, L. G., 1986, *SvA*, **30**, 720
- [54] Lukyanov, L. G., 1986, *AZh*, **63**, 1222
- [55] Maindl, T. I. 1996, *Completing the Inventory of the Solar System*, Astronomical Society of the Pacific Conference Proceedings, volume 107, T.W. Rettig and J.M. Hahn, Eds., 147
- [56] Maindl, Th. I., & Hagel, H., 1996, "Proceedings of the 3rd International Workshop on Positional Astronomy and Celestial Mechanics", held in Cuenca, Spain, October 17 - 21, 1994, Valencia: Universitat, Observatorio Astronomico, edited by Eleonora I. Yagudina, Maria J. Martinez Uso, and Alicia Cordero Barbero, 183
- [57] Markellos, V. V., & Moran, P. E., & Black, W., 1975, *Ap&SS*, **33**, 385
- [58] McKenzie, R., & Szebehely, V., 1981, *CeMec*, **23**, 223
- [59] Melita, M. D., & Brunini, A., 2001, *MNRAS*, **311**, L12
- [60] Murray, C. D., 1994, *Icarus*, **112**, 465
- [61] Murray, C. D., & Dermott, S. F., 1999, "Solar System Dynamics", Cambridge Univ. Press
- [62] Palutan, F., 1994, *JBIS*, **47**, 497
- [63] Papadakis, K. E., 1998, *CeMDA*, **72**, 235

- [64] Perdios, E., & Zagouras, C. G., 1991, CeMDA, **51**, 75
- [65] Quinlan, G .D., 1999, astro-ph/9901136
- [66] Quinlan, G. D., and Tremaine, S., 1990, AJ, **100**, 1694
- [67] Ragos, O., & Zagouras, C., 1991, CeMDA, **50**, 325
- [68] Roels, J., 1975, CeMec, **12**, 327
- [69] Richardson, D. L., 1980, CeMec, **22**, 231
- [70] Selaru, D., & Cucu-Dimitrescu, C., 1995, CeMDA, **61**, 333
- [71] Shrivastava, A. K. & Garain, D., 1991, **51**, 67
- [72] Simmons, J. F. L., McDonald, A. J. C., & Brown, J. C., 1985, CeMec, **35**, 145
- [73] Singh, J., & Ishwar, B., 1984, CeMec, **32**, 297
- [74] Skokos, Ch. & Dokoumetzidis, A., 2001, A&A, **367**, 729
- [75] Sharma, R. K., & Rao, P. V. S., 1975, CeMec, **12**, 189
- [76] Sharma, R. K., & Rao, P. V. S., 1978, CeMec, **18**, 185
- [77] Stiefel, E. L., 1970, CeMec, **2**, 274
- [78] Stiefel, E., Rössler, Waldvogel, J., & Burdet, C. A., 1967, Method of Regularization for computing Orbits in Celestial Mechanics, NASA Constractor Report NASA CR-769
- [79] Stiefel, E., L., & Sheifele, G., 1971, Linear and Regular Celestial Mecchanics, Springer
- [80] Todoran, I., 1994, Ap&SS, **215**, 237
- [81] Todoran, I., & Roman, R., 1992, Astron. Nachr., **313**, 315
- [82] Tsiganis, K., Dvorak, R., & Pilat-Lohinger, E., 2000, A&A, **354**, 1091
- [83] Weibel, W. M., Kaula, W. M., & Newman, W. I., 1990, Icarus, **83**, 382
- [84] Whipple, A. L., 1983, CeMec, **30**, 385
- [85] Whipple, A. L., 1984, CeMec, **33**, 271
- [86] Yoshida, H., 1993, CeMDA, **56**, 27
- [87] Zagouras, C. G., 1991, CeMDA, **51**, 331
- [88] Zagouras, C. G., & Kazantzis, P. G., 1979, Ap&SS, **61**, 389
- [89] Zhang, S.-P., & Innanen, K. A., 1988, AJ, **96**,1983
- [90] Zhang, S.-P., & Innanen, K. A., 1988, AJ, **96**, 1988
- [91] Zhang, S.-P., & Innanen, K. A., 1988, AJ, **96**, 1995

# 木星と平均運動共鳴またはその近傍にある小惑星の軌道

## Minor planet's orbits in or near mean motion resonances with Jupiter

中井宏、木下宙 (国立天文台)

H. Nakai and H. Kinoshita

National Astronomical Observatory

### Abstract

We investigate the characters of minor planet's orbits in or near mean motion resonances with Jupiter. Using semi-analytical methods and numerical integrations which take account of the perturbation of Venus, Earth(and Moon), Mars, Jupiter, Saturn, Uranus, and Neptune, then we get following results.

- If critical arguments( $\sigma$ ) and  $\theta$  which is the angle between the longitude of perihelion of asteroid and Jupiter( $\theta = \varpi_A - \varpi_J$ ) librate at the same time, the eccentricities of the asteroids become very large( $0.24 \leq e \leq 0.58$ , as for 3:2 mean motion resonance and more larger for inner mean motion resonances) except the 1:1 mean motion resonance. But numbered asteroids that have orbital elements with these conditions do not exist.
- When  $\sigma$  circulates and  $\theta$  librates coincidentally, the eccentricities of asteroids become very small( $e_{max} \simeq 0.08$ ) except 1:1 mean motion resonance.
- When  $\sigma$  and  $\theta$  librate in 1:1 mean motion resonance coincidentally, the eccentricities of asteroids become small( $0 \leq e \leq 0.09$ )

### 1. はじめに

軌道が確定し番号が付けられた小惑星は2001年12月現在32729個ある。観測が少なく軌道が確定していない小惑星は約12万個ある。その中には木星との平均運動共鳴関係になっている小惑星も数多く存在する。共鳴と軌道の安定性については興味ある問題として、古くから多くの人が研究している。メインベルトにおける共鳴については Yoshikawa (1989,1990,1991)、カイパーベルト領域については Fuse (1999)、太陽系外惑星系(GJ876)の安定性については、Marcy *et al.*(2001)、Laughlin *et al.*(2001)、Kinoshita *et al.*(2001)などの研究もある。我々は番号の付いた小惑星の中でギャップ、群に属する小惑星の分類を行い、各共鳴領域の軌道の特徴を調べている。同時に、小惑星の安定性やギャップ、群の形成に木星や他の惑星がどのように影響するかを検討し、メインベルトとカイパーベルトにおける共鳴や分布構造の類似点、相違点について調べている。また、太陽系内小惑星と太陽系外惑星系の安定性に種々の共鳴が果たす役割についても調べようとしている。ここでは、木星と平均運動共鳴の関係にある小惑星の軌道について検討した。

## 2 数値積分

小惑星の初期値は MPC(2001 NOV. 30) の軌道要素を用いた。木星と平均運動共鳴になる軌道長半径の範囲は、臨界引数 ( $\sigma$ ) が秤動する時の軌道長半径の変動幅を見越した範囲とした。この範囲は表 1 に示している。軌道長半径がこの範囲内にある小惑星を各共鳴の候補として数値積分を行った。しかし、平均運動共鳴の小惑星で離心率が非常に大きい小惑星では、接触軌道要素の軌道長半径が表 1 に示した範囲外にあることも考えられるため、今回数値積分を行った小惑星以外にも平均運動共鳴にある小惑星が存在する可能性はある。特にギャップについてはこの範囲を見直す必要がある。数値積分法は外挿法を用い、摂動天体は金星、地球 (月を含む)、火星、木星、土星、天王星、海王星の 7 天体と摂動天体として木星、土星の 2 天体だけを考慮した 2 つの場合の数値積分を行った。計算期間は 7 摂動天体で 10 万年、2 摂動天体で 100 万年の場合について行った。小惑星の質量は 0、摂動天体の初期値は惑星暦 DE405 から求めた。

### 3.1 摂動関数

摂動天体として木星 (離心率:  $e_J$ ) 1 天体を考える。その質量を  $m_J$ 、小惑星と木星の相互距離を  $\Delta$  とすると、

小惑星におよぼす木星の摂動関数は

$$R = k^2 m_J \left( \frac{1}{\Delta} - \frac{r_A}{r_J^2} \cos S \right)$$

で表される。但し  $r_A, r_J$  は小惑星、木星の日心距離、 $S$  は両天体間の角距離である。小惑星、木星の近日点経度を  $\varpi_A, \varpi_J$ 、真近点離角を  $f_A, f_J$ 、とし、小惑星の近日点引数を  $\omega_A$ 、木星軌道に対する小惑星の軌道傾斜角を  $i_A$  とすると、

$$\begin{aligned} \Delta^2 &= r_A^2 - 2r_A r_J \cos S + r_J^2, \\ \cos S &= \frac{1}{2}(1 + \cos i_A) \cos\{f_A - f_J + (\varpi_A - \varpi_J)\} \\ &\quad + \frac{1}{2}(1 - \cos i_A) \cos\{f_A + f_J - (\varpi_A - \varpi_J) + 2\omega_A\} \end{aligned}$$

となる。添え字の  $A$  は小惑星、 $J$  は木星を表す。 $\omega_A$  の周期は  $\varpi_A - \varpi_J$  に比べて短周期と仮定して、数値的に平均操作を行い消去する。同様に、真近点離角は短周期であるので、これも平均操作を行い消去すると、摂動関数  $R$  は  $e_A, i_A, \varpi_A - \varpi_J$  だけの関数となる。角運動量  $\Theta = \sqrt{1 - e_A^2} \cos i_A$  を決めると、系の自由度は 1 となり、等エネルギーカーブから  $\varpi_A - \varpi_J$  と  $e_A$  の関係が分かる。

### 3.2 等エネルギーカーブ

上述の方法で求めた小惑星と木星の近日点経度の差 ( $\theta = \varpi_A - \varpi_J$ ) と離心率 ( $e_A$ ) の関係を図 1、図 2、図 3、図 4 に示す。図 1 は 3:2 平均運動共鳴領域の  $\theta: e_A$  の等エネルギー曲線の例である。共鳴領域で、 $\theta$  が秤動するときの最大離心率を  $e_{Rmax}$ 、最小離心率を  $e_{Rmin}$  とする。添え字 R は共鳴にあることを示す。図 1 から  $\theta$  が 180 度のとき、軌

道の離心率が最小離心率 ( $e_{Rmin} = 0.24$ ) 以上または最大離心率 ( $e_{Rmax} = 0.58$ ) 以下の場合、 $\theta$  は180度の回りを稗動する。平均運動共鳴の状態、 $\theta (= \varpi_A - \varpi_J)$  が稗動すれば、離心率は非常に大きくなる可能性がある。また、 $\theta$  が稗動しなくても、境界付近では離心率は大きく変化する。図1中黒丸は、木星と3:2平均運動共鳴にある小惑星(4446) Carolyn の10万年間の数値積分の結果である。この小惑星は今回調べた中では10万年間の平均の離心率が一番大きな小惑星である。しかし、 $\theta$  が180度のときの離心率(約0.2)は  $e_{Rmin}(0.24)$  より小さいために  $\theta$  は常に回転している。一方、この小惑星の臨界引数 ( $\sigma = 3\lambda_J - 2\lambda_A - \varpi_A$ ) は0度の回りを稗動していて、小惑星が木星を追い越すのは小惑星の近日点付近になっている。このような大きな離心率の小惑星は木星との平均運動共鳴関係にないと、木星に大接近し軌道は不安定になる。

図2は3:2平均運動共鳴付近で共鳴から外れ、 $\sigma$  が回転する場合の  $\theta : e_A$  の等エネルギー曲線である。この場合、小惑星の軌道長半径は平均運動共鳴に相当する軌道長半径 ( $a=3.97$ ) とは異なる値 ( $a=3.76$ ) とした。しかし、軌道長半径を変化させても等エネルギーカーブの様子は本質的な差がなかった。平均操作は木星と小惑星の真近点離角で独立に行っている。図2から、 $\theta$  が0度のとき、離心率が最大離心率 ( $e_{Nmax} \simeq 0.08$ ) 以下の軌道は  $\theta$  が0度の回りを稗動する。添え字Nは平均運動から外れていることを示す。このように平均運動共鳴から外れると、永年共鳴ではあるが離心率はそれ程大きくならない。図2中の黒丸は小惑星(1144) Oda の10万年間の数値積分の結果である。

図3は1:1,4:3,3:2,2:1,7:3,5:2,3:1,4:1平均運動共鳴にある( $\sigma$  が稗動)場合で、 $\theta$  が稗動する場合の離心率の範囲を示している。図中+印は稗動の中心を表し、 $e_{Rcen}$  の添え字のRは共鳴にあることを示す。\*印はエネルギーの高い所で散逸過程では不安定になる場所である。群である4:3,3:2共鳴の  $e_{Rcen}$  は0.3~0.4、ギャップである2:1,7:3,5:2,3:1,4:1共鳴の  $e_{Rcen}$  は0.5~0.9と非常に大きいために、 $\sigma$  と  $\theta$  が同時に稗動する軌道は長期間安定に存在できない。しかし、1:1共鳴では  $e_{Rcen}$  が約0.05、 $e_{Rmax}$  が約0.09と小さい値のために、 $\sigma$  と  $\theta$  が同時に稗動する軌道がある。Yoshikawa(1989,1990,1991)によると、 $\sigma$  の変動幅により  $e_{Rcen}$  は変化するが、ここでは  $\sigma$  の変動幅を0とした。

図4には図3と同じ共鳴付近で、平均運動共鳴から外れ  $\sigma$  が回転する場合の  $\theta : e_A$  の等エネルギー曲線である。 $\theta$  が稗動するための中心の離心率 ( $e_{Ncen}$ ) は約0.02~0.05と非常に小さく、軌道長半径の増加につれて僅かに増加する。1:1共鳴付近は平均運動共鳴が崩れると軌道が不安定になるので、図は省略した。1:1共鳴の  $e_{Rmin}$ ,  $e_{Rcen}$ ,  $e_{Rmax}$  とその他の共鳴付近の  $e_{Nmin}$ ,  $e_{Ncen}$ ,  $e_{Nmax}$  とは殆ど同じ値で、等エネルギー曲線は1:1共鳴と他の共鳴付近で平均運動共鳴から外れた状態が同じようである。

#### 4. 考察

木星と平均運動共鳴にある小惑星の個数を表1に示す。トロヤ群(1:1共鳴)、ヒルダ群(3:2共鳴)では、共鳴の候補として選んだ小惑星の殆どが木星と平均運動共鳴の関係になっている。1:1共鳴では、10万年間軌道が安定なものは、全積分期間中全ての小惑星の  $\sigma$  が稗動している。3:2共鳴では、軌道が安定なものは、小惑星の  $\sigma$  が稗動しているか(図5-1)、 $\sigma$  が0度の回りを稗動しているとき  $\theta$  が回転し、 $\sigma$  が回転しているとき  $\theta$  が

秤動する状態を交互に繰り返す "Pericentric librator" である (図 5-2)。チューレ群の領域には現在 3 個の小惑星がある。(279) Thule は一時期  $\sigma$  が回転する "Pericentric librator" (図 6) で、その軌道は積分期間内では (7 摂動天体の場合 10 万年、2 摂動天体の場合 100 万年) 安定である。しかし、他の 2 個の小惑星、(3552) Don Quixote は約 800 年後、(20898) Fountainhills は約 1 万年後に 4:3 共鳴が崩れ軌道が不安定になる。

4:3, 3:2 共鳴では、 $\sigma$  が 180 度の回りを秤動しているとき  $\theta$  が回転し、 $\sigma$  が回転しているとき  $\theta$  が秤動する状態を交互に繰り返す "Apocentric librator" は見つからなかった。これらの共鳴になる軌道長半径は木星の軌道長半径に近いので、小惑星の遠日点で木星を追い越す "Apocentric librator" は木星と接近する可能性があるためと考えられる。

ギャップでは、平均運動共鳴の関係になる小惑星は少ない。しかし、2:1 共鳴領域では、 $\sigma$  が全期間完全に秤動する "librator" が 31 個、"Pericentric librator" が 131 個、"Apocentric librator" が 34 個ある。このように 2:1 共鳴では 3 種類の "librator" がある。それぞれの例として、図 7-1 に "librator" (1362) Griqua、図 7-2 に "Pericentric librator" (300) Geraldina、図 7-3 に "Apocentric librator" (528) Rezia の軌道要素を示す。

7:3 平均運動共鳴の小惑星は (5324) Lyapunov で 10 万年間の最大の離心率は 0.67、平均の離心率は 0.56 である。また、3:1 平均運動共鳴の小惑星は (6318) Cronkite で 10 万年間の最大の離心率は 0.73、平均の離心率は 0.54 である。この他に、3:1 共鳴付近には  $\sigma$  が 180° の回りの秤動と回転を繰り返す 3 個の小惑星 (2608), (6491), (19356) がある。しかし、これら 3 個の小惑星の  $\theta$  は常に回転していて、"Apocentric librator" のように  $\sigma$  が回転しているとき  $\theta$  が秤動することはない。また、これら小惑星の離心率は 0.65 以上なので、軌道は地球の軌道の内側に入り込んでいる。このように大きな離心率の軌道が長期間安定かどうかは今後検討する必要がある。

5:2、4:1 平均運動共鳴領域では  $\sigma$  が秤動する小惑星は見つからなかった。

平均運動共鳴から外れた場合、図 4 に示したように、 $\theta$  が秤動しているときの離心率は小さくなる。平均運動共鳴ではないが、7:3、4:1 共鳴付近に  $\theta$  が秤動する小惑星が存在する。これら  $\theta$  が秤動する全ての小惑星は離心率が最大になった時でも 3.2 節で説明した  $e_{Nmax}$  より小さい離心率である。

数値積分で求めた 10 万年間の平均の軌道長半径と平均の離心率の関係を図 8 (1:1)、図 9 (3:2)、図 10 (2:1) に示す。軌道が不安定になった小惑星は除外している。図 8、トロヤ群 (1:1) では、軌道が安定なものは全て木星と平均運動共鳴の関係になっている。その中で離心率の大きさにより 3 グループに分類される。(1) 離心率の小さい小惑星は  $\theta$  が常に秤動 (図 8 白丸) する。(2) 離心率の大きい小惑星は  $\theta$  が常に回転 (図 8 黒四角) する。(3) 中間の離心率の小惑星は  $\theta$  が秤動と回転を繰り返す (図 8 +印)。各グループの個数はそれぞれ 93 個、210 個、191 個で、各グループに属する代表的な小惑星はそれぞれ (1871), (617), (1208) である。

図 9、ヒルダ群 (3:2) では、積分期間中完全に平均運動共鳴にあるものを丸印、"Pericentric librator" を三角印で表している。軌道が安定なものはこの 2 種類だけで、離心率の小さいものが "Pericentric librator" となる。一方、図 10 に示す 2:1 共鳴では、 $\sigma$  が秤動し完全に、木星と平均運動共鳴の関係にある小惑星 (丸印)、"Pericentric librator" (三

角印)、“Apocentric libration”(四角印)がある。黒小点は平均運動共鳴でない小惑星を示している。図8,9,10には $\theta$ が秤動するための離心率の最大値も示している。離心率の最大値が $e_{Rmax}$ (1:1共鳴), $e_{Nmax}$ (3:2,2:1共鳴)より小さい軌道は $\theta$ が秤動する。また、図9,10中には、小惑星の軌道が火星軌道、木星軌道に交差するための離心率も示している。これより大きな離心率の軌道は火星や木星の軌道と交差することになる。平均運動共鳴付近で離心率の大きな小惑星は全て $\sigma$ が0度の回りを秤動し、木星との大接近を回避している。これら小惑星の軌道傾斜角は20度~30度と大きいので、このことが火星との接近を回避している可能性があり、検討しなければならない。

摂動天体を金星から海王星までの7天体と木星・土星の2天体とした差は、軌道長半径が小さい共鳴、5:2、3:1共鳴などで大きく現れる。例えば3:1共鳴付近の小惑星(887)では、 $\sigma$ は摂動天体が2天体の場合100万年間安定に秤動を繰り返すが、摂動天体が金星から海王星の7天体になると、6万年後には $\sigma$ は秤動から回転になり、平均運動共鳴が崩れる。このように、内側の共鳴領域では内惑星の存在が小惑星の分布構造に影響してくる。一方、木星に近い共鳴(1:1,3:2,2:1)の個数は摂動天体の数では変化がなかった。

## 5. まとめ

半解析的摂動手法によれば、木星と平均運動共鳴の関係にあり( $\sigma$ が秤動)、同時に、小惑星と木星の近日点経度の差( $\theta$ )が永年共鳴で秤動すると、

- 1:1共鳴以外では小惑星の離心率は非常に大きくなる可能性がある(3:2共鳴での離心率 $\sim 0.6$ )。そのような小惑星が安定に存在するためには、火星をはじめ他の大惑星との大接近を回避する機構が必要になるが今回の調査ではそのような小惑星は見つからなかった。
- 1:1共鳴では小惑星の離心率は小さい領域に限定( $e_{Rmax} \simeq 0.09$ )される。実際、93個の小惑星がこの様な軌道であった。

木星と平均運動共鳴の関係にない( $\sigma$ が回転)場合は $\theta$ が秤動すると、

- 1:1共鳴以外では小惑星の離心率は最大でも $e_{Nmax}(\simeq 0.08)$ 以下である。
- 1:1共鳴付近では小惑星の軌道は不安定で存在できない。

半解析的摂動手法によるこれらの結果は数値積分の結果と良く一致する。

1:1共鳴では $\sigma$ は $\pm 60$ 度の回りを秤動する。そのとき、 $\theta$ が $\mp 60$ 度の回りを秤動する小惑星がある。この場合永年共鳴であるが離心率はそれ程大きくならない。

$\sigma$ が0度の回りを秤動する小惑星は4:3,3:2,2:1共鳴付近に存在し、180度の回りを秤動する小惑星は2:1,7:3,3:1共鳴付近に存在する。

## 6. おわりに

離心率が大きく火星や地球の軌道と交差する軌道が安定かどうかの検討と、摂動天体が外側と内側に存在するときの解析的解法も検討する必要がある。また、メインベルト小惑星の軌道の特徴がカイパーベルト小天体の軌道ではどう変化するかを調べることは今後の課題である。

## 7. 参考文献



Fuse, T. : 1999, Dynamical Structure of Edgeworth-Kuiper Belt Objects in/around Mean Motion Resonances with Neptune, thesis for a doctorate at the Graduate University for Advanced Studies.

Kinoshita, H., & Nakai, H. : 2001, Stability of the GJ876 Planetary System, *Publ. Astron. Soc. Japan*, **53**, pp.L25-L26.

Laughlin, G., & Chambers, J. E. : 2001, Short-Term Dynamical Interactions among Extrasolar Planets, *ApJ*, **551**, pp.L109-L113.

Marcy, G.W., Butler, R.P., Fischer, D., Vogt, S.S., Lissauer, J.J., & Rivera, E.J. : 2001, A Pair of Resonant Planets Orbiting GJ876, *ApJ*, **556**, pp.296-301.

MPC :The Minor Planet Circulars/Minor Planets and Comets

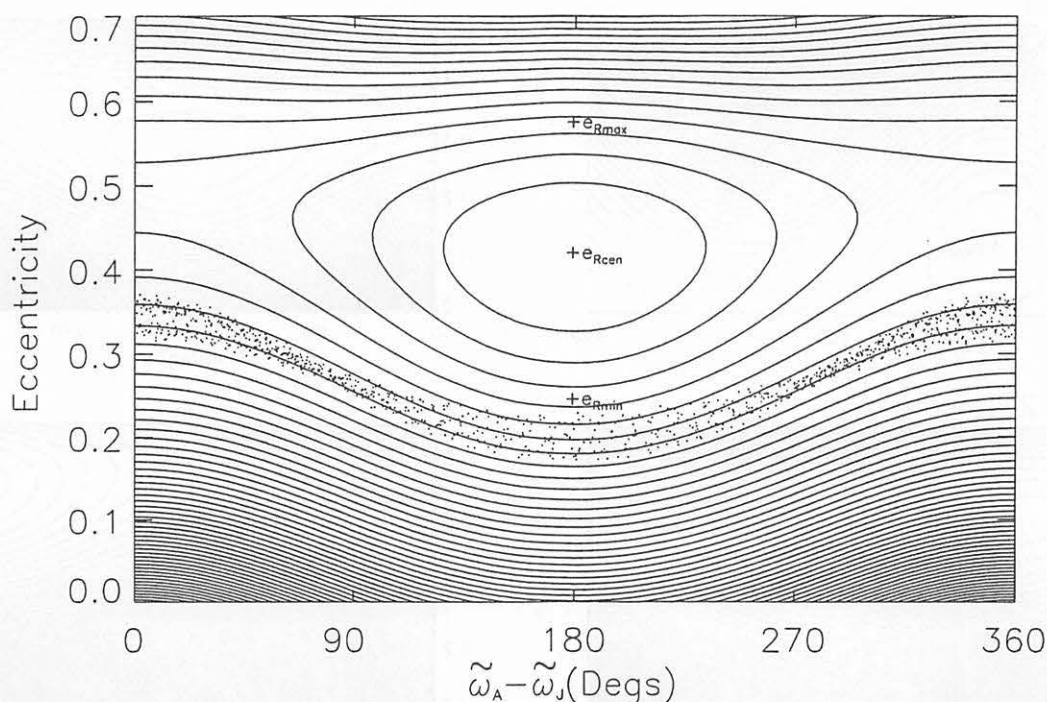
Yoshikawa, M. : 1989, A survey of the motions of asteroids in the commensurabilities with Jupiter, *Astron. Astrophys.*, **213**, pp.436-458.

Yoshikawa, M. : 1990, Motions of Asteroids at the Kirkwood Gaps  
I. On the 3:1 Resonance with Jupiter , *Icarus*, **87**, pp.78-102.

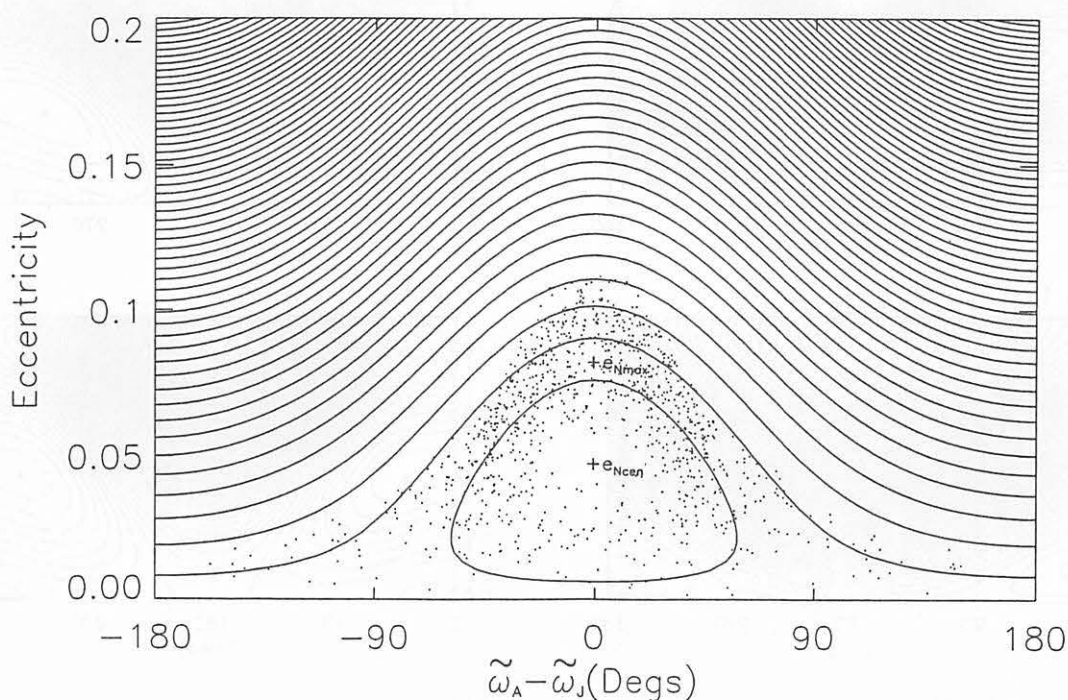
Yoshikawa, M. : 1991, Motions of Asteroids at the Kirkwood Gaps  
II. On the 5:2, 7:3, and 2:1 Resonances with Jupiter , *Icarus*, **92**, pp.94-117.

			survey			resonance		out resonance
		a(AU)	region(AU)		numbers	$\sigma$ (Lib.)	$\sigma$ (L<>C)	$\theta$ (Lib.)
Trojans	1:1	5.20	5.0	5.5	495	494	0	0
Thule	4:3	4.29	4.1	4.5	3	0	1(P)	0
Hilda group	3:2	3.97	3.7	4.2	184	166	15(P)	0
gap	2:1	3.28	3.20	3.40	836	31	131(P) 34(A)	0
	7:3	2.96	2.94	2.98	494	1	0	4
	5:2	2.82	2.80	2.84	380	0	0	0
	3:1	2.50	2.48	2.52	189	1	3	0
	4:1	2.06	1.94	2.14	179	0	0	1

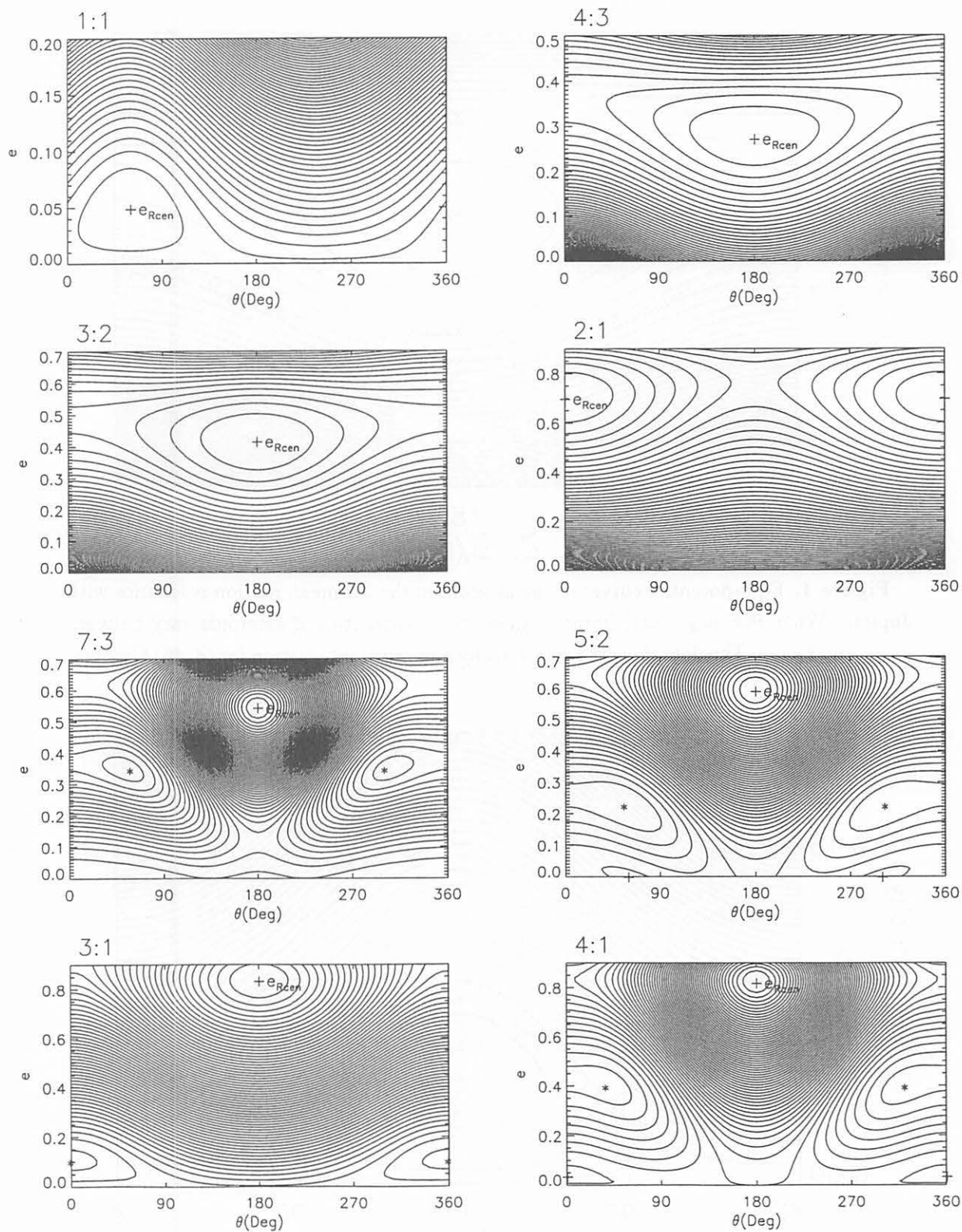
**Table 1.** The numbers of asteroids in the mean motion resonances with Jupiter.  $\sigma$  means critical argument.  $\theta = \varpi_A - \varpi_J$ , where  $\varpi_A$  and  $\varpi_J$  mean the longitude of perihelion of the asteroid and Jupiter respectively.  $\sigma$ (L<>C) indicates that  $\sigma$  alternates libration and circulation. P,A mean pericentric liblator and apocentric liblator, respectively.



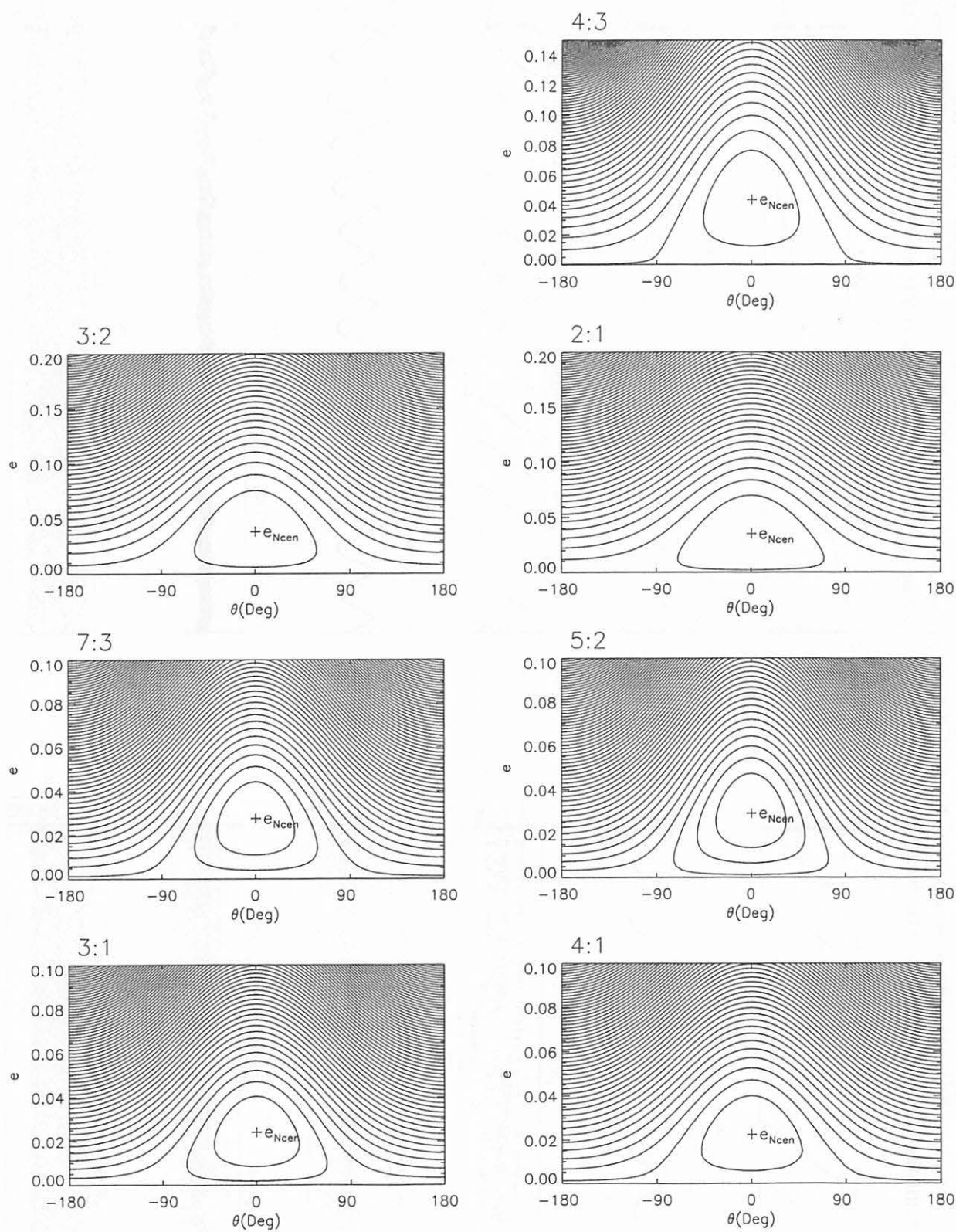
**Figure 1.** Equi-potential curves of the asteroid in the 3:2 mean motion resonance with Jupiter. When  $\theta(= \varpi_A - \varpi_J)$  librates, then the eccentricities of asteroids vary between  $e_{Rmin}$  and  $e_{Rmax}$ . The dots show the solutions by numerical integration for (4446) Carolyn.



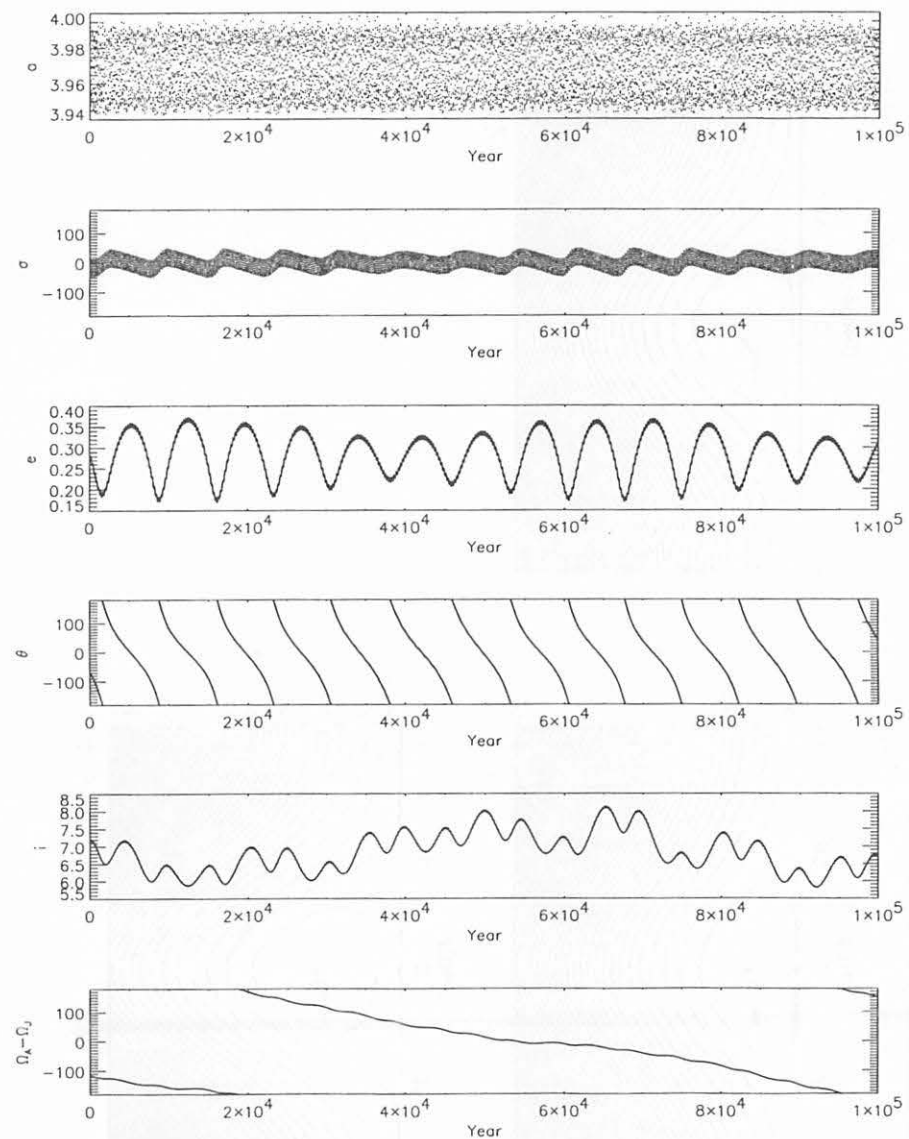
**Figure 2.** Equi-potential curves of the asteroid out of the 3:2 mean motion resonance with Jupiter. When  $\theta$  librates, then the eccentricities of asteroids vary between 0 and  $e_{Nmax}$ . The dots show the solutions by numerical integration for (1144) Oda.



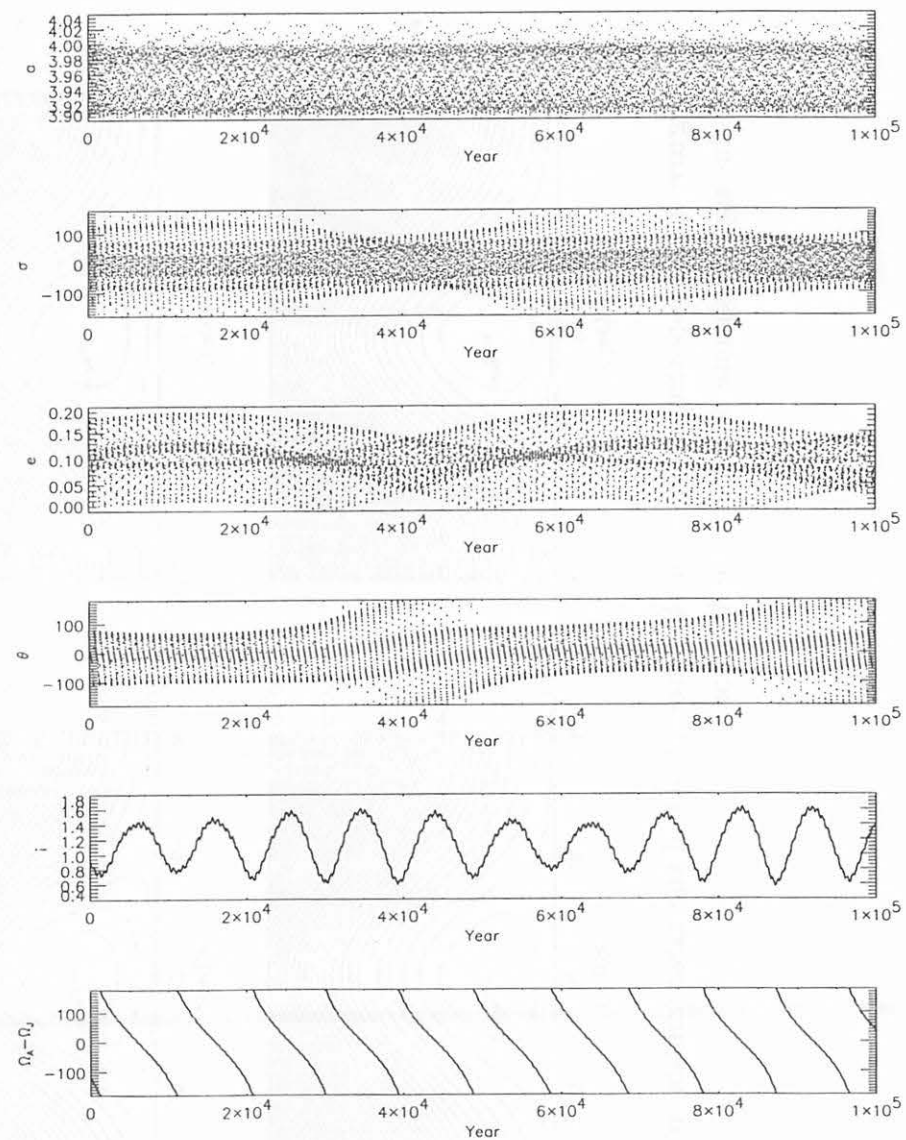
**Figure 3.** Equi-potential curves of the asteroid in the mean motion resonances. The plus signs designate centers of the libration and asterisk signs designate the peaks of energy.



**Figure 4.** Equi-potential curves of the asteroid out of the mean motion resonances. The plus signs designate centers of the libration.

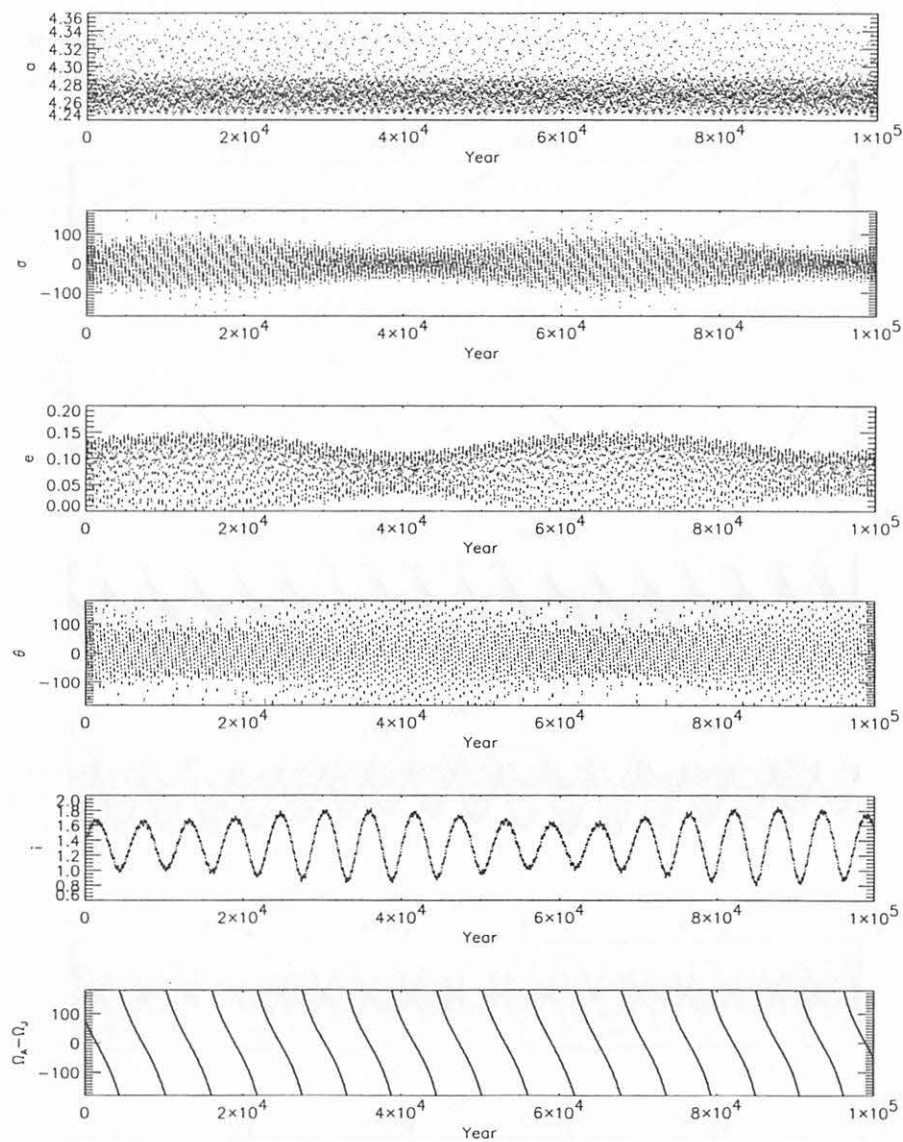


**Figure 5-1.** The orbital elements of (4446)Carolyn : 3:2 librador. The six panels represent the semi-major axis,  $\sigma = 3\lambda_J - 2\lambda_A - \varpi_A$ , the eccentricity,  $\theta = \varpi_A - \varpi_J$ , the inclination, and  $\Omega_A - \Omega_J$ , respectively.

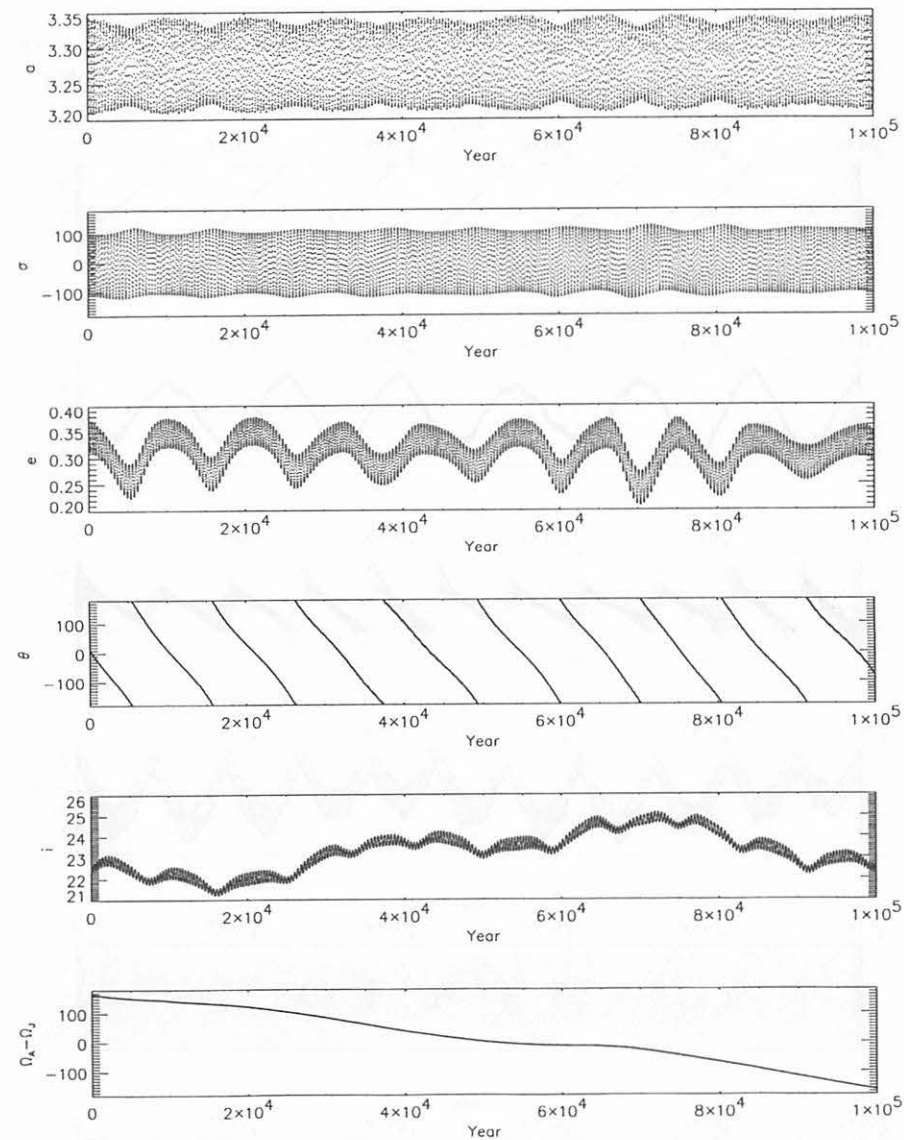


**Figure 5-2.** The orbital elements of (15626)2000 HR50 : 3:2 pericentric librador. The explanations for the vertical axes are same as for Figure 5-1.

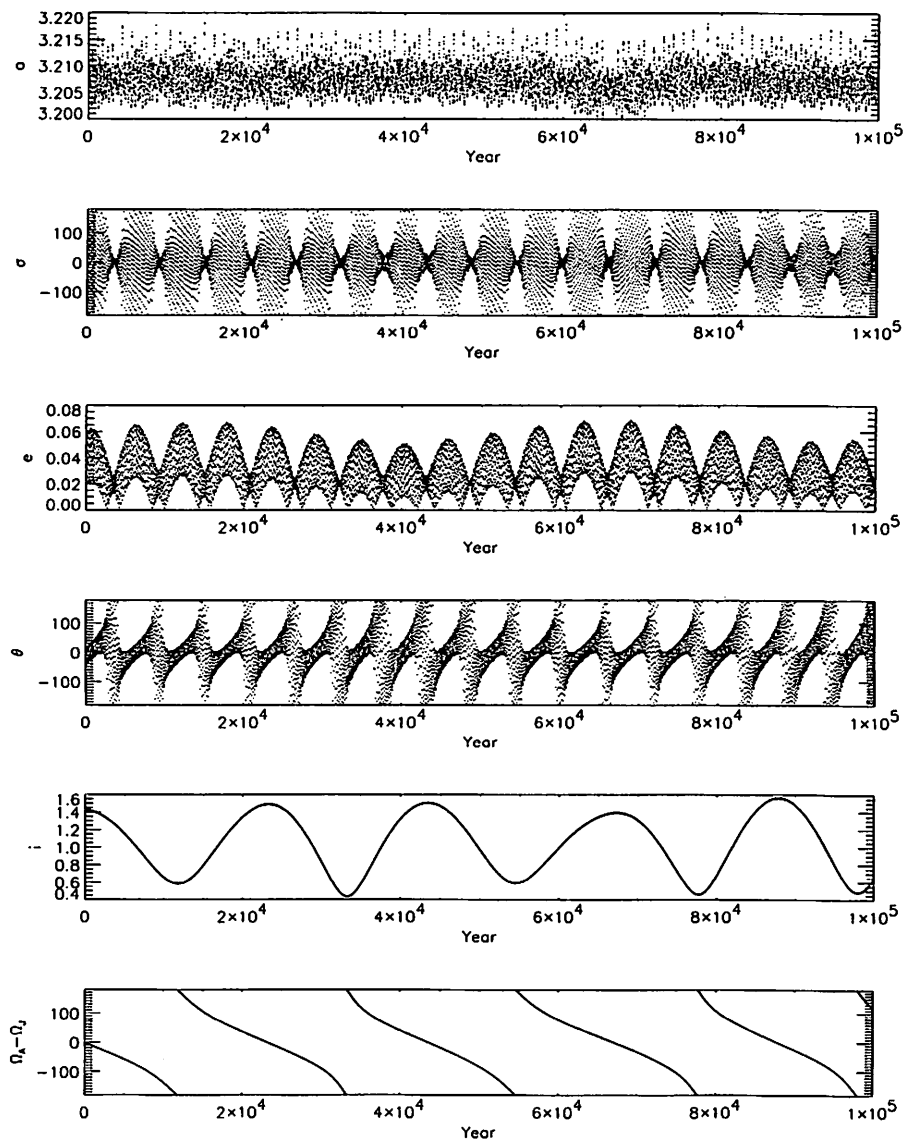




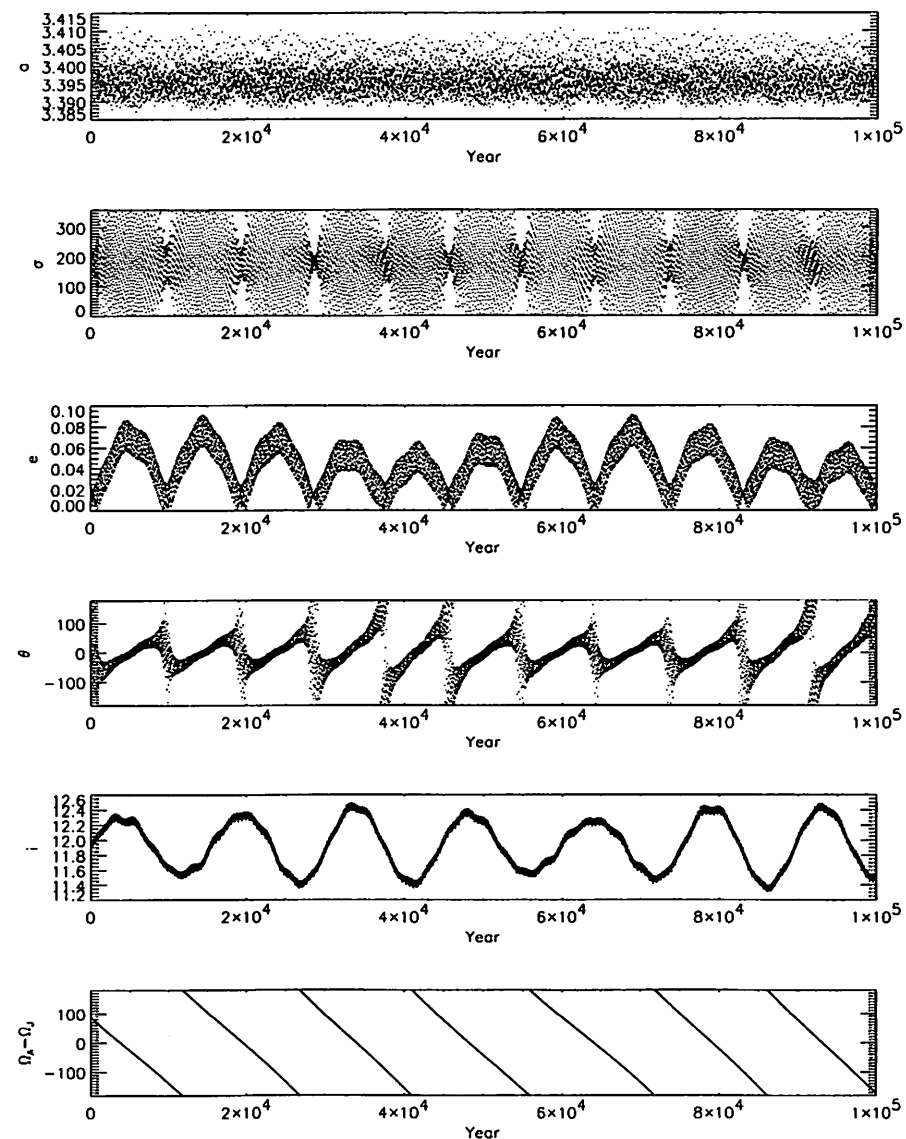
**Figure 6.** The orbital elements of (279)Thule : 4:3 pericentric libration. The six panels represent the semi-major axis,  $\sigma = 4\lambda_J - 3\lambda_A - \varpi_A$ , the eccentricity,  $\theta = \varpi_A - \varpi_J$ , the inclination, and  $\Omega_A - \Omega_J$ , respectively.



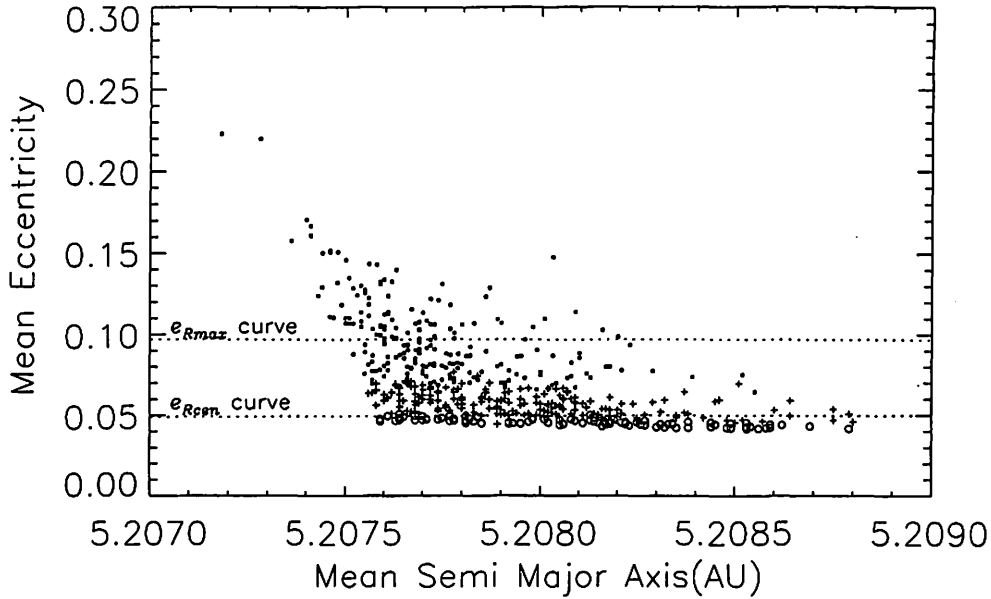
**Figure 7-1.** The orbital elements of (1362)Griqua : 2:1 libration. The six panels represent the semi-major axis,  $\sigma = 2\lambda_J - \lambda_A - \varpi_A$ , the eccentricity,  $\theta = \varpi_A - \varpi_J$ , the inclination, and  $\Omega_A - \Omega_J$ , respectively.



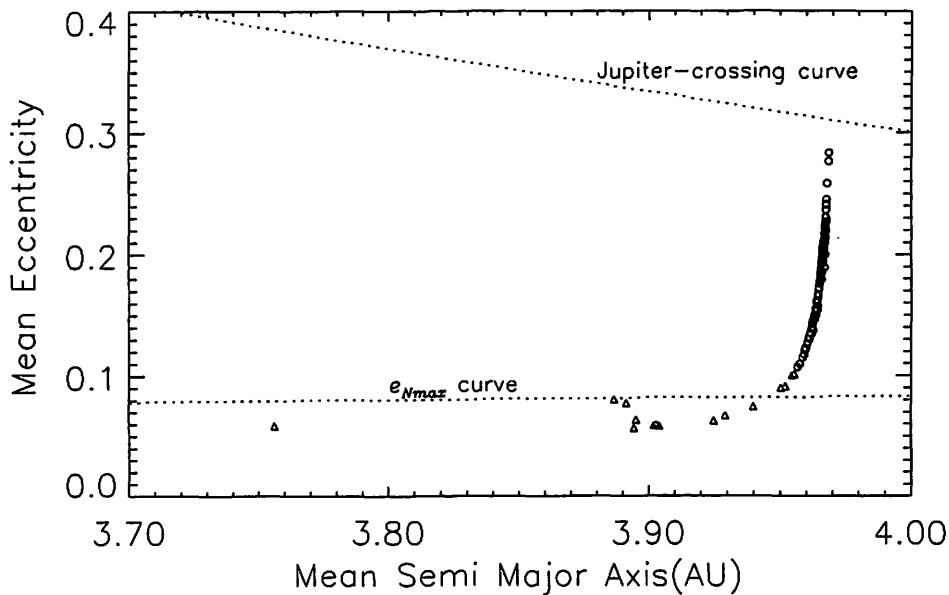
**Figure 7-2.** The orbital elements of (300)Geraldina : 2:1 pericentric libration. The explanations for the vertical axes are same as for Figure 7-1.



**Figure 7-3.** The orbital elements of (528)Rezia : 2:1 apocentric libration. The explanations for the vertical axes are same as for Figure 7-1.

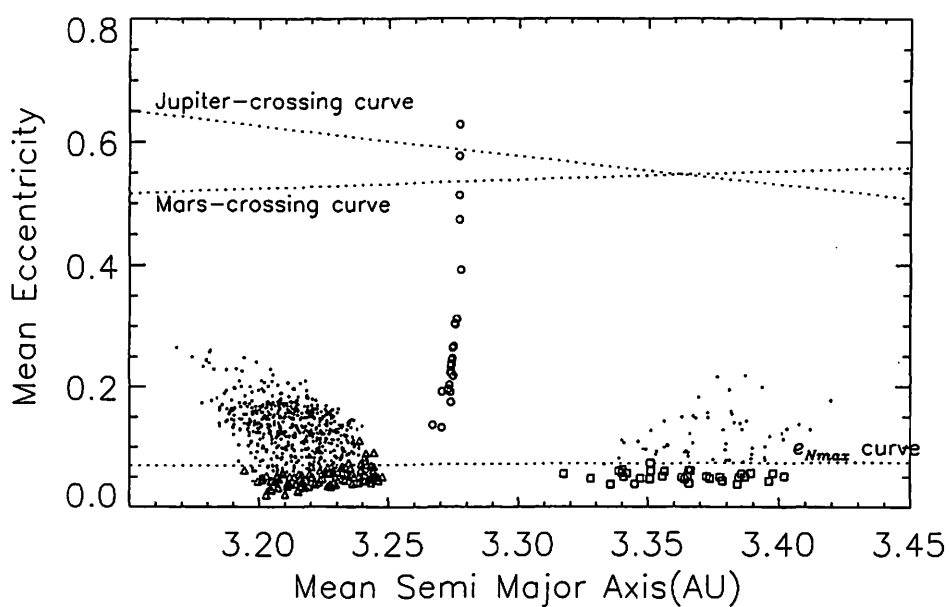


**Figure 8.** Mean semi-major axes versus mean eccentricities of the asteroids in the 1:1 mean motion resonance with Jupiter. Mean values of semi-major axes and eccentricities are calculated by numerical integrations for  $10^5$  years. The critical arguments ( $\sigma$ ) of all asteroids in this figure librate around  $\pm 60^\circ$  perfectly. The open circles indicate  $\theta$  librates always. The plus signs indicate that  $\theta$  alternates libration and circulation. The filled box signs indicate  $\theta$  circulate always.



**Figure 9.** Mean semi-major axes versus mean eccentricities of the asteroids in the 3:2 mean motion resonance with Jupiter. The circles indicate  $\sigma$  librates always. The triangles indicate  $\sigma$  interchanges libration with circulation alternately, then  $\theta$  contrariwise interchanges circulation with libration. When  $\sigma$  librates, it's librational center is  $0^\circ$ .





**Figure 10.** Mean semi-major axes versus mean eccentricities of the asteroids in the 2:1 mean motion resonance with Jupiter.

The circles and triangles are same as figure 9. The boxes indicate  $\sigma$  interchanges libration with circulation alternately, then  $\theta$  contrariwise interchanges circulation with libration. When  $\sigma$  librates, it's librational center is  $180^\circ$ . The dots indicate non-resonant asteroids:  $\sigma$  and  $\theta$  do not librate always.

# Orbital Theory of a Highly Eccentric Satellite Disturbed by a Massive Inner Satellite

Yoshimitsu MASAKI<sup>1</sup> and Hiroshi KINOSHITA<sup>2</sup>

<sup>1</sup> The Graduate University for Advanced Studies (National Astronomical Observatory, Japan)

masakiys@cc.nao.ac.jp

2-21-1, Osawa, Mitaka, Tokyo, 181-8588, Japan

<sup>2</sup> National Astronomical Observatory, Japan

Kinoshita@nao.ac.jp

2-21-1, Osawa, Mitaka, Tokyo, 181-8588, Japan

## ABSTRACT

We have developed an analytical theory for a celestial body orbiting in a highly eccentric orbit under the perturbational influence of an inner body which revolves in a circular orbit around a central body (a restricted three-body problem). We made the Hamiltonian closed in form to the orbital eccentricity.

We confirmed that our theory is highly accurate by comparing numerically integrated results. However, the theory loses its high accuracy when the eccentricity of the outer body is very large.

Our theory can be applied to some celestial bodies. The motion of the Neptunian satellite Nereid orbiting in a highly eccentric orbit ( $e = 0.75$ ) perturbed by Triton is one example. Our theory provides a degree of accuracy, with results generally much better than 30Km in the (osculating) semimajor axis of Nereid.

## 1 Introduction

Astronomical ephemerides provide precise positions of celestial bodies. Today, numerically integrated ephemerides are widely used in the world. DE (Development Ephemeris) series compiled by JPL is one example.

When we construct an analytical ephemeris, we usually handle a perturbation theory. Results are expressed as osculating elements which are functions of time. Since most planets or satellites revolve in nearly circular orbits, we expand a perturbing function in terms of powers of eccentricity. However, when we construct the orbital theory of a highly eccentric body, the power series of eccentricity converges quite slowly.

For example, a Neptunian satellite Nereid revolves on a highly eccentric orbit ( $e = 0.75$ ). Mignard(1975)'s study is the pioneered work on the motion of Nereid. Saad(2000) studied

the motion of Nereid using a canonical perturbation method of Hori type. The inner orbiting satellite, Triton, is not taken into consideration in these studies because its perturbational effect is weaker than the Sun's and is not detected by astrometric observations from ground telescopes.

Oberti (1990), Segerman and Richardson (1997) developed the analytical theory of Nereid under the perturbational influence of Triton. In their work, the motion of Nereid is described in the barycentric coordinate system of Neptune and Triton to express the problem in simple form. Oberti (1990) expanded the Hamiltonian in eccentric anomaly, thereafter Segerman and Richardson (1997) expanded the Hamiltonian in eccentricity. However, it is suspected that these theories can provide good accuracy because they did not show any experimental check in their papers.

In this study, we have developed an orbital theory of a highly eccentric body under the perturbative influence of an inner revolving body. It is one kind of restricted three-body problem. In the previous paper (Masaki and Kinoshita (2001)), we proposed an orbital theory of the planar restricted problem.

We construct an analytical theory using a Lie-type canonical perturbation method, proposed by Hori (1966). This theory can be applied not only to a planar restricted problem but also to an inclined case. We made the Hamiltonian closed in form to the eccentricity in order to apply it to any highly eccentric orbit.

## 2 Analytical formulation

We describe the motion of a celestial body (hereafter, we call it the 'outer') moving in a highly eccentric orbit around a pair consisting of the primary body and an inner revolving body (called the 'primary' and the 'inner', respectively). See Figure 1. The mass of the inner body is small enough compared to the primary, and that of the outer can be neglected (i.e. mass-less particle). For brevity's sake, we can say that the inner body orbits around the primary star in a circular motion.

Hereafter, we designate masses of the Primary, Inner and Outer as  $M, m', m$ , respectively. The symbols used in this paper are listed in Appendix pages. The universal gravity constant is written by  $k^2$ .

When we consider the motion of the outer orbiting body, it is preferable to introduce the barycentric (Jacobi) coordinate system of the primary and the inner, because the time-variation in osculating elements is limited in the small magnitude as we see in Brouwer and Clemence (1961), Oberti (1990) and Segerman and Richardson (1997).

The (perturbed) Hamiltonian for Outer becomes,

$$F = \mu \left[ \frac{1}{2a} + \frac{Mm'}{(M+m')^2} \frac{r'^2}{r^3} P_2(\cos S) + \frac{Mm'(M-m')}{(M+m')^3} \frac{r'^3}{r^4} P_3(\cos S) + \dots \right],$$

where  $S$  is the elongation between Inner and Outer, and  $P_i$  is a Legendre polynomial of degree  $i$ . Using spherical trigonometry,  $S$  can be expressed by the angular orbital elements,  $f, \omega, \Omega$  and

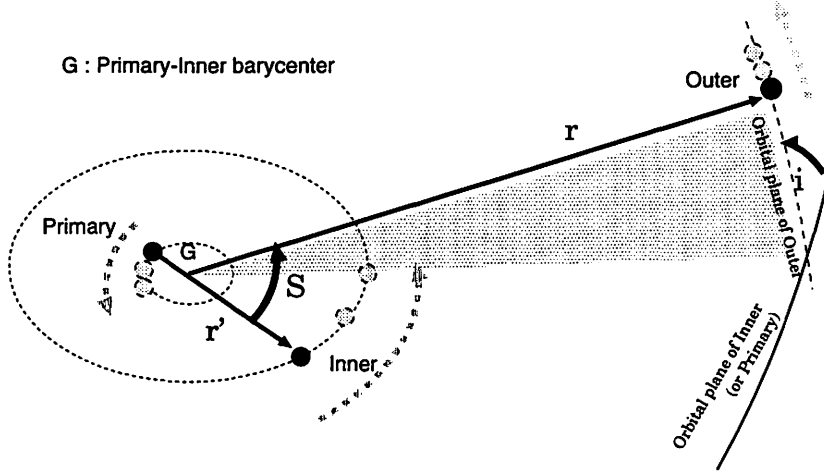


Figure 1: A schematically illustrated model of an inclined restricted problem

$\lambda'$  (See Figure 2) as follows:

$$\cos S = \cos(f + \omega) \cos(\lambda' - \Omega) + \sin(f + \omega) \sin(\lambda' - \Omega) \cos I$$

From now on, we simplify the problem: (1) Inner moves in a circular orbit. (2) Outer's perturbation does not affect the motion of the Primary or Inner (i.e. we neglect the mass of Outer). In other words, we consider a circular restricted three-body problem.

The elongation  $S$  contains a variable  $\lambda'$ , which depends on time:

$$\lambda' \equiv k = n't + \text{const..}$$

To make a Hamiltonian independent of time, we introduce a canonical conjugate action variable,  $K$ . The term  $-n'K$  has to be added to the Hamiltonian.

There are three independent angular variables in this system:  $f, g$  and  $h - k$ , i.e. there are three degrees of freedom. Therefore, we can deduce the Hamiltonian including only three sets of canonical variables,  $(y_1, x_1)$ ,  $(y_2, x_2)$  and  $(y_3, x_3)$ , after a suitable canonical transformation.

$$F(l, g, h, k, L, G, H, K) \longrightarrow F(y_1, y_2, y_3, x_1, x_2, x_3).$$

One example is:

$$y_1 = l, \quad y_2 = g, \quad y_3 = h - k, \quad y_4 = k$$

$$x_1 = L, \quad x_2 = G, \quad x_3 = H, \quad x_4 = K + H.$$

Here,  $(l, g, h)$  and  $(L, G, H)$  are the canonical set of Delaunay variables. Since  $y_4$  does not depend on  $F$  anymore, we can eliminate  $x_4$  from the Hamiltonian. Finally, we obtain  $F$ :

$$F = \mu \frac{1}{2a} + n'x_3$$

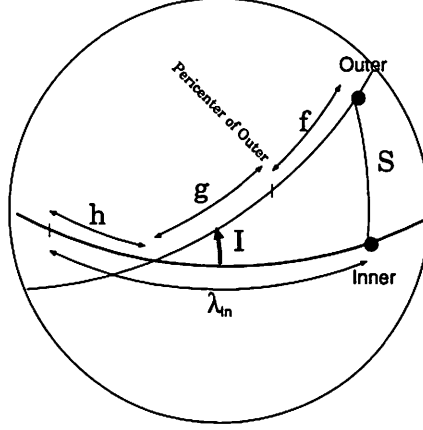


Figure 2: An angular distance  $S$  expressed by angular orbital elements.

$$\begin{aligned}
 & +\mu \frac{M m'}{(M+m')^2} \frac{a^3}{r^3} \frac{a'^2}{a^3} \frac{1}{4} (3 \cos(2S) + 1) \\
 & +\mu \frac{M m'(M-m')}{(M+m')^3} \frac{a^4}{r^4} \frac{a'^3}{a^4} \frac{1}{8} (5 \cos(3S) + 3 \cos S) + \dots
 \end{aligned}$$

### 3 Hori's canonical perturbation theory

In this section, we briefly discuss Hori's canonical perturbation theory. See Hori(1966)'s work in detail.

Suppose Hamiltonians  $F(x, y)$  and  $F^*(x^*, y^*)$  are expanded in a small parameter  $\epsilon$ , i.e.,

$$\begin{aligned}
 F &= F_0 + F_1 + F_2 + \dots \\
 F^* &= F_0^* + F_1^* + F_2^* + \dots,
 \end{aligned}$$

where subscripts mean powers of  $\epsilon$ .

An arbitrary function  $f$  of canonical variables  $(x, y)$  can be developed in a converged series of  $\epsilon$  using the Lie theorem:

$$f(x, y) = \sum_{n=0}^{\infty} \frac{\epsilon^n}{n!} D_s^n f(x^*, y^*),$$

where  $D_s$  is an operator of  $n$ -times Poisson bracket with  $S$ , i.e.:

$$\begin{aligned}
 D_s^0 &= f \\
 D_s^1 &= \{f, S\} \\
 D_s^2 &= \{\{f, S\}, S\} \\
 \dots &= \dots
 \end{aligned}$$

Substituting a Hamiltonian  $F$  in  $f$ , we obtain  $F$  in a series of  $\epsilon$  with variables  $(x^*, y^*)$ . Comparing this expression with  $F^*(x^*, y^*)$ , we obtain the following equivalences for each power

of  $\epsilon$ :

$$\begin{aligned}
F_{0*} &= F_0 \\
F_{1*} &= \{F_0, S_1\} + F_1 \\
F_{2*} &= \{F_0, S_2\} + \{F_1, S_1\} + \frac{1}{2}\{\{F_0, S_1\}, S_1\} + F_2 \\
F_{3*} &= \{F_0, S_3\} + \{F_1, S_2\} + \{F_2, S_1\} \\
&\quad + \frac{1}{2}\{\{F_0, S_2\}, S_1\} + \frac{1}{2}\{\{F_0, S_1\}, S_2\} + \frac{1}{2}\{\{F_1, S_1\}, S_1\} \\
&\quad + \frac{1}{6}\{\{\{F_0, S_1\}, S_1\}, S_1\} + F_3 \\
\ldots &= \ldots
\end{aligned}$$

Variables before transformation,  $(x, y)$  require expression in variables  $(x^*, y^*)$ . With Hori's theory, they are written:

$$\begin{aligned}
x &= x^* + \sum_{n=1}^{\infty} \frac{\epsilon^n}{n!} D_S^{n-1} \frac{\partial S}{\partial y^*} \\
&= x^* + \epsilon \frac{\partial S}{\partial y^*} + \frac{1}{2} \epsilon^2 \left\{ \frac{\partial S}{\partial y^*}, S \right\} + \ldots
\end{aligned}$$

and

$$\begin{aligned}
y &= y^* - \sum_{n=1}^{\infty} \frac{\epsilon^n}{n!} D_S^{n-1} \frac{\partial S}{\partial x^*} \\
&= y^* - \epsilon \frac{\partial S}{\partial x^*} - \frac{1}{2} \epsilon^2 \left\{ \frac{\partial S}{\partial x^*}, S \right\} - \ldots
\end{aligned}$$

## 4 Building an analytical theory

We decompose perturbations into four parts according to their periodicities. They are:

- Short periodic perturbation (caused by the revolution of the inner body)
- Intermediate periodic perturbation (caused by the revolution of the outer body)
- Long periodic perturbation (caused by the circulation of the pericenter of the outer body)
- Secular perturbation .

In this chapter, we mainly describe the equations for the perturbed Hamiltonian only up to  $P_2$  terms of Legendre polynomials for the short or intermediate periodic perturbations, and up to  $P_6$  terms for the long periodic ones.

### 4.1 Short periodic perturbation: Elimination of short periodic terms

We average the Hamiltonian over the short periodic variable,  $y_3$  and decrease the number of degrees of freedom by one. In other words, we let the new Hamiltonian  $F^*$  be free from  $y_3$ . The Hamiltonian is transformed into:

$$F(y_1, y_2, y_3, x_1, x_2, x_3) \longrightarrow F^*(y_1, y_2, x_1, x_2).$$

The original Hamiltonian (before transformation) is written as:

$$F = F_0 + F_1 + F_2$$

where,

$$\begin{aligned} F_0 &= n' x_3 \\ F_1 &= \frac{\mu}{2a} \\ F_2 &= \mu \frac{Mm'}{(M+m')^2} \frac{a^3}{r^3} \frac{a'^2}{a^3} \left[ \frac{1}{8}(-1 + 3\theta^2) \right. \\ &\quad + \frac{3}{8}(1 - \theta^2) \cos(2f + 2y_2) \\ &\quad + \frac{3}{16}(1 - \theta^2)^2 \cos(2f + 2y_2 - 2y_3) \\ &\quad + \frac{3}{8}(1 - \theta^2) \cos(2y_3) \\ &\quad \left. + \frac{3}{16}(1 + \theta)^2 \cos(2f + 2y_2 + 2y_3) \right], \end{aligned}$$

and

$$\theta \equiv \cos I.$$

Subscripts are approximate orders of a small parameter  $\frac{n}{n'}$ , and from here on, the expressions are neglected  $(\frac{m'}{M+m'})^2$  or higher order terms.

A new Hamiltonian  $F^* = \sum_i F_i^*$  and a generating function  $S = \sum_i S_i$  are:

$$\begin{aligned} F_0^* &= F_0(\text{UNPERTURBED}) \\ F_1^* &= [F_1]_{sec} \\ &= \frac{\mu^2}{2x_1^2} \\ S_1 &= \int [F_1]_{per} dt^* \\ &= 0 \\ F_2^* &= [\{F_1, S_1\} + F_2]_{sec} \\ &= \frac{1}{8} \mu \frac{Mm'}{(M+m')^2} \frac{a^3}{r^3} \frac{a'^2}{a^3} \left[ (-1 + 3\theta^2) + 3(1 - \theta^2) \cos(2f + 2y_2) \right] \\ S_2 &= \int [\{F_1, S_1\} + F_2]_{per} dt^* \\ &= -\frac{\mu}{n'} \frac{Mm'}{(M+m')^2} \frac{a^3}{r^3} \frac{a'^2}{a^3} \left[ -\frac{3}{32}(1 - \theta^2)^2 \sin(2f + 2y_2 - 2y_3) \right. \\ &\quad + \frac{3}{16}(1 - \theta^2) \sin(2y_3) \\ &\quad \left. + \frac{3}{32}(1 + \theta)^2 \sin(2f + 2y_2 + 2y_3) \right] \\ F_3^* &= [\{F_1, S_2\}]_{sec} \\ &= 0 \\ S_3 &= \int [\{F_1, S_2\}]_{per} dt^* \end{aligned}$$

$$\begin{aligned}
F_4^* &= [\{F_1, S_3\}]_{sec} + O((\frac{m}{M+m'})^2) \\
&= 0 \\
\dots &= \dots \\
F_i^* &= [\{F_1, S_{i-1}\}]_{sec} + O((\frac{m}{M+m'})^2) \\
S_i &= \int [\{F_1, S_{i-1}\}]_{per} dt^* + O((\frac{m}{M+m'})^2) \\
\dots &= \dots
\end{aligned}$$

Here,  $\{X, Y\}$  is an operation of the Poisson bracket of  $X$  and  $Y$ .  $[Q]_{sec}$  and  $[Q]_{per}$  are operations of getting secular and periodic parts of  $Q$ , respectively. An artificial time  $t^*$  satisfies the following relations:

$$\begin{aligned}
\frac{dx^*}{dt^*} &= \frac{\partial F_0}{\partial y^*} \\
\frac{dy^*}{dt^*} &= -\frac{\partial F_0}{\partial x^*}.
\end{aligned}$$

## 4.2 Intermediate periodic perturbation: Elimination of intermediate periodic terms

Next, we eliminate an intermediate periodic perturbation by using a canonical transformation

$$F^*(y_1, y_2, x_1, x_2) \longrightarrow F^{**}(y_2, x_2).$$

we have to add additional terms for  $S_1^*$  to avoid contaminating secular trends in angular variables.

We obtain the Hamiltonian  $F^{**}$  and a generating function  $S^*$  as follows:

$$\begin{aligned}
F_0^{**} &= F_0^* \\
F_1^{**} &= F_1^*(\text{UNPERTURBED}) \\
F_2^{**} &= \frac{1}{8}\mu \frac{Mm'}{(M+m')^2} \frac{1}{\eta^3} \frac{a'^2}{a^3} (-1 + 3\theta^2) \\
S_1^* &= \mu \frac{Mm'}{(M+m')^2} \frac{1}{n\eta^3} \frac{a'^2}{a^3} \left[ \frac{1}{8}(-1 + 3\theta^2)(f + e \sin f - y_1) \right. \\
&\quad \left. + \frac{3}{8}(1 - \theta^2) \left\{ \frac{1}{2}e \sin(f + 2y_2) + \frac{1}{2} \sin(2f + 2y_2) + \frac{1}{6}e \sin(3f + 2y_2) \right\} \right. \\
&\quad \left. - \frac{1}{16}(1 - \theta^2) \frac{1}{e^2} (2 - 3e^2 - 2\eta^3) \sin(2y_2) \right] \\
F_3^{**} &= O((\frac{m}{M+m'})^2) \\
S_2^* &= O((\frac{m}{M+m'})^2) \\
\dots &= \dots
\end{aligned}$$



### 4.3 Long periodic perturbation: Elimination of long periodic terms

Finally, we eliminate the long periodic variable  $y_2$  from the Hamiltonian and obtain a new Hamiltonian  $F^{***}$  free from any angular variables  $y$ .

However, we have seen in the previous section, the Hamiltonian  $F^{**}$  does not contain an angular variable. This is attributed to the fact that we take only  $P_2$  perturbational contributions into consideration. In other words, if we include higher  $P_i$  terms (in practice, only the even numbers of  $i$  contribute intermediate or long periodic perturbations),  $F^{**}$  contains trigonometric functions of  $y_2$ . For example, if we follow up till  $P_6$  terms,

$$\begin{aligned}
 F_0^{**} &= n'x_3 \\
 F_1^{**} &= \frac{\mu^2}{2x_1^2} \\
 F_2^{**}(P_2) &= \mu \frac{Mm'}{(M+m')^2} \frac{1}{\eta^3} \frac{a'^2}{a^3} \frac{1}{8} (-1 + 3\theta^2) \\
 F_2^{**}(P_4) &= \mu C_4 \frac{1}{\eta^7} \frac{a'^4}{a^5} \left[ \frac{9}{1024} (3 - 30\theta^2 + 35\theta^4)(2 + 3e^2) \right. \\
 &\quad \left. - \frac{45}{512} (1 - \theta^2)(1 - 7\theta^2)e^2 \cos(2y_2) \right] \\
 F_2^{**}(P_6) &= \mu C_6 \frac{1}{\eta^{11}} \frac{a'^6}{a^7} \left[ \frac{25}{32768} (-5 + 105\theta^2 - 315\theta^4 + 231\theta^6)(8 + 40e^2 + 15e^4) \right. \\
 &\quad + \frac{2625}{32768} (1 - \theta^2)(1 - 18\theta^2 + 33\theta^4)e^2(2 + e^2) \cos(2y_2) \\
 &\quad \left. + \frac{1575}{65536} (1 - \theta^2)^2 (-1 + 11\theta^2)e^4 \cos(4y_2) \right],
 \end{aligned}$$

where,

$$\begin{aligned}
 C_4 &\equiv \frac{Mm'(M^3 + m'^3)}{(M + m')^5} \\
 C_6 &\equiv \frac{Mm'(M^5 + m'^5)}{(M + m')^7}
 \end{aligned}$$

From now on, we redefine  $F_2^{**}$  and  $F_3^{**}$ :

$$\begin{aligned}
 F_2^{**} &= F_2^{**}(P_2) \\
 F_3^{**} &= F_2^{**}(P_4) + F_2^{**}(P_6).
 \end{aligned}$$

We consider the new  $F_2^{**}$  as an unperturbed Hamiltonian for a canonical transformation

$$F^{**}(y_2, x_2) \longrightarrow F^{***},$$

and eliminate  $y_2$  from  $F^{**}$ .

The new Hamiltonian  $F^{***} = \sum_i F_i^{***}$  and the generating function  $S^{**} = \sum_i S_i^{**}$  are:

$$\begin{aligned}
 F_0^{***} &= F_0^{**} \\
 F_1^{***} &= F_1^{**} \\
 F_2^{***} &= F_2^{**}(\text{UNPERTURBED})
 \end{aligned}$$

$$\begin{aligned}
F_3^{***} &= [F_3^{**}]_{sec} \\
&= \mu \left[ C_4 \frac{1}{\eta^7} \frac{a'^4}{a^5} \frac{9}{1024} (3 - 30\theta^2 + 35\theta^4)(2 + 3e^2) \right. \\
&\quad \left. + C_6 \frac{1}{\eta^{11}} \frac{a'^6}{a^7} \frac{25}{32768} (-5 + 105\theta^2 - 315\theta^4 + 231\theta^6)(8 + 40e^2 + 15e^4) \right] \\
S_1^{**} &= \int [F_3^{**}]_{per} dt^{***} \\
&= -\frac{15}{128} B_4 \frac{a'^2}{\eta^3} \frac{n}{(-1 + 5\theta^2)} e^2 (1 - \theta^2)(1 - 7\theta^2) \sin(2y_2) \\
&\quad - \frac{175}{32768} B_6 \frac{a'^4}{a^2 \eta^7} \frac{n}{(-1 + 5\theta^2)} \left[ -20e^2 (2 + e^2)(1 - \theta^2)(1 - 18\theta^2 + 33\theta^4) \sin(2y_2) \right. \\
&\quad \left. 3e^4 (1 - \theta^2)^2 (-1 + 11\theta^2) \sin(4y_2) \right] \\
F_4^{***} &= [\{F_3^{**}, S_1^{**}\}]_{sec} \\
S_2^{**} &= \int [\{F_3^{**}, S_1^{**}\}]_{per} dt^{***} \\
\dots &= \dots,
\end{aligned}$$

where,

$$\begin{aligned}
B_4 &\equiv \frac{M^3 + m'^3}{(M + m')^3} \\
B_6 &\equiv \frac{M^5 + m'^5}{(M + m')^5}.
\end{aligned}$$

It is noted that  $S_1^{**}$  is  $O((\frac{m'}{M+m'})^0)$  because in integrating  $\int Q dt^{***}$ , a quantity  $Q$  is divided by a factor of  $O((\frac{m'}{M+m'})^1)$ .

#### 4.4 Secular perturbation

We have obtained a Hamiltonian  $F^{***}$  which does not depend on any angular variables.

The equations of motion are:

$$\begin{aligned}
\frac{dx^{***}}{dt} &= \frac{\partial F^{***}}{\partial y^{***}} (\equiv 0) \\
\frac{dy^{***}}{dt} &= -\frac{\partial F^{***}}{\partial x^{***}}
\end{aligned}$$

From them, we obtain

$$\begin{aligned}
x^{***} &= const. \\
y^{***} &= \left(-\frac{\partial F^{***}}{\partial x^{***}}\right)t + const..
\end{aligned}$$

That is, the action variables (the semimajor axis  $a^{***}$ , the eccentricity  $e^{***}$  and the inclination  $I^{***}$ ) are constants, while the angular variables (the mean anomaly  $l^{***}$ , the argument of the perihelion  $\omega^{***}$  and the longitude of the ascending node  $\Omega^{***}$ ) increase (or decrease) linearly with time  $t$ .

Orbital elements  $a^{***}, e^{***}, I^{***}, l^{***}, \omega^{***}, \Omega^{***}$  deduced from  $x^{***}$  and  $y^{***}$  are mean orbital elements.

#### 4.5 Deriving osculating elements

Osculating elements  $E$  for canonical variables are summed up by the following contributions:

- Mean elements  $E^{***}$
- Contribution from long periodic perturbation  $\delta E^{**}$
- Contribution from intermediate periodic perturbation  $\delta E^*$
- Contribution from short periodic perturbation  $\delta E$ .

Therefore, we evaluate the following procedure: first,

$$E^{**} = E^{***} + \delta E^{**}(E^{***}),$$

then,

$$E^* = E^{**} + \delta E^*(E^{**}),$$

and finally

$$E = E^* + \delta E(E^*).$$

Neglecting  $O((\frac{m}{M+m'})^2)$  terms, then,

$$\begin{aligned} \delta E^{**} &= \{E^{**}, S^{**}\} \\ \delta E^* &= \{E^*, S^*\} \\ \delta E &= \{E, S\}. \end{aligned}$$

### 5 Checking accuracy (comparison with numerical integrations)

The accuracy of ephemerides is assessed by comparing the results to observational data, i.e. calculating residuals of  $O - C$  (observed) – (calculated)). In general, observations are contaminated by noise, the constants are fit by methods of least squares.

In this study, we check the accuracy of our theory by comparing numerical integrations and calculating residuals of ((analyticals) – (numericals)). See Figure 3. Numerical integration is performed by Bulirsch-Stoer (“extrapolation method”) code in double-precision accuracy. This code provides high accuracy in results, suitable for our aim. We start integration with position and velocity values that are converted into Cartesian coordinates from a set of analytical osculating elements at the initial time. Then, at a time  $T$ , residuals (analytical results minus numerically integrated ones) are calculated as:

$$(\text{Residuals}) = (\text{Analytical results}) - (\text{Numerical results}).$$

When we calculate residuals in angular variables,  $y_1$ ,  $y_2$  and  $y_3$ , secular trends (the slopes of regression line for raw data) are subtracted, not to hide fine structures in residuals for output figures.

Comparison between analytical and numerical results

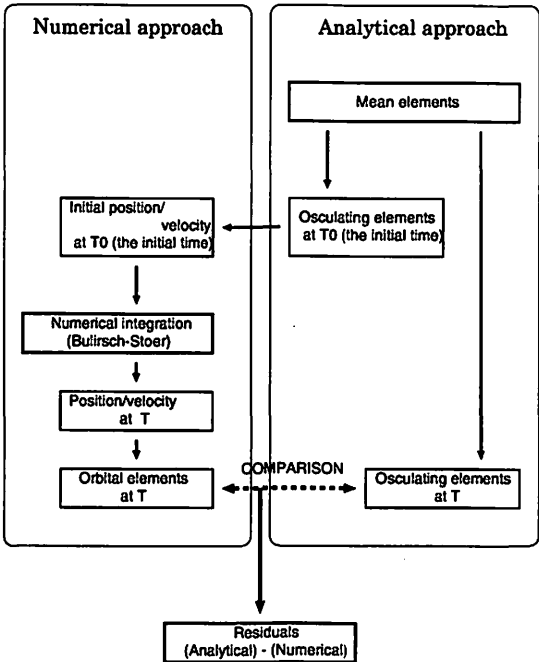


Figure 3: Flowchart for calculating residuals

## 6 Accuracy of our analytical theory

### 6.1 Results and residuals

Now we check residuals between analytical and numerical results. The former contains up to  $P_5$  perturbations (and up to  $S_5, S_1^*$  and  $S_2^{**}$  terms for each periodic perturbations), while the latter contains the full perturbing force caused by the inner body in the equations of motion. We use the same parameters as in the problem of Nereid. (See Table 1, except for  $a = 5.5 \times 10^6$ [Km] and  $e = 0.75$ .) The residuals are shown in Fig. 4.

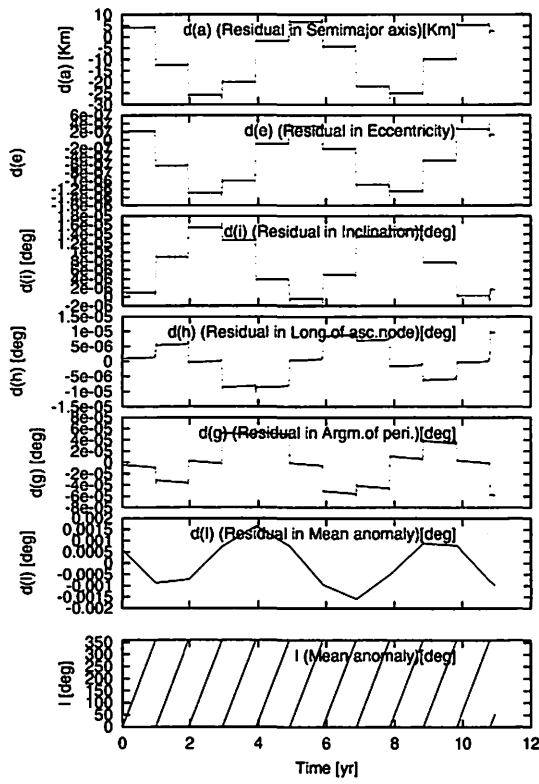
Short periodic variations, which are expressed by trigonometric functions of  $y_3$  (the periodicity is  $\sim 6$ [days]), are not seen in the residuals. That is, they are almost perfectly calculated by our theory.

However, a step-shaped discrepancies are seen in the residuals. The maximum divergence of our analytical solution from the numerical one is the level of  $\sim 30$ [Km] in the semimajor axis.

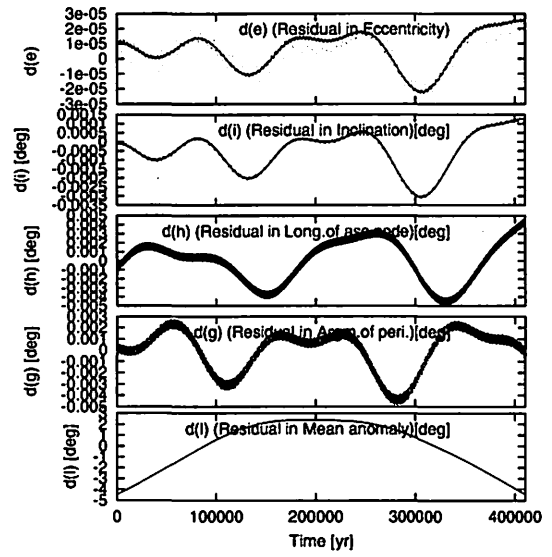
Next, we check residuals for a long timespan. In  $e, I$  and  $\omega$ , misfits of  $\cos 4\omega$  or  $\sin 4\omega$  terms dominate. ( $\omega$  has the periodicity of  $\sim 3.4 \times 10^5$ [yr].) They come from the higher order terms of  $(\frac{m'}{M+m'})^2$  that we have neglected in our construction of a theory.

Residuals in mean anomaly  $l$  remain quadratic trends. They are due to growth of numerical round-off errors in longitude (or position) with time. I.e., they are attributed to the errors in numerical integration, not to those in our analytical theory.

Therefore, we conclude that our analytical theory represents the true orbital motion of an



(a)



(b)

Figure 4: Residuals of numerical results minus analytical ones (a) for a short timespan and (b) for a long timespan. Secular trends in the angular variables are subtracted.

Table 1: Parameters and (mean) orbital elements used in this study. The inclination of the outer body refers to the orbital plane of the inner one. We use the same parameters as in the Nereid problem except for those of the semimajor axis and the eccentricity.

Item	Model
Mass of the primary body [ $M_{\odot}$ ]	$5.1514 \times 10^{-5}$
Mass of the inner body	$(2.89 \times 10^{-4}) \times (\text{Mass of Primary})$
Mass of the outer body	0. (test particle)
Semimajor axis [Km]	variable (Integer multiples of $5.5 \times 10^6$ )
Eccentricity	variable
Inclination [deg]	132.4
Longitude of ascending node [deg]	0.0
Argument of pericenter [deg]	0.0
Initial longitude of the outer body [deg]	0.0
Semimajor axis of the inner body [Km]	$14.15 \times 24764$
Eccentricity of the inner body	0.0
Inclination of the inner body [deg]	0.0

eccentric celestial body fairly well, except for step-shaped errors. Our theory maintains the accuracy of within 30[Km] in the osculating semimajor axis.

## 6.2 Normalized residuals

In this section, we define the term ‘normalized’ residual range:

$$(\text{normalized residual range}) = \frac{\text{residual range}}{\text{magnitude of perturbations}}$$

The reason for introducing the normalized residual range is the following: if we change orbital parameters of the outer body, such as  $a$  or  $e$ , the magnitudes of perturbation are also changed.

We adopted the initial values as shown in Table 1. We vary the values of  $a$  and  $e$  and check the normalized residual ranges. We used the same values as in the problem of Nereid for other parameters.

The growth of the normalized residual ranges with the eccentricity is shown in Fig. 5. Our theory maintains a high degree of accuracy for a wide range of eccentricity, especially for a larger semimajor axis. The normalized residual range on the order of  $\sim 10^{-5}$  is a machine precision limit for calculating residuals. Therefore, our theory perfectly agrees with numerical results.

However, for a larger eccentricity or for a smaller semimajor axis, the combined analytical model degrades accurately. This is due to the following:

- Offset growth

According to the results of numerical integration (Fig. 6), step-shaped abrupt changes in the orbital elements are observed when Nereid passes its pericenter. We define ‘offset’ as a

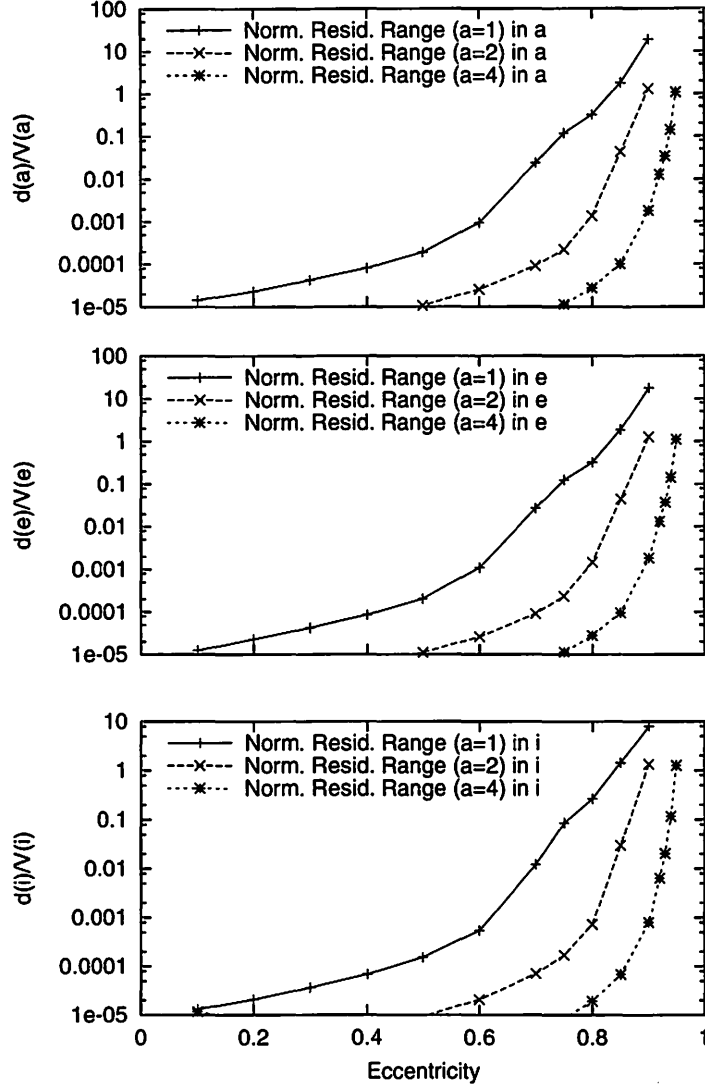


Figure 5: The growth of the normalized residual ranges with the eccentricity. Each curve shows the values for the same semimajor axis. (From the upper panel to the lower) Normalized residual ranges in the semimajor axis, the eccentricity and the inclination.

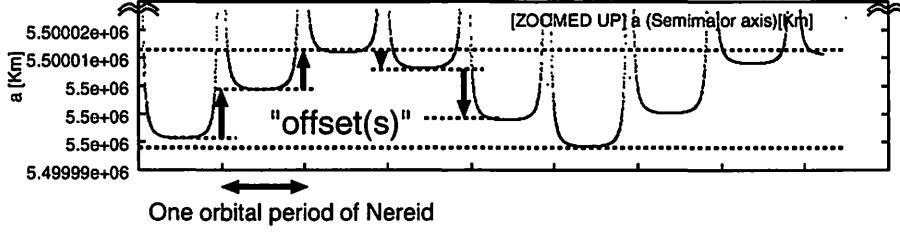


Figure 6: Definition of offset in this study.

difference of the values of an orbital element between the successive revolutionary periods. Our analytical theory does not describe these phenomena.

- Ill-convergent series of  $S_i$

For short periodic perturbations, the series of  $S_i$  converges by a factor of  $\frac{n}{n'}$ . However,  $S_i$  is calculated as follows:

$$\begin{aligned}
 S_i &= \int \{F_1, S_{i-1}\} dt^* \\
 &= -\frac{n}{n'} \int \frac{\partial S_{i-1}}{\partial y_1} dy_2 \\
 &= -\frac{n}{n'} \frac{(1 + e \cos f)^2}{\eta^3} \int \frac{\partial S_{i-1}}{\partial f} dy_2.
 \end{aligned}$$

The factor of  $\frac{(1+e \cos f)^2}{\eta^3}$  becomes a larger value for a large eccentricity, especially at the pericenter,  $\frac{(1+e)^2}{\eta^3}$ . Therefore it prevents  $S_i$  from converging.

- Truncational error of Legendre polynomials

Our theory is truncated by the  $P_5$  terms of Legendre polynomials. Thus, the accuracy of our theory is decreased when the outer body approaches the pericenter.

- Truncational error of canonical transformation

Our theory is truncated by the  $S_5$  terms for short periodic perturbations and neglects terms of the order of  $(\frac{m'}{M+m'})$ . It cannot fully explain all of the perturbations.

## 7 Discussion

### 7.1 Validity of neglecting higher order terms of $(\frac{m'}{M+m'})$

In this study, we deal with the small factor of the order  $(\frac{m'}{M+m'})^2$  or higher terms as negligible parameters.

In practice, this simplification is valid for the system that the theory is applicable to. For the Nereid system,

$$\frac{m'}{M} \sim 2.89 \times 10^{-4},$$



the squared value is approximately of the order of  $10^{-7}$ . The value is small enough to neglect the terms of the order of  $(\frac{m'}{M+m'})^2$ .

Similarly, for the extrasolar planetary system, it depends on each system. The upper limit of the planetary mass is bounded by 13 Jupiter masses. (This is the lower limit of the mass for brown dwarfs. Recently, this classification has come to be accepted among astronomers. See Martin et.al.(1999) or Oppenheimer et.al.(2000).) Therefore, for an extrasolar planet orbiting around a sun-like star,  $(\frac{m'}{M+m'})^2$  is smaller than the order of  $10^{-4}$ .

## 7.2 Reliability of Numerical Integration

We have checked the accuracy of our analytical theory by comparing it with the numerical results integrated by Bulirsch-Stoer. Results of numerical integrations are degraded by numerical errors, like a 'round-off'. They occur in a round-off process at the smallest digit throughout calculation. For a long-interval calculation, the result suffers severely from these errors; therefore, 'good' integration codes are for the purpose.

The Bulirsch-Stoer code used in this study is widely admitted as a highly accurate integration code for a relatively short interval (which means that it is not suited for the age of the Solar system). Murison (1989) discussed the usefulness of the code for keeping a Jacobi integral value, which is the integral for the restricted three-body problem in the corotating coordinate system with a perturbing body, throughout his numerical integration.

## 8 Further Application of This Study

### 8.1 Application to the Nereid system

Nereid, a satellite of Neptune, moves on the most eccentric orbit of known satellites in the Solar system. It was discovered by Kuiper in 1949. The semimajor axis of the orbit is about  $5.5 \times 10^6$  [Km] ( $\sim 220$  radii of Neptune). However, its orbital eccentricity reaches 0.75, Nereid approaches Neptune in  $1.4 \times 10^6$  [Km]. Before investigation by the spacecraft Voyager II, only two satellites were known to be orbiting Neptune. Triton, the first one discovered, by Lassell in 1846, orbits in a nearly circular but retrograde orbit once every six days.

After the discovery of Nereid, many astronomers reported its osculating elements, such as van Biesbroeck (1951,1957), Rose(1974) and Veillet(1982). Due to its long orbital period (nearly 1 year), astrometric observations for a long span are required to obtain orbital elements accurately. Besides, its faintness (19th magnitude) has obstructed the acquisition of clear images of Nereid from ground-based telescopes. The accuracy of ground-based astrometrical observations is about 0.1[arcsec]. It corresponds to the length of about 2200[Km] at the mean distance of the Neptune.

Nereid's Kepler motion is disturbed incessantly by the Sun, Triton and other disturbers. The Sun plays an important role in the time-variation of orbital elements of Nereid.

Mignard (1975) first studied the motion of Nereid, and built an ephemeris (Mignard (1981)) analytically. He used canonical transformations and took only the solar perturbation ( $P_2$  and  $P_3$  terms) into account.

For the Voyager II mission's flight program, Jacobson(1990, 1991) constructed a precise numerical ephemeris of Neptunian satellites.

In an analytical approach, Oberti (1990) showed periodic and secular perturbation terms of Nereid using a canonical perturbation method of Deprit (1968) type. In this work, the solar perturbation ( $P_2$  to  $P_4$  terms) and that of Triton ( $P_2$  and  $P_3$  terms) were included in his Hamiltonian and set the origin at the Neptune-Triton barycenter. Segerman and Richardson (1997) also studied the motion of Nereid. They took the solar perturbation ( $P_2$  and  $P_3$  terms), that of Triton ( $P_2$  and  $P_3$  terms) and the  $J_2$  effect of Neptune into consideration.

Saad(2000) studied the motion of Nereid using a canonical perturbation method of Hori type. He considered only the solar perturbation ( $P_2$  term only). The secular perturbation he solved analytically, based on Kinoshita and Nakai(1999)'s work.

## 8.2 Application to newly discovered outer satellites

Gladman et.al. (1997) discovered two satellites of Uranus on CCD-images using the Hale 5-meter telescope at Mount Palomar. Some successive observations revealed that these satellites orbit far from Uranus in retrograde. Today, they are named Caliban (semimajor axis  $\sim 280$  radii of Uranus) and Sycorax ( $\sim 480$  radii of Uranus). After this discovery, satellites orbiting the outer region have been reported one after another for Jupiter, Saturn and Uranus.

These newly discovered satellites tend to have slightly more eccentric orbits (eccentricity  $\sim 0.5$ ). For Jovian and Saturnian outer satellites, perturbation by the Sun is predominant. However, satellites of the farther revolving planet from the Sun have the same spatial configuration (the planet + the inner orbiting satellite + the outer orbiting satellite) as those in this study.

## 8.3 Application to extrasolar planetary systems

In the mid-1990's, Mayor and Queloz (1995) reported the presence of a Jupiter-mass companion to 51 Pegasi through observations of its radial velocity. On the study of extrasolar planets, some review papers have already been published (for example, Marcy and Butler(1998) or Marcy et.al.(2000)). Today's findings are mainly due to radial velocity observation. Slightly more massive planets than Jupiter are orbiting around the sun-like stars. Some extrasolar giant planets have been detected in the neighborhood of the primary stars. They are called "hot Jupiters."

This property supports the expectation that a more highly eccentric planet revolves further from the known extrasolar planets. Our theory can be applied to such a system.

Our theory can also be applied to the motion of a planet around a binary pair, when we take the higher order terms of  $(\frac{m'}{M+m'})$  into consideration.

# 9 Conclusions

We have developed an analytical theory on the motion of a celestial body orbiting in a highly eccentric orbit. The body is perturbed by an inner celestial body which revolves in a circular

orbit around the main star. Our theory is constructed using Hori's canonical perturbation method without expanding the Hamiltonian in eccentricity. In order to check the accuracy of our theory, we compared the analytical results with numerically integrated ones.

We ascertained that our theory provides the orbital elements with high accuracy. The semi-major axis ratio of the outer body to the inner one is larger, and our theory maintains its high degree of accuracy in the case of a higher eccentricity.

We also found that both results diverge in the case of a very large eccentricity. This is mainly due to the following: (1) The series of the generating function  $S_i$  for the short periodic perturbations becomes less convergent, or more diverse. (2) The offset phenomena: abrupt changes in the orbital elements take place when the outer body passes through its pericenter. They are not represented in the analytical theory.

We tried to apply our theory to the Neptunian satellite Nereid that orbits in a highly eccentric orbit ( $e = 0.75$ ) perturbed by the inner revolving satellite, Triton. Our theory maintains a good degree of accuracy, yielding results better than 30Km in the osculating semimajor axis of Nereid.

Our analytical theory can also be applied to other highly eccentric orbits. Some extrasolar planets are known as "Hot Jupiters", which revolve around their primary stars in circular orbits at small distances from them. Another new planet may exist in the outer field. The motion of such a new planet can be described by our theory.

## 10 References

- Brouwer, D. and Clemence, G.M., 1961, *Methods of Celestial Mechanics*, Academic Press.
- Deprit, A., 1968, Canonical Transformations Depending on a Small Parameter, *Cele. Mech. Dyn. Astron.*, **1**, 12.
- Gladman, B.J., Nicholson, P.D., Burns, J.A. and Kavelaars, J.J., 1997, reported in IAUC 6764 (Marsden, B.G. ed.).
- Hori, G., 1966, Theory of General Perturbations with Unspecified Canonical Variables, *Publ. Astron. Soc. Japan*, **18**, 287.
- Jacobson, R.A., 1990, The orbits of the satellites of Neptune, *Astron. Astrophys.*, **231**, 241.
- Jacobson, R.A., 1991, Triton and Nereid astrographic observations from Voyager 2, *Astron. Astrophys. Suppl. Ser.*, **90**, 541.
- Kinoshita, H. and Nakai, H., 1999, Analytical Solution of the Kozai Resonance and its Application, *Cele. Mech. Dyn. Astron.*, **75**, 125.
- Marcy, G.W. and Butler, R.P., 1998, Detection of Extrasolar Giant Planets, *Annu. Rev. Astron. Astrophys.*, **36**, 57.

Marcy, G.W., Cochran, W.D. and Mayor, M., 2000, Extrasolar Planets around Main-sequence Stars, in *Protostars and Planets IV* (Mannings, V., Boss, A.P. and Russell, S.S. eds.), The University of Arizona Press, 1285.

Martin, E.L., Brandner, W. and Basri, G., 1999, A Search for Companions to Nearby Brown Dwarfs: The Binary DENIS-P J1228.2-1547, *Science*, **283**, 1718.

Masaki, Y. and Kinoshita, H., 2001, Construction of Analytical Expressions of Nereid's Motion Perturbed by Triton (I) Planar Problem, in *Proceedings of the 33rd Symposium on Celestial Mechanics* (Kokubo, E., Ito, T. and Arakida, H. eds.), 189.

Mayor, M. and Queloz, D., 1995, A Jupiter-mass Companion to a Solar-type Star, *Nature*, **378**, 355.

Mignard, F., 1975, Satellite à forte excentricité. Application à Néréide, *Astron. Astrophys.*, **43**, 359.

Mignard, F., 1981, The Mean Elements of Nereid, *Astron. J.*, **86**, 1728.

Murison, M.A., 1989, On an Efficient and Accurate Method to Integrate Restricted Three-body Orbits, *Astron. J.*, **97**, 1496.

Oberti, P., 1990, An accurate solution for Nereid's motion I. Analytical modeling, *Astron. Astrophys.*, **239**, 381.

Oppenheimer, B.R., Kulkarni, S.R. and Stauffer, J.R., 2000, Brown Dwarfs, in *Protostars and Planets IV* (Mannings, V., Boss, A.P. and Russell, S.S. eds.), The University of Arizona Press, 1313.

Rose, L.E., 1974, Orbit of Nereid and the mass of Neptune, *Astron. J.*, **79**, 489.

Russell, S.S. and Boss, A.P., 1998, Protostars and Planets, *Science*, **281**, 932.

Saad, A.S., 2000, The Theory of Motion and Ephemerides of the Second Neptunian Satellite Nereid, Doctoral thesis, The Graduate University for Advanced Studies.

Segerman, A.M. and Richardson, D.L., 1997, An Analytical Theory for the Orbit of Nereid, *Cele. Mech. Dyn. Astron.*, **66**, 321.

van Biesbroeck, 1951, The Orbit of Nereid, Neptune's Second Satellite, *Astron. J.*, **56**, 110.

van Biesbroeck, 1957, The Mass of Neptune from a New Orbit of its Second Satellite Nereid, *Astron. J.*, **62**, 272.

Veillet, C., 1982, Orbital Elements of Nereid from New Observations, *Astron. Astrophys.*, **112**, 277.

## A Appendix — Analytical Expressions (The Inclined Problem)

We use the following descriptions in this paper.

$M$	Mass of the primary body
$m'$	Mass of the inner body
$m$	( $\equiv 0$ ) Mass of the outer body
$\mu$	$k^2(M + m' + m) \equiv n^2 a^3$
$a$	Semimajor axis of the outer body
$a'$	Semimajor axis of the inner body
$n$	Mean motion of the outer body
$n'$	Mean motion of the inner body
$e$	Eccentricity of the outer body
$\eta$	$\sqrt{1 - e^2}$
$I$	Inclination of the outer body
$\theta$	$\equiv \cos I$
$f$	True anomaly of the outer body
$r$	Radius of the outer body
$\lambda'$	Longitude of the inner body
$y_2$	( $\equiv \omega$ ) Argument of pericenter of the outer body
$y_3$	( $\equiv h - \lambda'$ )

$$\begin{aligned}
 C_2 &\equiv \frac{M m'}{(M + m')^2} \\
 C_3 &\equiv \frac{M m' (M^2 - m'^2)}{(M + m')^4} \\
 C_4 &\equiv \frac{M m' (M^3 + m'^3)}{(M + m')^5} \\
 C_5 &\equiv \frac{M m' (M^4 - m'^4)}{(M + m')^6} \\
 C_6 &\equiv \frac{M m' (M^5 + m'^5)}{(M + m')^7}
 \end{aligned}$$

Hereafter, we neglect  $O((\frac{m}{M+m'})^2)$  or higher order terms.

### A.1 Hamiltonians

Original Hamiltonian

$$F = n' x_3 + \frac{\mu^2}{2 x_1^2}$$

$$\begin{aligned}
& +\mu C_2 \frac{a'^2}{a^3} \left(\frac{a}{r}\right)^3 \left[ \frac{1}{8}(-1+3\theta^2) \right. \\
& \quad +\frac{3}{8}(1-\theta^2)\cos(2f+2y_2) \\
& \quad +\frac{3}{16}(1-\theta^2)^2\cos(2f+2y_2-2y_3) \\
& \quad +\frac{3}{8}(1-\theta^2)\cos(2y_3) \\
& \quad \left. +\frac{3}{16}(1+\theta)^2\cos(2f+2y_2+2y_3) \right] \\
& +\mu C_3 \frac{a'^3}{a^4} \left(\frac{a}{r}\right)^4 \left[ \frac{15}{64}(1-\theta-\theta^2+\theta^3)\cos(f+y_2-3y_3) \right. \\
& \quad +\frac{5}{64}(1-3\theta+3\theta^2-\theta^3)\cos(3f+3y_2-3y_3) \\
& \quad +\frac{15}{64}(1+\theta-\theta^2-\theta^3)\cos(f+y_2+3y_3) \\
& \quad +\frac{5}{64}(1+3\theta+3\theta^2+\theta^3)\cos(3f+3y_2+3y_3) \\
& \quad +\frac{3}{64}(-1+11\theta+5\theta^2-15\theta^3)\cos(f+y_2-y_3) \\
& \quad +\frac{15}{64}(1-\theta-\theta^2+\theta^3)\cos(3f+3y_2-y_3) \\
& \quad +\frac{3}{64}(-1-11\theta+5\theta^2+15\theta^3)\cos(f+y_2+y_3) \\
& \quad \left. +\frac{15}{64}(1+\theta-\theta^2-\theta^3)\cos(3f+3y_2+y_3) \right] \\
& +\dots
\end{aligned}$$

**Hamiltonians  $F^*$ ,  $F^{**}$  and  $F^{***}$**

$$\begin{aligned}
F^* &= n'x_3 \\
& \quad +\frac{\mu^2}{2x_1^2} \\
& \quad +\frac{1}{8}\mu C_2 \frac{a'^2}{r^3} [(-1+3\theta^2)+3(1-\theta^2)\cos(2f+2y_2)] \\
& \quad +\frac{3}{512}\mu C_4 \frac{a'^4}{r^5} [3(3-30\theta^2+35\theta^4) \\
& \quad \quad -20(1-8\theta^2+7\theta^4)\cos(2f+2y_2) \\
& \quad \quad +35(1-\theta^2)^2\cos(4f+4y_2)] \\
& \quad +\frac{5}{8192}\mu C_6 \frac{a'^6}{r^7} [10(-5+105\theta^2-315\theta^4+231\theta^6) \\
& \quad \quad -105(-1+19\theta^2-51\theta^4+33\theta^6)\cos(2f+2y_2) \\
& \quad \quad +126(-1+\theta^2)^2(-1+11\theta^2)\cos(4f+4y_2) \\
& \quad \quad -231(-1+\theta^2)^3\cos(6f+6y_2)] \\
& \quad +\dots
\end{aligned}$$

$$\begin{aligned}
F^{**} &= n'x_3 \\
& \quad +\frac{\mu^2}{2x_1^2} \\
& \quad +\frac{1}{8}\mu C_2 \frac{a'^2}{a^3} \frac{1}{\eta^3} (-1+3\theta^2) \\
& \quad +\frac{9}{1024}\mu C_4 \frac{a'^4}{a^5} \frac{1}{\eta^7} [(3-30\theta^2+35\theta^4)(2+3e^2)-10(1-8\theta^2+7\theta^4)e^2\cos(2y_2)]
\end{aligned}$$

$$\begin{aligned}
& + \frac{25}{65536} \mu C_6 \frac{a'^6}{a^7} \frac{1}{\eta^{11}} [ + 2(-5 + 105\theta^2 - 315\theta^4 + 231\theta^6)(8 + 40e^2 + 15e^4) \\
& \quad - 210(-1 + 19\theta^2 - 51\theta^4 + 33\theta^6)e^2(2 + e^2) \cos(2y_2) \\
& \quad + 63(-1 + \theta^2)^2(-1 + 11\theta^2)e^4 \cos(4y_2)] \\
& + \dots \\
F^{****} & = n' x_3 \\
& + \frac{\mu^2}{2x_1^2} \\
& + \frac{1}{8} \mu C_2 \frac{a'^2}{a^3} \frac{1}{\eta^3} (-1 + 3\theta^2) \\
& + \frac{9}{1024} \mu C_4 \frac{a'^4}{a^5} \frac{1}{\eta^7} (3 - 30\theta^2 + 35\theta^4)(2 + 3e^2) \\
& + \frac{25}{32768} \mu C_6 \frac{a'^6}{a^7} \frac{1}{\eta^{11}} (-5 + 105\theta^2 - 315\theta^4 + 231\theta^6)(8 + 40e^2 + 15e^4) \\
& + \dots
\end{aligned}$$

## A.2 $P_2$ -limited generating functions

For short periodic terms

$$\begin{aligned}
S_1 & = 0 \\
S_2 & = -\frac{3}{16} \mu C_2 \frac{1}{n'} \frac{a'^2}{r^3} \\
& \quad \left[ [2\theta \sin(2f + 2y_2)] \cos(2y_3) \right. \\
& \quad \quad + [(1 - \theta^2) \\
& \quad \quad \quad + (1 + \theta^2) \cos(2f + 2y_2)] \sin(2y_3) \left. \right] \\
S_3 & = -\frac{3}{64} \mu C_2 \frac{na}{\eta n'^2} \frac{a'^2}{r^4} \\
& \quad \left[ [-6(-1 + \theta^2)e \sin(f) \right. \\
& \quad \quad + (1 + \theta^2)\{-e \sin(f + 2y_2) \\
& \quad \quad \quad + 4 \sin(2f + 2y_2) \\
& \quad \quad \quad + 5e \sin(3f + 2y_2)\}] \cos(2y_3) \\
& \quad \quad + [2\theta\{-e \cos(f + 2y_2) \\
& \quad \quad \quad + 4 \cos(2f + 2y_2) \\
& \quad \quad \quad + 5e \cos(3f + 2y_2)\}] \sin(2y_3) \left. \right] \\
S_4 & = -\frac{3}{256} \mu C_2 \frac{(na)^2}{\eta^2 n'^3} \frac{a'^2}{r^5} \\
& \quad \left[ [2\theta\{3e^2 \sin(2y_2) \right. \\
& \quad \quad - 10e \sin(f + 2y_2) \\
& \quad \quad \quad + 2(8 - 5e^2) \sin(2f + 2y_2) \\
& \quad \quad \quad + 54e \sin(3f + 2y_2) \\
& \quad \quad \quad + 35e^2 \sin(4f + 2y_2)\}] \cos(2y_3) \\
& \quad \quad + [-18e^2(1 - \theta^2) \\
& \quad \quad \quad - 12(-1 + \theta^2)e \cos(f) \left. \right]
\end{aligned}$$

$$\begin{aligned}
& -30(-1 + \theta^2)e^2 \cos(2f) \\
& +3(1 + \theta^2)e^2 \cos(2y_2) \\
& -10(1 + \theta^2)e \cos(f + 2y_2) \\
& -2(1 + \theta^2)(-8 + 5e^2) \cos(2f + 2y_2) \\
& +54(1 + \theta^2)e \cos(3f + 2y_2) \\
& +35(1 + \theta^2)e^2 \cos(4f + 2y_2)] \sin(2y_3) \Big] \\
S_5 = & -\frac{3}{1024} \mu C_2 \frac{(na)^3 a'^2}{\eta^3 n'^4 r^6} \\
& \Big[ [(-1 + \theta^2)\{6e(-4 + 45e^2) \sin(f) \\
& -192e^2 \sin(2f) \\
& -210e^3 \sin(3f)\} \\
& + (1 + \theta^2)\{15e^3 \sin(f - 2y_2) \\
& +40e^2 \sin(2y_2) \\
& +e(-68 + 45e^2) \sin(f + 2y_2) \\
& -16(-4 + 13e^2) \sin(2f + 2y_2) \\
& -e(-436 + 105e^2) \sin(3f + 2y_2) \\
& +712e^2 \sin(4f + 2y_2) \\
& +315e^3 \sin(5f + 2y_2)\} \} \cos(2y_3) \\
& + [2\theta\{-15e^3 \cos(f - 2y_2) \\
& +40e^2 \cos(2y_2) \\
& +e(-68 + 45e^2) \cos(f + 2y_2) \\
& -16(-4 + 13e^2) \cos(2f + 2y_2) \\
& -e(-436 + 105e^2) \cos(3f + 2y_2) \\
& +712e^2 \cos(4f + 2y_2) \\
& +315e^3 \cos(5f + 2y_2)\} \} \sin(2y_3) \Big]
\end{aligned}$$

**For intermediate periodic terms**

$$\begin{aligned}
S_1^* &= \frac{1}{16} \mu C_2 \frac{a'^2}{a^3} \frac{1}{\eta^3 n} [2(-1 + 3\theta^2)\{(f - l) + e \sin f\} \\
& + (1 - \theta^2)\{3e \sin(f + 2y_2) + 3 \sin(2f + 2y_2) + e \sin(3f + 2y_2)\} \\
& - (1 - \theta^2) \frac{1}{e^2} \{2 - 3e^2 - 2\eta(1 - e^2)\} \sin(2y_2)] \\
S_2^* &= O\left(\left(\frac{m}{M + m'}\right)^2\right)
\end{aligned}$$

**For long periodic terms**

$$S^{**} = 0$$

### A.3 $P_3$ -limited generating functions

**For short periodic terms**

$$\begin{aligned}
S_1 &= 0 \\
S_2 &= -\frac{5}{96} \mu C_3 \frac{1}{n'} \frac{a'^3}{r^4}
\end{aligned}$$



$$\begin{aligned}
& \left[ [\theta \{-3(-1 + \theta^2) \sin(f + y_2) \right. \\
& \quad + (3 + \theta^2) \sin(3f + 3y_2)\}] \cos(3y_3) \\
& \quad + [-3(-1 + \theta^2) \cos(f + y_2) \\
& \quad \left. + (1 + 3\theta^2) \cos(3f + 3y_2)\}] \sin(3y_3) \right] \\
& - \frac{3}{32} \mu C_3 \frac{1}{n'} \frac{a'^3}{r^4} \\
& \left[ [\theta \{(-11 + 15\theta^2) \sin(f + y_2) \right. \\
& \quad - 5(-1 + \theta^2) \sin(3f + 3y_2)\}] \cos(y_3) \\
& \quad + [(-1 + 5\theta^2) \cos(f + y_2) \\
& \quad \left. - 5(-1 + \theta^2) \cos(3f + 3y_2)\}] \sin(y_3) \right] \\
S_3 = & - \frac{5}{576} \mu C_3 \frac{na}{\eta n'^2} \frac{a'^3}{r^5} \\
& \left[ [3(-1 + \theta^2) \{3e \sin(y_2) \right. \\
& \quad - 2 \sin(f + y_2) \\
& \quad - 5e \sin(2f + y_2)\} \\
& \quad + (1 + 3\theta^2) \{-e \sin(2f + 3y_2) \\
& \quad + 6 \sin(3f + 3y_2) \\
& \quad + 7e \sin(4f + 3y_2)\}] \cos(3y_3) \\
& \quad + [3\theta(-1 + \theta^2) \{3e \cos(y_2) \\
& \quad - 2 \cos(f + y_2) \\
& \quad - 5e \cos(2f + y_2)\} \\
& \quad + \theta(3 + \theta^2) \{-e \cos(2f + 3y_2) \\
& \quad + 6 \cos(3f + 3y_2) \\
& \quad \left. + 7e \cos(4f + 3y_2)\}] \sin(3y_3) \right] \\
& - \frac{3}{64} \mu C_3 \frac{na}{\eta n'^2} \frac{a'^3}{r^5} \\
& \left[ [(-1 + 5\theta^2) \{-3e \sin(y_2) \right. \\
& \quad + 2 \sin(f + y_2) \\
& \quad + 5e \sin(2f + y_2)\} \\
& \quad - 5(-1 + \theta^2) \{-e \sin(2f + 3y_2) \\
& \quad + 6 \sin(3f + 3y_2) \\
& \quad + 7e \sin(4f + 3y_2)\}] \cos(y_3) \\
& \quad + [\theta(-11 + 15\theta^2) \{-3e \cos(y_2) \\
& \quad + 2 \cos(f + y_2) \\
& \quad + 5e \cos(2f + y_2)\} \\
& \quad - 5\theta(-1 + \theta^2) \{-e \cos(2f + 3y_2) \\
& \quad + 6 \cos(3f + 3y_2) \\
& \quad \left. + 7e \cos(4f + 3y_2)\}] \sin(y_3) \right] \\
S_4 = & - \frac{5}{3456} \mu C_3 \frac{(na)^2}{\eta^2 n'^3} \frac{a'^3}{r^6} \\
& \left[ [3\theta(-1 + \theta^2) \{15e^2 \sin(f - y_2) \right.
\end{aligned}$$

$$\begin{aligned}
& +8e \sin(y_2) \\
& +2(-2+15e^2) \sin(f+y_2) \\
& -32e \sin(2f+y_2) \\
& -35e^2 \sin(3f+y_2)\} \\
& +\theta(3+\theta^2)\{3e^2 \sin(f+3y_2) \\
& -16e \sin(2f+3y_2) \\
& -2(-18+7e^2) \sin(3f+3y_2) \\
& +104e \sin(4f+3y_2) \\
& +63e^2 \sin(5f+3y_2)\} \cos(3y_3) \\
& +[3(-1+\theta^2)\{-15e^2 \cos(f-y_2) \\
& +8e \cos(y_2) \\
& +2(-2+15e^2) \cos(f+y_2) \\
& -32e \cos(2f+y_2) \\
& -35e^2 \cos(3f+y_2)\} \\
& +(1+3\theta^2)\{3e^2 \cos(f+3y_2) \\
& -16e \cos(2f+3y_2) \\
& -2(-18+7e^2) \cos(3f+3y_2) \\
& +104e \cos(4f+3y_2) \\
& +63e^2 \cos(5f+3y_2)\} \sin(3y_3)] \\
& -\frac{3}{128} \mu C_3 \frac{(na)^2}{\eta^2 n'^3} \frac{a'^3}{r^6} \\
& \left[ [\theta(-11+15\theta^2)\{-15e^2 \sin(f-y_2) \right. \\
& \quad -8e \sin(y_2) \\
& \quad -2(-2+15e^2) \sin(f+y_2) \\
& \quad +32e \sin(2f+y_2) \\
& \quad +35e^2 \sin(3f+y_2)\} \\
& \quad +5\theta(-1+\theta^2)\{-3e^2 \sin(f+3y_2) \\
& \quad +16e \sin(2f+3y_2) \\
& \quad +2(-18+7e^2) \sin(3f+3y_2) \\
& \quad -104e \sin(4f+3y_2) \\
& \quad -63e^2 \sin(5f+3y_2)\} \cos(y_3) \\
& \quad +[(-1+5\theta^2)\{15e^2 \cos(f-y_2) \\
& \quad -8e \cos(y_2) \\
& \quad -2(-2+15e^2) \cos(f+y_2) \\
& \quad +32e \cos(2f+y_2) \\
& \quad +35e^2 \cos(3f+y_2)\} \\
& \quad +5(-1+\theta^2)\{-3e^2 \cos(f+3y_2) \\
& \quad +16e \cos(2f+3y_2) \\
& \quad +2(-18+7e^2) \cos(3f+3y_2) \\
& \quad -104e \cos(4f+3y_2) \\
& \quad \left. \left. -63e^2 \cos(5f+3y_2)\} \sin(y_3) \right] \right]
\end{aligned}$$

$$S_5 = -\frac{5}{20736} \mu C_3 \frac{(na)^3}{\eta^3 n'^4} \frac{a'^3}{r^7}$$

$$\begin{aligned}
& \left[ 3(-1 + \theta^2) \{ -105e^3 \sin(2f - y_2) \right. \\
& \quad + 18e^2 \sin(f - y_2) \\
& \quad - 5e(-4 + 45e^2) \sin(y_2) \\
& \quad + 4(-2 + 59e^2) \sin(f + y_2) \\
& \quad + 3e(-52 + 105e^2) \sin(2f + y_2) \\
& \quad - 466e^2 \sin(3f + y_2) \\
& \quad - 315e^3 \sin(4f + y_2) \} \\
& \quad + (1 + 3\theta^2) \{ -15e^3 \sin(3y_2) \\
& \quad + 70e^2 \sin(f + 3y_2) \\
& \quad + e(-172 + 63e^2) \sin(2f + 3y_2) \\
& \quad - 12(-18 + 35e^2) \sin(3f + 3y_2) \\
& \quad - e(-1156 + 189e^2) \sin(4f + 3y_2) \\
& \quad + 1670e^2 \sin(5f + 3y_2) \\
& \quad + 693e^3 \sin(6f + 3y_2) \} \} \cos(3y_3) \\
& \quad + [3\theta(-1 + \theta^2) \{ 105e^3 \cos(2f - y_2) \\
& \quad - 18e^2 \cos(f - y_2) \\
& \quad - 5e(-4 + 45e^2) \cos(y_2) \\
& \quad + 4(-2 + 59e^2) \cos(f + y_2) \\
& \quad + 3e(-52 + 105e^2) \cos(2f + y_2) \\
& \quad - 466e^2 \cos(3f + y_2) \\
& \quad - 315e^3 \cos(4f + y_2) \} \\
& \quad + \theta(3 + \theta^2) \{ -15e^3 \cos(3y_2) \\
& \quad + 70e^2 \cos(f + 3y_2) \\
& \quad + e(-172 + 63e^2) \cos(2f + 3y_2) \\
& \quad - 12(-18 + 35e^2) \cos(3f + 3y_2) \\
& \quad - e(-1156 + 189e^2) \cos(4f + 3y_2) \\
& \quad + 1670e^2 \cos(5f + 3y_2) \\
& \quad + 693e^3 \cos(6f + 3y_2) \} \} \sin(3y_3) \Big] \\
& - \frac{3}{256} \mu C_3 \frac{(na)^3 a'^3}{\eta^3 n'^4 r^7} \\
& \left[ [ -(-1 + 5\theta^2) \{ -105e^3 \sin(2f - y_2) \right. \\
& \quad + 18e^2 \sin(f - y_2) \\
& \quad - 5e(-4 + 45e^2) \sin(y_2) \\
& \quad + 4(-2 + 59e^2) \sin(f + y_2) \\
& \quad + 3e(-52 + 105e^2) \sin(2f + y_2) \\
& \quad - 466e^2 \sin(3f + y_2) \\
& \quad - 315e^3 \sin(4f + y_2) \} \\
& \quad - 5(-1 + \theta^2) \{ -15e^3 \sin(3y_2) \\
& \quad + 70e^2 \sin(f + 3y_2) \\
& \quad + e(-172 + 63e^2) \sin(2f + 3y_2) \\
& \quad - 12(-18 + 35e^2) \sin(3f + 3y_2) \\
& \quad - e(-1156 + 189e^2) \sin(4f + 3y_2) \\
& \quad + 1670e^2 \sin(5f + 3y_2)
\end{aligned}$$

$$\begin{aligned}
& +693e^3 \sin(6f + 3y_2)\} \cos(y_3) \\
& +[-\theta(-11 + 15\theta^2)\{105e^3 \cos(2f - y_2) \\
& -18e^2 \cos(f - y_2) \\
& -5e(-4 + 45e^2) \cos(y_2) \\
& +4(-2 + 59e^2) \cos(f + y_2) \\
& +3e(-52 + 105e^2) \cos(2f + y_2) \\
& -466e^2 \cos(3f + y_2) \\
& -315e^3 \cos(4f + y_2)\} \\
& -5\theta(-1 + \theta^2)\{-15e^3 \cos(3y_2) \\
& +70e^2 \cos(f + 3y_2) \\
& +e(-172 + 63e^2) \cos(2f + 3y_2) \\
& -12(-18 + 35e^2) \cos(3f + 3y_2) \\
& -e(-1156 + 189e^2) \cos(4f + 3y_2) \\
& +1670e^2 \cos(5f + 3y_2) \\
& +693e^3 \cos(6f + 3y_2)\} \sin(y_3) \Big]
\end{aligned}$$

**For intermediate periodic terms**

$$S^* = 0$$

**For long periodic terms**

$$S^{**} = 0$$

# Size and Spatial Distributions of Sub-km Main-Belt Asteroids

Yoshida, Fumi & Nakamura, Tsuko  
National Astronomical Observatory of Japan,  
2-21-1, Osawa, Mitaka, Tokyo, 181-8588, Japan  
Tel 81-422-34-3627, Fax 81-422-34-3627  
E-mail [yoshdafm@cc.nao.ac.jp](mailto:yoshdafm@cc.nao.ac.jp)

## Abstract

This paper is the results of the first systematic investigation of very small Main-belt Asteroids (sub-km MBAs) using the Subaru Prime-Focus Camera (Suprime-Cam), which has an  $8K \times 10K$  mosaic CCD array on the 8.2m Subaru telescope atop Mauna Kea, Hawaii. We call this survey SMBAS (Sub-km Main-Belt Asteroid Survey). Observations were carried out on February 22 and 25, 2001 (HST) and the  $\sim 3.0\text{deg}^2$  sky area near opposition and near the ecliptic was searched. We detected 1111 moving objects down to  $R \sim 26$  mag (including very slow Trans Neptunian Objects). In this survey, we could not determine the exact orbit of each moving object, because of its short observational arc, which is only 2 hours. Instead we estimated statistically the semi-major axis ( $a$ ) and inclination ( $I$ ) of each moving object from its apparent sky motion vector, and then obtained the size and spatial distributions of sub-km MBAs. The main results of SMBAS are summarized as follows: (1) the sky number density of MBAs is found to be  $\sim 290$  per  $\text{deg}^2$  down to  $R \sim 24.4$  mag (for MBAs) near opposition and near the ecliptic. (2) the slope of the cumulative size distribution for sub-km MBAs ranging from 0.5 km to 1 km in diameter is fairly shallower ( $\sim 1.2$ ) than that for large MBAs obtained from the past asteroid surveys ( $\sim 1.8$ ). This means that the number of sub-km MBAs is much more depleted than a result extrapolated from the size distributions for large asteroids. (3) the depletion of sub-km MBAs is clearer in the outer main-belt than in the inner main-belt. (4) the spatial distribution of the smaller asteroids indicates a wider  $I$ -distribution near the mean motion resonances ( $2.8 \sim 3.1\text{AU}$ ) in the outer zone of the main-belt.

## 1. Introduction

The current size, spatial and compositional distributions of the Main-Belt Asteroids (MBAs) have been believed to reflect a long-term history of collisional evolution (e.g.,

Wetherill, 1989). Good knowledge of the Cumulative Size Distribution (hereafter CSD) of asteroids in the main-belt brings an insight into collisions between MBAs, the production rate of Near-Earth Asteroids (NEAs) and meteoroids, the cratering rate on the surfaces of the inner planets, the impact strengths of asteroids and so on. It may also allow us to infer the accretion process in the main-belt region in the initial stage of our solar system and the original mass of the main-belt (e.g. Kuiper *et al*, 1958, Anders, 1965, Jedicke &

---

\* Based on data collected at Subaru Telescope, which is operated by the National Astronomical Observatory of Japan.

\* Send offprint requests to [yoshdafm@cc.nao.ac.jp](mailto:yoshdafm@cc.nao.ac.jp)

MetCalfel, 1998).

From such motivations, some systematic investigation of MBAs as summarized in Table I have so far been done and the CSDs of MBAs have been revealed down to a few km in diameter ( $D$ ). However, we emphasize here the importance of sub-km MBAs whose sizes are a few hundred meters in  $D$  from the two view points as follows; 1) the majority (about 70~80%) of NEAs are sub-km-sized (<http://cfa-www.harvard.edu/iau/mpc.html>), and are widely supposed to originate from sub-km MBAs and 2) this size region lies near the border-line size separating two typical catastrophic impact mechanisms, namely those in the strength regime and the gravity regime (e.g., Melosh & Ryan, 1997, Durda *et al.*, 1998). Concerning the first point, it is generally accepted that NEAs originated from MBAs through collision processes between asteroids in the main-belt and the subsequent gravitational perturbations associated with the Kirkwood gaps (e.g., Wisdom, 1983, Morbidelli & Moons, 1995). However, this dynamical conjecture has never been confirmed observationally, because of the faintness of sub-km MBAs. Hence, in this respect our SMBAS may shed direct light on physical relations between NEAs and sub-km MBAs.

And if the above second point is correct, there may be difference between the CSD slope for sub-km MBAs and that for known large MBAs, and it may be able to be interpreted as the difference in the collisional nature between the strength regime and the gravity regime. For those reasons, we consider that the observational study of sub-km MBAs is very crucial in the collisional history of the main-belt and performed the following observations.

In this paper Section 2 deals with SMBAS observations and data reduction and Section 3 explains our detection technique of moving objects. We described in Section 4 positional and photometric measurements of asteroids, including the determination of the detection limiting magnitude. In Section 5 and 6, the method of statistical estimations of the semi-major axis and inclination for each asteroid and a observational bias correction method are treated, respectively. Section 7 mentions the main results derived from our SMBAS, namely the size and the spatial distributions for sub-km MBAs. And finally in Section 8 and 9, we discuss physical implication for our obtained results and future prospect.

**TABLE I**  
**The Previous Asteroid Surveys and SMBAS**

Survey	Observation time	Telescope size(m)	Limiting magnitude	Sky coverage (deg <sup>2</sup> )	Number of detected asteroids	Slope of the CSD for detected asteroids
YMS	1950-52	0.25	$V_p < 14-16$	14,400	1,550	2.4 (for $D = 30\sim 300\text{km}$ )
PLS	1960	1.25	$V_p < 20-21$	216	$> 2,000$	1.8 (for $D > 5\text{km}$ )
Spacewatch	1992-95	0.90	$V < 21$	3,740	59,226	1.8 (for $D > 5\text{km}$ )
SDSS	1998-2000	2.5	$r^* < 21.5$	500	$\sim 13,000$	1.3 (for $D = 1\sim 5\text{km}$ )
SMBAS	2001	8.2	$R \sim 24.4$	3.26	1,111	1.2 (for $D = 0.5\sim 1\text{km}$ )

$V_p$  : photographic magnitude.

$r^*$  : R-band (the effective wavelengths is 6280 Å) in SDSS.

## 2. Observations and Data reductions

### 2.1 Observations

Observations were carried out on February 22 and 25, 2001 (HST) by using the 8.2m Subaru telescope atop Mauna Kea, Hawaii. We used the 8K×10K wide-field mosaic camera acronymic as Suprime-Cam (Subaru Prime-Focus Camera) (Korniyama *et al.*, 2000), attached at the prime focus of the telescope. Suprime-Cam covers the field of view of  $\sim 34' \times \sim 27'$  at the prime focus (F/2.0) and it consists of ten CCD chips ( $2048 \times 4096$  pixels for each chip; the pixel size is  $0.2''$ ). However, since one CCD chip did not work in our observations, we actually used nine CCD chips. The field of view was  $\sim 0.22 \text{ deg}^2$  with nine CCDs. The searched sky includes the ecliptic area near opposition at RA.  $= 10^{\text{h}} 22^{\text{m}}$ , DEC.  $= +10^\circ 20'$ , which was within an ecliptic latitude  $\pm 1^\circ$ . The seven sky fields were selected carefully so that they are relatively star-free and did not include bright stars. The *R*-band filter used, which is most efficient in terms of both quantum efficiency of CCD and the peak intensity of solar spectra. Each exposure time was 7 min. The seeing size was 0.8–1.0 arcsec on Feb. 22 and 0.6–0.9 arcsec on Feb. 25. The same fields were taken in two nights. The total surveyed area during the two nights was  $2.97 \text{ deg}^2$ .

Two observational modes were performed : 1) Wide Field (WF) survey mode and 2) Deep Field (DF) survey mode. In WF survey, we took three images of the same field with a time interval of about 55 min. In DF survey, eleven images of the same field were taken every 11 min in succession. In both modes, however, the observational arc for each moving object is about two hours. We also observed six Landolt standard stars at different airmass for photometric calibrations (Landolt, 1992).

The observational data described here are actually the same as those of the Wide-Field Survey of Edgeworth-Kuiper Belt Objects that

have already been reported by Kinoshita *et al.* (2002). However, since the purpose of observations and the method of data analysis are quite different between Kinoshita's and this work. We distinguish our survey from that by Kinoshita *et al.* for Edgeworth-Kuiper Belt Objects (EKBOs) by calling this survey the Sub-km Main-belt Asteroids Survey (hereafter SMBAS). The observational mode and exposure time were optimized for detection of EKBOs.

### 2.2 Data reductions

Image reduction was carried out on a chip-by-chip basis using the standard method with NOAO IRAF. First, the averaged output value for the overscan region of each CCD was subtracted from each CCD image data. Second, the overscan region was trimmed and then the image consisting only of the effective area was made. Third, in order to correct the two-dimensional bias pattern of each CCD, the bias image was subtracted from each CCD image. The bias image was produced by averaging a few raw bias frames which were taken every night. Next, we made corrections of difference in pixel-sensitivity over a CCD-chip, namely a traditional flat-field calibration. For that purpose, we took several images of the twilight sky with the field-centers offset slightly from each other. Then, a median flat-field image was constructed from them, by which each CCD image was divided to get uniform sensitivity.

## 3. Detection of moving objects

There are two approaches to detect moving objects in observed images, that is, detection by visual inspection and that by computer software. We adopted here the former approach, whereas most large-scale survey programs conducted in the last decade for NEAs and EKBOs relied on the latter one (e.g., SDSS, LINEAR). However, both approaches have their own merits and

demerits. Software detection is believed to be objective, free from careless mistakes made by the human interface and appropriate for handling large amount of data. But a properly designed procedure of visual detection can also be as objective as the software approach.

On the other hand, there seems to be a tendency that use of only software detection gives the limiting magnitude of roughly at least 1.0 ~ 1.5 magnitude shallower than that for visual detection; this is reasonable because at critical signal levels or in blended images, even sophisticated detection algorithms can never surpass the overall judging ability of the human eye and brain. As a result, software detection generally needs some help of more or less visual confirmation. We therefore believe, as shown later, that our technique of visual detection gives reliable results compatible with software approach, especially for medium-sized data of less than a few thousand objects. Regarding this, it is worthwhile to cite the recent survey observation of EKBOs by Millis *et al.* (2002), who, after comparing software detection and, they report that their technique works well even with only two-exposure pseudo-colored images.

After the basic reduction mentioned in Section 2 were applied to all object frames, we made combined images to recognize and count moving objects easily. Concretely, for all object frames in the WF survey which consist of three exposures, we subtracted the first image from the second one, and added the resulting image to the third image. An example of such new images made by the above operation is shown in Fig.1a. One can see moving objects as trains of separated black-white-black dots. This technique is basically the same as that proposed in Yoshida *et al.* (2001). Then we counted the number of moving objects by careful eye-inspection.

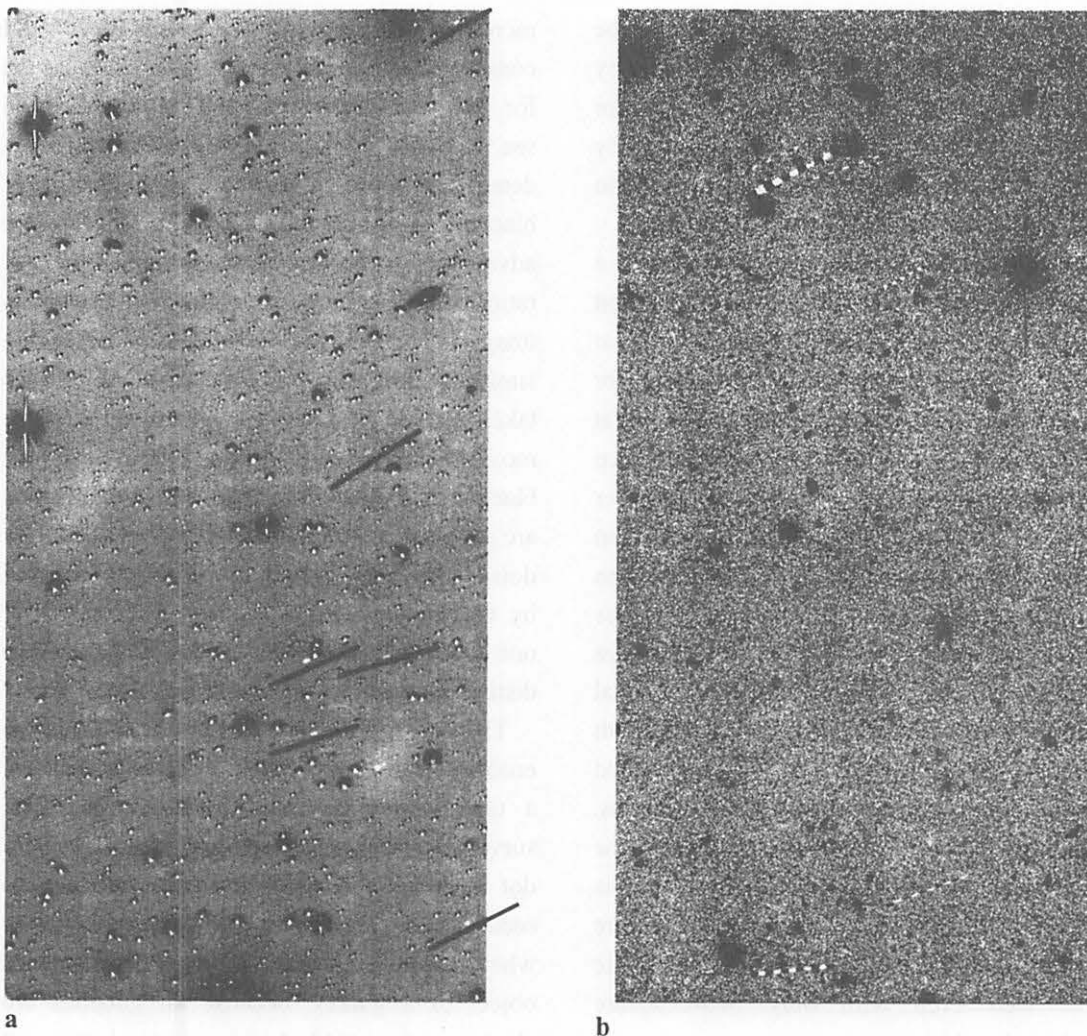
When multiple images were combined, the absolute value of the sky mean level for positive and negative images was equaled within a certain error level. This operation

increases the dispersion of the sky level compared to that for single exposure images (as for quantitative aspect of sky level fluctuation see Section 4.2). However, the easy detectability of asteroids shown as trains of black and white dots was much more advantageous than some degradation of S/N ratio in composite images. Fig.1b shows an image to which the above operations were similarly applied using a series of 11 images taken in the DF survey. One can recognize moving objects as a sequence of black-and-white striped bars. In Fig.1a, stars are generally seen as grouping black-and-white dots slightly shifted each other; this was caused by the telescopic pointing error during about one hour exposure interval. This also helps us distinguish stars from moving objects.

This black-and-white image technique also enables us to easily identify moving objects as a time sequence. For example, in the WF survey, we can surely confirm that the white dot corresponds to the image taken at the second exposure. It may also help us judge whether one elongated object is either a moving object or a galaxy, because all galaxies are always seen as black images, and we can confidently distinguish all white images as moving objects. Our this technique is also useful to confirm that a moving object is the same one on the neighboring CCD chips.

In actual detection of moving objects, we divided all of the processed images into partially overlapped small sub-frames the separated parts of  $100'' \times 100''$  (about  $500 \times 500$  pixel), so that its size can cover sufficiently the motion of asteroids in inner main-belt during two hours. And then we magnified all separated images and checked by careful eye-inspection twice separated by a few days. As a result we detected 1194 moving objects. Then, after removing the same moving objects that strode over neighboring CCDs, we eventually recognized 1111 moving objects.





**FIG.1**—The moving objects detected in each one of the CCD images in WF survey and DF surveys. (a) only some trains of black-and-white dots as identified asteroids are marked by lines for clarity. Fifteen moving objects are detected in  $2K \times 4K$  chip image altogether. Black-white-black dots appear fairly separated because of long exposure intervals ( $\sim 55$  min.). Field stars and galaxies appear as slightly shifted groupings of black-white-black dots, due to the telescope guiding error during the three exposures. (b) twenty-three detected moving objects are included in this image. They appear as black-and-white straight bars because of short exposure intervals ( $\sim 11$  min.). Field stars and galaxies appear as black images. Up is north and left is east in these images. All moving objects moved from left to right (due to retrograde motions near opposition).

#### 4. Photometry and Measurement of positions

##### 4.1 Photometry

We carried out aperture photometric measurements of detected moving objects using

IRAF-APPHOT. Then we added to the measured brightness of each moving object the following two corrections. First, we made the correction of difference in the sensitivity between CCD chips. The relative response for each CCD to the incoming radiation on the Suprime-Cam was calibrated comparing the mean count of sky background brightness

between CCDs. Second, the correction of the atmospheric extinction arising from the variation of airmass was made by using the extinction coefficients obtained from several Landolt photometric standard stars observed at some different airmasses on each night.

From measurements of the brightness of moving objects at each exposure time, we found that the mean amplitude of the intrinsic light variations of the moving objects (caused by their rotation) is  $\sim 0.25$  mag. This value is about ten times larger than the measuring photometric error ( $\sim 0.03$  mag) of each object. Therefore, the absolute magnitudes of all detected moving objects may include the error of  $\sim 0.25$  mag. However we emphasize that this kind of error can accordingly be averaged out when we construct the size or spatial distributions from many objects.

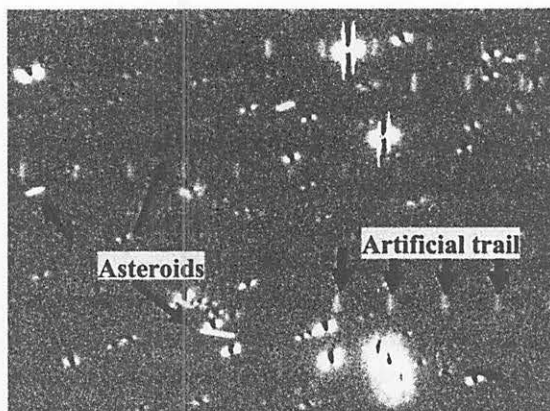
#### 4.2 Determination of the limiting magnitude

According to Kinoshita *et al.* (2002), the *R*-band limiting magnitude of point sources for our observing run is 26.1 mag. However, it is necessary for us to examine independently the limiting magnitude of moving objects by our own method, because our moving objects were much more trailed than EKBOs and we used a special technique, that is, the black and white image method. Since we combined multiple images to detect moving object, the mean sky fluctuation was increased to  $1.8\sim 2.0\sigma$  ( $\sigma$ : standard deviation of the sky brightness variation for a 1-exposure image) for 3-exposure composites and to  $2.9\sim 3.3\sigma$  for 11-exposure ones. Considering that the variation of the sky brightness follows the photon (namely, Poisson or Gaussian) statistics, the detection probability of a star with its peak intensity of  $1\sigma$  is calculated to be 68.3 %, and 95.5 % for  $2\sigma$ -peak, 99.7 % for  $3\sigma$ -peak, respectively (Meyer, 1975) for single exposure images. Hence, this can be interpreted as that the objects barely observed in 3-exposure composite images with  $1.8\sim 2.0\sigma$  have

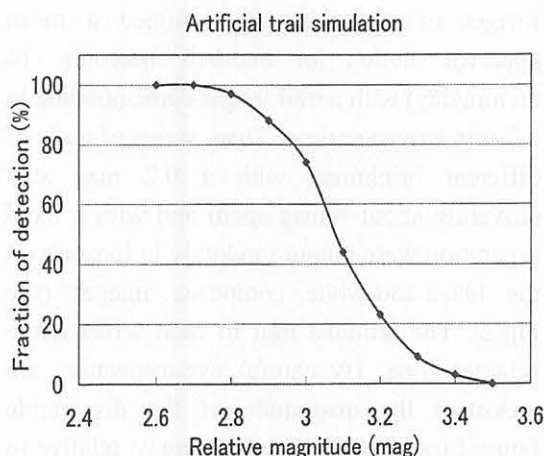
92.8~95.5 % detection probability, and 99.6~99.9 % probability for those in 11-exposure composites with  $2.9\sim 3.3\sigma$ . In other words, we may safely say that all of the detected asteroids in our composite images can be found in single exposure images with probabilities higher than 90%. This is a quantitative basis for our detecting moving objects in SMBAS.

The next step is to determine the limiting magnitude for trailed asteroids in black-and-white composite images that we made. For the purpose, we conducted simulation experiments using the task IRAF-MKOBJECTS. We first produced a stellar image with the FWHM of an average seeing-size for the observed night, then we made its slightly shifted image, added the two, and continue the process, to give a train of superimposed images; this is to mimic trailed images of asteroids. We assumed a mean apparent motion for mid-belt asteroids (14 arcmin/day) with a trail length corresponding to a 7-min exposure time. Then, series of trails of different brightness with a 0.2 mag step (covering about 4 mag span) and with a fixed separation were output randomly in location on the black-and-white composite images (see Fig.2). The leftmost trail in each series is the brightest one. By careful eye-inspection, we measured the magnitude of the discernible faintest trail with 0.1 mag accuracy, relative to the brightest one.

In practice, overlapping of some part of the trails with background stars and galaxies often occurred, so that we had to attempt many series of trails. Among them, we picked up 87 cases in which the brightest and faintest discernible trails could safely be measured, and plotted a percentage detection frequency of the faintest trails as a function of magnitude (in Fig.3). The origin of the magnitude in abscissa is arbitrary. By measuring magnitude differences between the brightest trail and some nearly photometric standard stars using asteroid trails as a mediator, we connected the abscissa to the standard



**FIG.2**—An example of simulated asteroid trails superimposed on a composite Suprime-Cam image. The magnitude difference between adjacent artificial trails is kept to be 0.2 mag. The length of trails is characterized by an average motion for mid-belt asteroids during 7 min exposure time and by a mean seeing size for the observed night.



**FIG.3**—Trail detectability as a function of magnitude. The ordinate shows what fraction of trailed asteroids can be visually detected at a given magnitude, relative to the brightest trail. The origin of abscissa is visually arbitrary.

magnitude system. One can see that the detection probability in Fig.3 changes from 100% to 0% over a magnitude range of 0.8–0.9. This is good agreement with Fig.4 in Millis *et al.* (2002).

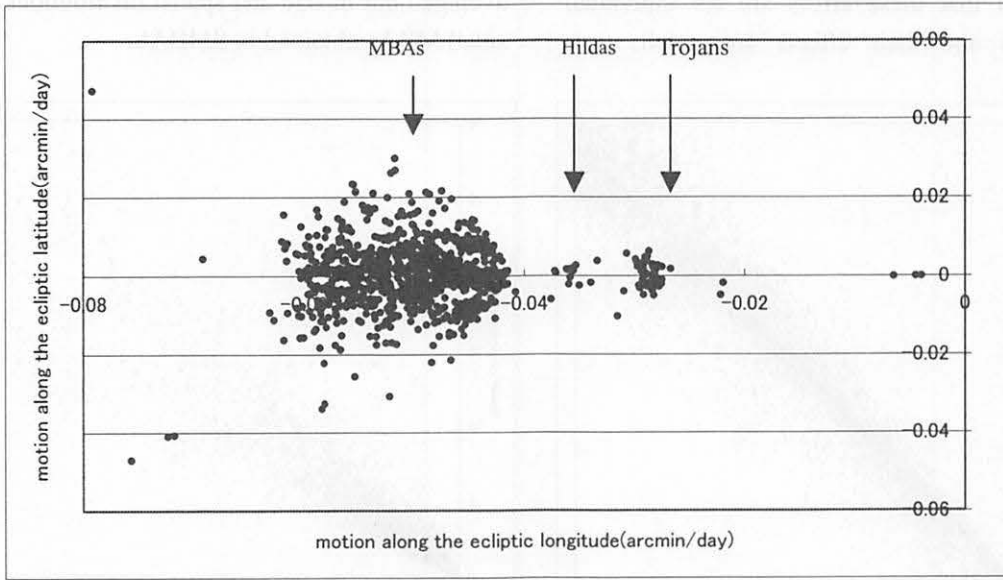
Here we adopted a 90 % perfect detection level in Fig.3 as the limiting magnitude in SMBAS. The magnitude corresponds to 24.4 mag in *R*-band. Note that this limiting magnitude is for mid-belt MBAs with their typical motions, but not for stars (namely point sources).

#### 4.3 Positional measurements of moving objects

The position for each moving object was measured again with IRAF-APPHOT relative to about ten USNO-A2 stars (<http://tdc-www.harvard.edu/software/catalogs/ua2.html>) that we picked up on the same frame. The apparent velocity for each object was calculated from its positions corresponding to all the exposure images. Fig.4 shows the apparent daily motions along the ecliptic longitude and the ecliptic latitude for each moving object detected in SMBAS. From Fig.4, we can easily distinguish between MBAs and the other groups of moving objects by their motions. We discuss only MBAs in the next section, because our interest focuses on small asteroids in the main-belt in this paper.

### 5. Estimates of semi-major axis and inclination for asteroids

Since the observational arc for each asteroid detected in SMBAS is only two hours, we can not determine its exact orbital elements. Thus instead we adopted a method to derive approximate semi-major axis ( $a$ ) and inclination ( $I$ ) from the sky motion vector of each asteroid under the assumption that its orbital eccentricity ( $e$ ) is zero. This method is based on geometrical and kinematical relations in the two-body problem, which was initially proposed by *Bowell et al.* in 1990. We call it *Bowell's method* in this paper. Since, however, the  $e$ -values of the typical MBAs lie actually in the range from 0 to  $\sim 0.2$ , we had to estimate in a statistical sense by Monte Carlo simulations the possible errors between the  $a$  and  $I$  obtained



**FIG.4---Apparent motions along the ecliptic longitude and the latitude of moving objects detected in SMBAS.** One can also see a considerable number of Hilda and Trojan asteroids.

by Bowell's method and true orbital elements for each asteroids. We hereafter denote the former as  $a'$  and  $I'$  and the latter as  $a$  and  $I$ . The following is the outline results from Nakamura & Yoshida (2001) and Yoshida & Nakamura (2001a).

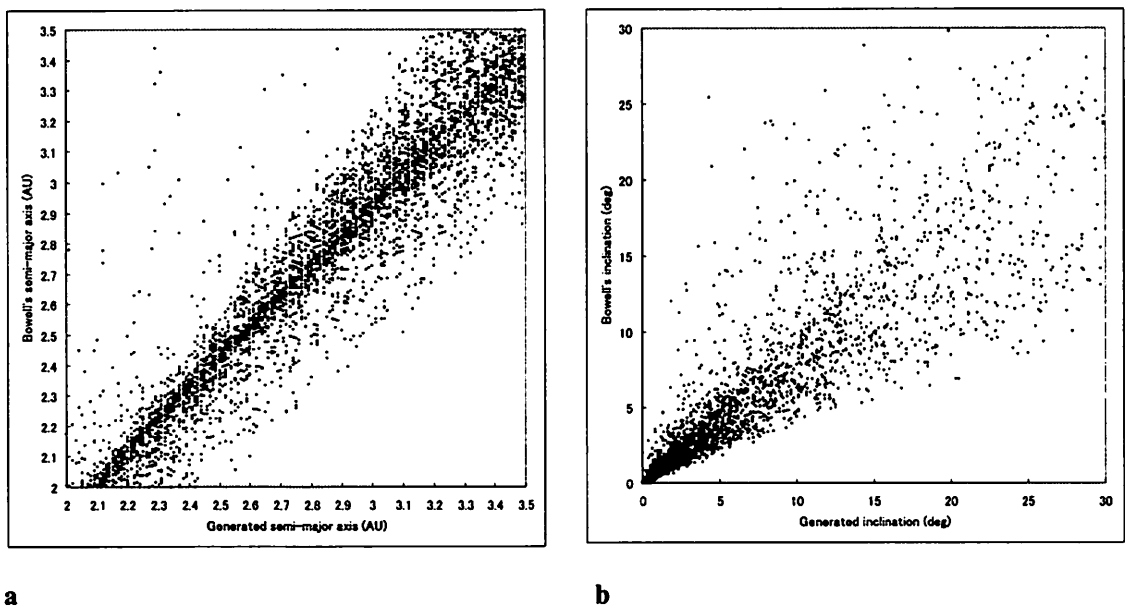
First, we generated many hypothetical asteroids with various orbital elements in a computer and picked up a few thousand asteroids that entered the observational window in SMBAS. We selected the ranges in orbital elements of the generated asteroids to be slightly wider than the ranges of known MBAs. The observational window area was set to be nearly the same as the actual observational window in SMBAS. The orbital elements ranges of the generated asteroids and the observational windows are listed in Table II. Next, we calculated their daily motions at opposition using a two-body ephemeris generator. We then compared the  $a$  and  $I$  for each hypothetical asteroid with the  $a'$  and  $I'$  calculated from its motion vector.

Fig.5a is a reproduction from Nakamura and Yoshida (2001) which shows the  $a$  and  $a'$  calculated by the above simulation. It seems that there is little systematic difference between

the  $a$  and the  $a'$ , though the scattering roughly attains to  $\sim 0.1$  AU. Fig.5b shows the  $I$  and  $I'$  calculated by the same simulation. The difference between the  $I$  and the  $I'$  is seem to be considerably large, especially for  $I > 10$  deg. We summarized more quantitative results in Table III a and Table III b.

After we calculated the  $a'$  and  $I'$  of asteroids detected in SMBAS using Bowell's method, we obtained systematic errors in their  $a'$  and  $I'$  based on the mean values of  $(a-a')$  and  $(I-I')$  shown in Table III a and Table III b. In order to give random errors in estimating errors of the  $a'$  and  $I'$ , we also calculated standard deviations, namely  $SD(a-a')$  and  $SD(I-I')$ , which were shown in the fourth column in Table III a and Table III b. From Table III a, we can see that random errors, namely  $SD(a-a')$  exceed systematic errors, namely  $mean(a-a')$ , for each zone, so that correction of the systematic error may bring about little important in  $a$ -estimate. On the other hand, Table III b shows that the  $mean(I-I')$  for the high-inclination MBA zone is larger than  $SD(I-I')$ , and hence the correction of the systematic error is essential, though random errors in  $I$  are fairly large for three zones. We

comment that these errors are for individual asteroids and their effects are much more averaged out in size and spatial distributions of small MBAs obtained in SMBAS.



**FIG.5 (a) and (b)—Generated  $a$  vs. Bowell's  $a'$  plot and Generated  $I$  vs. Bowell's  $I'$  plot, respectively, which is based on Nakamura and Yoshida (2001). The horizontal axis : the value of  $a$  or  $I$  generated in a computer run. The vertical axis: the value of  $a'$  or  $I'$  calculated from the sky motion vector of each asteroid generated in a computer run. The input parameters are from Table II .**

**TABLE II ---The Orbital-element Ranges of Simulation-generated Asteroids and Assumed Observational Window**

$a$ (AU)	$2.75 \pm 0.75$
$I$ (deg)	$15 \pm 15$
$e$	$0.2 \pm 0.2$
Angular elements	uniform over $0\sim 360^\circ$
Observational window	$2^\circ \times 2^\circ$

**TABLE III a---Errors of Semi-major Axis Obtained from Bowell’s Orbit**

Zone	range (AU)	mean ( $a-a'$ )	SD ( $a-a'$ )
Inner-belt	$2.0 < a < 2.6$	0.075	0.14
Mid-belt	$2.6 < a < 3.0$	0.070	0.13
Outer-belt	$3.0 < a < 3.5$	0.083	0.15

SD : the standard deviation

**TABLE III b---Errors of Inclination Obtained form Bowell’s Orbit**

Zone	range (deg)	mean ( $I-I'$ )	SD ( $I-I'$ )
Low-incl.	$0 < I < 10$	0.27	1.5
Medium-incl.	$10 < I < 20$	2.35	4.4
High-incl.	$20 < I < 30$	6.37	5.8

SD : the standard deviation

## 6. Observational bias corrections

SMBAS was conducted in a very small sky area only near opposition and near the ecliptic. For such observational conditions, we must consider some specific observational biases. Nakamura & Yoshida (2001) have already estimated observational biases for a small area near opposition and near the ecliptic (see Fig.6a, b). They calculated the observational biases as functions of  $a$  and  $I$  for an assumed observational field of view ( $5^\circ \times 4^\circ$ , a little wider than in actual SMBAS), centered at opposition and on the ecliptic. Fig.6a shows the *relative bias* as a function of  $a$ . The relative bias is defined here as the number ratio between near-ecliptic distant asteroids with, say,  $r \sim 6$  AU ( $r$ : heliocentric distance) and those with  $r=a$  (AU). Three relative bias curves are calculated for circular, near-circular, and elliptic orbits. Fig.6b shows the relative bias as a function of  $I$ . The relative bias is defined here to be the number ratio between ecliptic asteroids ( $I \sim 0$ ) and those with  $I$ . The relative bias curves were calculated for the inner- ( $a=2.3$  AU), middle- ( $a=2.7$  AU), and outer-MBAs ( $a=3.1$  AU), respectively.

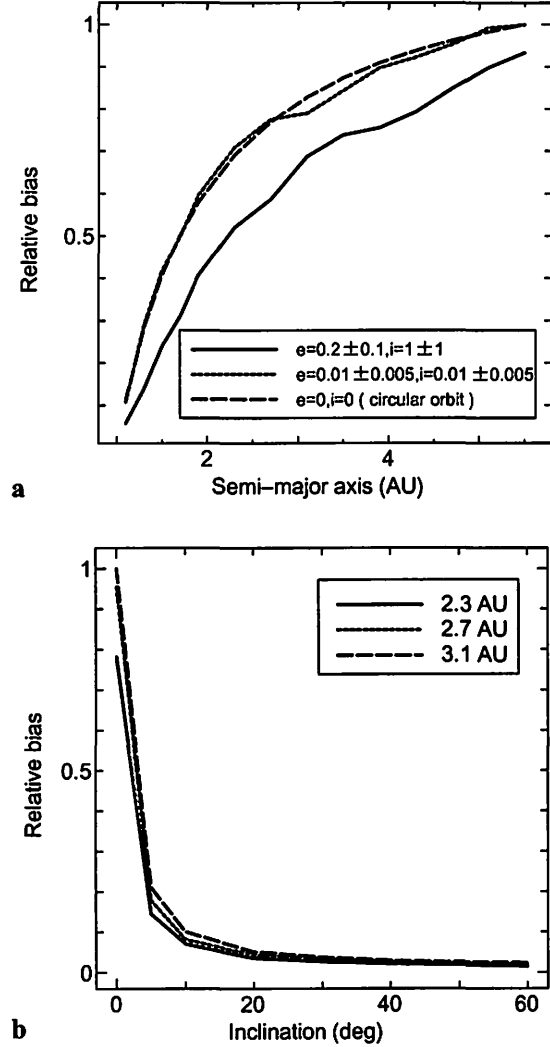
The  $a$ -distribution and  $I$ -distribution for MBAs detected in SMBAS, which we describe in Section 7, were corrected by using the bias corrections shown in Fig.6a and b. Since we are mainly interested in profiles of size and spatial distributions, the relative bias corrections are sufficient for our purpose and absolute bias corrections are unnecessary.

## 7. Results

### 7.1 Overview of results

In this paper, we defined the main-belt zone as  $a=2\sim 3.5$  AU in accordance with the past surveys. After the corrections of the  $a$  and  $I$  described in Section 5, we regarded all the moving objects that fell in  $a=$  between 2 and

3.5 AU as MBAs. We found 861 MBAs in SMBAS. From this, we estimated that the sky number density of MBAs down to  $R=24.4$  mag is  $\sim 290/\text{deg}^2$  near opposition and near the ecliptic.



**FIG.6 (a) —The relative bias as a function of  $a$ .** The relative bias is the number ratio between near-ecliptic asteroids in distant orbits ( $r \sim 6$  AU) and those at  $r=a$  (AU). The relative bias curves are calculated for circular, near-circular and elliptic orbits. **(b) The relative bias as a function of  $I$ .** The relative bias is the number ratio between ecliptic asteroids ( $I \sim 0$ ) and those with  $I$ . Three relative bias curves are calculated for the inner-, middle-, and outer-MBAs. These figures are taken from Nakamura and Yoshida (2001).



Next we calculated the absolute magnitudes ( $H_R$ ) of each MBA by Eq.(1) (e.g., see Ephemerides of Minor Planets for 2001) :

$$H_R = R - 5 \log(\Delta \cdot r) - p(\alpha) - \delta V, \quad (1)$$

where  $R$  is the apparent magnitude in  $R$ -band of an asteroid in question,  $\Delta$  and  $r$  stand for the geocentric and heliocentric distances (in AU), respectively,  $p(\alpha)$  is the phase function ( $\alpha$ : phase angle, namely Earth-asteroid-Sun angle), and  $\delta V$  is the light variation due to rotation. Since SMBAS is observations near opposition the value of  $p(\alpha \sim 0^\circ)$  is negligibly small. Though  $r = a(1 - e^2) / [1 + e \cos(f + \omega)]$  near the ecliptic, where  $\omega$  is argument of perihelion and  $f$  is true anomaly, we can regard  $r \sim a$  for each asteroids, because of our assumption that  $e = 0$  (see Section 5). We can not estimate the  $\delta V$  for individual asteroids, because we didn't observe their lightcurve over all phase. However the phase of every asteroids are random. So their  $\delta V$  will be averaged out when we construct size distributions from their observations. So we put here  $\delta V = 0$  for convenience.

If the albedo ( $p$ ) of an asteroid is known or assumed, its diameter ( $D$ ) can be obtained approximately from its  $H_R$  (except for color effects) by Eq. (2) (Bowell *et al.* 1989), which is a modified version of the formula by Bowell and Lumme (1979):

$$\log D = 3.1295 - 0.5 \log p - 0.2 H_R. \quad (2)$$

We have used the averaged albedo of C- and S-type asteroids, that is, using an empirical formula  $\log D = 3.65 - 0.2 H$ , because we could not measure albedos of asteroids in SMBAS. Fig.7 shows a comparison between the  $H_R$ -distribution of MBAs detected in SMBAS and the  $H$ -distribution of 85,150 known MBAs September 2000 data of ([ftp://ftp.lowell.edu/pub/elgb/astorb.html](http://ftp.lowell.edu/pub/elgb/astorb.html)).

Since the  $(V-R)$  color for asteroids with typical taxonomic types is known to range between  $-0.05$  and  $+0.25$ , we need not distinguish  $H_R$ - and  $H$ -magnitudes within accuracy of our diameter estimate and such a comparison is therefore justified. The approximate  $D$  corresponding to the  $H_R$  was also indicated by arrows on the upper horizontal axis in Fig.7. From Fig.7, we see that the peak of the  $H_R$ -distribution for known MBAs is  $\sim 15$  mag (corresponding to  $D \sim 4.5$  km), whereas the peak for MBAs obtained with SMBAS is  $\sim 20$  mag (corresponding to  $D \sim 450$  m). Therefore, one can understand that our SMBAS could substantially observe sub-km MBAs whose size region is *one order of magnitude smaller* than that of known asteroids; this is really an unknown world. The range of MBAs in SMBAS for the minimum and the maximum size covers from  $\sim 0.1$  to  $\sim 10$  km.

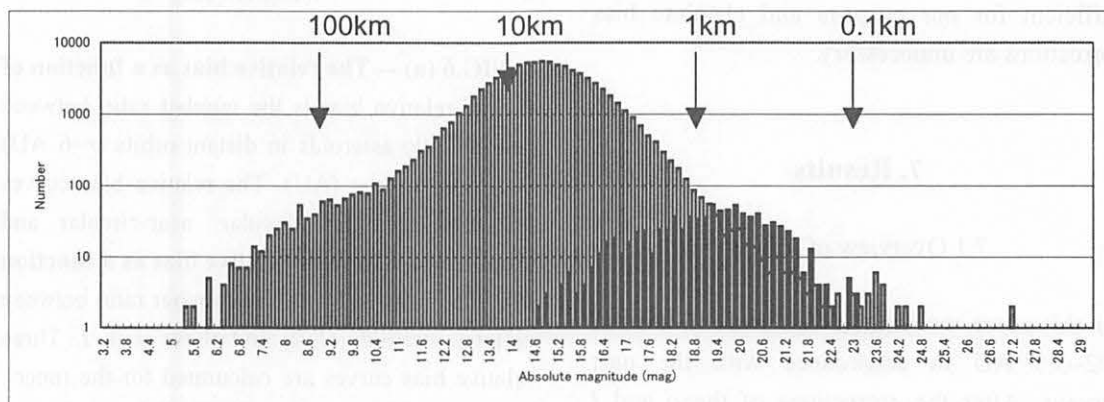


FIG.7---The  $H_R$ -distribution of MBAs detected in our observations and that for known MBAs (solid boxes: MBAs detected in SMBAS; open boxes: known MBAs).

## 7.2 $a$ - and $I$ -distributions

Fig.8a shows the  $a$ -distribution of 861 MBAs detected in SMBAS. The black, white, and gray boxes respectively indicate, the raw  $a$ -distribution of 861 MBAs, the  $a$ -distribution of MBAs corrected by the relative bias calculated for elliptic orbits, and that of MBAs corrected by the relative bias calculated for circular orbits given in Fig.6a. In Section 5, we said that the  $a$  of each asteroid in our estimation has a mean error of 0.13~0.15 AU. However, the  $a$ 's error in the histogram of Fig.8a should be generally smaller, because the statistical error in an  $a$ -bin is improved by the amount of  $\Delta a/\sqrt{n}$  ( $\Delta a$ : the  $a$ -error for a single asteroid,  $n$ : the data number in that  $a$ -bin), unless the data number is too small. Hence the histogram, especially the one corrected for elliptic orbit is worth comparing with that for known asteroids.

Fig.8b shows the  $a$ -distribution of 85,150 known MBAs discussed in Section 7.1 and Fig.7. In order to fairly compare the result of known MBAs with that for the MBAs from SMBAS, we drew Fig.8b by *intentionally degrading the resolution* for the  $a$ . In the  $a$ -distribution of known MBAs, it is well known that the diminution of asteroids near 2.1, 2.5 and 2.9 AU is reflected by the existence of the Kirkwood resonant gaps. One can see the depression near  $a \sim 2.1$  and 2.5 AU which is same as that of the known MBAs in the  $a$ -distribution of the small MBAs discovered with SMBAS. This might imply that, even for sub-km MBAs, the 3:1 Kirkwood gap still gives strong dynamical effects. On the other hand, the depression of the small MBAs near  $a \sim 2.9$  AU is so not conspicuous as that of the known MBAs; on the contrary the number seems to increase in this region in Fig.8a.

Fig.9a shows the  $I$ -distribution of MBAs detected in SMBAS. The black and white boxes represent the raw  $I$ -distribution of 861 MBAs from SMBAS and the  $I$ -distribution of MBAs corrected by the relative bias calculated for middle-belt asteroids in Fig.6b, respectively.

According to the statistical consideration already mentioned, a mean error of 1.5~5.8 deg in the  $I$  for individual asteroids in our estimation is also improved by the amount of  $\Delta I/\sqrt{n}$  ( $\Delta I$ : the  $I$ -error for a single asteroid,  $n$ : the data number in a  $I$ -bin), unless the data number is too small. Fig.9b shows the  $I$ -distribution of the 85,150 known MBAs. Though the  $I$ -distribution of known MBAs steeply decreases along with the increase of  $I$ , that of the relative-bias-corrected MBAs is seen to be comparatively uniform over the  $I$ -range shown in Fig.9a, except the lack of asteroids near  $I \sim 12^\circ$ . It is not clear now whether the lack is real or an artifact caused by statistics of small sample number, and will be a target for future exploration.

## 7.3 Size distribution in the whole main-belt

Next we discuss the cumulative size distribution (hereafter CSD) of sub-km MBAs detected in SMBAS. It is well known that a cumulative number distribution for MBAs brighter than a certain magnitude  $H_R$  is expressed as

$$\log N(<H_R) = C + \alpha H_R. \quad (3)$$

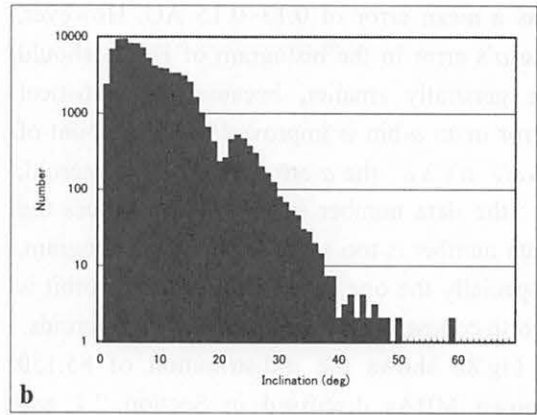
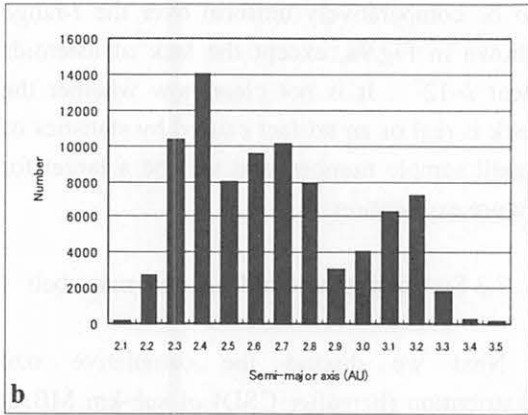
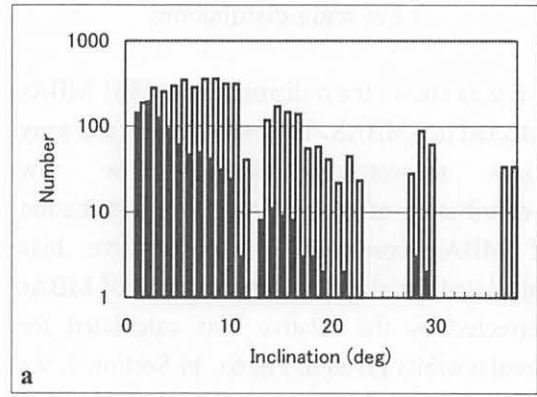
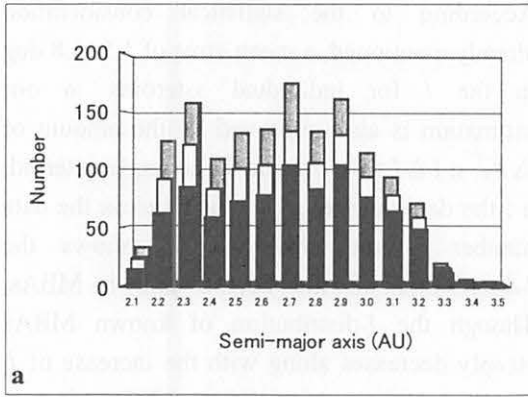
where  $\alpha$  and  $C$  are constants. The value of  $\alpha$  is often referred to as the slope for  $\log N$  vs.  $H_R$  plot. If we rewrite Eq. (3) with the help of Eq. (2), the result is equal to

$$N(>D) \propto D^{-b}. \quad (4)$$

The power-law index ( $b$ ) in Eq. (4), which corresponds to the slope for  $\log N$  vs.  $\log D$  plot, is connected to  $\alpha$  by  $b = 5\alpha$ . In this paper, we use  $b$  to express the slope of the CSD of asteroids.

Fig.10 shows the differential (white-box histogram) and cumulative (filled dots)  $H_R$ -magnitude distributions for 861 MBAs detected in SMBAS. The  $D$  corresponding to the  $H_R$  was also shown on the upper horizontal





**FIG.8---(a) The  $a$ -distribution of MBAs detected in SMBAS.** The black, white, and gray boxes respectively indicate the  $a$ -distribution of 861 MBAs from SMBAS, that of MBAs corrected by the relative bias calculated for ecliptic orbit, and that of MBAs corrected by the relative bias calculated for circular orbit in Fig.6a. **(b) The  $a$ -distribution of 85,150 known MBAs.**

**FIG.9---(a) The  $I$ -distribution of MBAs detected in SMBAS.** The black and white boxes indicate the  $I$ -distribution of 861 MBAs from SMBAS and one of MBAs corrected by the relative bias calculated for middle-belt asteroids in Fig.6b. Note that there are two detected asteroids near  $I = 37$  and  $38$  deg. **(b) The  $I$  distribution of 85,150 known MBAs.**

axis. The  $H_R$  of each asteroid was estimated by using Eq. (1) with its  $a$  (see Section 7.1). Therefore, the  $H_R$  of each asteroid includes the  $a$ -error caused by the assumption that  $e=0$ . The  $H_R$ -error caused by the  $a$ -error is  $0.25\sim 0.38$  mag. Again statistically, however, the  $H_R$ -error in a  $H_R$ -bin decreases by the amount of  $\Delta H_R / \sqrt{n}$  ( $\Delta H_R$ : the mean error derived from the mean  $a$ -error,  $n$ : data number in a  $H_R$ -bin). We showed the error bar only for the cumulative  $H_R$ -magnitude in Fig.10 for clarity.

The solid line was drawn to compare with the slope ( $b\sim 1.75$ ) for the CSD from PLS and Spacewatch survey. It seems that the slope of

the CSD for asteroids brighter than  $H_R\sim 17$  is a little steeper than  $\sim 1.75$ , and that of asteroids fainter than  $H_R\sim 18$  is much more gentle. Spacewatch (Jedicke & Mefcalfe, 1998) and SDSS (Ivezic *et al.*, 2001) ascertained that the slope of the CSD for small MBAs is shallower than previous estimates and cannot be represented by a single power-law. We can also confirm their results in Fig.10. Nevertheless, in order to compare our results with the  $b$  obtained from the past survey, we attempted to obtain the  $b$  by fitting with the least squares method. The best-fit value of the  $b$  for the asteroids with  $18.3 < H(\text{mag}) < 19.7$  (corresponding to  $0.5 < D$

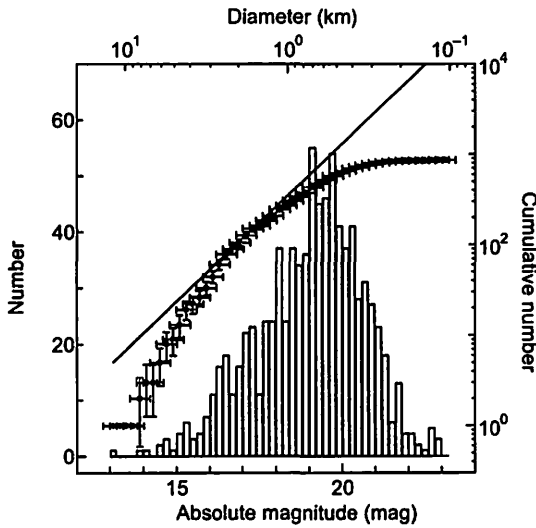


FIG.10—The  $H_R$ -distribution of MBAs detected with SMBAS in the whole main-belt.

(km)<1, assuming a mean albedo for known C- and S-type MBAs in Eq. (2)) is found to be  $1.18 \pm 0.03$ . Hence we attain the same conclusion with that of SDSS for asteroids conducted by Ivezić *et al.* (2001) that the number of small MBAs extrapolated so far based on the number of large MBAs ( $D>5$ km) was overestimated.

#### 7.4 Size distribution in three zones of the main-belt

Furthermore, we partitioned the main-belt into the inner, middle and outer zones defined by  $2.0 < a < 2.6$ ,  $2.6 < a < 3.0$ , and  $3.0 < a < 3.5$  AU, respectively. This division is conformable to the previous surveys: YMS, PLS and Spacewatch survey. Since the limit of the detectable magnitude becomes brighter along with the increase of asteroid's heliocentric distance, it is important to take such a distance effect into account when we examine the CSD of MBAs for each zone of the main-belt. The detectable  $H_R$ -magnitude for each zone were here calculated at the farthest position of each zone, namely 2.6, 3.0, and 3.5 AU by using Eq. (1) and the limiting magnitude ( $R=24.4$  mag) discussed in Section 4.2, and they were listed in Table IV. We can say that SMBAS detected

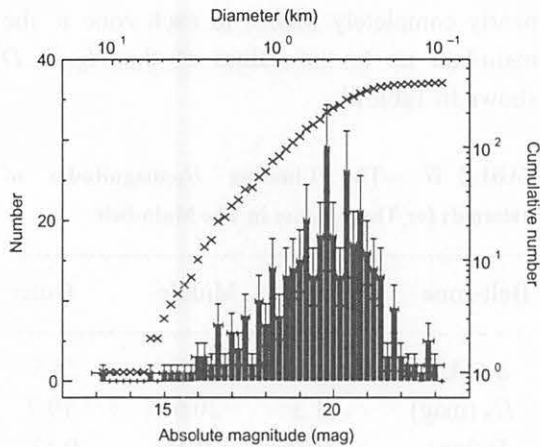
nearly completely MBAs in each zone in the main-belt up to the values of the  $H_R$  or  $D$  shown in Table IV.

TABLE IV ---The Limiting  $H_R$ -magnitudes of Asteroids for Three Zones in The Main-belt

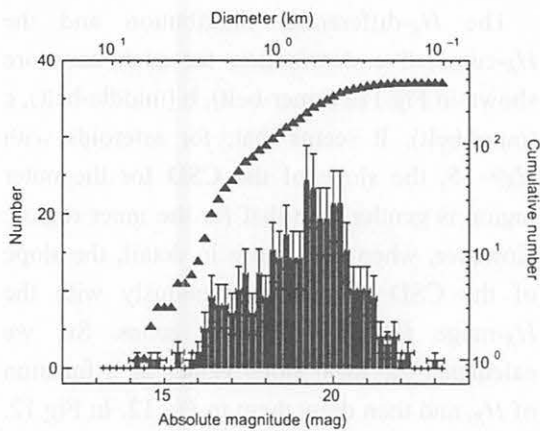
Belt-zone	Inner	Middle	Outer
$a$ (AU)	2.6	3.0	3.5
$H_R$ (mag)	21.3	20.5	19.7
$D$ (km)	0.25	0.35	0.51

The  $H_R$ -differential distribution and the  $H_R$ -cumulative distribution for each zone are shown in Fig.11a (inner-belt), b (middle-belt), c (outer-belt). It seems that, for asteroids with  $H_R \sim 15$ , the slope of the CSD for the outer region is gentler than that for the inner region. However, when seen more in detail, the slope of the CSD changes continuously with the  $H_R$ -range for any of three zones. So, we calculated the local slope values as a function of  $H_R$ , and then drew them in Fig.12. In Fig.12, the crosses, triangles, and open circles connected with curves show the changes in the CSD slopes for the inner-, mid-, and outer-MBAs, respectively. For asteroids with  $H_R > 17.5$ , it is seen that the slopes of the CSDs becomes gentler in order of the inner-, middle-, and outer-belt. Namely, it is obvious that the  $b$  of the CSD for outer-belt in Fig.12 is smaller than that for inner-belt. On the other hand, for asteroids with  $H_R < 17.5$ , the slope is steepest in the middle-belt, though it is possible that this trend may be an artifact caused by small sample statistics.

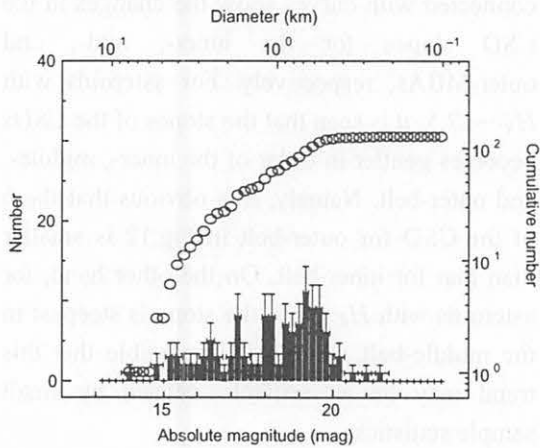
Now we focus on the CSD slopes of sub-km MBAs for the size range from 0.5km to 1 km in diameter in each zone, since we are especially interested in the CSD for sub-km MBAs ( $D < 1$  km). The discussion below is meaningful, because the number of MBAs detected in SMBAS is large enough to deal with statistically in this size region, and SMBAS



a inner-belt

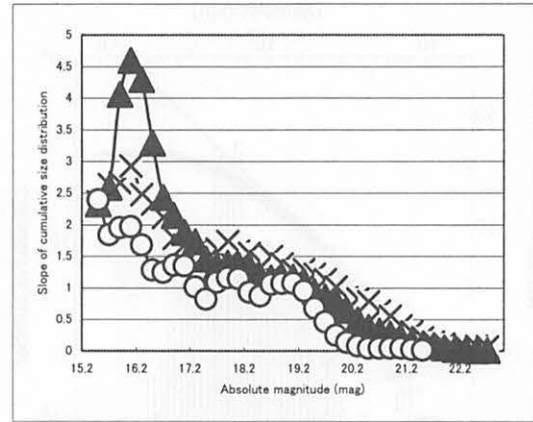


b middle-belt



c outer-belt

**FIG.11---The size distribution detected with SMBAS in three different zones of the main-belt.** The box histograms are the differential  $H_R$ -magnitude distribution of the MBAs (a)inner, (b)middle and (c)outer main-belt, respectively. The crosses, triangles, and circles show the cumulative  $H_R$ -magnitude distribution.



**FIG.12---Local slopes of the CSDs shown as a function of the absolute magnitude.** The crosses, triangles , and circles show the CSD slope curves of MBAs in the inner-, middle- and outer-zone, respectively. A conspicuous peak near  $H=16.2$  in the middle-belt curve is likely to be due to an artifact caused by small number statistics.

detection of MBAs is nearly complete down to  $D \sim 0.5$  km in the whole main-belt (see Table IV). Since the slope change is not so clear within the narrow sized-range, we examined the  $b$  of the CSD for MBAs with  $0.5 < D(\text{km}) < 1$  by fitting with the least squares method. The slopes of the CSDs thus calculated in the sized-range for the three main-belt zones are found to be  $1.37 \pm 0.04$ ,  $1.13 \pm 0.03$ , and  $1.10 \pm 0.03$ , respectively (see Table V). From the results, it is likely that the CSD slopes for sub-km MBAs are systematically different.

**TABLE V ---The Slopes of CSDs of MBAs with  $0.5 < D(\text{km}) < 1$**

Belt-zone	Inner	Middle	Outer
slope ( $b$ )	$1.37 \pm 0.04$	$1.13 \pm 0.03$	$1.01 \pm 0.03$

### 7.5 Spatial distribution in the whole main-belt

Fig.13 shows the spatial distribution of MBAs detected in SMBAS. Note that each data point includes the errors for the  $a$  and  $I$

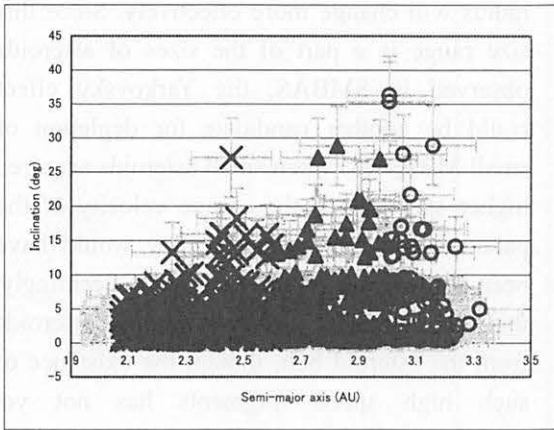


FIG.13---the spatial distribution as the  $a$  vs.  $I$  plot for MBAs detected in the SMBAS. Symbols correspond to the marks given in Fig.12.

mentioned in Section 5, namely the error  $\sim 0.1$  AU in the  $a$  and the error of  $2^\circ$  (for high inclination  $\sim 6^\circ$ ) in the  $I$  for each asteroid (also see Table III a, Table III b). Probably due to these errors, one cannot confirm the well-known Kirkwood gaps in Fig.13. Instead, however, there seems to be several vague gaps near  $I \sim 12^\circ$  and  $\sim 23^\circ$  over the whole main-belt. The cause of the depletion of asteroids with  $I \sim 23^\circ$  is known theoretically. It is due to the secular perturbations by Jupiter. On the other hand, the depletion of asteroids with  $I \sim 12^\circ$  has not been clear so far (nobody has ever pointed out that there is a weak depletion of asteroids with  $I = 10 \sim 12^\circ$  seen in PLS (Van Houten *et al.*, 1970)). Furthermore, the as of the high-inclination asteroids seem to be located near the heliocentric distances corresponding to the mean motion resonances of Jupiter, namely 2.5, 2.8 and 3.0 AU. This may imply that those asteroids are under the transport process in which asteroids are exported through the resonances from the main-belt to the other place in the solar system. Since it is likely that there is an inverse correlation between the size and velocity for fragments produced in collisions among asteroids (e.g., Nakamura & Fujiwara, 1991), smaller fragments should have faster release velocity. If so, it is reasonable that

the  $I$ -distribution of smaller MBAs has a wider range than that of large MBA. There seems to be a trend that asteroids with larger  $a$  show a higher  $I$ .

## 8 Discussion

### 8.1 Size distribution of sub-km MBAs

First we comment on the sky number density of small MBAs,  $\sim 290/\text{deg}^2$  ( $R < 24.4$  mag) which is given in Section 7.1. Although Poisson statistics teaches us that a formal error of the sky density is about 10 %, we suspect that the actual error will probably be much higher, perhaps by several ten percent. In fact, such a high variation of the number density can be seen depending on time, if plot on the sky positions for known asteroids of about a hundred thousand (<http://ftp.lowell.edu/pub/elgb/astorb.html>).

Hence, we consider that the sky number density is a less stable quantity to characterize the size distribution of small MBAs than their slope, and the latter is more important.

We have already seen in Section 7 that the CSD slope for sub-km MBAs with 0.5~1 km in diameter detected from SMBAS is obviously shallower than that for multi-km asteroids estimated with YMS, PLS, and Spacewatch surveys. Yoshida *et al.* (2001b) also found that the CSD slope of sub-km MBAs is  $\sim 1.0$  from their preliminary survey that was similar in principle to SMBAS but adopted a more simplified approach than SMBAS. Ivezić *et al.* (2001) found recently that the slope is 1.3 for MBAs with  $1 < D(\text{km}) < 5$  in SDSS project. Considering the above three results, it seems highly certain that small asteroids are not so plentiful as had been expected from past observations of larger ones.

However, when seen in more detail, our CSD slopes of sub-km MBAs with  $D = 0.5 \sim 1 \text{ km}$  are shallower than that obtained in SDSS. This implies that the number of sub-km MBAs

detectable only in SMBAS is much more depleted compared with the prediction by Ivezić *et al.* (2001). Thus we infer that there is size dependence in the mechanism to remove asteroids from the main-belt. For example, Nakamura (1994) found that smaller asteroids more abundantly exist more closely to the centers of the Kirkwood gaps, from where NEAs are believed to be supplied.

As possible causes of the depletion of small asteroids in the main-belt, researchers have so far proposed three physical processes as follows: (1) small asteroids may become a part of large asteroids that have the structure of strengthless “rubble-piles”. Likely existence of rubble-pile asteroids which consist of re-accumulated impact fragments was theoretically predicted for the first time by Weidenschilling (1981). In fact, the asteroid (253) Mathilde observed by the NEAR spacecraft (Veverka *et al.*, 1999), and (216) Kleopatra observed by radar (Ostro *et al.*, 2000) have been regarded as having the rubble pile structure, because of their observed low bulk density. Recent collisional theories and experiments suggest that the impact energy needed to disperse an asteroid is greater than that to thoroughly shatter it, for asteroids larger than a few km to sub-km in size. This means that it is more difficult to disperse collisional fragments for asteroids in such sizes. If this is the case, the number of small MBAs should be depleted, because they will be incorporated into a part of the large rubble-piled asteroids. (2) small asteroids would have been thrown into the Kirkwood gaps by the Yarkovsky effect (anisotropic repulsion force due to thermal emission), and then they would have been removed from the main-belt (e.g. Farinella & Vokrouhlický, 1999). According to calculations by Farinella & Vokrouhlický (1999), the  $a$  of the asteroids with 1~10 km in radius can be moved by a few hundredths of AU by the Yarkovsky effect during their collisional lifetimes (10~1000 million years). In particular, the  $a$  of small asteroids with 10~100 m in

radius will change more effectively. Since this size range is a part of the sizes of asteroids observed in SMBAS, the Yarkovsky effect could be another candidate for depletion of small MBAs. (3) when small asteroids acquired higher speeds than the escape velocity of the parent body in a collision, they would have been thrown out of the main-belt. Seemingly, this is the simplest process to remove asteroids from the asteroid belt, though the existence of such high speed fragments has not yet established in laboratory experiments. Presently we cannot say from only our SMBAS, which process among the above three candidates is more plausible. To solve the problem, we may need detailed observations for the CSD of NEAs and/or for the CSDs of the craters on the surfaces of inner planets and satellites, in addition to the CSD of small MBAs. In any case, we would say that our investigation of the CSD for sub-km MBAs (namely NEA-sized asteroids) can provide an important step in estimating both the supply rate of NEAs and the formation rate of rubble pile asteroids.

Next, we discuss the CSDs of sub-km MBAs investigated for three zones of the main-belt. As we have already shown in Section 7.3, it is fairly certain that there is a difference in the slopes between the inner- and outer-zones for asteroids with 0.5~1 km in diameter. The slope in the inner-zone is relatively steep ( $\sim 1.4$ ), while that in the outer-zone is shallow ( $\sim 1.0$ ). However, we must remember here that transformation from the brightness ( $H_R$ -mag) of asteroids to the size considerably depends on their albedo. For well-observed MBAs, we know that S-type asteroids with a high albedo are abundant in the inner main-belt while C-type ones with a low albedo are dominant in the outer region of the main-belt; namely the number ratio of the S-type and C-type asteroids varies with the heliocentric distance (<http://pdssbn.astro.umd.edu>). Note that a C-type asteroid is about twice as large as a S-type one with the same absolute magnitude, due to the difference in albedo. Xu *et al.* (1995)

**TABLE VI**  
**Two kinds of Mean Slopes of the CSDs for Asteroids with  $0.5 < D(\text{km}) < 1$**

Belt-zone	Inner-belt $2.0 < a < 2.6$	Middle-belt $2.6 < a < 3.0$	Outer-belt $3.0 < a < 3.5$
mean slope 1	$1.37 \pm 0.04$	$1.13 \pm 0.03$	$1.01 \pm 0.03$
mean slope 2	$1.58 \pm 0.04$	$1.13 \pm 0.03$	$0.83 \pm 0.07$

Note : mean slope 1 ; the mean slope of the CSD for asteroids of the size-range estimated based on the assumption that any asteroids have the mean albedo of well-known MBAs, mean slope 2 ; the mean slope of the CSD for asteroids of the size-range estimated by considering the abundance ratio of the S-type and C-type asteroids in the main-belt.

showed, in the Small Main-belt Asteroid Spectroscopic Survey (SMASS), that the majority of the small main-belt asteroids ( $D < 20$  km) are C- and S-type asteroids, and their distributions are similar to the one of large asteroids. Therefore, we assumed that all of the inner-belt asteroids have the albedo of S-type asteroids, the middle-belt ones have a mean albedo between S-type and C-type asteroids, and the outer-belt ones have the albedo of C-type asteroids. Then we re-estimated the CSD slopes for the three main-belt zones by taking the albedo effects into account and showed the results in Table VI. In Table VI, the mean slope 1 indicates the slope of the CSD for asteroids in the sized-range estimated based on the mean albedo of well-known MBAs, which has been already shown in the Table VI. The mean slope 2 represents the slope estimated by considering the ratio mentioned above of the S-type and C-type asteroids in the main-belt. In Table VI, except for the mean slope 2 of outer-belt asteroids, all slopes were calculated for asteroids with  $0.5 < D(\text{km}) < 1$ . The mean slope 2 of the outer-belt was estimated for asteroids with  $0.7 < D(\text{km}) < 1$ , because completeness of asteroid detection is only up to  $H_R = 19.7$  mag ( $D = 0.7$  km, assuming the albedo of C-type asteroids) in this region (refer to Table IV). The above result indicates that consideration of the S/C number ratio for

MBAs makes clearer, the difference in the slopes of CSD between inner- and outer-MBAs.

In the PLS, the slope of CSD in the inner-belt was also somewhat steeper than those in the other zone of the main-belt (Van Houten *et al.* 1970). We here again propose that there really exists a difference in the CSD depending upon the location of sub-km MBAs, namely it is comparatively steep in the inner-belt, and shallow in the outer-belt.

Then what is the cause of this difference ? Here are some considerations. We may infer that some specific mechanism to remove a large number of small asteroids had worked in a distant past, or rubble-pile asteroids have been produced more effectively in the outer-belt rather than in the inner-belt. We may also be able to consider that it is a result of the difference in the distributions between S- and C-type asteroids. For larger asteroids, in this respect, Anders (1965) had suggested that the frequency of collisions is different between the inner-belt and the outer-belt : in the inner belt, impact frequency is only a few times throughout its history, while collisions are severe and asteroids are highly fragmented in the outer-belt due to proximity to Jupiter. We know that C-type asteroids and S-type asteroids are like carbonaceous chondrites and silicate rocks, respectively. Furthermore, by recent

space probe investigations, we know that the bulk densities of (243) Ida (S-type asteroid) and (253) Mathilda (C-type asteroid) are  $\sim 2.6 \text{ g/cm}^3$  and  $\sim 1.3 \text{ g/cm}^3$ , respectively. It is therefore likely that the different outcomes would occur in collisions of bodies that have the different material and density.

On the other hand, Nolan *et al.* (2001) found by their numerical simulations that since a shock wave after a collision fractures an asteroid in advance of crater excavation flow, impact results are controlled by gravity ; the tensile strength is unimportant whether asteroids are initially intact or rubble-piles. If this is the case, it means that the dispersion of fragments after a collision is independent from the tensile strength of parent bodies. Hence, in short, whether a correlation exists or not between the size distribution of the collisional fragments and the tensile strength of the parent bodies is not yet known. Therefore, in order to pursue the cause of the difference of slopes of the CSDs in three zones of the main-belt that we found, it is necessary to investigate separately the size distributions of C-type and S-type asteroids. And the point will be discussed a little more in detail in 9.2.

## 8.2 Spatial distribution of sub-km MBAs

Finally, we discuss the spatial distribution of sub-km MBAs. Fig.13 indicates that the spatial distribution of the smaller-sized asteroids has a wider  $I$ -distribution in the outer-part of the main-belt. We also pointed out that a small number of asteroids with high inclination were seen in the neighborhood of the distinct Kirkwood gaps at 2.5, 2.8, and 3.0 AU. Dynamical consideration generally shows that, once an asteroid gets trapped into a gap, it undergoes a chaotic orbital transition, its  $e$  grows unexpectedly to a very high value, and finally it is ejected from the main-belt to the near-Earth region or to the other regions in the solar system. In some cases where a mean-motion resonance is coupled with a

secular resonance or the Kozai resonance, it is shown that the inclination also pumps up to a high level (Morbidelli and Moons 1995). So it is possible that the high-inclination asteroids near the Kirkwood gaps in Fig.13 might correspond to such chaotic asteroids. In the Fig.13, asteroids near the gaps (at 2.5, 2.8, and 3.0 AU) may be in the process of delivery outside the main-belt.

The depletion of asteroids with  $I \sim 12^\circ$  over the whole main-belt may also be significant, which it is seen in Fig.13, because the recent studies of numerical direct orbital integration and the analytical method showed that the pumping-up of  $I$ s and  $e$ s of asteroids caused by sweeping secular resonances in the main-belt (Nagasawa *et al.* (2000)).

# 9. Conclusions and future prospect

## 9.1 Summary and conclusions

We detected 1111 moving objects down to  $R=24.4 \text{ mag}$  in the sky area of  $2.97 \text{ deg}^2$  near opposition and near the ecliptic in SMBAS. Then, we identified 861 MBAs by estimating the  $a$  of each moving object from its sky motion vector. The sky number density of MBAs was  $\sim 290$  per  $\text{deg}^2$  down to  $R=24.4 \text{ mag}$  near opposition and near the ecliptic. We found that the slope of the CSD for small MBAs ranging from a few km to sub-km is fairly shallower ( $\sim 0.8$ - $1.6$ , depending upon locations in the main-belt) than that for large MBAs ( $\sim 1.8$ ) obtained from the past asteroid surveys. This means that the number of sub-km MBAs is much more depleted than a result extrapolated from the size distribution for large asteroids. The CSD slope of the inner sub-km MBAs was somewhat steeper than that of the outer sub-km MBAs. The investigation of the spatial distribution suggests that there seems to be a trend that asteroids with larger  $a$  show a higher  $I$  and there seems to be a new gap near  $I \sim 12^\circ$  over the whole main-belt.



From the above mentioned results, we conclude that overall size and spatial distributions of very small asteroids can fairly be different from those for large asteroids. However, one possible weak point of our survey described in this paper might be smallness of the survey area, only covering about  $3 \text{ deg}^2$ . In this respect, we plan to widen the survey area in near future observations, in order to make our conclusions more reliable.

### 9.2 Future prospect

As a next step, we must clarify and well interpret the difference for the slopes of the CSD in the individual regions of the main-belt. For the purpose, it is necessary for us to investigate each size distribution for C-type and S-type asteroids because these two taxonomic types are major component of MBAs. Their two types can be discriminated from  $(B-V)$  or  $(V-R)$  color observations. We performed such observations in late October, 2001, and data reductions are now progressing. As above mentioned, since it is believed that there is some correspondence between the asteroid material and taxonomic types, such a survey observation conducted by us will also allow us to argue interrelations between the collisional processes, orbital evolution and material distribution of MBAs.

### Acknowledgement

We gratefully thank Dr. Yutaka Komiyama (support astronomer), Mr. Gray Fujiwara, Mr. Bob Potter, and Miss Sumiko Harasawa (night operators), Suprime-Cam team, and Subaru telescope supporting staff for their kind support in our observations. We would also like to express our thanks to Dr. Takashi Ito and staff of the Astronomical Data Analysis Center, NAOJ whose advice in data reduction was very useful.

This paper is based on the first author's

dissertation, submitted to Kobe University, in partial fulfillment of requirements for the doctorate.

### References

- Anders, E. 1965. Fragmentation history of asteroids. *Icarus* **4**, 398-408.
- Bowell, E., B. A. Skiff, L. H. Wasserman, and K. S. Russell. 1990. Orbital information from asteroid motion vectors, in *ACM- III*, pp.19-24. Uppsala Univ.
- Bowell, E., B. Hapke, D. Domingue, K. Lumme, J. Peltoniemi and A. W. Harris 1989. Application of photometric models to asteroids. in *Asteroids II*, (R. P. Binzel, T. Gehrels, and M. S. Matthews, Eds.), pp.524-556. the Univ. of Arizona press.
- Bowell, E. and Lumme, K., (1979) Colorimetry and magnitudes of asteroids. in *Asteroids*, (Gehrels, T. Eds.), pp.132-169. the Univ. of Arizona press.
- Donnison, J. R., and M. P. Wiper 1999. Bayesian statistical analysis of asteroid rotation rates. *Mon. Not. R. Astron. Soc.* **302**, 75-80.
- Durda, D.D., R. Greenberg, R. Jedicke 1998. Collisional models and scaling laws: A new interpretation of the shape of the main-belt asteroid size distribution. *Icarus*, **135**, 431-440.
- Farinella, P. and D. Vokrouhlicky 1999. Semimajor axis mobility of asteroidal fragments. *Science* **283**, 1507-1510.
- Harris, A. W., and J. A. Burns 1979a. Asteroid rotation I. Tabulation and analysis of rates, pole positions and shapes. *Icarus* **40**, 115-144.
- Harris, A. W. 1979b. Asteroid rotation rates II. A theory for the collisional evolution of rotation rates. *Icarus* **40**, 145-153.
- Ivezic, Z., and 31 colleagues (for THE SDSS COLLABORATION) 2001. Solar system objects observed in the Sloan Digital Sky Survey commissioning data. *Astron. J.* **122**,



- 2749-2784.
- Jedicke, R. and T. S. Metcalfe 1998. The orbital and absolute magnitude distributions of main belt asteroids. *Icarus* **131**, 245-260.
- Kinoshita, D., N. Yamamoto, J. Watanabe, T. Fuse, O. Hainaut, H. Boehnhardt, S. Ida, M. Nagasawa, F. Yoshida 2002. Wide-field survey near the ecliptic with Subaru telescope, submitted to *AJ*.
- Komiyama, Y., and 32 colleagues 2000. High-resolution images of the ring nebula taken with the Subaru telescope. *Pub. Astron. Soc. Japan* **52**, 93-98.
- Kuiper, G. P., Y. Fujita, I. Gehrels, J. K. Groeneveld, G. Van Biesbroeck, and C. J. Van Houten 1958. Survey of asteroids. *Astrophys. J. Suppl.* **3**, 289-427.
- Landolt, A. U. 1992. UBVRI photometric standard stars in the magnitude range  $11.5 < V < 16.0$  around the celestial equator. *Astron. J.* **104**, 340-371, 436-491.
- Melosh, H. J. and E. V. Ryan 1997. Asteroids: shattered but not dispersed. *Icaurs* **129**, 562-564.
- Meyer, S.L. 1975. *Data Analysis for Scientists and Engineers*. Chap. 25, John Wiley and Sons, New York.
- Millis, R. L., Buie, M. W., Wasserman, L. H., Elliot, J. E., Kern, S. D., and Wagner, R.M 2002. The deep ecliptic survey : A search for Kuiper belt objects and Centaurs. I. Description of methods and initial results. *AJ* **123**, 2083-2109.
- Morbidelli, A. and M. Moons 1995. Numerical evidence on the chaotic nature of the 3/1 mean motion commensurability. *Icarus* **115**, 60-65.
- Nagasawa, M., H. Tanaka, and S. Ida 2000. Orbital evolution of asteroids during depletion of the solar nebula. *AJ* **119**, 1480-1497.
- Nakamura, A. and Fujiwara, A. 1991. Velocity distribution of fragments formed in a simulated collisional disruption *Icarus* **92**, 132-146.
- Nakamura, T. and F. Yoshida 2001. Statistical method for deriving spatial and size distribution of sub-km main belt asteroids from their sky motions. submitted to *Pub. Astron. Soc. Japan*.
- Nakamura, T. 1994. Size dependence of asteroid belt structure. *Seventy-five years of Hirayama asteroid families ASP conference series, Vol. 63* (T. Kozai, R. P. Binzel, and T. Hirayama, Eds.), pp.52-61.
- Nolan, M. C., E. Asphaug, R. Greenberg, and H. J. Melosh 2001. Impacts on asteroids : fragmentation, regolith transport, and disruption. *Icarus* **153**, 1-15.
- Ostro, S. J., R. S. Hudson, M. C. Nolan, J. I. Margot, D. J. Scheeres, D. B. Campbell, C. Magri, J. D. Giorgini, and D. K. Yeomans 2000. Radar Observations of Asteroid 216 Kleopatra. *Science* **288**, 836-839.
- Van Houten, C. J., G. I. Van Houten, P. Herget, T. Gehrels 1970. The Palomar-Leiden survey of faint minor planets. *Astr. Astrophys. Suppl.*, **2**, 339-448.
- Veverka J., P. Thomas, A. Harch, B. Clark, J. F., Bell III, B. Careich, J. Joseph, S. Murchie, N. Izenberg, C. Chapman, W. Merline, M. Malin, L. McFadden, and M. Robinson 1999. NEAR encounter with asteroid 253 Mathilda : overview. *Icarus* **140**, 3-16.
- Wetherill, G. W. 1989. Origin of the asteroid belt. in *Asteroids II* (R. P. Binzel, T. Gehrels, and M. S. Matthews, Eds.), *the Univ. of Arizona press*, pp.661-680.
- Weidenschilling, S. J. 1981. How fast can an asteroid spin ?. *Icarus* **46**, 124-126.
- Wisdom, J. 1983. Chaotic behavior and the origin of the 3/1 Kirkwood gap. *Icarus* **56**, 51-74.
- Xu, S., R. P., Binzel, T. H. Burbine, and S. J. Bus 1995. Small main-belt asteroid spectroscopic survey : initial results. *Icarus* **115**, 1-35.
- Yoshida, F. and T. Nakamura 2001a. Sub-km Main Belt Asteroid Survey (SMBAS) using SUBARU telescope — The Estimation Methods of Size Distributions — . *Proceedings of the 33rd symposium on Celestial Mechanics*, 269-276.

- Yoshida, F., T. Nakamura, T. Fuse, Y. Komiyama, M. Yagi, S. Miyazaki, S. Okamura, M. Ouchi, and M. Miyazaki 2001b. First Subaru Observations of Sub-km Main Belt Asteroids. *Pub. Astron. Soc. Japan Letter* **53**, L13-L16.
- Yoshida, F. and T. Nakamura 2000. SUBARU sub-km main belt asteroid survey plan — Statistical estimation of size distribution—. *Proc. of 33rd ISAS Lunar and Planetary Sympo.*, 21-24.

# Sub-km小惑星のデータに力学的摩擦の徴候は見られるか?

吉田 二美(国立天文台・神戸大)

箱根天体力学N体力学研究会 Mar.11-13,2002

微小メインベルト小惑星のサイズ分布・空間分布を調べるために、すばる望遠鏡を用いて観測を行い、861個の小惑星帯小惑星の軌道長半径( $a$ )、軌道傾斜角( $i$ )、絶対等級( $H$ )のデータを得た。

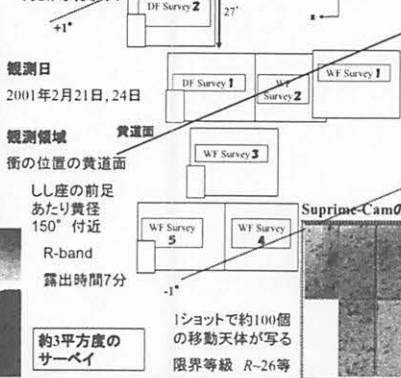
このデータをもとに、 $2 < a < 3.5$  AUの小惑星帯小惑星において、力学的摩擦が有効であるかどうかを検討した。

## すばる望遠鏡によるsub-km小惑星の観測 Sub-km Main-Belt Asteroid Survey (SMBAS)

### 観測装置



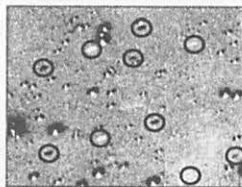
### 観測領域



### 小惑星探し

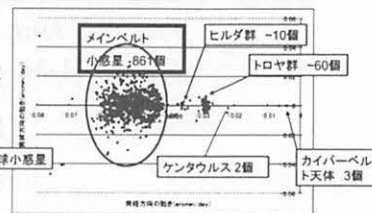
同じ領域を撮ったイメージをスケールレベルを合わせて重ねると移動天体は3個1組の白黒の点か、白黒棒として見える。これらの天体の数を注意深く目で数えた。

全部で1111個の移動天体を発見



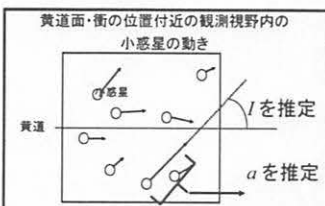
### 検出された天体

- 検出された1111個の天体のうち、既知の天体は60個程度であった。
- 小惑星帯小惑星を各移動天体の移動速度の違いを利用して分離した(右図)。
- 検出された小惑星のサイズ範囲は0.1~10km(過半数は1km以下)であった。



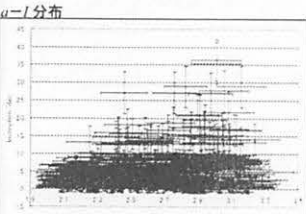
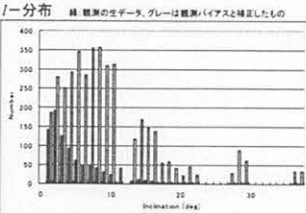
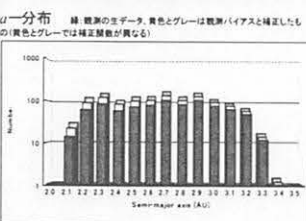
### 各小惑星帯小惑星の軌道決定の精度

我々の観測では1個の小惑星に対して2時間の観測arcしかないの、正確な軌道決定はほぼ不可能である。そこで各小惑星の軌道離心率をゼロと仮定して、各小惑星の見かけの移動速度と方向から、 $a$ (軌道長半径)と $i$ (軌道傾斜角)を推定した。各小惑星の真の軌道要素と見かけの運動ベクトルから得られる軌道要素を比較することにより、小惑星帯小惑星に限れば下記の誤差内で $a, i$ が推定できることをシミュレーションにより確かめた。

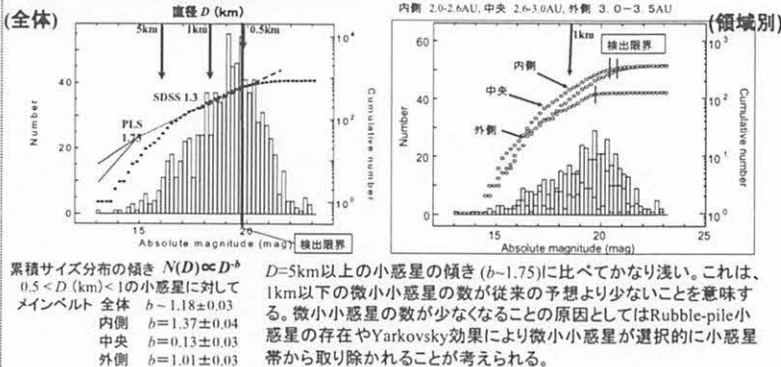


- 軌道長半径( $a$ ) 推定誤差  $\Delta a \sim \pm 0.1$  AU
- 軌道傾斜角( $i$ )  $\Delta i_1 \sim \pm 1.5^\circ$  ( $0 < i < 10$  の小惑星に対して)  $\Delta i_2 \sim \pm 4.4^\circ$  ( $10 < i < 20$  の小惑星)  $\Delta i_3 \sim \pm 5.8^\circ$  ( $20 < i < 30$  の小惑星)

### Sub-km小惑星の空間分布



### Sub-km小惑星のサイズ分布



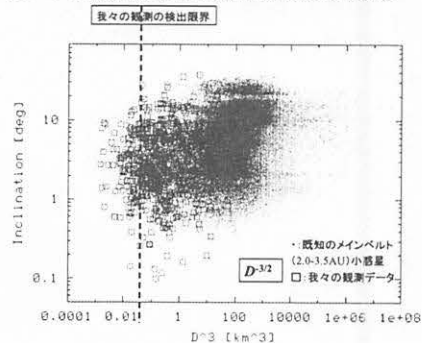
Sub-km小惑星の $i$ -分布において、 $i \sim 12^\circ$  付近にギャップがあるように見える。 $a$ - $i$ 分布でこのギャップはメインベルト全体にわたって存在するのがわかる。1960年に行われたPLS(直径5km以上の小惑星帯小惑星のサイズ分布及び空間分布を調べた)でもこのギャップは見えていた。このギャップが力学的成因(惑星との永年共鳴など)によるものか、 $i \sim 12^\circ$  を挟んでファミリーが存在するために相対的に $i \sim 12^\circ$  の小惑星が少なく見えるのか、今のところわからない。今回の我々の観測はほぼ黄道面上のサーベイであるため軌道傾斜角の大きい小惑星が検出される可能性は少ない。我々は黄道面に垂直な方向にサーベイ観測を行ってこのギャップが確かに存在するかどうか確かめる必要がある。

さらに小惑星帯の外側で微小小惑星の枯渇は顕著である。小惑星帯の内と外でのサイズ分布の違いは、小惑星の構成物質が異なるため(実際に小惑星帯の内と外で小惑星のタイプは異なる)か、もしくは木星の影響の違いを反映するものと思われる。外側ほど木星の摂動を強く受けるので、小惑星帯の外側の小惑星は内側の小惑星よりも強い衝突を何度も経験した可能性がある。次期観測では小惑星のタイプごとに小惑星帯の外側と内側のサイズ分布を調べる予定である。

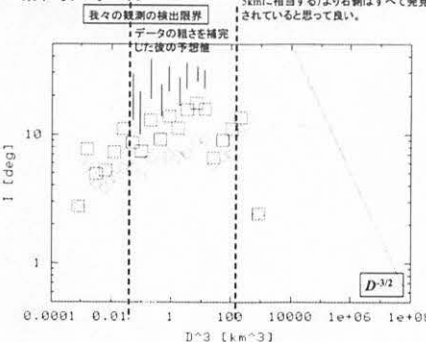
## Sub-km小惑星のデータに力学的摩擦の徴候は見られるか?

今回の箱根N体力学研究会のテーマは『力学的摩擦』であった。サイズの異なる粒子の集団の中で、系の運動状態を成る種の熱平衡状態へと近付ける効果が力学的摩擦と言える。これを小惑星帯で言えば、質量の大きな小惑星のランダム速度を小さくし、質量の小さな小惑星のランダム速度を大きくする効果を力学的摩擦は持っている可能性がある。勿論、小惑星帯の天体の個数密度が低くて相互作用が小さすぎる場合には、力学的摩擦は効果を発揮しない。我々の観測により発見された非常に小さな小惑星のランダム速度は、既知の大きな小惑星との相互作用によって増大されるのか? 軌道傾斜角のデータをアルベドを仮定して得た小惑星の質量のデータと比較し、以下のように検証してみた。

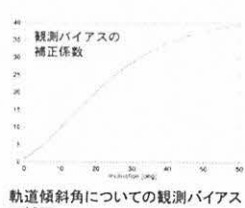
生データ 固有軌道傾斜角に直してある(cf. Brower & Clemence 1961)。



二乗平均データ  $< i^2 >^{1/2}$



既知の小惑星について線形補正( $D=5$  kmに相当する)より右側はすべて発見されていると想定して良い。



軌道傾斜角についての観測バイアスの補正

● 我々の観測データ  
● 既知の小惑星の二乗平均データ  
● 我々の観測データの二乗平均データ  
● 我々の観測データの $i$ についての観測バイアスを補正した二乗平均データ

### 現時点での結論

- 小惑星帯では力学的摩擦は弱いように見えるが、sub-km小惑星の観測データを増やしていくことで、力学的摩擦が見えてくる可能性がある。
- $D > 5$  kmでは力学的摩擦は全く見られない。

すばる望遠鏡による観測は継続中。  
2001年10月 データ解析中  
2002年9月に再度観測予定

謝辞: 観測生データから、固有軌道傾斜角への変換は国立天文台・データ解析センターの伊藤幸士さんにやりました。

Mar.10, 2002 Junji Yoshida

# 離散力学系におけるレヴィフライト Dynamical systems which produce the Lévy flights

Tomoshige Miyaguchi\* and Yoji Aizawa†

Department of Applied Physics, Faculty of Science and Engineering,  
Waseda University, Tokyo 169-8555, Japan

## Abstract

We introduce a one-dimensional map producing flights of arbitrary length and explain that the orbits and the density functions that evolve under this map have the same properties as Lévy flight. We derive an approximated Frobenius-Perron equation and prove that this equation converges to the Lévy diffusion equation.

## 1 Introduction

Gaussian statistics, which are characterized by exponentially decaying tails of the density function, provide an important framework of statistical mechanics and thermodynamics. But physical systems also frequently exhibit scaling behavior, e.g., Hamiltonian systems[1], fluid mechanical systems [2], and economics systems[3]. Lévy stable distributions are believed to be important for the statistical treatment of such systems. The symmetric Lévy stable distribution is defined by its characteristic function,

$$\hat{p}(k) = \exp(-b|k|^\alpha), \quad (1)$$

where  $\alpha$  ( $0 < \alpha \leq 2$ ) is the Lévy index and  $b$  ( $b > 0$ ) is a scale factor. For  $\alpha = 2$  and  $\alpha = 1$ , the corresponding distributions are Gaussian and Cauchy distributions, respectively. Except for the case  $\alpha = 2$ , the variance of a stable distribution is infinite, and for  $\alpha < 1$ , the mean is also infinite.

Lévy stable distributions were defined by Paul Lévy[4] as rescaled sums of independent, identically distributed random variables. This definition naturally leads to the concept of a diffusion process called ‘Lévy flight’, whose density function is a Lévy stable distribution at any instant if the initial conditions are given by a delta function  $\delta(x)$ . Lévy flight is defined by the Lévy diffusion equation[5],

$$\frac{\partial \rho_t(x)}{\partial t} = -b \int_{-\infty}^{\infty} D^\alpha(x - x') \rho_t(x') dx', \quad (2)$$

where  $D^\alpha(x)$  is the Fourier transform of the generalized function  $|k|^\alpha$ :

$$D^\alpha(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} |k|^\alpha \exp(-ikx) dk. \quad (3)$$

---

\*E-mail: tomo@aizawa.phys.waseda.ac.jp

†E-mail: aizawa@aizawa.phys.waseda.ac.jp

Integrating the Fourier transform of Eq. (2), it is easily shown that the density function  $\rho_t(x)$  is a Lévy stable distribution whose Lévy index is  $\alpha$ . For  $\alpha = 2$ , this equation is a normal diffusion equation, whose space derivative is of order 2. Contrastingly, for  $\alpha < 2$ , this equation becomes non-local.

Although Lévy flights are defined in a purely probabilistic manner, we would like to know the dynamical structures generating Lévy flights. Recently, deterministic models of normal and anomalous diffusion have been introduced by many people [7, 8]. In those models, the flight lengths are finite. For Lévy flight, however, arbitrarily long flights are significant[9]. For this reason, in this article, we propose a one-dimensional map producing flights of arbitrary length and discuss its convergence to Lévy flight.

In §2, we introduce a model consisting of a one-dimensional map. In §3, we derive the relation between the approximated Frobenius-Perron equation (FPE) of this map and Lévy diffusion equation and the last section is devoted to a summary.

## 2 One-dimensional maps

Consider the one-dimensional map  $f_\gamma(x)$  defined by

$$f_\gamma(x) = \begin{cases} -\frac{a\epsilon^\gamma}{(a+x)^\gamma}, & -a \leq x \leq -(a-\epsilon) \\ \frac{a}{a-\epsilon}x, & -(a-\epsilon) \leq x \leq (a-\epsilon) \\ \frac{a\epsilon^\gamma}{(a-x)^\gamma}, & a-\epsilon \leq x \leq a \end{cases} \quad (4)$$

which satisfies the discrete translational symmetry

$$f_\gamma(x + 2na) = f_\gamma(x) + 2na. \quad (n = 0, \pm 1, \pm 2, \dots) \quad (5)$$

In Eq. (4),  $a$ ,  $\gamma$  and  $\epsilon$  are parameters obeying the restrictions  $a > 0$ ,  $\gamma > 0$  and  $0 < \epsilon < a$ . As shown in Fig. 1, this map diverges at lattice points,  $x = (2n+1)a$  ( $n = 0, \pm 1, \pm 2, \dots$ ), i.e., the mapping generates arbitrarily long flights.

We define the orbit  $x_t$  by successive iterations of  $f_\gamma(x)$ :

$$x_{t+1} = f_\gamma(x_t). \quad (6)$$

Three different realizations of the orbit  $x_t$  are shown in Figs. 2(a)-(c). In Fig. 2(a), small fluctuations are dominant, and the orbit resembles Brownian motion. On the other hand, for large  $\gamma$ , long flights are generated, and the orbits display behavior similar to Lévy flights.

Next, we calculate and approximate the FPE of our map  $f_\gamma(x)$ . In what follows, we assume  $-a \leq x \leq a$ , and the result will be extended to the entire real axis, due to the translational symmetry of our system. FPE is given by

$$\rho_{t+1}(x) = \sum_{n=-\infty}^{\infty} \frac{\rho_t(x_n)}{|f'_\gamma(x_n)|} \quad (7)$$

$$= \frac{a-\epsilon}{a} \rho_t\left(\frac{a-\epsilon}{a}x\right) + \frac{1}{\gamma} \sum_{n \neq 0} \frac{(a\epsilon^\gamma)^{\frac{1}{\gamma}}}{|x - 2na|^{1+\frac{1}{\gamma}}} \rho_t(x_n), \quad (8)$$

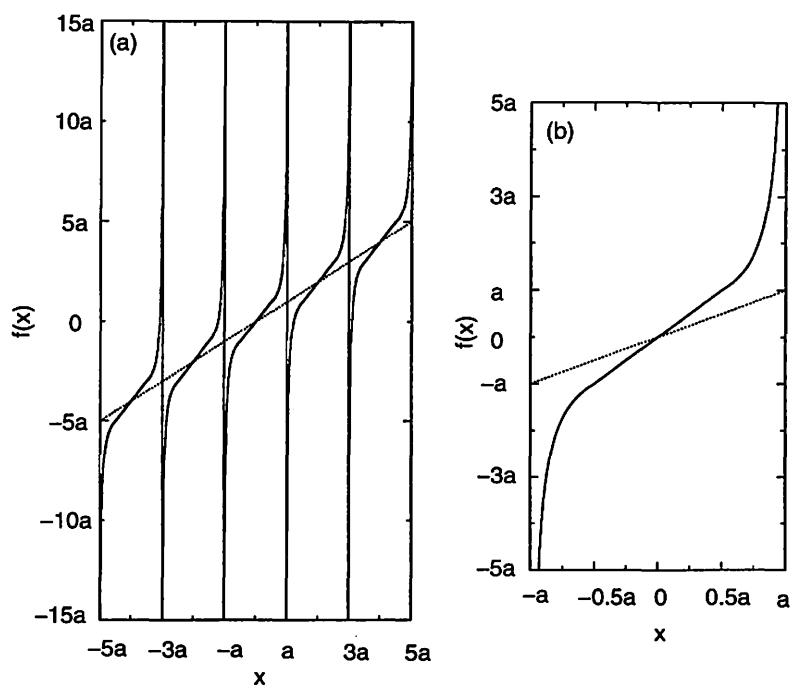


Figure 1: An example of the map  $f_\gamma(x)$ , for  $\epsilon = 0.5a$  and  $\gamma = 1.0$ . (a)  $f_\gamma(x)$  is presented over five periods. (b) The magnification over one period. At every lattice point, the value of the map diverges. The dashed lines are linear functions of slope 1.0.

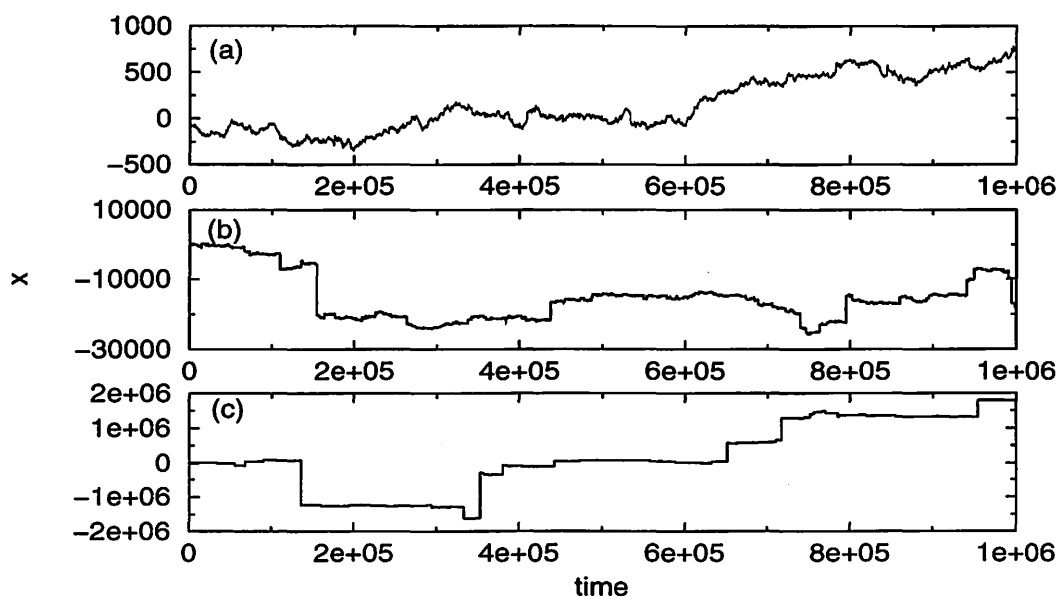


Figure 2: Numerical simulations of the orbits  $x_t$ . Each graph corresponds to a different value of  $\gamma$ : (a) 0.4, (b) 0.8, (c) 1.2, at  $a = 1$  and  $\epsilon = 0.1026$ .

where  $x_n$  ( $n = 0, \pm 1, \pm 2, \dots$ ) are the inverse images  $f_\gamma^{-1}$  of  $x$  ( $x_n = f_\gamma^{-1}(x)$ ), which are given by

$$x_n = \begin{cases} (2n+1)a - \epsilon \left( \frac{a}{x-2na} \right)^{\frac{1}{\gamma}} & \text{for } n \leq -1, \\ \frac{a-\epsilon}{a}x & \text{for } n = 0, \\ (2n-1)a + \epsilon \left( \frac{a}{2na-x} \right)^{\frac{1}{\gamma}} & \text{for } n \geq 1. \end{cases} \quad (9)$$

If the density  $\rho_t(x)$  is a smooth function of  $x$ , we can approximate the FPE by setting  $x_0 \cong x$  ( $n = 0$ ) and  $x_n \cong 2na$  ( $n \neq 0$ ). Then, rewriting the infinite sum as a Riemann integral and using the continuous time approximation, we finally get

$$\frac{\partial \rho_t(x)}{\partial t} \cong \frac{\epsilon a^{\frac{1}{\gamma}-1}}{2\gamma} \left[ \int_{|x-x'| \geq a} \frac{\rho_t(x')}{|x-x'|^{1+\frac{1}{\gamma}}} dx' - \frac{2\gamma}{a^{\frac{1}{\gamma}}} \rho_t(x) \right]. \quad (10)$$

The first term on the right-hand side of Eq. (10) is the density, which flows into the interval  $-a \leq x \leq a$ , and the second term is the flow from the same interval. As mentioned above, this equation is extended to the whole real axis,  $-\infty < x < \infty$ . Integrating Eq. (10) over  $x$ , one can show that this equation satisfies the conservation of probability.

### 3 Analytical results

The purpose of this section is to derive the relation between the approximated FPE of Eq. (10) and Lévy diffusion equation given in Eq. (2).

We begin by transforming the Lévy diffusion equation into a form comparable with Eq. (10). Here we assume  $0 < \alpha < 2$ . The function  $D^\alpha(x)$  [Eq. (3)] can be expressed as [10]

$$D^\alpha(x) = \lim_{\epsilon \rightarrow 0} E^\alpha(x, \epsilon), \quad (11)$$

where the function  $E^\alpha(x, \epsilon)$  is given by

$$E^\alpha(x, \epsilon) = -\frac{\Gamma(\alpha+1)}{2\pi i} \left[ \frac{(-1)^{\frac{\alpha}{2}}}{(x+i\epsilon)^{\alpha+1}} - \frac{(-1)^{-\frac{\alpha}{2}}}{(x-i\epsilon)^{\alpha+1}} \right]. \quad (12)$$

Here we have assumed  $\arg(x \pm i\epsilon) \in (-\pi, \pi)$  and  $(-1)^\alpha = \exp(i\pi\alpha)$ . Combining Eqs. (2) and (11), let us exchange the order of the integral and the limit:

$$\begin{aligned} \frac{\partial \rho_t(x)}{\partial t} &= -b \lim_{\epsilon \rightarrow 0} \int_{-\infty}^{+\infty} E^\alpha(x-x', \epsilon) \rho_t(x') dx' \\ &= b \lim_{\epsilon \rightarrow 0} \frac{\Gamma(\alpha+1)}{2\pi i} \left[ \int_{-\infty}^{+\infty} \frac{(-1)^{\frac{\alpha}{2}} \rho_t(x')}{(x-x'+i\epsilon)^{\alpha+1}} dx' - \int_{-\infty}^{+\infty} \frac{(-1)^{-\frac{\alpha}{2}} \rho_t(x')}{(x-x'-i\epsilon)^{\alpha+1}} dx' \right]. \end{aligned}$$

After substituting  $x' = z + i\epsilon$  into the first integral and  $x' = z - i\epsilon$  into the

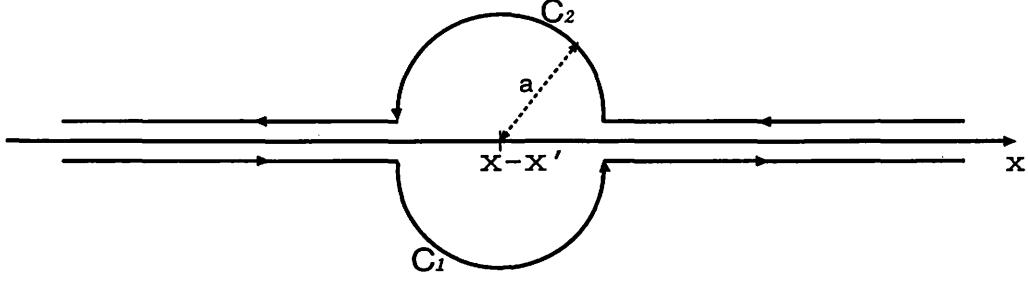


Figure 3: Path of integration.  $C_1$  is a semicircle in the lower half plane and  $C_2$  is a semicircle in the upper half plane. The radius of  $C_1$  and  $C_2$  is  $a$ .

second, we obtain

$$\frac{\partial \rho_t(x)}{\partial t} = b \lim_{\epsilon \rightarrow 0} \frac{\Gamma(\alpha + 1)}{2\pi i} \left[ \int_{-\infty - i\epsilon}^{+\infty - i\epsilon} \frac{(-1)^{\frac{\alpha}{2}} \rho_t(z)}{(x - z)^{\alpha+1}} dz - \int_{-\infty + i\epsilon}^{+\infty + i\epsilon} \frac{(-1)^{-\frac{\alpha}{2}} \rho_t(z)}{(x - z)^{\alpha+1}} dz \right],$$

where we have assumed that  $\rho_t(x)$  is an analytic function of  $z$ . Furthermore, we choose the path of integration as shown in Fig. 3:

$$\begin{aligned} \frac{\partial \rho_t(x)}{\partial t} &= b \frac{\Gamma(\alpha + 1) \sin(\frac{\pi}{2}\alpha)}{\pi} \int_{|x-x'| \geq a} \frac{\rho_t(x')}{|x - x'|^{\alpha+1}} dx' \\ &\quad + b \lim_{a \rightarrow 0} \frac{\Gamma(\alpha + 1)}{2\pi i} \left[ \int_{C_1} \frac{(-1)^{\frac{\alpha}{2}} \rho_t(z)}{(x - z)^{\alpha+1}} dz + \int_{C_2} \frac{(-1)^{-\frac{\alpha}{2}} \rho_t(z)}{(x - z)^{\alpha+1}} dz \right]. \end{aligned}$$

Finally, after evaluation of the integrals, we find

$$\frac{\partial \rho_t(x)}{\partial t} = \lim_{a \rightarrow 0} b \frac{\Gamma(\alpha + 1) \sin(\frac{\pi}{2}\alpha)}{\pi} \left[ \int_{|x-x'| \geq a} \frac{\rho_t(x')}{|x - x'|^{1+\alpha}} dx' - \frac{2\rho_t(x)}{\alpha a^\alpha} \right] \quad (13)$$

Comparing Eqs. (10) and (13), we find

$$\alpha = \frac{1}{\gamma}, \quad (14)$$

and

$$b = \frac{\epsilon a^{\frac{1}{\gamma}-1}}{2\gamma} \frac{\pi}{\Gamma(\frac{1}{\gamma} + 1) \sin(\frac{\pi}{2\gamma})}. \quad (15)$$

That is to say, the approximated FPE [Eq.(10)] converges to the Lévy diffusion equation in the limit  $a \rightarrow 0$ .

As seen from Eq. (14), the Lévy index  $\alpha$  is characterized by only the parameter  $\gamma$  in the map  $f_\gamma(x)$ , and the scale factor  $b$  is also determined by Eq. (15). These results are valid only to the case  $0 < \alpha < 2$  (i.e.,  $\gamma > 0.5$ ).

## 4 Summary

Here we introduced a one-dimensional map that generates arbitrarily long flights and found relation between this map and symmetric Lévy flight. We derived an approximated Frobenius-Perron equation and showed that this equation converges to the Lévy diffusion equation for  $\gamma > 0.5$ .



## References

- [1] M. F. Shlesinger, G. M. Zaslavsky and J. Klafter, *Nature (London)* **363** (1993), 31.
- [2] H. Takayasu, *Prog. Theor. Phys.* **72** (1984), 471.
- [3] R. N. Mantegna and H. E. Stanley, in *Lévy flights and related topics in physics*, ed. M.F. Shlesinger et al. (Springer-Verlag, New York, 1995), p. 300.
- [4] P. Lévy, *Théory de l'Addition des Variables Aléatoires* (Gauthier-Villars, Paris, 1937).
- [5] V. Seshadri and B. J. West, *Proc. Natl. Acad. Sci. USA.* **79** (1982), 4501.
- [6] W. Feller, *An Introduction to Probability Theory and its Applications*, 2nd ed. (John Wiley & Sons, New York, 1971), vol.II.
- [7] T. Geisel and S. Thomae, *Phys. Rev. Lett.* **52** (1984), 1936.
- [8] S. Grossmann and H. Fujisaka, *Phys. Rev. A* **26** (1982), 1779.
- [9] J. P. Bouchaud and A. Georges, *Phys. Rep.* **195** (1990), 127.
- [10] I. M. Gel'fand, *Generalized functions* (Academic Press, New York, 1964).
- [11] T. Miyaguchi and Y. Aizawa, *Prog. Theor. Phys.* **106** (2001), 697.

# Dynamical Ordering of Non-Birkhoff Orbits and Topological Entropy in the Standard Mapping

Yoshihiro YAMAGUCHI<sup>1</sup> and Kiyotaka TANIKAWA<sup>2</sup>

<sup>1</sup> Teikyo Heisei University, Ichihara, Chiba 290-0193, Japan.

<sup>2</sup> National Astronomical Observatory, Mitaka, Tokyo 181-8588, Japan.

## Abstract

The standard mapping is an analytical, reversible monotone twist mapping. The appearance ordering, i.e., the so called dynamical ordering, of the symmetric non-Birkhoff periodic orbits (SNBO) in the standard mapping is derived. Essential use is made of the reversibility. After the establishment of various properties of the symmetry axes under the mapping, two theorems on the dynamical ordering are proved. Then the braids for SNBOs are constructed with the aid of techniques developed in the braid group theory. The lower bound of the topological entropy of the system possessing an SNBO is estimated by the eigenvalue of the reduced Burau matrix representation of the braid constructed from the SNBO. Behavior of the topological entropy in the integrable limit is discussed.

# 1 Introduction

The standard mapping  $T$  is defined in cylinder.

$$y_{n+1} = y_n + af(x_n), \quad (1)$$

$$x_{n+1} = x_n + y_{n+1} \pmod{2\pi}, \quad (2)$$

where  $a$  is a positive parameter and  $f(x) = \sin x$ . There are two fixed points  $P = (0, 0)$  and  $Q = (\pi, 0)$ , where  $P$  is a saddle and  $Q$  is an elliptic point ( $0 < a < 4$ ) or a saddle with reflection ( $a > 4$ ). For convenience, we call a point  $(2\pi, 0)$  a saddle  $P'$ .

As we have repeatedly exploited its property in various occasions,<sup>1,2,14,20,22,23</sup> the standard mapping belongs to a class of systems possessing reversibility (see §2.3). All orbits are classified either into symmetric or non-symmetric ones. The standard mapping belongs to a class of systems which are called monotone twist. All orbits are classified either into monotone or non-monotone ones (§2.2). Symmetric monotone periodic orbits exist down to the integrable limit. In a sense, these objects are not interesting. In order to study the chaotic behavior of the system or how the system becomes more chaotic with parameter, it seems necessary to study other types of periodic orbits. In the standard mapping, symmetric non-monotone periodic orbits are quite suitable for this purpose. There exist non-symmetric periodic orbits bifurcated from symmetric ones.<sup>22</sup>

The appearance order of periodic orbits is called the *dynamical ordering*. The most famous ordering has been proved by Sharkovskii<sup>3</sup>) in one-dimensional mappings. Extensions to two-dimensional mappings and to systems described by ordinary differential equations have been carried out by many authors.<sup>4)-10)</sup> In the preceding papers,<sup>1),2)</sup> we studied the dynamical ordering of the symmetric non-Birkhoff orbits in the standard mapping and its family mappings, and in the forced oscillator. In this paper, we extend the dynamical ordering derived in these papers, construct braids for periodic orbits in the dynamical ordering, and estimate the topological entropy of the standard mapping at parameter values for which there are periodic orbits whose braid types are determined. The appearance of non-Birkhoff orbits is related to the non-integrability of systems,<sup>1),2)</sup> and to the breakup of KAM (Kolmogorov–Arnold–Moser) invariant curves.<sup>11),12)</sup>

In §2, we provide several useful concepts and notation used in the following sections. In §3, the properties needed in the proof of theorems are proved. The theorems on the dynamical orderings are proved in §4. The braid for non-Birkhoff orbit listed in the dynamical ordering is constructed in §5 and the topological entropy is estimated in §6. In the final section, we point out several problems to be solved.

## 2 Preliminaries

### 2.1 Notation

We have defined the mapping in cylinder. In sections 3 and 4, we almost always work in universal cover  $\mathbf{R}^2$  and use the lift mapping. We move from cylinder to universal cover

and vice versa when we construct braids of the non-Birkhoff periodic orbits in section 5. We lift the mapping  $T$  to universal cover in such a way that the fixed points of  $T$  are also fixed under the lift. Then the lift is uniquely determined. To avoid the notational complexity, we use the same notation  $T$  for the lift mapping and coordinates  $(x, y)$  for universal cover. We will make a note if there may be a possibility of confusion.

The orbit of a point  $z \in \mathbf{R}^2$  is denoted  $o(T, z) = \{\dots, T^{-1}z, z, Tz, \dots\}$ . Following Boyland and Hall,<sup>11)</sup> we define the extended orbit of a point  $z \in \mathbf{R}^2$  by

$$eo(T, z) = \{T^k z + (2\pi l, 0) : k, l \in \mathbf{Z}\}. \quad (3)$$

We usually abbreviate  $o(T, z)$  and  $eo(T, z)$  as  $o(z)$  and  $eo(z)$ . Let  $\pi_1(z)$  (resp.  $\pi_2(z)$ ) be projection to the  $x$ -coordinate (resp.  $y$ -coordinate) of  $z$ . The rotation number  $\nu$  of an orbit of  $z \in \mathbf{R}^2$  is defined as

$$\nu = \limsup_{n \rightarrow \infty} \frac{\pi_1(T^n z) - \pi_1(z)}{n}. \quad (4)$$

## 2.2 Birkhoff and Non-Birkhoff periodic orbits

A point  $z \in \mathbf{R}^2$  is called a  $p/q$ -periodic point for the standard mapping  $T : \mathbf{R}^2 \rightarrow \mathbf{R}^2$  if

$$T^q z - (2\pi p, 0) = z.$$

A  $p/q$ -periodic point  $z \in \mathbf{R}^2$  is called Birkhoff by Hall<sup>15)</sup> if for any  $r, s \in eo(z)$

$$\pi_1(r) < \pi_1(s) \Rightarrow \pi_1(Tr) < \pi_1(Ts). \quad (5)$$

Otherwise, the point is said to be non-Birkhoff. A Birkhoff (resp. non-Birkhoff) periodic point is abbreviated as a BP (resp. NBP). Corresponding orbits are denoted as a BO and an NBO.

We give a little bit more precise definition for non-Birkhoff periodic points or orbits. If we lift a  $p/q$ -periodic orbit in cylinder to universal cover, there are  $p$  different orbits corresponding to the original one. In fact, a  $p/q$ -periodic orbit in cylinder has  $q$  points in one period. In universal cover, the orbit cover  $p$  copies of cylinder in one period. In these copies, there are  $q \times p$  points. Therefore, we need  $p$  different  $p/q$ -periodic orbits. Taking this into account, let us define two types of non-Birkhoff periodic points (NBP).

**Definition.** Suppose a  $p/q$ -periodic point  $z \in \mathbf{R}^2$  is given.

(1) If for some  $r, s \in eo(z)$  with  $r \in o(s)$ ,

$$\pi_1(r) < \pi_1(s) \Rightarrow \pi_1(Tr) \geq \pi_1(Ts), \quad (6)$$

then, point  $z$  is called a non-Birkhoff periodic point of Type I.

(2) If for any  $r, s \in eo(z)$  with  $r \in o(s)$ ,

$$\pi_1(r) < \pi_1(s) \Rightarrow \pi_1(Tr) < \pi_1(Ts) \quad (7)$$

and if for some  $r', s' \in eo(z)$  with  $r' \notin o(s')$ ,

$$\pi_1(r') < \pi_1(s') \Rightarrow \pi_1(Tr') \geq \pi_1(Ts'), \quad (8)$$

then, point  $z$  is called a non-Birkhoff periodic point of Type II.

**Remarks.**

1) Obviously, if a  $1/q$ -periodic orbit is non-Birkhoff, it is always of Type I. Both Types I and II are possible for  $p/q$ -NBOs ( $p \geq 2, q \geq 2$ ).

2) A point  $r \in eo(z)$  satisfying

$$\pi_1(T^{-1}r) < \pi_1(r), \pi_1(r) \geq \pi_1(Tr), \text{ or} \quad (9)$$

$$\pi_1(T^{-1}r) > \pi_1(r), \pi_1(r) \leq \pi_1(Tr), \quad (10)$$

is called a turning-back point or simply a turning point. There is an even number of turning-back points in universal cover in an NBO of Type I. There are no turning-back points in an NBO of Type II. Points monotonically proceed to the same direction in the  $x$  coordinate under the mapping.

In this paper, we restrict our attention to the first type. An example of NBOs of the second type has been displayed in Refs. 6) and 8).

## 2.3 Symmetry axes and symmetric periodic orbits

A mapping is reversible if it is decomposed into a product of two involutions.<sup>16)</sup> Since the standard mapping is doubly reversible,<sup>14)</sup> there are two forms of the product. Roughly speaking, the first one represents the left-right symmetry, i.e., symmetric points are disposed rather horizontally, and the second one the up-down symmetry, i.e., symmetric points are disposed rather vertically. In this paper, we mainly use the first form. Thus  $T$  is expressed as the product of involutions  $H$  and  $G$ .

$$T = H \circ G, \quad (11)$$

$$G : (x, y) \leftrightarrow (-x, y + af(x)) \pmod{2\pi}, \quad (12)$$

$$H : (x, y) \leftrightarrow (-x + y, y) \pmod{2\pi}. \quad (13)$$

where  $G^2 = id = H^2$  and  $\det \nabla G = \det \nabla H = -1$ . The sets of fixed points of  $G$  and  $H$  are the symmetry axes. In universal cover, the symmetry axes are expressed in the forms

$$S_1^{(m)} : x = 2\pi m, \quad (14)$$

$$S_2^{(m)} : x = 2\pi m + \pi, \quad (15)$$

$$S_3^{(m)} : y = 2(x - 2m\pi), \quad (16)$$

$$S_4^{(m)} : y = 2(x - (2m + 1)\pi) \quad (17)$$

where  $m(-\infty < m < \infty)$  is an integer. We denote  $S_1 = S_1^0$ ,  $S_2 = S_2^0$ ,  $S_3 = S_3^0$ , and  $S_4 = S_4^0$ .  $S_1$  and  $S_2$  are symmetry axes of  $G$ , whereas  $S_3$  and  $S_4$  are symmetry axes of  $H$ . In order to specify a branch of symmetry axis with  $y > 0$  (resp.,  $y < 0$ ), we add a suffix  $+$  ( $-$ ) to the expression of axis.

A periodic orbit which has points in symmetry axes is called symmetric. A  $p/q$ -periodic symmetric Birkhoff (resp. non-Birkhoff) orbit is denoted by an  $p/q$ -SBO (resp. an  $p/q$ -SNBO). A point of an SBO (resp. an SNBO) is denoted by an SBP (resp. an SNBP). Let  $\{p_0, p_1, \dots, p_{q-1}\}$  be a set of points from one period of an  $p/q$ -SBO or an  $p/q$ -SNBO. We summarize the relation of these points and the symmetry axes in Table I.<sup>(12), (16)</sup>

Table I. The relation of the symmetry axes and the symmetric periodic orbits.

$p(\geq 0)$	$q$	$p_0$	Transit	$p(\geq 0)$	$q$	$p_0$	Transit
Odd	$2k$	$S_1$	$p_k \in S_2$	Odd	$2k+1$	$S_1$	$p_{k+1} \in S_4$
Odd	$2k$	$S_2$	$p_k \in S_1$	Odd	$2k+1$	$S_2$	$p_{k+1} \in S_3$
Odd	$2k$	$S_3$	$p_k \in S_4$	Odd	$2k+1$	$S_3$	$p_k \in S_2$
Odd	$2k$	$S_4$	$p_k \in S_3$	Odd	$2k+1$	$S_4$	$p_k \in S_1$
Even	$2k$	$S_1$	$p_k \in S_1$	Even	$2k+1$	$S_1$	$p_{k+1} \in S_3$
Even	$2k$	$S_2$	$p_k \in S_2$	Even	$2k+1$	$S_2$	$p_{k+1} \in S_4$
Even	$2k$	$S_3$	$p_k \in S_3$	Even	$2k+1$	$S_3$	$p_k \in S_1$
Even	$2k$	$S_4$	$p_k \in S_4$	Even	$2k+1$	$S_4$	$p_k \in S_2$

Here we present an easy way to distinguish SNBPs from SBPs and distinguish two types of SNBPs. Let us suppose that a  $p/q$ -periodic point is bifurcated in a symmetry axis. We claim that this is a non-Birkhoff periodic point. Indeed, let us suppose the point is bifurcated in either  $S_1$  or  $S_2$ . We already have an  $p/q$ -SBP in the axis. This means that the whole  $p/q$ -periodic orbits are not in the graph of a Lipschitz function,<sup>(4), (6)</sup> which implies that new born point is a non-Birkhoff point. Next, let us take a positive rotation number and suppose a point  $r$  is bifurcated in either  $S_{3+}$  or  $S_{4+}$ . Let  $s$  be the Birkhoff point in the axis. Assume that  $r$  and  $s$  satisfy the relation  $\pi_1(r) < \pi_1(s)$ . One easily confirms that  $\pi_1(T^{-1}r) > \pi_1(T^{-1}s)$ . This means either  $r$  or  $s$  is non-Birkhoff. Then  $r$  is non-Birkhoff. For a negative rotation number we take points in  $S_{3-}$  or  $S_{4-}$  and argue in a similar manner.

Thus we have proved

**Proposition 2.1.** Symmetric periodic points bifurcated in symmetry axes and stay there are SNBPs.

Next let us see, using examples, the geometrical difference of SNBPs of Types I and II. Suppose that  $T^2r \in T^2S_{1+}^{(0)} \cap S_{4-}^{(0)}$ . Then  $r \in S_{1+}$  is an  $1/3$ -SNBP of Type I, since  $\pi_1(r) < \pi_1(T^2r)$  and  $\pi_1(T^2r) > \pi_1(r)$ . Suppose that  $T^2r \in T^2S_{1+}^{(0)} \cap S_{3+}^{(1)}$  and this point is bifurcated via saddle-node bifurcation. Then  $r \in S_{1+}^0$  is an  $2/3$ -SNBP of Type II, since  $\pi_2(r), \pi_2(T^2r), \pi_2(T^2r)$  are all positive and there are no turning-back points.

### 2.4 Stable and unstable manifolds

Two saddles  $P$  and  $P'$  have stable and unstable manifolds. Let  $W_u^1$  be a branch of unstable manifold starting at  $P$  toward the right direction and  $W_s^1$  be a branch of stable manifold going toward  $P'$  from the left side in universal cover. Let  $W_u^2$  be a branch of unstable manifold starting at  $P'$  and  $W_s^2$  be a branch of stable manifold going toward  $P$  (see Fig. 1).

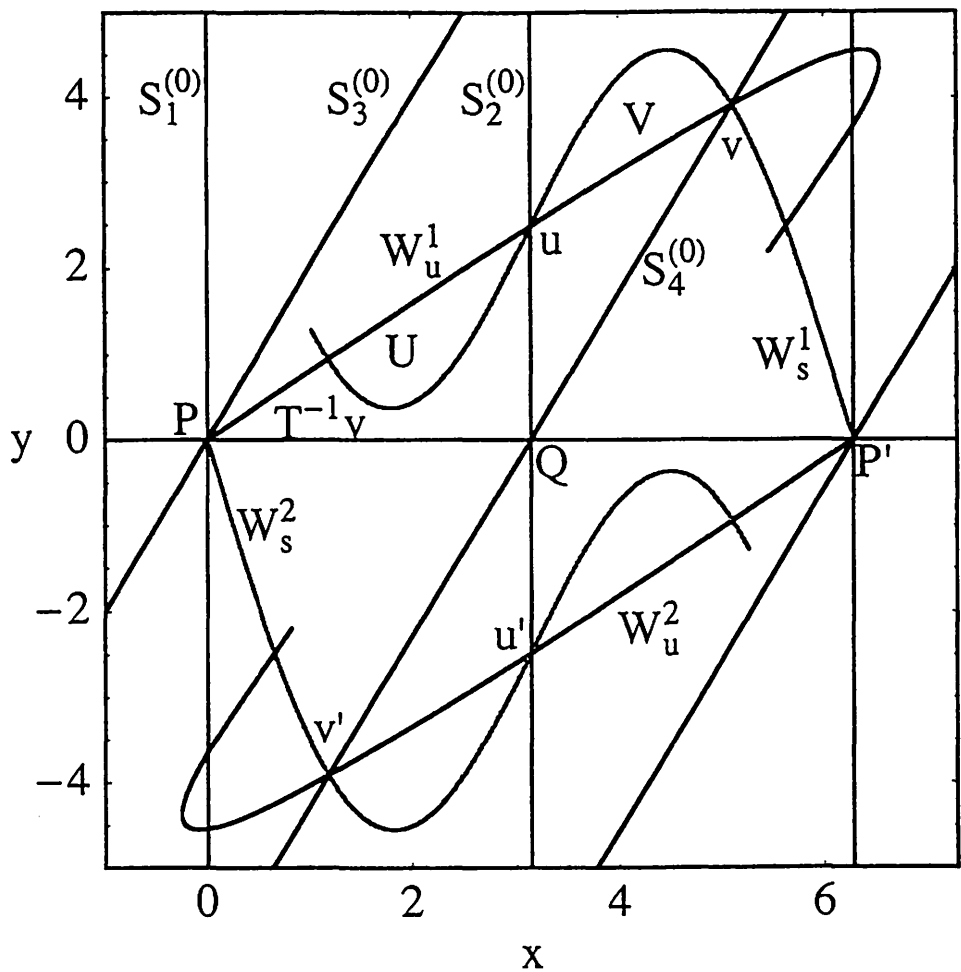


Fig. 1. Structure of stable and unstable manifolds in the standard mapping at  $a = 3.2$ . The symmetry axes  $S_i^{(0)}$  ( $i = 1, \dots, 4$ ) are also displayed.

The existence of transverse intersection of  $W_u^1$  and  $W_s^1$  at  $u$  has been proved,<sup>18)</sup> where  $u$  is an intersection point of  $W_u^1$  and  $S_{2+}^{(0)}$ . The slope of  $W_s$  at  $u$  is larger than that of  $W_u$ . Due to the up-down symmetry, two unstable manifolds  $W_u^1$  and  $W_u^2$  are symmetric with respect to  $Q$ . This is the case also for  $W_s^1$  and  $W_s^2$ . Hence  $W_u^2$  and  $W_s^2$  intersect transversely at  $u'$  where  $u'$  is the symmetrical point of  $u$  with respect to  $Q$ . Due to the

reversibility, the stable and unstable manifolds are related by

$$GW_u^i = W_s^i, \quad (18)$$

$$HW_u^i = W_s^i, \quad (19)$$

with  $i = 1, 2$ .

Let  $r, s \in W_u$  where  $r$  is closer to  $P$  than  $s$  is along  $W_s$ . A closed arc of  $W_u$  between  $r$  and  $s$  will be written as  $[r, s]_{W_u}$ . Open and semi-open arcs are defined similarly. Arcs of other manifolds are defined in a similar manner. Points  $u$  and  $v$ , and their forward and backward iterates are the primary homoclinic points.<sup>19)</sup> Let  $\gamma_u = [u, v]_{W_u}$  and  $\gamma_s = [u, v]_{W_s}$ . Let  $V$  be an open region bounded by  $\gamma_u$  and  $\gamma_s$  and  $U$  an open region bounded by  $[T^{-1}v, u]_{W_u}$  and  $[T^{-1}v, u]_{W_s}$ . These are primary homoclinic lobes. Due to Eq. (18), two lobes  $U$  and  $V$  are related by

$$U = GV. \quad (20)$$

We define the open intervals in the symmetry axes.

$$I_n^{(-m)} = T^{-n}V \cap S_{1+}^{(-m)} \quad (m \geq 0, n \geq 1), \quad (21)$$

$$J_n^{(-m)} = T^{-n}V \cap S_{2+}^{(-m)} \quad (m \geq 1, n \geq 1), \quad (22)$$

$$K_n^{(-m)} = T^{-n}V \cap S_{3+}^{(-m)} \quad (m \geq 0, n \geq 0), \quad (23)$$

$$L_n^{(-m)} = T^{-n}V \cap S_{4+}^{(-m)} \quad (m \geq 1, n \geq 0). \quad (24)$$

These may be empty for a certain range of parameter values. In Fig. 2, several intervals are displayed. We further introduce

$$I^{(-m)} = \cup_{n \geq 1} I_n^{(-m)}, \quad J^{(-m)} = \cup_{n \geq 1} J_n^{(-m)}, \quad (25)$$

$$K^{(-m)} = \cup_{n \geq 0} K_n^{(-m)}, \quad L^{(-m)} = \cup_{n \geq 0} L_n^{(-m)}. \quad (26)$$



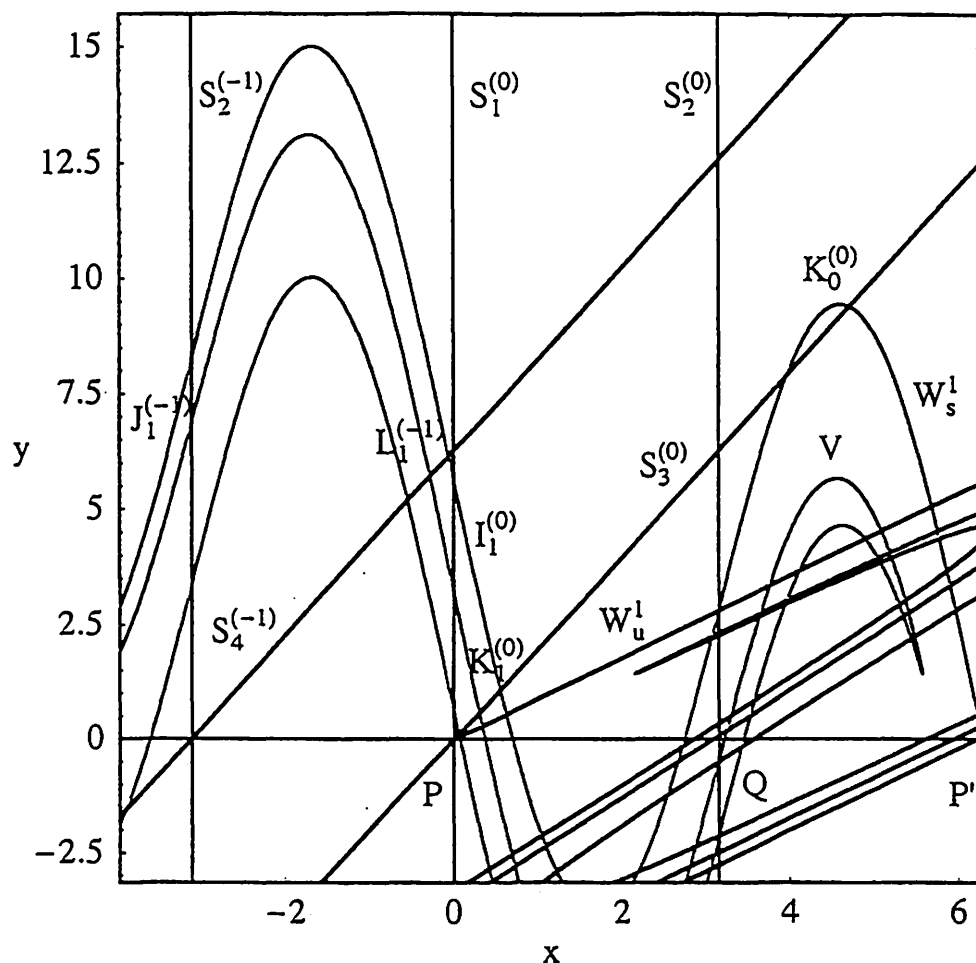


Fig. 2. Intervals  $I_1^{(0)}$ ,  $J_1^{(-1)}$ ,  $K_0^{(0)}$ ,  $K_1^{(0)}$  and  $L_1^{(-1)}$  are observed at  $a = 8$ .

We write

$$p/q \in I_n^{(-m)} \quad (27)$$

if an  $p/q$ -SNBP exists in  $I_n^{(-m)}$ . We use the same notation for other intervals. We define the critical parameter values.

[1]  $a_c(p/q \in I_n^{(-m)}) = \inf\{a > 0 : p/q \in I_n^{(-m)}\}$ .

[2]  $a_c(I_n^{(-m)}) = \inf\{a > 0 : I_n^{(-m)} \neq \emptyset\}$ .

The critical values for SNBPs in other intervals are similarly defined.

### 3 Properties

**Property 3.1.** Let  $r \in V$  and  $s \in U$ . Then we have

$$\pi < \pi_1(r) < 2\pi, \quad (28)$$

$$0 < \pi_1(s) < \pi. \quad (29)$$

**Proof.** We have  $\pi_1(u) = \pi$  and  $\pi < \pi_1(v) < 2\pi$ . Thus  $V$  is located in the region between  $x = \pi$  and  $x = 2\pi$ . The first relation is proved. Equation (20) implies the second one. Q.E.D.

Let  $D$  be an open region bounded by  $B_1 = [Q, u]_{S_2^{(0)}}$ ,  $B_2 = [u, v]_{W_1^1}$ ,  $B_3 = [v, P']_{W_1^1}$ , and  $B_4 = [Q, P']$  in the  $x$  axis, and let  $Z$  be an open region outside  $D$  satisfying  $y > 0$  and  $y < 2(x - \pi)$ .

**Property 3.2.**

$$TD \cap Z = \emptyset. \quad (30)$$

**Proof.**  $TB_1$  is a line segment of  $y = x - \pi$  in  $D$  and then  $TB_1 \cap Z = \emptyset$ .  $TB_3$  is  $[Tv, P']_{W_1^1}$ . One easily confirms  $\pi_2(Ts) < 0$  for any  $s \in B_4 \setminus \{Q, P'\}$ . So  $TB_4 \cap Z = \emptyset$ . Since  $TB_2 \cap TB_4 = \emptyset$  holds,  $TV$  is located in the region surrounded by  $TB_1$ ,  $[Tu, P']_{W_1^1}$ , and  $TB_4$ . Thus we have  $TB_2 \cap Z = \emptyset$ . Q.E.D.

**Property 3.3.** Any periodic point in primary homoclinic lobe  $V$  is an NBP of Type I.

**Proof.** Let  $z \in V$  be a  $p/q$ -periodic point. The relation  $\pi_1(z) - \pi_1(T^{-1}z) = \pi_2(z) > 0$  holds. The orbit of  $z$  goes forward from  $T^{-1}z$  to  $z$ . Thus the existence of turning-back is obvious if  $p \leq 0$ . So we take  $p > 0$ . A point of  $TV$  is either in  $D$  or below the  $x$ -axis. A point below the  $x$ -axis is a turned-back point. A point of  $TV \subset D$  under one more iterate decreases its  $y$ -coordinate but never go into  $Z$  by Property 3.2. Since  $z$  is a periodic point, its finite iterate goes out of  $D$ . It necessarily goes below the  $x$ -axis. Therefore in any case, the orbit of  $z$  has a turning-back point. Q.E.D.

**Property 3.4.**  $I_n^{(-m)}$  consists of a unique component. The same is true for other intervals.

**Proof.** Let  $\gamma_s = [u, v]_{W_1^1}$ . The slope of the graph of  $\gamma_s$  strictly decreases as  $x$  increases.<sup>20)</sup> The largest slope is at  $u$  and is greater than 1. So there is a unique point  $w \in \gamma_s$  at which the slope is 1. Then  $T^{-1}\gamma_s$  has a unique point  $T^{-1}w$  at which the slope diverges. This implies that there is only one component in  $I_1^{(0)}$  if  $I_1^{(0)}$  exists at all.

If  $I_2^{(0)}$  has two components, there exist at least three points in the arc  $T^{-2}\gamma_s$  at which the slope diverges. Then  $T^{-1}\gamma_s$  has at least three points at which the slope is equal to 1, and thus it has at least three points at which the slope diverges. This contradicts the above mentioned property of  $T^{-1}\gamma_s$ . Repeating this procedure, the statement is proved. Q.E.D.

**Property 3.5.** The critical values satisfy the following relations.

$$a_c(I_{n+1}^{(-m)}) < a_c(I_n^{(-m)}), \quad a_c(J_{n+1}^{(-m)}) < a_c(J_n^{(-m)}) \quad (m \geq 0, n \geq 1), \quad (31)$$

$$a_c(K_{n+1}^{(-m)}) < a_c(K_n^{(-m)}), \quad a_c(L_{n+1}^{(-m)}) < a_c(L_n^{(-m)}) \quad (m \geq 0, n \geq 0), \quad (32)$$

$$a_c(I_n^{(-m)}) < a_c(I_n^{(-m-1)}), \quad a_c(J_n^{(-m)}) < a_c(J_n^{(-m-1)}) \quad (m \geq 0, n \geq 1), \quad (33)$$

$$a_c(K_n^{(-m)}) < a_c(K_n^{(-m-1)}), \quad a_c(L_n^{(-m)}) < a_c(L_n^{(-m-1)}) \quad (m \geq 0, n \geq 0), \quad (34)$$

$$a_c(J_n^{(-m-1)}) > a_c(L_n^{(-m-1)}) > a_c(I_n^{(-m)}) > a_c(K_n^{(-m)}) \quad (m \geq 0), \quad (35)$$

$$a_c(K_n^{(-m)}) > a_c(I_{n+1}^{(-m)}) \quad (m \geq 0), \quad (36)$$

$$a_c(L_n^{(-m)}) > a_c(J_{n+1}^{(-m)}) \quad (m \geq 0). \quad (37)$$

**Proof.** Equations (31)–(35) follow from the construction of intervals and the lambda lemma.<sup>21)</sup> We shall give a proof of Eq. (36). Let us increase the parameter and suppose  $T^{-1}[u, v]_{W_1^1}$  touches  $S_{3+}^{(0)}$  for the first time at some  $a$ . There exist an interval in  $y = x(x \geq 0)$  such that two end points are intersection points of  $T^{-1}[u, v]_{W_1^1}$  and  $y = x(x \geq 0)$ . By definition, the backward image of this interval is  $I_2^{(0)}$ . Thus we have the relation  $a_c(K_1^{(0)}) > a_c(I_2^{(0)})$ . Repeating this, Eq. (36) is proved. The proof for Eq. (37) is similar. Q.E.D.

**Property 3.6.**

$$\lim_{n \rightarrow \infty} a_c(R_n^{(-m)}) = 0 \quad (38)$$

where  $R = \{I, J\}$  ( $m \geq 0, n \geq 1$ ) and  $R = \{K, L\}$  ( $m \geq 0, n \geq 0$ ).

**Proof.** According to the lambda lemma, the intersection of  $W_{u(s)}^1$  and  $S_{1+}^{(-m)}$  exists for any small  $a > 0$ . This implies Eq. (38). Q.E.D.

**Property 3.7.**

$$\lim_{n \rightarrow \infty} T^n R_n^{(-m)} = \gamma_u \quad (39)$$

where  $R$  is defined in Property 3.6.

**Proof.** For simplicity let us only consider the case  $R = I$  and  $m = 0$ . The proof for other  $R$  and  $m$  is basically the same. Let  $\Gamma = [v, Tu]_{W_u^1}$  and let  $\Gamma^{(-k)}$  be the set of points of  $\Gamma$  shifted to the left by  $2k\pi, k \geq 1$ . Obviously  $\Gamma^{(-k)}$  is an arc of the unstable manifold emanating from a fixed point at  $(-2k\pi, 0)$ . For a given  $a > 0$ , there is a minimum  $k_0 \geq 1$  such that  $\Gamma^{(-k_0)} \cap S_{1+}^0 = \emptyset$ . Let  $\hat{\Gamma} = \Gamma^{(-k_0)}$ . By the lambda lemma,  $T^n \hat{\Gamma}$  has an arc arbitrarily close to  $[P, v]_{W_u^1}$  for sufficiently large  $n > 0$ .  $T^n I_n^0$  is sandwiched by  $T^n \hat{\Gamma}$  and  $[P, v]_{W_u^1}$ . Q.E.D.

**Property 3.8.** If  $I_n^{(-m)}$  exists at all, then it contains SNBPs of Type I. The same property holds for other intervals  $J_n^{(-m)}, K_n^{(-m)}$  and  $L_n^{(-m)}$ .

**Proof.** In view of Property 3.3, we need only to show that  $I_n^{(-m)}$  contains a periodic point, i.e.,  $T^k I_n^{(-m)}$  intersects another symmetry axis for some  $k > 0$ . By the lambda lemma, there exists a positive  $j_0$  such that  $T^j \gamma_u \cap S_{2-}^{(0)} \neq \emptyset$  for  $j \geq j_0$  and in addition, intersections are transverse. By Property 3.7,  $T^n I_n^{(-m)}$  is arbitrarily close to  $\gamma_u$  as a whole for sufficiently large  $n$ . Then one confirms that  $T^k I_n^{(-m)} \cap S_{2-}^{(0)} \neq \emptyset$  for  $k = n + j$  for sufficiently large  $j \geq j_0$ . The proofs for other cases are similar and omitted. Q.E.D.

Property 3.8 gives two relations for the critical values.

$$a_c(p/q \in R_n^{(-m)}) > a_c(R_n^{(-m)}), \quad (40)$$

$$\lim_{q \rightarrow \infty} a_c(p/q \in R_n^{(-m)}) = a_c(R_n^{(-m)}) \quad (41)$$

where  $R$  is defined in Property 3.6.

**Property 3.9.** The minimum period of SNBPs in  $I_n^{(-m)}$  or  $J_n^{(-m)}$  is  $2n + 1$  where  $n \geq 1$ . The minimum period of SNBPs in  $K_n^{(-m)}$  or  $L_n^{(-m)}$  is  $2n + 2$  where  $n \geq 0$ .

**Proof.** We know that  $T^n I_n^{(-m)} \subset V$ . At least one iteration is needed to go from  $V$  to  $U$ . Due to the symmetry with respect to  $G$ ,  $n$  iterations are needed to go from  $U$  to  $I_n^{(-m+1)}$ , which is  $I_n^{(-m)}$  shifted by  $2\pi$  toward the right direction. Thus the first statement is proved.

We know that  $T^n K_n^{(-m)} \subset V$ . At least one iteration is needed to go from  $V$  to  $U$ , and  $(n + 1)$  iterations are needed to go from  $U$  to  $K_n^{(-m+1)}$  due to the symmetry with respect to  $H$ . The second statement is proved. Q.E.D.

## 4 Dynamical ordering

There are many types of SNBOs with  $2n(n \geq 1)$  turning points if its period is large enough. For example, if the period is five, SNBOs with two and four turning points are possible. In the following, we restrict our attention to  $p/q$ -SNBOs ( $p \geq 0, q \geq 2$ ) with two turning points.

We introduce two symbols  $\rightarrow (\downarrow)$  and  $\leftrightarrow$ . Let  $R, R'$  be one of  $I, J, K$ , or  $L$ . If the existence of an  $p/q$ -SNBP in  $R_n^{(-m)}$  implies the existence of an  $p'/q'$ -SNBP in  $R_{n'}^{(-m')}$ , then we write as

$$\frac{p}{q} \in R_n^{(-m)} \rightarrow \frac{p'}{q'} \in R_{n'}^{(-m')}, \quad (42)$$

If both an  $p/q$ -SNBP in  $R_n^{(-m)}$  and an  $p'/q'$ -SNBP in  $R_{n'}^{(-m')}$  appear at the same value of  $a$ , then

$$\frac{p}{q} \in R_n^{(-m)} \leftrightarrow \frac{p'}{q'} \in R_{n'}^{(-m')}. \quad (43)$$

### 4.1 Dynamical ordering for $p/2$ - and $p/3$ -SNBOs

In the dynamical ordering derived in the next subsection,  $p/2$ - or  $p/3$ -SNBOs occupy special positions. We summarize their properties in this section.

Property 3.9 implies that an  $p/2$ -SNBP appears in  $K_0^{(-m)}$  and  $L_0^{(-m)}$ . A dynamical ordering among them proved in Ref. 22) is reproduced in Table II. For completeness sake, we cite the short proof.

Table II. Dynamical ordering for  $p/2$ -SNBPs.

$K_0^{(0)}$	$\leftarrow$	$L_0^{(-1)}$	$\leftarrow$	$K_0^{(-1)}$	$\leftarrow$	$L_0^{(-2)}$	$\leftarrow$
						0/2	$\leftarrow$
						$\downarrow$	$\nearrow$
				0/2	$\leftarrow$	1/2	$\leftarrow$
				$\downarrow$	$\nearrow$	$\downarrow$	$\nearrow$
		0/2	$\leftarrow$	1/2	$\leftarrow$	2/2	$\leftarrow$
		$\downarrow$	$\nearrow$	$\downarrow$	$\nearrow$	$\downarrow$	$\nearrow$
0/2	$\leftarrow$	1/2	$\leftarrow$	2/2	$\leftarrow$	3/2	$\leftarrow$
$\downarrow$	$\nearrow$	$\downarrow$	$\nearrow$	$\downarrow$	$\nearrow$	$\downarrow$	$\nearrow$
1/2	$\leftarrow$	2/2	$\leftarrow$	3/2	$\leftarrow$	4/2	$\leftarrow$

**Proof of Table II.** We estimate the critical values  $a_c(p/2 \in K_0^{(-m)})$  and  $a_c(p/2 \in L_0^{(-m)})$ , at which an  $p/2$ -SNBP appears in the corresponding interval.

$$a_c(p/2 \in K_0^{(-m)}) \approx 2\pi(3 + 4m - p) \quad (n \geq 0, 2m + 1 \geq p \geq 0), \quad (44)$$

$$a_c(p/2 \in L_0^{(-m)}) \approx 2\pi(1 + 4m - p) \quad (n \geq 1, 2m \geq p \geq 0). \quad (45)$$

The fact that the  $x$  coordinate of  $p/2$ -SNBP is approximately equal to  $3\pi/2$  is used. Equations (44) and (45) determine the dynamical ordering in Table II. Q.E.D.

We consider  $p/3$ -SNBPs appearing in  $I_0^{(-m)}$  and  $J_0^{(-m)}$ . The detailed discussion has been done in Ref. 23). We show in Table III the extended ordering including the case with  $p = 0$ .

Table III. Dynamical ordering for  $p/3$ -SNBPs.

$I_1^{(0)}$	$\leftarrow$	$J_1^{(-1)}$	$\leftarrow$	$I_1^{(-1)}$	$\leftarrow$	$J_1^{(-2)}$	$\leftarrow$
						0/3	$\leftarrow$
						$\downarrow$	
				0/3	$\leftarrow$	1/3	$\leftarrow$
				$\downarrow$		$\downarrow$	
		0/3	$\leftarrow$	1/3	$\leftarrow$	2/3	$\leftarrow$
		$\downarrow$		$\downarrow$		$\downarrow$	
0/3	$\leftarrow$	1/3	$\leftarrow$	2/3	$\leftarrow$	3/3	$\leftarrow$
$\downarrow$		$\downarrow$		$\downarrow$		$\downarrow$	
1/3	$\leftarrow$	2/3	$\leftarrow$	3/3	$\leftarrow$	4/3	$\leftarrow$

**Proof of Table III** We estimate the critical values  $a_c(p/3 \in I_1^{(-m)})$  and  $a_c(p/3 \in J_1^{(-m)})$ , at which an  $p/3$ -SNBP appears in the corresponding interval.

$$a_c(p/3 \in I_1^{(-m)}) \approx 2(3m - p + 9/4)\pi \quad (i \geq 0, 2m + 1 \geq p \geq 0), \quad (46)$$

$$a_c(p/3 \in J_1^{(-m)}) \approx 2(3m - p + 3/4)\pi \quad (i \geq 1, 2m \geq p \geq 0). \quad (47)$$

The fact that the  $x$  coordinate of  $p/3$ -SNBP is approximately equal to  $3\pi/2$  is used. Equations (46) and (47) determine the dynamical ordering in Table III. Q.E.D.

## 4.2 Two theorems

**Theorem 1.** The following dynamical ordering for  $p/q$ -SNBPs in  $I_n^{(-m)}$  and  $J_n^{(-m)}$  holds with  $0 \leq p \leq (2m + 1)$  for  $I_n^{(-m)}$ , and  $0 \leq p \leq 2m$  for  $J_n^{(-m)}$ .

$I_0^{(-m)}, J_0^{(-m)}$ :	$p/3$	$\rightarrow$	$p/4$	$\rightarrow$	$p/5$	$\rightarrow$	$p/6$	$\rightarrow$
	$\downarrow$		$\downarrow$		$\downarrow$		$\downarrow$	
$I_1^{(-m)}, J_1^{(-m)}$ :	$p/5$	$\rightarrow$	$p/6$	$\rightarrow$	$p/7$	$\rightarrow$	$p/8$	$\rightarrow$
	$\downarrow$		$\downarrow$		$\downarrow$		$\downarrow$	
$I_2^{(-m)}, J_2^{(-m)}$ :	$p/7$	$\rightarrow$	$p/8$	$\rightarrow$	$p/9$	$\rightarrow$	$p/10$	$\rightarrow$
	$\downarrow$		$\downarrow$		$\downarrow$		$\downarrow$	

**Proof.** The conditions for  $p$  are determined by Eqs (46) and (47). We prove the cases for SNBPs in  $I_n^{(-m)}$ . The proof for SNBPs in  $J_n^{(-m)}$  is similar, and thus is omitted.

(1) Proof of  $p/k \in I_n^{(-m)} \rightarrow p/(k+1) \in I_n^{(-m)} (k \geq 2n+1)$ .

(1-1) Both  $p$  and  $k$  are odd.

The assumption  $p/k \in I_n^{(-m)}$  implies the relation  $T^{(k+1)/2} I_n^{(-m)} \cap S_{4-}^{-(m-(p-1)/2)} \neq \emptyset$ . Relation  $T^{(k+1)/2} I_n^{(-m)} \cap S_{2-}^{-(m-(p-1)/2)} \neq \emptyset$  follows from the relative positions of  $S_{2-}^{-(m-(p-1)/2)}$  and  $S_{4-}^{-(m-(p-1)/2)}$ . In fact,  $S_{4-}^{-(m-(p-1)/2)}$  is located to the left side of  $S_{2-}^{-(m-(p-1)/2)}$ . The intersection points are those of SNBO with  $\nu = p/(k+1)$  starting from  $I_n^{(-m)}$ .

(1-2)  $p$  is odd and  $k$  is even.

The assumption implies the relation  $T^{k/2} I_n^{(-m)} \cap S_{2-}^{-(m-(p-1)/2)} \neq \emptyset$ . The intersection points are mapped under  $T$  to the left of  $S_{4-}^{-(m-(p-1)/2)}$  (see Fig. 3). This implies  $T^{(k+2)/2} I_n^{(-m)} \cap S_{4-}^{-(m-(p-1)/2)} \neq \emptyset$ .

(1-3)  $p$  is even and  $k$  is odd.

The proof is similar to (1-1), and thus is omitted.

(1-4) Both  $p$  and  $k$  are even.

The proof is similar to (1-2), and thus is omitted.

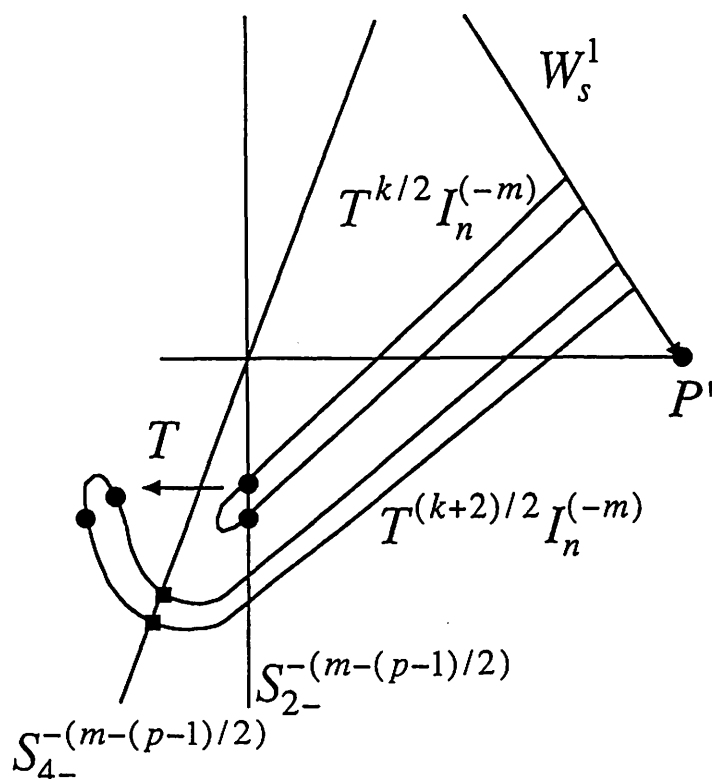


Fig. 3. Disposition of  $S_{2-}^{-(m-(p-1)/2)}$  and  $S_{4-}^{-(m-(p-1)/2)}$  and that of  $T^{k/2} I_n^{(-m)}$  and  $T^{(k+2)/2} I_n^{(-m)}$ .

(2) Proof of  $p/k \in I_n^{(-m)} \rightarrow p/(k+2) \in I_{n+1}^{(-m)}$  ( $k \geq 2n+1$ ).

(2-1) Both  $p$  and  $k$  are odd.

The assumption  $p/k \in I_n^{(-m)}$  implies  $T^{(k+1)/2} I_n^{(-m)} \cap S_{4-}^{-(m-(p-1)/2)} \neq \emptyset$ . The relative positions of  $T^{(k+1)/2} I_n^{(-m)}$  and  $T^{(k+3)/2} I_{n+1}^{(-m)}$  are displayed in Fig. 4.  $T^{(k+3)/2} I_{n+1}^{(-m)}$  is located outside the closed area bounded by  $T^{(k+1)/2} I_n^{(-m)}$  and an arc of  $W_s^1$ . We express this configuration  $T^{(k+3)/2} I_{n+1}^{(-m)} \cap T^{(k+1)/2} I_n^{(-m)} = \emptyset$ . This implies  $T^{(k+3)/2} I_{n+1}^{(-m)} \cap S_{4-}^{-(m-(p-1)/2)} \neq \emptyset$ . The intersection points are those of  $p/(k+2)$ -SNBOs starting in  $I_{n+1}^{(-m)}$ .

For the other cases, the proofs are similar to (2-1), and thus are omitted. Q.E.D.

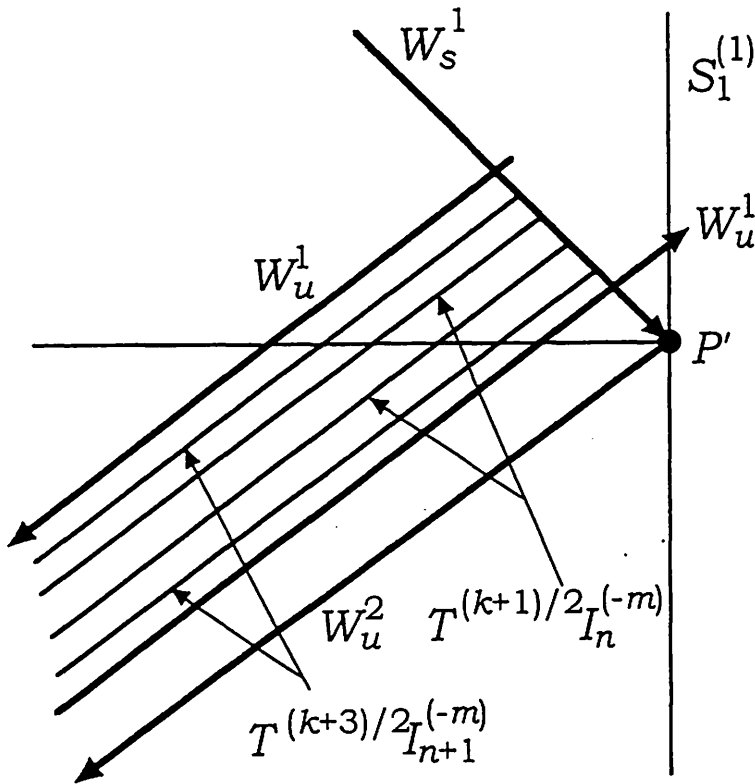


Fig. 4. Relation of  $T^{(k+1)/2} I_n^{(-m)}$  and  $T^{(k+3)/2} I_{n+1}^{(-m)}$ .



**Theorem 2.** The following dynamical ordering for  $p/q$ -SNBPs in  $K_n^{(-m)}$  and  $L_n^{(-m)}$  holds where  $0 \leq p \leq 2m + 1$  for  $K_n^{(-m)}$  and  $0 \leq p \leq 2m$  for  $L_n^{(-m)}$ .

$K_0^{(-m)}, L_0^{(-m)}:$	$p/2$	$\rightarrow$	$p/3$	$\rightarrow$	$p/4$	$\rightarrow$	$p/5$	$\rightarrow$
	$\downarrow$		$\downarrow$		$\downarrow$		$\downarrow$	
$K_1^{(-m)}, L_1^{(-m)}:$	$p/4$	$\rightarrow$	$p/5$	$\rightarrow$	$p/6$	$\rightarrow$	$p/7$	$\rightarrow$
	$\downarrow$		$\downarrow$		$\downarrow$		$\downarrow$	
$K_2^{(-m)}, L_2^{(-m)}:$	$p/6$	$\rightarrow$	$p/7$	$\rightarrow$	$p/8$	$\rightarrow$	$p/9$	$\rightarrow$
	$\downarrow$		$\downarrow$		$\downarrow$		$\downarrow$	

**Proof.** Equations (44) and (45) determine the conditions for  $p$ . We prove the cases for SNBPs in  $K_n^{(-m)}$ . The proof for  $L_n^{(-m)}$  is similar.

(1) Proof of  $p/k \in K_n^{(-m)} \rightarrow p/(k+1) \in K_n^{(-m)}$  ( $k \geq 2n+2$ ).

(1-1) Both  $p$  and  $k$  are odd.

The assumption  $p/k \in K_n^{(-m)}$  implies  $T^{(k-1)/2}K_n^{(-m)} \cap S_{2-}^{-(m-(p-1)/2)} \neq \emptyset$ . The intersection points are mapped to the left of  $S_{4-}^{-(m-(p-1)/2)}$ . This implies  $T^{(k+1)/2}K_n^{(-m)} \cap S_{4-}^{-(m-(p-1)/2)} \neq \emptyset$ .

(1-2)  $p$  is odd and  $k$  is even.

The assumption implies  $T^{k/2}K_n^{(-m)} \cap S_{4-}^{-(m-(p-1)/2)} \neq \emptyset$ .  $T^{k/2}K_n^{(-m)} \cap S_{2-}^{-(m-(p-1)/2)} \neq \emptyset$  follows from the disposition of  $S_{2-}^{-(m-(p-1)/2)}$  and  $S_{4-}^{-(m-(p-1)/2)}$ .

(1-3)  $p$  is even and  $k$  is odd.

The proof is similar to (1-1), and thus is omitted.

(1-4) Both  $p$  and  $k$  are even.

The proof is similar to (1-2), and thus is omitted.

(2) Proof of  $p/k \in K_n^{(-m)} \rightarrow p/(k+2) \in K_{n+1}^{(-m)}$  ( $k \geq 2n+2$ ).

(2-1) Both  $p$  and  $k$  are odd.

The assumption  $p/k \in K_n^{(-m)}$  implies  $T^{(k-1)/2}K_n^{(-m)} \cap S_{2-}^{-(m-(p-1)/2)} \neq \emptyset$  and  $T^{(k+1)/2}K_{n+1}^{(-m)} \cap T^{(k-1)/2}K_n^{(-m)}$ . This implies  $T^{(k+1)/2}K_{n+1}^{(-m)} \cap S_{2-}^{-(m-(p-1)/2)} \neq \emptyset$ .

For the other cases, the proofs are similar to (2-1), and thus are omitted. Q.E.D.

### 4.3 Accumulation of critical values

We now describe the behavior of critical values  $a_c(1/(2i+j) \in I_i^{(0)})$  and  $a_c(1/(2i+j) \in K_i^{(0)})$  as functions of  $i$  and  $j$ . We express the numerical results in a three-dimensional plot (see Figs. 5(a) and (b)). The maximum is  $a_c(1/3 \in I_1^{(0)})$  in Fig. 5(a) and  $a_c(1/2 \in K_0^{(0)})$  in Fig. 5(b). For a fixed  $i$ , the critical values decrease as  $j$  increases. The accumulation value is  $a_c(I_i^{(0)})$  (resp.,  $a_c(K_i^{(0)})$ ).

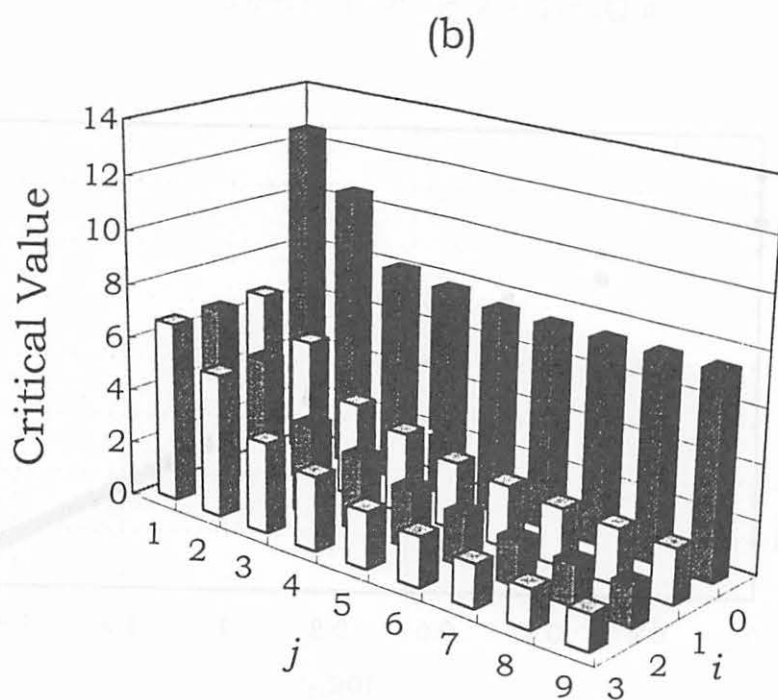
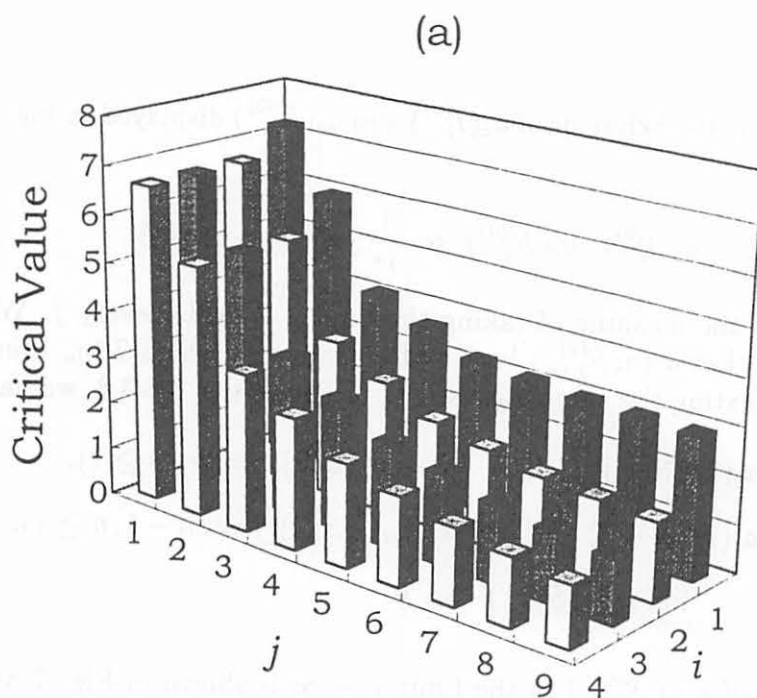


Fig. 5. Three-dimensional plot of (a)  $a_c(1/(2i+j)) \in I_i^{(0)}$  and (b)  $a_c(1/(2i+j)) \in K_i^{(0)}$

In the limit  $i \gg 1$ , the behavior of  $a_c(I_i^{(0)})$  and  $a_c(K_i^{(0)})$  displayed in Fig. 6 is modeled by

$$a_c(I_i^{(0)}), a_c(K_i^{(0)}) \propto \frac{1}{i^\alpha} \quad \text{with } \alpha \simeq 0.93 \quad (48)$$

Here we discuss the meaning of taking the limit  $i \rightarrow \infty$  for every  $j$ . We shall define two critical values. Let  $a_c(n; S_{2,4-}^{(0)})$  be a critical value at which  $T^n \gamma_u$  touches  $S_{2,4-}^{(0)}$  for the first time. Repeating the same discussions in Properties 3.6-3.8, we have

$$\lim_{i \rightarrow \infty} a_c(1/(2i+j) \in I_j^{(0)}) = a_c(j; S_{2-}^{(0)})(j = 2n, n \geq 1), \quad (49)$$

$$\lim_{i \rightarrow \infty} a_c(1/(2i+j) \in I_j^{(0)}) = a_c(j; S_{4-}^{(0)})(j = 2n-1, n \geq 1), \quad (50)$$

$$\lim_{j \rightarrow \infty} a_c(j; S_{2,4-}^{(0)}) = 0. \quad (51)$$

The accumulation of  $a_c(j; S_{2,4-}^{(0)})$  in the limit  $j \rightarrow \infty$  is shown in Fig. 7 and is modeled by

$$a_c(j; S_{2,4-}^{(0)}) \propto \frac{1}{j^\delta} \quad \text{with } \delta \simeq 0.87. \quad (52)$$

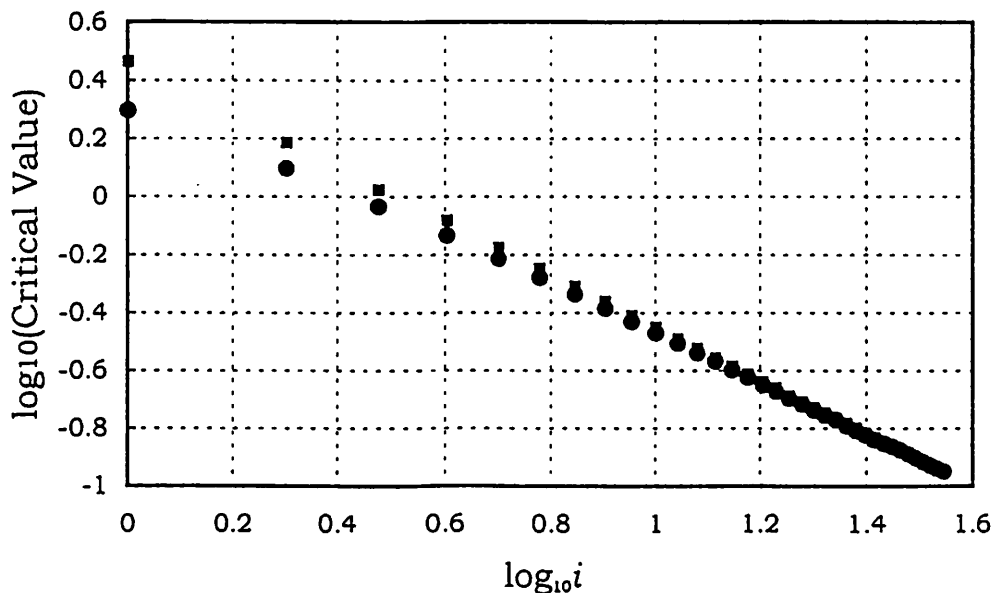


Fig. 6. Plot of  $a_c(I_i^{(0)})$ (square) and  $a_c(K_i^{(0)})$ (circle) as a function of  $i$ .

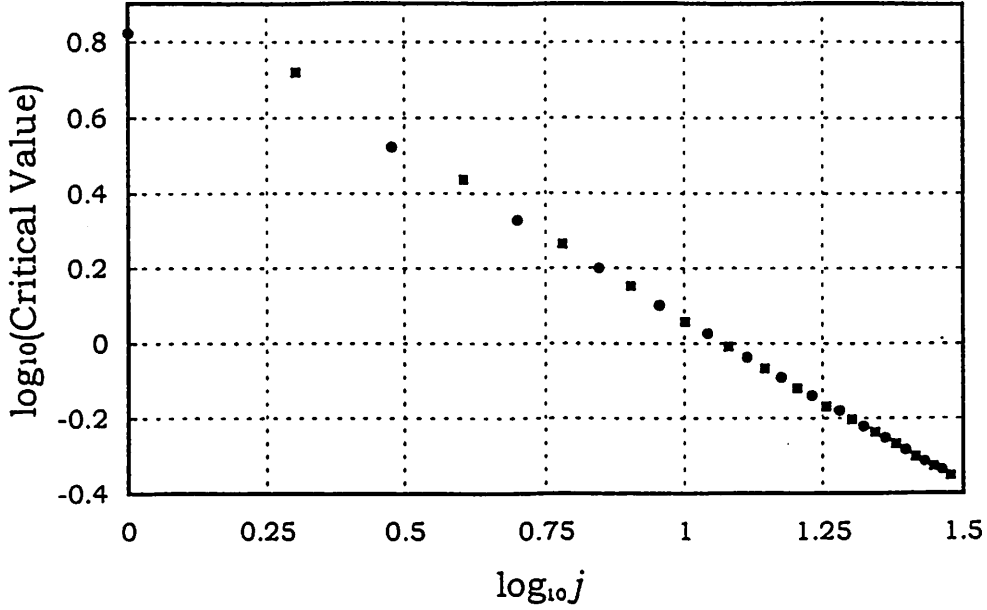


Fig. 7. Accumulation of critical values  $a_c(j; S_2^{(0)})$  (square) and  $a_c(j; S_4^{(0)})$  (circle).

## 5 Braid

We express the reason why we construct the braid of SNBO. If there exists a particular periodic orbit in the system and its braid is *pseudo-Anosov* (pA) type, the lower bound of topological entropy is positive and thus there exists chaos in this system. The concept of pA is an extension of Anosov property showing the hyperbolicity of system.<sup>28)</sup>

Let us introduce notation.  $\text{DO}_{R^{(-m)}}^p$  ( $R = I, J, K$ , or  $L$ ) represents the dynamical order itself realized in  $R$  where  $p$  shows the numerator of rotation number.  $\text{DO}_{R^{(-m)}}^p(i, j)$  ( $R = I, J, K$ , or  $L$ ) is the  $(i, j)$  element of  $\text{DO}_{R^{(-m)}}^p$  where  $i \geq 1$  and  $j \geq 1$  for  $I^{(-m)}$  and  $J^{(-m)}$ , and  $i \geq 0$  and  $j \geq 1$  for  $K^{(-m)}$  and  $L^{(-m)}$ .

## 5.1 Braid for SNBOs in $DO_{I^{(0)}}^1$

For an  $1/(2i+j)$ -SNBO with  $j = 2k + 1, k \geq 0, i \geq 1$ , we divide its orbital points  $\{p_0, p_1, \dots, p_{2i+2k}\}$  of one period into four groups  $\mathcal{A}, \mathcal{B}, \mathcal{C}$  and  $\mathcal{D}$  as

$$\begin{aligned} \mathcal{A} &= \{p_0, \dots, p_{i-1}\}, \quad \mathcal{B} = \{p_i, \dots, p_{i+k}\}, \\ \mathcal{C} &= \{p_{i+k+1}, \dots, p_{i+2k+1}\}, \quad \text{and} \quad \mathcal{D} = \{p_{i+2k+2}, \dots, p_{2i+2k}\} \end{aligned}$$

where  $p_0$  is always taken in  $I_i^{(0)}$ ,  $\mathcal{A} \setminus \{p_0\}$  and  $\mathcal{C}$  are located between  $S_1^{(0)}$  and  $S_2^{(0)}$ , and  $\mathcal{B}$  and  $\mathcal{D}$  between  $S_2^{(0)}$  and  $S_1^{(1)}$ . We have  $\mathcal{A} \setminus \{p_0\} = G\mathcal{D}$  and  $\mathcal{B} = G\mathcal{C}$  by reversibility.

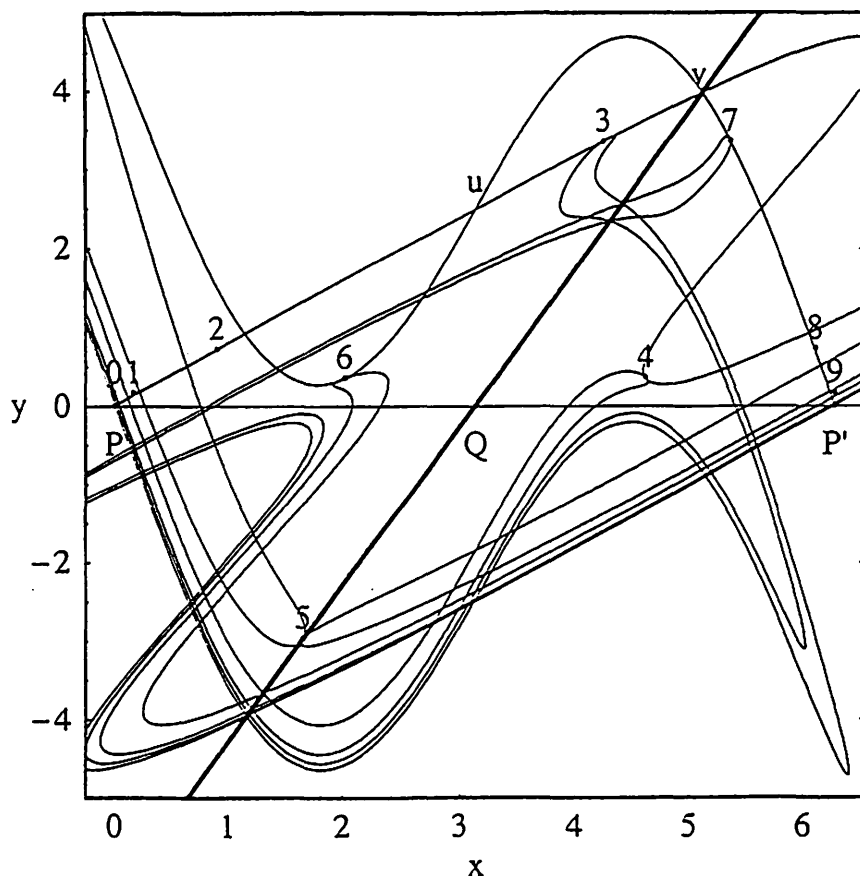


Fig. 8. The orbit of an  $1/9$ -SNBP in  $I_3^{(0)}$  where an integer  $k$  stands for  $p_k$ . The stable and unstable manifolds are also displayed.

Figure 8 displays the orbit of an  $1/9$ -SNBP in  $I_3^{(0)}$  where  $\mathcal{A} = \{p_0, p_1, p_2\}$ ,  $\mathcal{B} = \{p_3, p_4\}$ ,  $\mathcal{C} = \{p_5, p_6\}$  and  $\mathcal{D} = \{p_7, p_8\}$ . The schematic version is illustrated in Fig. 9(a). The group  $\mathcal{A}$  is located to the left of  $W_u^1$  and  $\mathcal{C}$  is to the right. Due to the reversibility with respect to  $S_2^{(0)}$ ,  $\mathcal{B}$  is located to the right of  $W_s^1$  and  $\mathcal{D}$  is to the left.

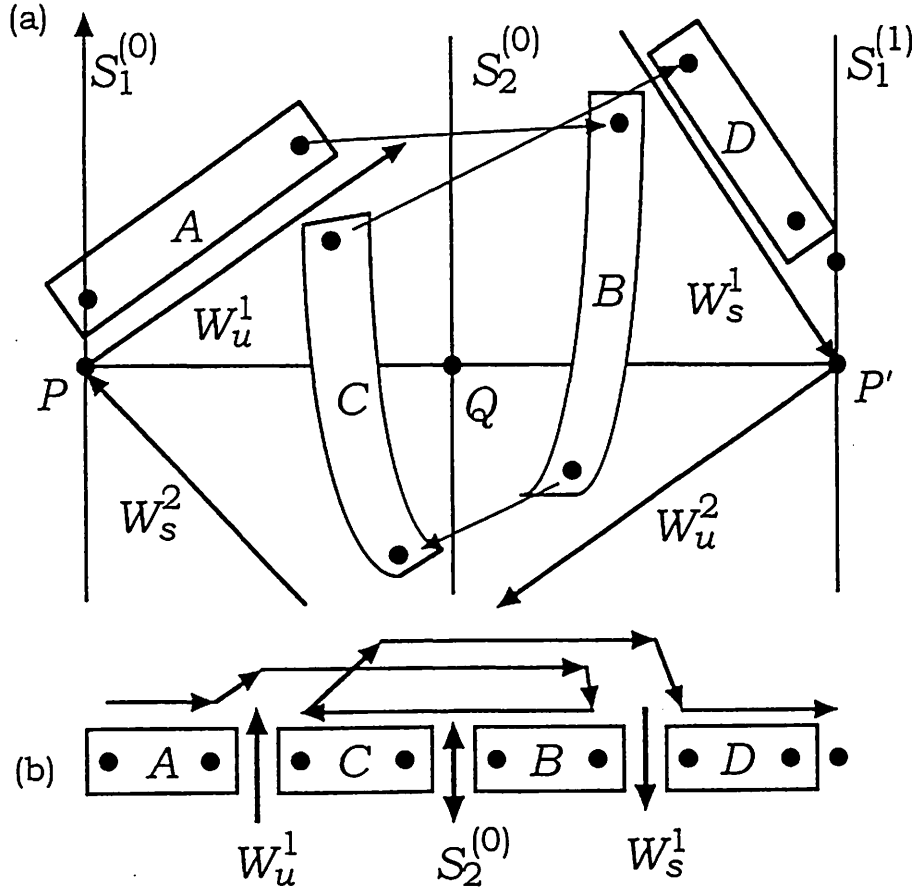


Fig. 9. (a) Configuration of four groups  $\mathcal{A}$ ,  $\mathcal{B}$ ,  $\mathcal{C}$  and  $\mathcal{D}$ , and (b) reconstruction of these groups into the fundamental orbital order.

For the behavior of the orbital points, we observe the following properties. These are justified next paragraph.

[P1] There is one turing point in  $\mathcal{B}$  and  $\mathcal{C}$ .

[P2] For the points in  $\mathcal{A}$  and  $\mathcal{D}$ , the following relations hold.

$$\pi_1(p_m) < \pi_1(p_{m'}) \quad (53)$$

where  $0 \leq m < m' \leq i-1$  in  $\mathcal{A}$  and  $i+2k+2 \leq m < m' \leq 2i+2k$  in  $\mathcal{D}$ .

[P3] For the points in  $\mathcal{A}$  and  $\mathcal{C}$ , the following relations hold if the corresponding points exist at all.

$$\pi_1(p_{i-1}) < \pi_1(p_{i+2k+1}), \quad (54)$$

$$\pi_1(p_{i-2}) < \pi_1(p_{i+2k}), \quad (55)$$

...

[P4] For the points in  $\mathcal{B}$  and  $\mathcal{D}$ , the following relations hold if the corresponding points exist at all.

$$\pi_1(p_i) < \pi_1(p_{i+2k+2}), \quad (56)$$

$$\pi_1(p_{i+1}) < \pi_1(p_{i+2k+3}), \quad (57)$$

...

[P5] The  $x$  coordinate of any point in  $\mathcal{B}$  is larger than those of points in  $\mathcal{C}$ .

Here [P1] means that one turning-back point is located in  $\mathcal{B}$  and one turning-forward point in  $\mathcal{C}$ . [P2] means that the orbital points in  $\mathcal{A}$  and  $\mathcal{D}$  are well-ordered. In fact, the turning back does not occur since the value of  $y$ -coordinate of orbital points is positive. Equation (56) in [P4] is true since  $p_i \in V$  and  $p_{i+2k+2} \in TU$ , and  $TU$  is located to the right of  $V$ . The other inequalities reflect the fact that  $p_{i+2k+2}$  in  $TU$  is mapped to the right of  $Tv$  and  $p_i$  in  $V$  is mapped to the left of  $Tv$ , and so on. Note that  $p_{2i+2k+1} \in T^iU$  is located in  $S_{1+}^{(1)}$  and  $p_{i+k+1} \in T^{k+1}V$  is located in  $S_{4-}^{(0)}$ . Operating  $G$  to the equations of [P4], we have [P3]. [P5] is derived from the reversibility with respect to  $G$ .

Taking into account [P1] through [P5], we place the four groups of points in a line as

$$O_o(i, j) = (\mathcal{A} \uparrow \mathcal{C} \downarrow \mathcal{B} \downarrow \mathcal{D}). \quad (58)$$

A schematic illustration is given in Fig. 9(b). We call the expression an *orbital order* of  $DO_{I^{(0)}}^1(i, j)$ . The suffix  $o$  stands for 'odd', i.e., Eq. (58) is for odd periodic orbits. Three symbols  $\uparrow, \downarrow$  and  $\downarrow$  in the orbital order represent  $W_u^1$ ,  $S_2^0$ , and  $W_s^1$ , respectively. These symbols stress that four groups are separated by three objects. Later, these will be frequently omitted. If in particular  $p_k$  and  $p_{i+k+1}$  are the two turning points and there are no other turning points, then we call the order the *fundamental orbital order* (FOO). The role of the FOO in constructing braids will be seen in what follows.

Let us consider  $DO_{I^{(0)}}^1(3, 3)$ . Its FOO(see Fig. 10(a)) is

$$O_o(3, 3) = (012654378). \quad (59)$$

There are two additional orbital orders which satisfy [P1-5] and which can be realized by actual orbits(Fig. 10(b) and (c)):

$$O'_o(3, 3) = (012563478), \quad (60)$$

$$O''_o(3, 3) = (015263748). \quad (61)$$

We can not exchange 1 and 5 due to [P3], and 4 and 8 due to [P4]. The orbital order directly determined by the orbit displayed in Fig. 8 is  $O'_o(3, 3)$ .

We describe now the construction of a braid of an  $1/q$ -SNBO using the information on the orbital order. A braid is constructed in two steps. The first step is a construction of  $(q-1)$  strings from 0 to 1, 1 to 2, ...,  $(q-2)$  to  $(q-1)$ , and the second step is that of the final string from  $(q-1)$  to 0.

Step [1]. A string corresponding to the forward motion goes behind a string corresponding to the backward motion when they cross.

Step [2]. The string from  $(q - 1)$  to 0 goes behind all other strings when they cross.

The first rule expresses the twist property and the second one the rigid rotation with  $\nu = 1/q$  revolving round cylinder.

Braid  $\beta(3, 3)$  constructed from  $O_o(3, 3)$  is shown in Fig. 10(a). Its expression in terms of generators is

$$\beta(3, 3) = \sigma_4^{-1} \sigma_5^{-1} \sigma_6^{-1} \sigma_3^{-1} \sigma_4^{-1} \sigma_5^{-1} \zeta_9, \quad (62)$$

where  $\zeta_9 = \sigma_8 \cdots \sigma_1$ . Braid  $\beta'(3, 3)$  constructed from  $O'_o(3, 3)$  and braid  $\beta''(3, 3)$  from  $O''_o(3, 3)$  (Figs. 10(b) and (c)) are expressed in terms of generators as

$$\beta'(3, 3) = \sigma_6^{-1} \sigma_5^{-1} \sigma_6^{-1} \sigma_3^{-1} \sigma_4^{-1} \sigma_3^{-1} \zeta_9, \quad (63)$$

$$\beta''(3, 3) = \sigma_7^{-1} \sigma_6^{-1} \sigma_5^{-1} \sigma_4^{-1} \sigma_3^{-1} \sigma_2^{-1} \zeta_9. \quad (64)$$

These braids are equivalent to  $\beta(3, 3)$  via the Markov move as

$$\beta'(3, 3) = \sigma_4 \sigma_6^{-1} \beta(3, 3) \sigma_4^{-1} \sigma_6, \quad (65)$$

$$\beta''(3, 3) = \sigma_3 \sigma_7^{-1} \beta'(3, 3) \sigma_3^{-1} \sigma_7. \quad (66)$$

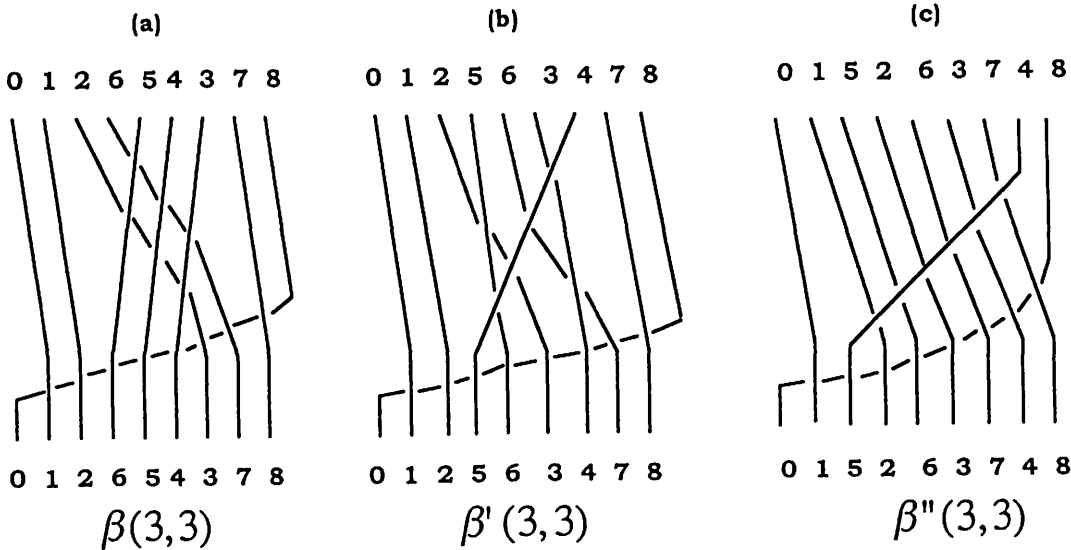


Fig. 10. (a) Braid constructed from  $O_o(3, 3)$ , (b) that from  $O'_o(3, 3)$  and (c) that from  $O''_o(3, 3)$  where an integer  $k$  stands for  $p_k$ .

To obtain  $\beta''(3, 3)$ , we exchange two strings 4 and 7 in  $\beta'(3, 3)$ . The string from 7 to 8 gets ahead of the string from 4 to 5 in the upper side. As a result, new intersection point appears. We add  $\sigma_7^{-1}$  in the left side of  $\beta'(3, 3)$ . In the bottom, the string from 3 to 4 and the string from 6 to 7 do not intersect each other. In order to untwist them,



we add  $\sigma_7$  in the right of  $\beta'(3, 3)$ . Due to symmetry, we add  $\sigma_3$  and  $\sigma_3^{-1}$  and thus have  $\beta''(3, 3)$ . It turns out that the exchange of points in the orbital order does not change the braid type if it does not violate [P1-5]. The construction of the FOO is simple compared with that of other orbital orders. This is the reason why we use the FOO to construct the braid.

Let us see what happens if we accomplish inhibited exchanges of two strings 3 and 7 and of strings 2 and 6 in Fig. 10(a). These exchanges do not satisfy [P3] and [P4]. We have an orbital order (016254738) and have a braid  $\sigma_7^{-1}\beta(3, 3)\sigma_3^{-1}$ . Since the number of intersection points of braid is different from that of  $\beta(3, 3)$ , the new braid is not equivalent to  $\beta(3, 3)$ .

Let us next consider periodic orbits of even period. We divide the orbital points (one period) of an  $1/(2i + j)$ -SNBO ( $j = 2k, k \geq 1$ ) into  $\{p_{i+k}\}$  and four sets  $\mathcal{A}, \mathcal{B}, \mathcal{C}$ , and  $\mathcal{D}$  as

$$\begin{aligned}\mathcal{A} &= \{p_0, \dots, p_{i-1}\}, \quad \mathcal{B} = \{p_i, \dots, p_{i+k-1}\}, \\ \mathcal{C} &= \{p_{i+k+1}, \dots, p_{i+2k}\}, \quad \text{and} \quad \mathcal{D} = \{p_{i+2k+1}, \dots, p_{i+2k-1}\}.\end{aligned}$$

and arrange them in a line as

$$O_e(i, j) = (\mathcal{A} \uparrow \widehat{\mathcal{C}(i+k)} \mathcal{B} \downarrow \mathcal{D}). \quad (67)$$

where  $(\widehat{i+k})$  stands for  $p_{i+k} \in S_2^{(0)}$ . The symbols have the same meaning as before.  $(\widehat{i+k})$  plays the role of  $\uparrow$ . The FOO is defined as the orbital order in which turning points are  $p_i$  and  $p_{i+k}$ .

Now the expressions of braids for  $O_o(2, 3) = (0154326)$  and  $O_e(2, 4) = (01654327)$  are determined.

$$\begin{aligned}\beta(2, 3) &= \sigma_3^{-1}\sigma_4^{-1}\sigma_5^{-1}\sigma_2^{-1}\sigma_3^{-1}\sigma_4^{-1}\zeta_7, \\ &= \sigma_1^{-1}\sigma_2^{-1}\sigma_3^{-1}\sigma_3^{-1}\sigma_2^{-1}\sigma_1^{-1}\zeta_7,\end{aligned} \quad (68)$$

$$\begin{aligned}\beta(2, 4) &= \sigma_3^{-1}\sigma_4^{-1}\sigma_5^{-1}\sigma_6^{-1}\sigma_2^{-1}\sigma_3^{-1}\sigma_4^{-1}\sigma_5^{-1}\zeta_8, \\ &= \sigma_1^{-1}\sigma_2^{-1}\sigma_3^{-1}\sigma_4^{-1}\sigma_4^{-1}\sigma_3^{-1}\sigma_2^{-1}\sigma_1^{-1}\zeta_8,\end{aligned} \quad (69)$$

where Reidemeister and Markov moves<sup>5),24)</sup> are operated to derive the second equation in Eqs. (68) and (69). As a result, the expression of braid for  $O_o(i, j)$  or  $O_e(i, j)$  is derived.

$$\beta(i, j) = \zeta_{j+1}^{-1}\rho_{j+1}^{-1}\zeta_{2i+j} \quad (70)$$

where  $\rho_k = \sigma_1 \cdots \sigma_{k-1}$  and  $\zeta_k = \sigma_{k-1} \cdots \sigma_1$ .

Here let us compare  $\beta(i, j)$  and braids investigated by Boyland.<sup>8)</sup> Consider two annuli and  $(2i+j)$  strings connecting them. The braid  $\zeta_{j+1}^{-1}\rho_{j+1}^{-1}$  means that the first string passes behind the second through  $(j+1)$ -th strings, and then passes in front of these strings. This twisting is caused by the non-Birkhoffness of the orbit. The braid  $\zeta_{2i+j}$  represents a rigid rotation of a  $1/(2i+j)$ -Birkhoff orbit. According to this geometrical interpretation, the braid  $\beta(i, 1)$  corresponds to the braid constructed by Boyland.

## 5.2 Braid for SNBOs in $DO_{K^{(0)}}^1$

### 5.2.1 Braid for SNBOs in $DO_{K^{(0)}}^1$ with $i \geq 1$

Figure 11 displays SNBOs starting in  $S_{3+}^{(0)}$ . The points of one period of an orbit are grouped into four sets  $\mathcal{A}, \mathcal{B}, \mathcal{C}$  and  $\mathcal{D}$ . For even period  $q = 2i + j + 1$  ( $j = 2k - 1, k \geq 1$ ), we define the orbital order.

$$\begin{aligned} O_e(i, j) &= (\mathcal{A} \uparrow \mathcal{C} \downarrow \mathcal{B} \downarrow \mathcal{D}), \\ \mathcal{A} &= \{p_0, \dots, p_{i-1}\}, \quad \mathcal{B} = \{p_i, \dots, p_{i+k-1}\}, \\ \mathcal{C} &= \{p_{i+k}, \dots, p_{i+2k-1}\}, \quad \text{and} \quad \mathcal{D} = \{p_{i+2k}, \dots, p_{2i+2k-1}\}. \end{aligned} \quad (71)$$

The FOO is such that  $p_i$  and  $p_{i+2k-1}$  are the turning points. For odd period  $q = 2i + j + 1$  ( $j = 2k, k \geq 1$ ), we define the orbital order.

$$\begin{aligned} O_o(i, j) &= (\mathcal{A} \uparrow \widehat{\mathcal{C}(i+k)} \mathcal{B} \downarrow \mathcal{D}), \\ \mathcal{A} &= \{p_0, \dots, p_{i-1}\}, \quad \mathcal{B} = \{p_i, \dots, p_{i+k-1}\}, \\ \mathcal{C} &= \{p_{i+k+1}, \dots, p_{i+2k}\}, \quad \text{and} \quad \mathcal{D} = \{p_{i+2k+1}, \dots, p_{2i+2k}\}. \end{aligned} \quad (72)$$

The FOO is such that  $p_i$  and  $p_{i+2k}$  are the turning points. Note that  $\mathcal{A} = \mathcal{G}\mathcal{D}$  and  $\mathcal{C} = \mathcal{G}\mathcal{B}$  hold in both  $O_e(i, j)$  and  $O_o(i, j)$ .

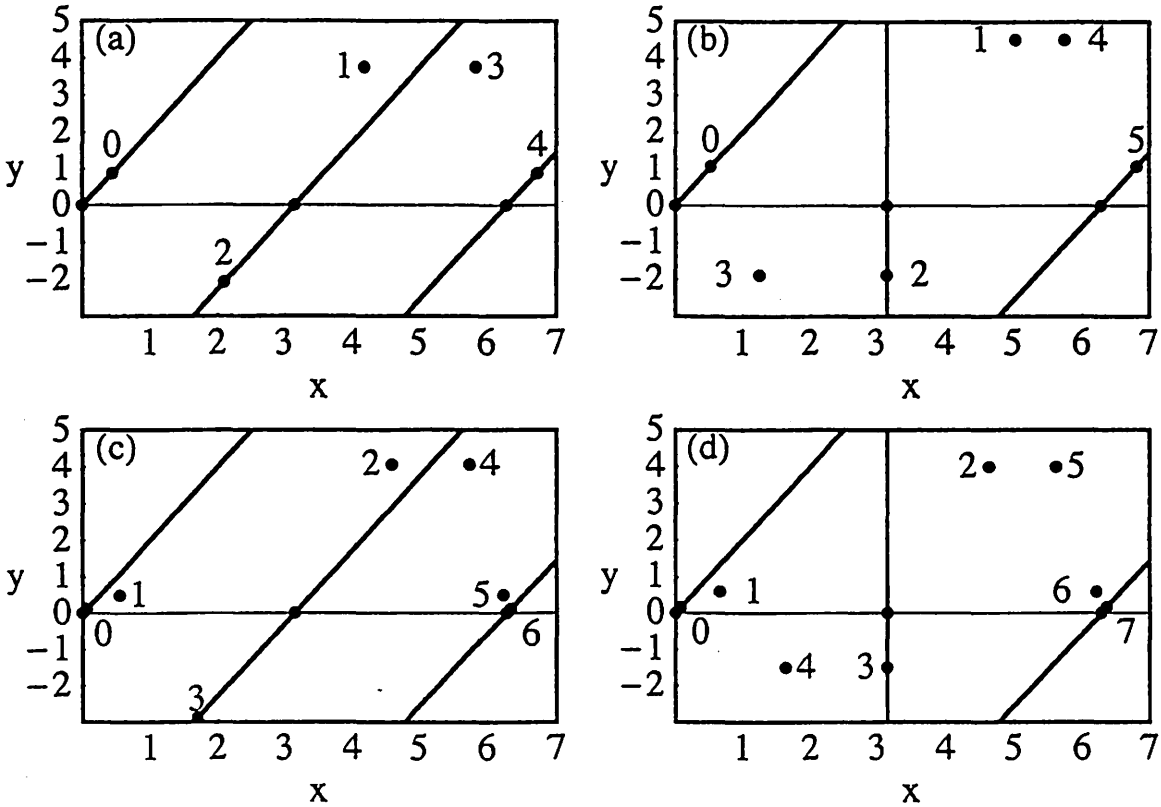


Fig. 11. (a)  $1/4 \in K_1^{(0)}$ , (b)  $1/5 \in K_1^{(0)}$ , (c)  $1/6 \in K_2^{(0)}$  and (d)  $1/7 \in K_2^{(0)}$ .

Using the same rules as in §5.1, the braid  $\beta(i, j)$  for  $DO_{K^{(0)}}^1(i, j)$  is derived.

$$\beta(i, j) = \zeta_{j+1}^{-1} \rho_{j+1}^{-1} \zeta_{2i+j+1}. \quad (73)$$

### 5.2.2 Braid for SNBOs in $DO_{K^{(0)}}^1$ with $i = 0$

In Fig. 12, two orbits with  $\nu = 1/2$  and  $1/3$  are displayed. For both orbits, the turning back of the orbit occurs at  $p_0$  (symbol 0 in the figure).

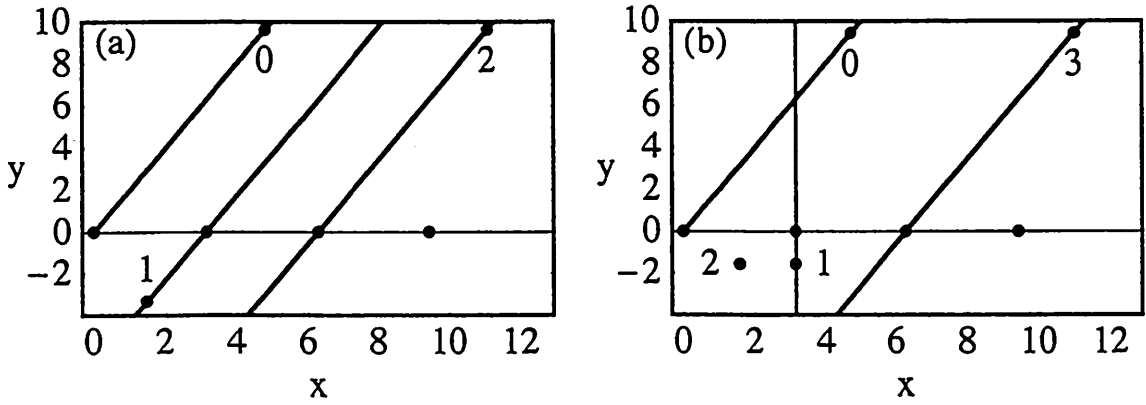


Fig. 12. (a)  $1/2 \in K_0^{(0)}$  and (b)  $1/3 \in K_0^{(0)}$ .

In order to have a braid of pA from an  $1/2$ -SNBO ( $DO_{K^{(0)}}(0, 1)$ ), we need at least one more string. To accomplish this, we use the information in the vicinity of the period-2 orbit. We first deform cylinder to annulus. Next we shrink the inner circle of annulus to a point (termed as  $c$ ). The point  $c$  stands for the fixed point at infinity. A string from  $c$  to  $c$  is the additional string. We construct the braid of three strings connecting two annuli. Let a string from  $p_0$  to  $p_1$  be  $A$ , a string from  $p_1$  to  $p_0$  be  $B$  and a string from  $c$  to  $c$  be  $C$ . Two strings  $A$  and  $B$  revolve round  $c$ . This gives how to intersect two strings ( $A$  and  $C$ ,  $B$  and  $C$ ). If  $A$  and  $B$  intersect each other, we draw the string of the forward motion behind the string of the backward motion. This intersection corresponds to the twist property. As a result, we have the expression of braid.

$$\beta(0, 1) = \sigma_1 \sigma_2^{-1} \sigma_1^{-1} \sigma_2 \sigma_2 = \sigma_1^2 \sigma_2^{-1}. \quad (74)$$

The braid  $\beta(0, 1)$  expresses the tangled structure of the braid of the period-2 orbit and that of  $c$ .

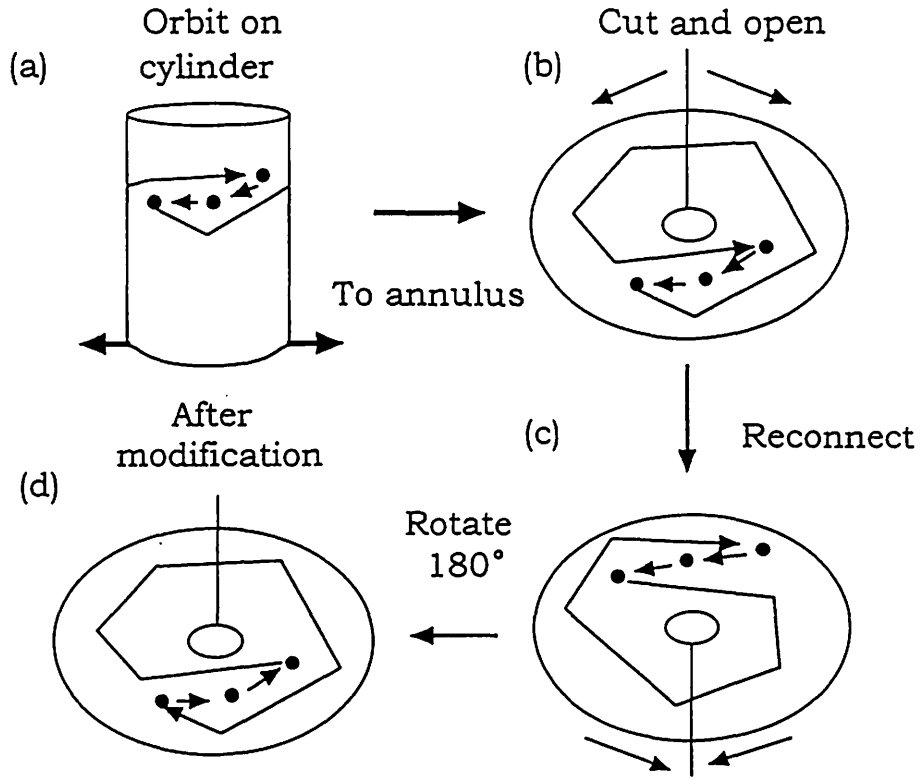


Fig. 13. Deformation of cylinder to annulus ((a)  $\rightarrow$  (b)). Cutting and recombination of annulus ((b)  $\rightarrow$  (c) ). After rotating (c) by  $180^\circ$ , we have (d)

Next we construct the braid for an  $1/3$ -SNBO. The orbital order is  $O_o = (\uparrow 2\hat{1}0 \downarrow)$ , and this gives a braid  $\sigma_1^{-1}\sigma_2^{-1}$ . But this braid is not pA. We apply the transformation illustrated in Fig. 13. The first step is the deformation of cylinder to annulus. We cut open annulus and glue both borders at the opposite side so that the inner and outer circles exchange their role. Thus we have new orbit in the modified annulus and have the expression of braid.

$$\begin{aligned}\beta(0,2) &= \sigma_2^{-1}\sigma_1^{-1}\sigma_1^{-1}\sigma_2^{-1}\sigma_2^{-1}\sigma_1^{-1}, \\ &= (\sigma_2^{-1}\sigma_1^{-1})^3\sigma_1\sigma_2^{-1}.\end{aligned}\quad (75)$$

Equation (75) implies that  $\beta(0,2)$  is pA.

The orbital orders for the elements with  $j(\geq 2)$  of  $q = j + 1$  are derived.

$$O_o(0,j) = (\uparrow B\hat{k}A \downarrow), \quad (76)$$

where  $A = (p_0, \dots, p_{k-1})$ ,  $B = (p_{k+1}, \dots, p_{2k})$  ( $j = 2k, k \geq 1$ ) and  $p_0$  and  $p_{2k}$  are the turning points, and

$$O_e(0,j) = (\uparrow B\uparrow A \downarrow), \quad (77)$$

where  $A = (p_0, \dots, p_k)$ ,  $B = (p_{k+1}, \dots, p_{2k+1})$  ( $j = 2k + 1, k \geq 1$ ) and  $p_0$  and  $p_{2k+1}$  are the turning points. In both expressions,  $B = GA$  holds.

These determine the braids  $\sigma_1^{-1} \dots \sigma_q^{-1}$ . Unfortunately these are not pA. Therefore, we apply the transformation mentioned above and then have the braids of pA.

$$\beta(0, j) = \rho_{j+1}^{-2} \zeta_{j+1}^{-1} \quad (j \geq 2). \quad (78)$$

## 6 Topological entropy

We have derived braids of SNBOs in §5. In this section, we use these to estimate the lower bound of the topological entropy of a system which possesses an SNBO with braid  $\beta$ . The lower bound  $h(\beta)$  of the entropy can be estimated by the maximum absolute eigenvalue ( $\lambda_{max} = \text{Max}(|\lambda_i|)$ ) of the reduced Burau matrix representation  $M_\beta(t)$ .<sup>25),26)</sup>

$$h(\beta) = \ln \lambda_{max}. \quad (79)$$

The reduced Burau matrix representation  $M_\beta(t)$  has a parameter  $t$  defined by  $t = \exp(i\theta)$  ( $0 \leq \theta < 2\pi$ ). In order to calculate  $h(\beta)$ , we need the value of  $\text{Max}(|\lambda_i|)$  as a function of  $t$  or  $\theta$ . For the braid of the third order,  $\lambda_{max}$  is obtained at  $t = -1$ .<sup>23)</sup> We need numerical calculations for the braids of arbitrary orders. In Appendix, we give a sample program written by MATHEMATICA<sup>27)</sup> to construct the reduced Burau matrix representation and to calculate its eigenvalues. Figure 14 displays the numerical results of  $\text{Max}(|\lambda_i|)$  as a function of  $\theta$ .

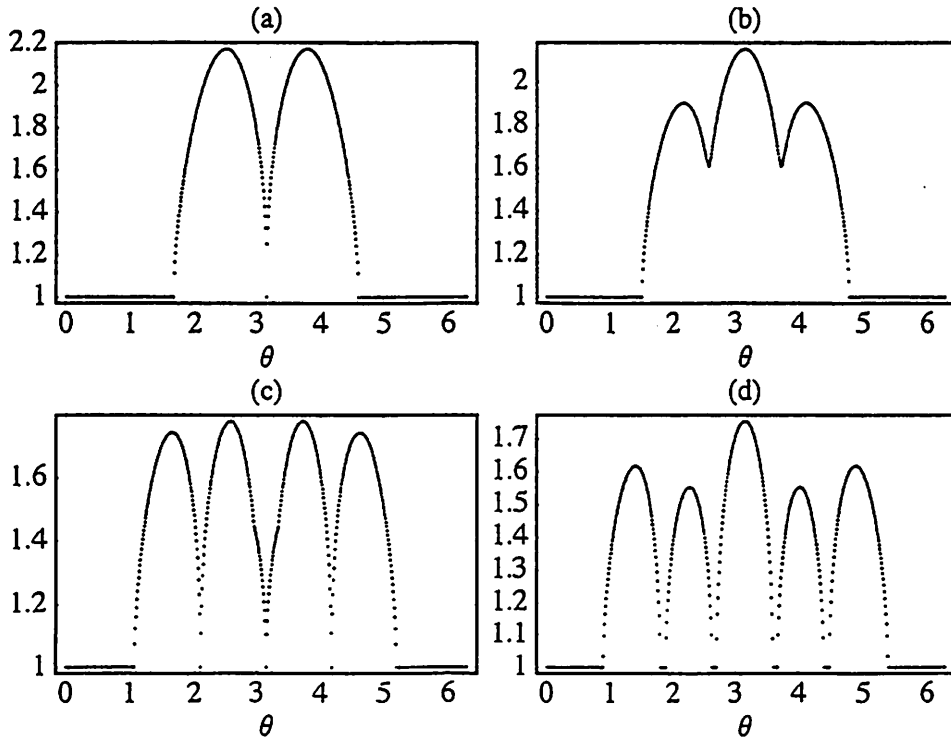


Fig. 14. The  $\theta$  dependence of  $\text{Max}(|\lambda_i|)$ . (a)  $1/4 \in I_1^{(0)}$ , (b)  $1/5 \in I_1^{(0)}$ , (c)  $1/6 \in I_2^{(0)}$  and (d)  $1/7 \in I_2^{(0)}$ .

## 6.1 Topological entropy estimated by using SNBO in $DO_{I(0)}^1$

We observe the following property for braids of odd orders through numerical calculations. Proof is in order. The property is generally not true for braids of even orders. Then we only calculate  $h(\beta)$  for SNBOs with odd periods.

**Observation.**  $\lambda_{max}$  is determined by the maximum absolute value of the real root ( $> 1$ ) of eigenfunction at  $t = -1$ .

We derive the eigenfunction of the reduced Burau matrix for braid  $\beta(1, j) = \zeta_{j+1}^{-1} \rho_{j+1}^{-1} \zeta_{j+2}$  in  $DO_{I(0)}^1(1, j)$  where  $j$  is assumed to be an odd integer satisfying  $j \geq 1$ . For example, Eq. (80) gives the reduced Burau matrix with  $t = -1$  for  $\beta(1, 3) = \sigma_3^{-1} \sigma_2^{-1} \sigma_1^{-1} \sigma_4$  of  $1/5 \in I_1^{(0)}$ , and Eq. (81) is its eigenfunction.

$$M_{\beta(1,3)}(-1) = \begin{pmatrix} 1 & -1 & 0 & 0 \\ 1 & 0 & -1 & 0 \\ 1 & 0 & 1 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix}, \quad (80)$$

$$\lambda^4 - 3\lambda^3 + 3\lambda^2 - 3\lambda + 1 = 0. \quad (81)$$

The eigenfunction for any  $j$  is obtained in Eq.(82).

$$\lambda^{j+1} + 3 \sum_{k=1}^j (-\lambda)^k + 1 = \frac{(\lambda - 2)(\lambda^{j+1} - 2) - 3}{\lambda + 1} = 0. \quad (82)$$

It is easy to see that the value of  $\lambda_{max}$  accumulates at 2 as  $j \rightarrow \infty$ . Equation (82) corresponds to the expression derived by Boyland.<sup>8)</sup>

In order to estimate  $\lambda_{max}$  for  $(i, j)$  element, we shall derive the eigenfunction for the cases with  $j > i$ .

$$\sum_{k=0}^{2i-2} (2k+1)(-\lambda)^k + (4i-1) \sum_{k=2i-1}^j (-\lambda)^k + \sum_{k=j+1}^{2i+j-1} (4i+2j-2k-1)(-\lambda)^k = 0. \quad (83)$$

Our main purpose is to estimate the topological entropy in the limit  $i, j \rightarrow \infty$ . Here we consider the limit  $j \rightarrow \infty$ . We keep the most divergent terms in Eq. (83), and let these be zero.

$$\lambda^{2i-1}(\lambda - 1) = 2. \quad (84)$$

Next we consider the cases with  $i > j$  and thus have

$$\sum_{k=0}^{j-1} (2k+1)(-\lambda)^k + (2j+1) \sum_{k=j}^{2i-1} (-\lambda)^k + \sum_{k=2i}^{2i+j-1} (4i+2j-2k-1)(-\lambda)^k = 0. \quad (85)$$

We also have

$$(-\lambda)^j(1 - \lambda) = 2. \quad (86)$$

For the cases with  $j = 2i - 1$ , the accumulation values  $\hat{h}$  of topological entropy  $h(\beta(i, j))$  at the  $i$ -th row and the  $j$ -th column are same. Let us consider the limit  $i, j \rightarrow \infty$ . Put  $\lambda = 1 + \epsilon (\epsilon > 0)$ . In the case with  $\epsilon \ll 1$ , we have the relation  $\hat{h} \approx \epsilon$ . Equations (84) is rewritten in the form:

$$k \ln(1 + \epsilon) = \ln(2/\epsilon) \quad (87)$$

where  $k = 2i - 1$ . In the limit  $k \rightarrow \infty$ , we have the topological entropy  $\hat{h}$ :

$$\hat{h} = \frac{\ln(2k)}{k} + O\left(\frac{\ln(\ln(2k))}{k}\right). \quad (88)$$

Numerical results of Eq. (87) are shown in Fig. 15. For large values of  $k$ , these results are in agreement with Eq. (88).

Equation (88) gives that the topological entropy is zero in the integrable limit and also means that the system with  $a > 0$  is pA. Combining Eqs. (48) and (88), we have the topological entropy  $\hat{h}$  as a function of  $a$  in the limit  $a \rightarrow 0$ .

$$\hat{h} \propto \frac{\ln(1/a)}{(1/a)^{1/\alpha}}. \quad (89)$$

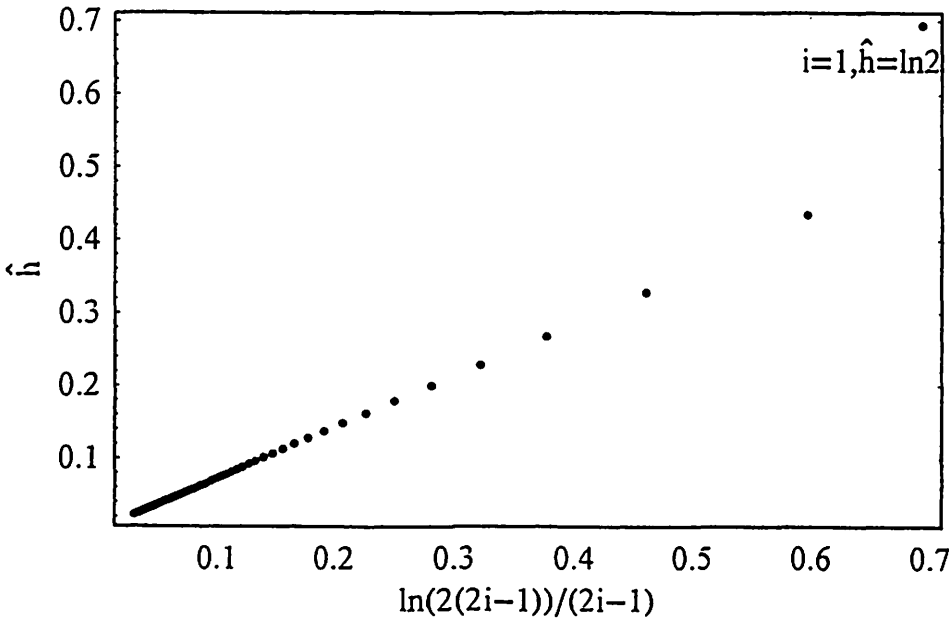


Fig. 15. Accumulation of topological entropy  $\hat{h}$ .

## 6.2 Topological entropy estimated by using SNBO in $\text{DO}_{K(0)}^1$

The topological entropy for  $\text{DO}_{K(0)}^1(0, 1)$  has been already estimated in Ref. 22).

$$h(\beta(0, 1)) = 2 + \sqrt{3}. \quad (90)$$

We derive the eigenfunction of  $\text{DO}_{K^{(0)}}^1(0, j)$  for even  $j$  since Observation of §6.1 is true in these cases.

$$(\lambda + 3)\lambda^{j+1} = 1 + 3\lambda. \quad (91)$$

In the limit  $j \rightarrow \infty$ , we have  $|\lambda| \rightarrow 3$ .

Next we study the topological entropy for  $\text{DO}_{K^{(0)}}^1(i, j) (i \geq 1)$ . In these cases, Observation of §6.1 is not true (see Fig. 16). Thus we calculate the topological entropy numerically. The accumulation value of the topological entropy of the limit  $j \rightarrow \infty$  at  $i = 1$  is  $\ln 1.63$ , and that of the limit  $i \rightarrow \infty$  at  $j = 1$  is  $\ln 2$ . The topological entropies  $\hat{h}_i (i \geq 1, j \rightarrow \infty)$  and  $\hat{h}_j (j \geq 1, i \rightarrow \infty)$  are displayed in Fig. 17, and accumulate to zero in the limit  $i, j \rightarrow \infty$  showing the integrable limit.

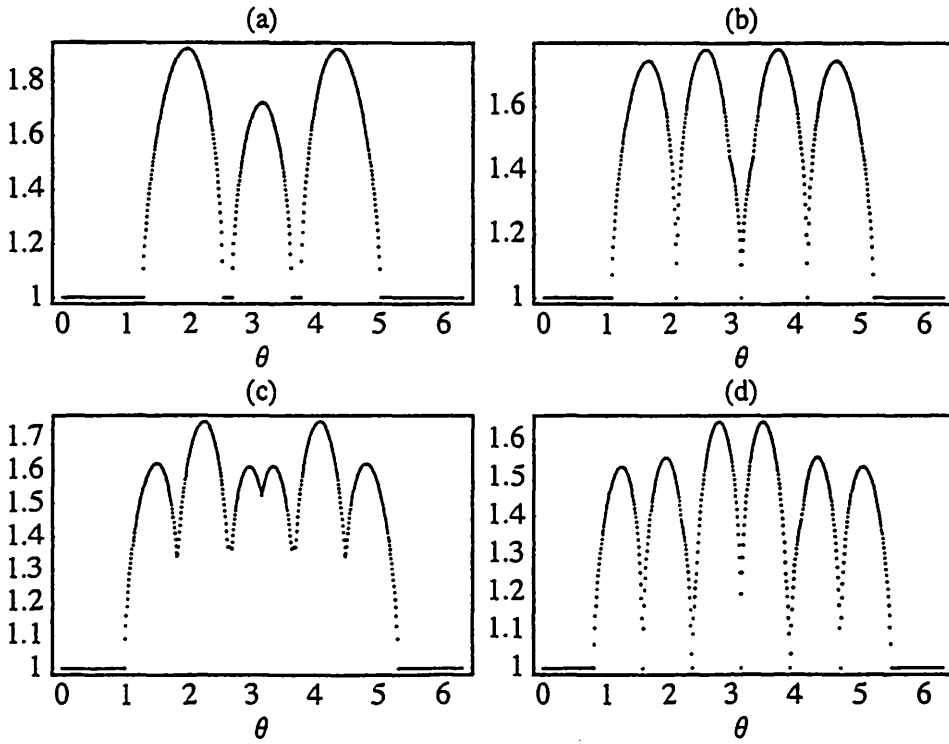


Fig. 16. The  $\theta$ -dependence of  $\text{Max}(|\lambda_i|)$ . (a)  $1/5 \in K_1^{(0)}$ , (b)  $1/6 \in K_1^{(0)}$ , (c)  $1/7 \in K_2^{(0)}$  and (d)  $1/8 \in K_2^{(0)}$ .



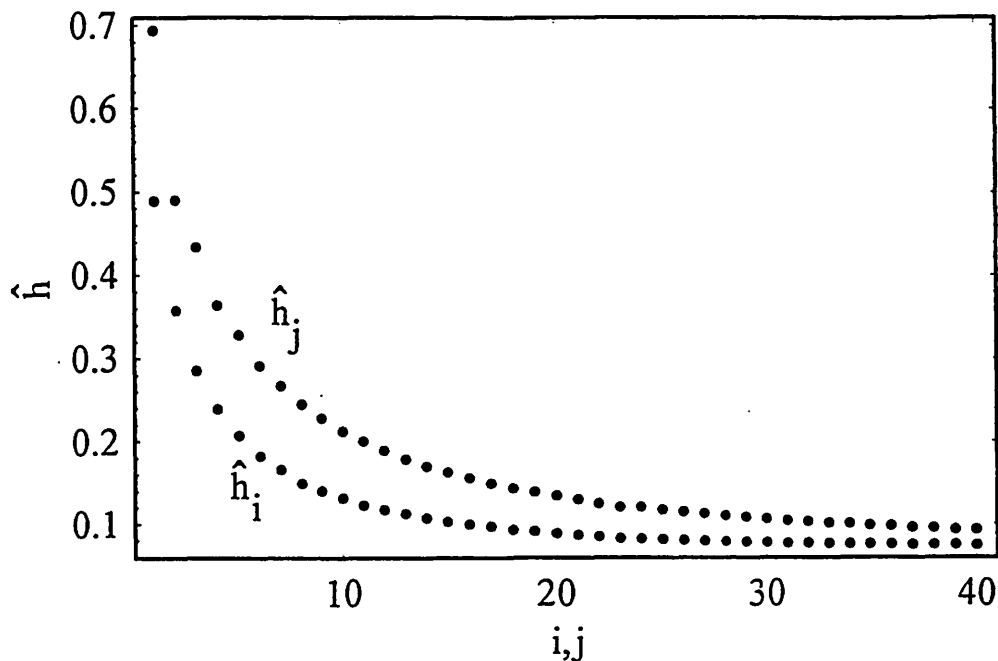


Fig. 17. Accumulation of topological entropy.

## 7 Remarks

We refer to several points to be pursued.

- [1] Forcing relation between different dynamical orderings, for example,  $DO_{I(0)}^p$  and  $DO_{K(0)}^{p'}$ .

To establish the relation,  $p/2$ - or  $p/3$ -SNBOs plays an important role.

- [2] Dynamical ordering for the SNBOs with  $2n$  ( $n \geq 2$ ) turning points and the estimation of topological entropies for these orbits.

- [3] Appearance order of NBOs not having points in symmetry axes.

We have a hypothesis for their appearance, i.e., all NBOs not having points in symmetry axes are bifurcated from mother SNBOs.

- [4] Dynamical ordering for SNBOs of Type II.

This will be useful to investigate the breakup of KAM curves.

- [5] Appearance of NBOs in the systems not having the left-right symmetry.

We believe that two theorems proved in §4 are true if the symmetry breaking perturbation is weak.

- [6] Forcing relation between the braids obtained in §5.

Can we derive the dynamical ordering of SNBOs by using it? This is a reverse approach.

- [7] Theoretical explanation of the power law decay of critical values obtained as Eqs. (48) and (52).

## Acknowledgements

The authors express their thanks to Drs. T. Matsuoka and E. Kin for discussions.

## References

- 1) K. Tanikawa and Y. Yamaguchi, *Chaos* **12** (2002), 33.
- 2) Y. Yamaguchi and K. Tanikawa, *Prog. Theor. Phys.* **106** (2001), 1097.
- 3) A. Sharkovskii, *Ukr. Mat. Z.* **16** (1964), 61.
- 4) L. Alsedà, J. Llibre and M. Misiurewicz, *Combinatorial Dynamics and Entropy in Dimension One* (World Scientific, 1993).
- 5) T. Matsuoka, in *Dynamical System 1* (World Scientific, 1986), p. 58; *Contemp. Math.* **152** (1993), 229. See also *Bussei Kenkyu*(Kyoto) **67** (1996), 1.
- 6) P. Boyland, *Contemp. Math.* **81** (1988), 119.
- 7) S. Baldwin, *Ergod. Th. & Dynam. Sys.* **11** (1991), 249.
- 8) P. Boyland, *Topology and its Appl.* **58** (1994), 223.
- 9) M. Handel, *Ergod. Th. & Dynam. Sys.* **17** (1997), 593.
- 10) J. Los, *I.H.E.S.* (1997), 5.
- 11) P. Boyland and G. R. Hall, *Topology* **26** (1987), 21.
- 12) I. Leage and R. S. Mackay, *Phys. Lett. A* **118** (1986), 274.
- 13) G. D. Birkhoff, *Acta. Math.* **43** (1920), 44.
- 14) K. Tanikawa and Y. Yamaguchi, *J. Math. Phys.* **28** (1987), 921; **30** (1989), 608.
- 15) G. R. Hall, *Ergod. Theor. & Dynam. Sys.* **4** (1984), 585.
- 16) R. de Vogelaere, in *Contribution to the Theory of Nonlinear Oscillations* Vol.IV (Princeton University Press, 1957).
- 17) P. Le Calvez, *Dynamical Properties of Diffeomorphisms of the Annulus and of the Torus* (AMS, 2000).
- 18) V. F. Lazutkin, I. G. Schachmannski and M. B. Tabanov, *Physica D*, **40** (1989), 235.
- 19) S. Wiggins, *Chaotic Transport in Dynamical Systems* (Springer-Verlag, 1991).
- 20) Y. Yamaguchi and K. Tanikawa, *Prog. Theor. Phys.* **103** (2000), 1127.
- 21) J. Palis and W. de Melo, *Geometric Theory of Dynamical Systems* (Springer, 1982).
- 22) Y. Yamaguchi and K. Tanikawa, *Prog. Theor. Phys.* **106** (2001), 691.
- 23) Y. Yamaguchi and K. Tanikawa, *Prog. Theor. Phys.* **104** (2000), 943.
- 24) S. Moran, *The Mathematical Theory of Knots and Braids* (North-Holland, 1983).
- 25) D. Fried, in *Geometric Dynamics*. ed. J. Palis Jr. *Lecture Notes in Mathematics* **1007** (Springer-Verlag, 1983). p. 261.
- 26) B. Kolev, *C. R. Acad. Sci. Paris*, **309**, Ser. I (1989), 835.
- 27) S. Wolfram, *THE MATHEMATICA BOOK*, Fourth Edition (Cambridge University Press, 1999).
- 28) A. J. Casson and S. A. Bleiber, *Automorphisms of Surface after Nielsen and Thurston* (Cambridge University Press, 1988).

## Appendix

A MATHEMATICA (Trademark of Wolfram Research) program which constructs the reduced Burau matrix representation and calculates its eigenvalues as a function of  $\theta$  is shown below.

Explanatory note.

[1] In the first block, input two values of  $n$  ( $\geq 4$ ) and  $ma$  where  $n$  is an order of braid and  $ma$  is a number of division of  $\theta$ .

[2] The second block is a preparation of generators.

[3] In the third block, input a braid by using  $s[k]$  and  $is[k]$  where  $s[k]$  is a generator  $\sigma_k$  and  $is[k]$  a generator  $\sigma_k^{-1}$ .

[4] In the fourth block, the eigenvalues are calculated and  $\text{Max}|\lambda_i|$  is displayed as a function of  $\theta$ . Final output is  $\lambda_{max}$ .

A sample program determines the eigenvalues of a braid  $\sigma_1\sigma_2^{-1}\sigma_3^{-1}$ .

### Sample program

```
(* 1st Block *)
(* Input an order of Braid *)
n = 4;
(* Input a number of division *)
ma = 360;

(* 2nd Block *)
(* Construction of generators s[1] - s[n - 1] and is[1] - is[n - 1] *)
Clear[t];
nn = n - 1; m = 1; v = {0};
Do[v = Append[v, 0], {k, 1, nn - 1}];
Do[d[i] = ReplacePart[v, 1, i], {i, 1, nn}];
d[m] = ReplacePart[d[m], -t, m];
d[m] = ReplacePart[d[m], 1, m + 1];
s[m] = Table[d[k], {k, 1, nn}];
m = 1; v = {0};
Do[v = Append[v, 0], {k, 1, nn - 1}];
Do[d[i] = ReplacePart[v, 1, i], {i, 1, nn}];
d[m] = ReplacePart[d[m], -1/t, m];
d[m] = ReplacePart[d[m], 1/t, m + 1];
is[m] = Table[d[k], {k, 1, nn}];
Do[v = {0};
  Do[v = Append[v, 0], {k, 1, nn - 1}];
  Do[d[i] = ReplacePart[v, 1, i], {i, 1, nn}];
  d[m] = ReplacePart[d[m], t, m - 1];
  d[m] = ReplacePart[d[m], -t, m];
```

```

    d[m] = ReplacePart[d[m], 1, m + 1];
    s[m] = Table[d[k], {k, 1, nn}], {m, 2, nn - 1}];
Do[v = {0};
  Do[v = Append[v, 0], {k, 1, nn - 1}];
  Do[d[i] = ReplacePart[v, 1, i], {i, 1, nn}];
  d[m] = ReplacePart[d[m], 1, m - 1];
  d[m] = ReplacePart[d[m], -1/t, m];
  d[m] = ReplacePart[d[m], 1/t, m + 1];
  is[m] = Table[d[k], {k, 1, nn}], {m, 2, nn - 1}];
m = nn; v = {0};
Do[v = Append[v, 0], {k, 1, nn - 1}];
Do[d[i] = ReplacePart[v, 1, i], {i, 1, nn}];
d[m] = ReplacePart[d[m], t, m - 1];
d[m] = ReplacePart[d[m], -t, m];
s[m] = Table[d[k], {k, 1, nn}];
m = nn; v = {0};
Do[v = Append[v, 0], {k, 1, nn - 1}];
Do[d[i] = ReplacePart[v, 1, i], {i, 1, nn}];
d[m] = ReplacePart[d[m], 1, m - 1];
d[m] = ReplacePart[d[m], -1/t, m];
is[m] = Table[d[k], {k, 1, nn}];
(* end *)

(* 3rd Block *)
(* Input a braidtype *)
b = s[1].is[2].is[3];

(* 4th Block *)
(* Calculation of Eigenvalues *)
Do[
  theta = 2Pi/ma*k;
  t = Cos[theta] + I*Sin[theta];
  gg = Eigenvalues[N[b]];
  Do[e[i] = Abs[Part[gg, i]], {i, 1, nn}];
  y[k] = Max[Table[e[k], {k, 1, nn}]];
  x[k] = N[theta], {k, 0, ma}];
(* Output *)
g1 = Table[{x[k], y[k]}, {k, 0, ma}];
ListPlot[g1, PlotStyle -> {RGBColor[1, 0, 0]}];
gg = Table[y[k], {k, 0, ma}];
Max[gg]

```

If the eigenfunction of the reduced Burau matrix with  $t = -1$  is needed, delete the fourth block and add the following statements.

```
(* Eigenfunction:  $f(x) = 0$ . Output is  $f(x)$ . *)  
t = -1;  
Det[b - x*IdentityMatrix[n - 1]]
```

# Non-Birkhoff Periodic Orbits in a Circle Mapping

Yoshihiro YAMAGUCHI<sup>1</sup> and Kiyotaka TANIKAWA<sup>2</sup>

<sup>1</sup> Teikyo Heisei University, Ichihara, Chiba 290-0193, Japan.

<sup>2</sup> National Astronomical Observatory, Mitaka, Tokyo 181-8588, Japan.

## Abstract

Birkhoff and Non-Birkhoff types of periodic orbits are defined in circle mappings. The dynamical order relation of non-Birkhoff periodic orbits (NBOs) with period longer than equal 3 is proved. The braids are constructed for NBOs and the topological entropy is estimated.

## 1 Introduction

Dynamics of one dimensional circle mapping offer useful information on periodic orbits, quasi-periodic and chaotic motions.<sup>1-3)</sup> In many cases, systems with one external parameter have been the target of research. The parameter region pertaining to the local motion<sup>4)</sup> or to the Arnold tongues<sup>3)</sup> has been investigated. Properties of systems in the parameter region where local and global motions mix are not made clear. By the global motion, we mean that of revolving the circle. In this situation, the mixing of the local and global motions (this will be called the mixed state) induces complicated phenomena. In this paper, we pay attention to the appearance of periodic orbits (called windows in the bifurcation diagram) in the parameter region of mixed state, and try to estimate the topological entropy of complicated motions.

We consider  $C^0$  mapping  $f$  on circle  $S^1$  defined by

$$\theta_{n+1} = f(\theta_n) \pmod{1}, \quad (1)$$

where  $f(\theta)$  is assumed to satisfy the following conditions.

- [1]  $f(\theta + 1) = f(\theta) + 1$ ,
- [2]  $f(\theta)$  has one local maximum point at  $\theta_{max} \in (0, \theta_c)$  and has one local minimum point at  $\theta_{min} \in (\theta_c, 1)$  for some  $0 < \theta_c < 1$ .
- [3] There exist two fixed points  $\theta_l$  and  $\theta_r$  satisfying  $\theta_c < \theta_l < \theta_r < 1$ . Note that  $\theta_r$  is an unstable fixed point.

Since we discuss periodic orbits revolving the circle, we work in universal cover  $\mathbf{R}^1$  of  $S^1$ . In  $\mathbf{R}^1$ , we use a lift  $\hat{f} : \mathbf{R}^1 \rightarrow \mathbf{R}^1$  of  $f : S^1 \rightarrow S^1$ . The lift  $\hat{f}$  is chosen to keep the fixed points for  $f$  fixed, so it is uniquely defined.

We address the following questions on the circle mappings satisfying [1] – [3], and answer partially to them. What periodic orbits exist in the mixed state, and what types of dynamical order relation between them hold? What is the topological entropy in the system? We use, in obtaining periodic orbits, the standard tools in the one-dimensional mappings such as primitive mappings, covering relations, oriented graphs

and so on. Construction of braids and estimation of topological entropy follow those of our preceding articles.<sup>12),13),14)</sup>

In §2, we introduce several definitions and notation, and prove the dynamical order relation for non-Birkhoff type periodic orbits (NBOs). In §3, the braids for NBOs are constructed, and the topological entropy is estimated. In §4, we give several remarks.

## 2 Dynamical order relation

### 2.1 Definition and notation

Birkhoffness or non-Birkhoffness of periodic orbits has been introduced in two-dimensional twist mappings.<sup>5)</sup> The notion is most relevant to circle mappings. A point  $\hat{\theta} \in \mathbf{R}^1$  is called a  $p/q$ -periodic point for  $\hat{f}$  if

$$\hat{f}^q(\hat{\theta}) = \hat{\theta} + p. \quad (2)$$

The orbit of  $\hat{\theta}$  is  $O(\hat{\theta}) = \{\dots, \hat{f}^{-1}(\hat{\theta}), \hat{\theta}, \hat{f}(\hat{\theta}), \dots\}$ . The extended orbit of  $\hat{\theta}$  is

$$EO(\hat{\theta}) = \{\hat{f}^k(\hat{\theta}) + m : k, m \in \mathbf{Z}\}. \quad (3)$$

A  $p/q$ -periodic point  $\hat{\theta}$  is called Birkhoff if for any  $\hat{r}, \hat{s} \in EO(\hat{\theta})$

$$\hat{r} < \hat{s} \Rightarrow \hat{f}(\hat{r}) < \hat{f}(\hat{s}). \quad (4)$$

If the extended orbit of a periodic point has a couple of points not satisfying Eq.(4), we call it the non-Birkhoff periodic point and its orbit the non-Birkhoff periodic orbit (NBO). From now on, we use the convention  $\hat{\theta}_k = \hat{f}^k(\hat{\theta}_0)$ .

Let us consider a  $p/q$ -periodic orbit  $O(\hat{\theta}_0)$ . If  $\hat{\theta}_k > \hat{\theta}_{k-1}$  and  $\hat{\theta}_k > \hat{\theta}_{k+1}$  hold at some  $k$  ( $1 \leq k < q$ ),  $\hat{\theta}_k$  is called a turning-back point. If  $\hat{\theta}_{k'} < \hat{\theta}_{k'-1}$  and  $\hat{\theta}_{k'} < \hat{\theta}_{k'+1}$  hold at some  $k'$  ( $0 \leq k' < q$ ),  $\hat{\theta}_{k'}$  is called a turning-forward point. We call these the *turning points*. In this paper, we consider NBOs with turning points. Note that there are NBOs with no turning points.<sup>12)</sup> We can choose  $\hat{\theta}_0$  as a starting point of the orbit so that  $\hat{\theta}_1$  be the first turning-back point. Let  $\hat{\theta}_{k_a+1}$  ( $k_a \geq 1$ ) be the first turning-forward point, and  $\hat{\theta}_{k_b+1}$  ( $k_b \geq 1$ ) be the last turning-forward point. If the orbit has only two turning points, then  $k_a = k_b$ . We restrict our attention to NBOs satisfying the condition

$$\hat{\theta}_0 < \hat{\theta}_{k_b+1} \leq \hat{\theta}_{k_a+1} < \hat{\theta}_1. \quad (5)$$

Later, we categorize NBOs with  $2n$  ( $n \geq 1$ ) turning points by  $k_a$  and the number of turning points.

If a closed interval  $I \subset \mathbf{R}^1$  contains one turning point, we denote it by  $\tilde{I}$ . Let  $I_1$  and  $I_2$  be two closed intervals satisfying  $\text{Int}(I_1) \cap \text{Int}(I_2) = \emptyset$  where  $\text{Int}(I)$  is the interior of  $I$ . If the relation  $I_2 \subset \hat{f}(I_1)$  holds, we write  $I_1 \succ I_2$ , and we say  $I_1$  covers  $I_2$ . We also call  $I_1 \succ I_2$  the oriented graph of intervals, or the covering relation.

The rotation number  $\nu$  of an orbit of  $\hat{\theta} \in \mathbf{R}$  is defined by

$$\nu = \lim_{n \rightarrow \infty} \sup \frac{\hat{f}^n(\hat{\theta}) - \hat{\theta}}{n}. \quad (6)$$

A  $1/q$ -NBO with  $k_a$  and two turning points is denoted by

$$\left(\frac{1}{q}\right)_{k_a}^2. \quad (7)$$

If the existence of a  $1/q$ -NBO with  $k_a$  and two turning points implies the existence of a  $1/q'$ -NBO with  $k'_a$  and two turning points, then we write as

$$\left(\frac{1}{q}\right)_{k_a}^2 \rightarrow \left(\frac{1}{q'}\right)_{k'_a}^2, \tag{8}$$

and we simply say that  $(1/q)_{k_a}^2$  implies  $(1/q')_{k'_a}^2$ . We also say that  $(1/q')_{k'_a}^2$  is dominated by  $(1/q)_{k_a}^2$ .

### 2.2 NBOs with period-3

We assume that there exists a  $1/3$ -NBO in the universal cover  $\mathbf{R}^1$  of  $\mathbf{S}^1$  in some parameter set satisfying the orbital order

$$0 \leq \hat{\theta}_0 < \hat{\theta}_2 < \hat{\theta}_1 < 1 \leq \hat{\theta}_3 = \hat{\theta}_0 + 1, \tag{9}$$

where  $\hat{\theta}_i = \hat{f}^i(\hat{\theta}_0)$  and  $k_a = k_b = 1$ . This orbit is  $(1/3)_1^2$  following the convention adopted in §2.1. Since the orbit turns back at  $\hat{\theta}_1$ , we have  $\hat{\theta}_l < \hat{\theta}_1 < \hat{\theta}_r$ . We define four intervals:

$$I_1 = [0, \hat{\theta}_{\max}], \tag{10}$$

$$I_2 = [\hat{\theta}_{\max}, \hat{\theta}_{\min}], \tag{11}$$

$$I_3 = [\hat{\theta}_{\min}, 1], \tag{12}$$

$$SI_1 = [1, \hat{\theta}_{\min} + 1]. \tag{13}$$

One observes that a  $(1/3)_1^2$  exists if the covering relation

$$I_1 \succ \tilde{I}_3 \succ \tilde{I}_2 \succ SI_1 \tag{14}$$

holds. In order to guarantee (14), we use three conditions.

$$\hat{f}(0) \leq \hat{\theta}_{\min}, \tag{15}$$

$$\hat{f}(\hat{\theta}_{\min}) \leq \hat{\theta}_{\max}, \tag{16}$$

$$\hat{f}(\hat{\theta}_{\max}) \geq 1 + \hat{\theta}_{\max}. \tag{17}$$

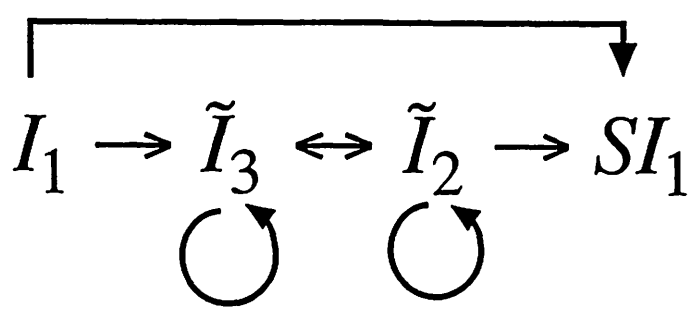


Fig. 1. The oriented graph of the intervals.



Equations (15)–(17) give an oriented graph of intervals shown in Fig. 1. This graph contains the covering relation (14) as a subgraph. Analysis of Fig. (1) gives Proposition 1.

**Proposition 1.**  $(1/3)_1^2$  implies  $1/n$ -NBOs ( $n \geq 4$ ).

**Proof.** A cycle  $I_1 \succ \tilde{I}_3 \succ I_3 \succ \tilde{I}_2 \succ SI_1$  gives a  $1/4$ -NBO. Similarly, we can construct a cycle with period longer than 4. (Q.E.D.)

## 2.3 Theorem

In this section, we elaborate the appearance order or the dynamical order of NBOs found in Proposition 1.

**Lemma 1.** For  $1/q$ -NBOs with  $q \geq 3$ ,  $k_a(\geq 1)$  and with two turning points for the circle mapping  $f$  satisfying [1]–[3], the following dynamical order relation holds.

$$\begin{array}{ccccccc}
 (1/3)_1^2 & \rightarrow & (1/4)_1^2 & \rightarrow & (1/5)_1^2 & \rightarrow & (1/6)_1^2 \rightarrow \cdots \\
 \downarrow & & \downarrow & & \downarrow & & \downarrow \\
 (1/5)_2^2 & \rightarrow & (1/6)_2^2 & \rightarrow & (1/7)_2^2 & \rightarrow & (1/8)_2^2 \rightarrow \cdots \\
 \downarrow & & \downarrow & & \downarrow & & \downarrow \\
 (1/7)_3^2 & \rightarrow & (1/8)_3^2 & \rightarrow & (1/9)_3^2 & \rightarrow & (1/10)_3^2 \rightarrow \cdots \\
 \downarrow & & \downarrow & & \downarrow & & \downarrow \\
 \vdots & & \vdots & & \vdots & & \vdots
 \end{array}$$

**Remarks.** We will use the matrix notation to specify the position of NBOs in Lemma 1. Regarding the above table as a matrix, we take  $(1/(2i+j))_i^j$  as the  $(i, j)$  ( $i, j \geq 1$ ) element. Thus, for example, we say  $(1, 1) \rightarrow (1, 2)$  if the forcing relation  $(1/3)_1^2 \rightarrow (1/4)_1^2$  holds.

**Proof:** In order to prove Lemma 1, we use *the primitive tight mapping*<sup>(2),6),7)</sup> which is the simplest piecewise linear mapping having a periodic orbit with the given orbital order. Using the information of the NBO defined by Eq. (9), we can construct the primitive tight mapping  $\hat{F}$ , shown in Fig. 2. In the figure, the relation  $SI_i = I_i + 3$  holds and the orbital order of the period-3 orbit is expressed by  $0 \rightarrow 2 \rightarrow 1 \rightarrow 3$ . Figure 3 is an oriented graph showing the covering relation between intervals. Each interval  $I_i$  in Fig. 2 has a unit length. We can change the length and use another continuous function connecting adjacent two points. However, new oriented graph for such mappings contains Fig. 3 as a subgraph. The oriented graph shown in Fig. 3 implies the existence of a cycle from  $I_1$  to  $SI_1$ .

Using the oriented graph, we can determine periodic orbits dominated by  $(1/3)_1^2$ . The following cycle gives a period-4 orbit.

$$I_1 \succ \tilde{I}_3 \succ \tilde{I}_2 \succ I_3 \succ SI_1. \quad (18)$$

The rotation number of the orbit is  $1/4$  since  $SI_1 \subset \hat{F}^4(I_1)$ . The orbit is non-Birkhoff because there are turning points in  $I_2$  and  $I_3$ . Obviously we have  $k_B = 1$ . Then the orbit is  $(1/4)_1^2$ . Thus  $(1, 1) \rightarrow (1, 2)$  is proved. It is to be noted that the orbital points of  $(1/4)_1^2$  are not at the endpoints of intervals since these are points of  $(1/3)_1^2$ . (This fact is true for cases treated below.)

We have two cycles for period-5 orbit satisfying  $SI_1 \subset \hat{F}^5(I_1)$ .

$$I_1 \succ \tilde{I}_3 \succ I_2 \succ \tilde{I}_2 \succ I_3 \succ SI_1, \tag{19}$$

$$I_1 \succ \tilde{I}_3 \succ I_3 \succ \tilde{I}_2 \succ I_2 \succ SI_1. \tag{20}$$

These give the same 1/5-NBOs with  $k_a = 2$ . Thus  $(1, 1) \rightarrow (2, 1)$  is proved.

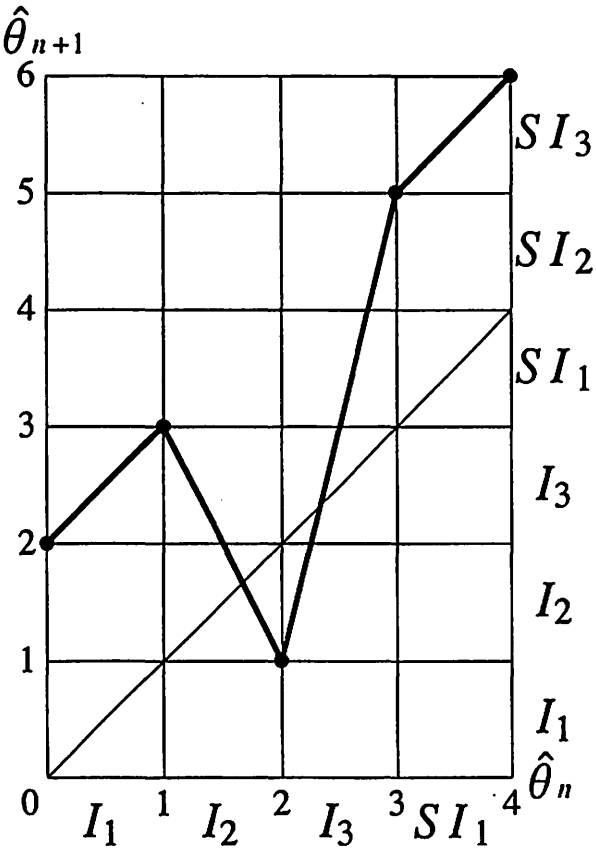


Fig. 2. Primitive tight mapping constructed by  $(1/3)_1^2$ .

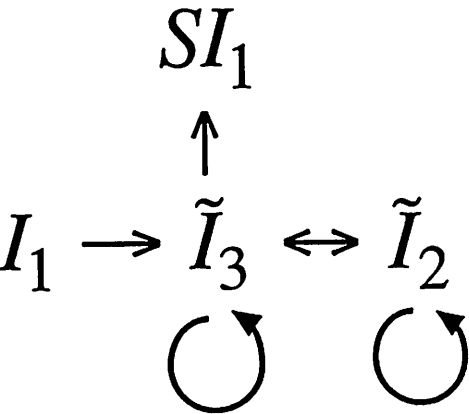


Fig. 3. The oriented graph of the intervals, and the unnecessary arrows are omitted.

Next we prove the general cases of  $(i, j) \rightarrow (i, j + 1)$  and  $(i, j) \rightarrow (i + 1, j)$ . The existence of  $(1/(2i + j))_i^2$  gives the primitive tight mapping shown in Fig. 4, in which

$\hat{F}(I_1^l) = I_{i+2}^l$  and  $I_{i+2}^r \subset \hat{F}(I_2^l)$  hold. The relation between the intervals is displayed in Fig. 5, where unnecessary arrows are omitted. We want to find the cycle from  $I_1^l$  to  $SI_1^l$ . The shortest orbit has a period  $(2i+j+1)$  and  $k_i = i$ . Thus the relation  $(i, j) \rightarrow (i, j+1)$  is proved. There are two ways to construct an orbit with period  $(2i+j+2)$  and  $k_B = i+1$ . If we use  $I_{i+2}^l$  twice, then we have  $I_1^l \succ \tilde{I}_{i+2}^l \succ I_{i+2}^l \succ I_{i+1} \succ \dots \succ SI_1^l$ . If we use  $I_2$  twice, we have  $I_1^l \succ \dots \succ I_3 \succ I_2 \succ \tilde{I}_2 \succ \dots \succ SI_1^l$ . These orbits are expressed by  $(1/(2i+j+2))_{i+1}^2$ . As a result,  $(i, j) \rightarrow (i+1, j)$  is proved. The proof completes. Q.E.D.

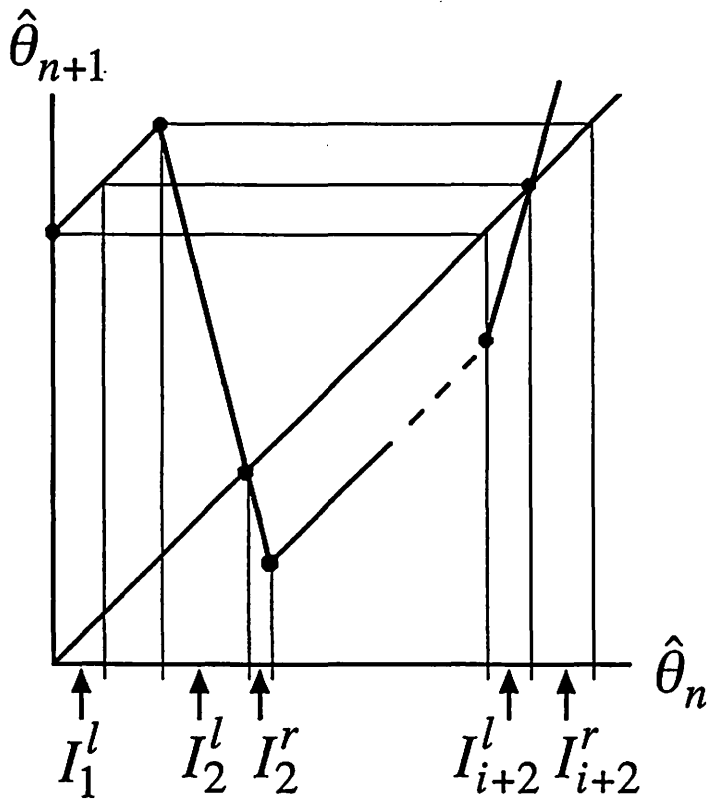


Fig. 4. Primitive tight mapping constructed by  $(1/(2i+j))_i^2$  where  $I_2 = I_2^l \cup I_2^r$  and  $I_{i+2} = I_{i+2}^l \cup I_{i+2}^r$ .

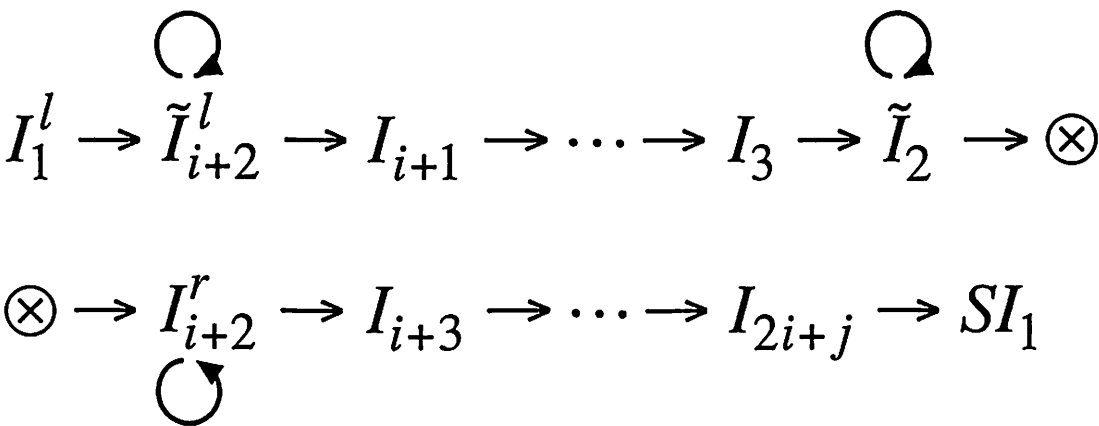


Fig. 5 The oriented graph of the intervals, and the unnecessary arrows are omitted.

We have two ways to obtain NBOs of higher periods. One observes this by looking at Fig.5 carefully. For example, the cycle

$$I_1 \succ \tilde{I}_{i+2}^l \succ \cdots \succ I_3 \succ \tilde{I}_2 \succ \tilde{I}_2 \succ I_{i+2}^r \succ \cdots \succ SI_1 \quad (21)$$

constructed by using the oriented graph shown in Fig. 5 gives an NBO with four turning points,  $k_a = i$  and  $k_b = i + 2$ . Using  $I_2$  repeatedly, we can prove the existence of NBOs with  $2n$  ( $n \geq 2$ ) turning points. Next using  $I_{i+2}^l$  or  $I_{i+2}^r$ , we construct cycles of longer periods without increasing the number of turning points. This property comes from the fact that either of  $I_2$  and  $I_{i+2}$  contains a fixed point.

In the oriented graph constructed by the primitive tight mapping for  $(1/q)_{k_a}^{2n}$ , the shortest cycle of an NBO from  $I_1$  to  $SI_1$  with  $2n$  turning points not using the edgepoints is  $(q + 1)$ . In fact, the orbit of  $(1/q)_{k_a}^{2n}$  passes one edgepoint of  $I_{i+2}$ , and the cycle not using edgepoints passes  $I_{i+2}$  twice such that it passes  $I_{i+2}^l$  to turn back and  $I_{i+2}^r$  to go out from the localized region. If we increase the number of turning points by 2, the period of a new cycle increases by 3. Summarizing these facts, we have Lemma 2.

**Lemma 2.** For  $i \geq 1$  and  $j \geq 1$ , the forcing relations hold.

$$\left(\frac{1}{2i+j}\right)_i^2 \rightarrow \left(\frac{1}{2i+j+3}\right)_i^4 \rightarrow \left(\frac{1}{2i+j+6}\right)_i^6 \rightarrow \cdots \quad (22)$$

Here we construct the order relation of NBOs with  $2m$  ( $m \geq 1$ ) turning points.

**Lemma 3.** For NBOs with  $\nu = 1/q$  ( $q \geq 3$ ),  $k_a (\geq 1)$  and  $2m$  turning points in the circle mapping  $f$  satisfying [1]-[3], the following dynamical order relation holds.

$$\begin{array}{ccccccc} (1/(3m))_1^{2m} & \rightarrow & (1/(3m+1))_1^{2m} & \rightarrow & (1/(3m+2))_1^{2m} & \rightarrow & \cdots \\ \downarrow & & \downarrow & & \downarrow & & \\ (1/(3m+2))_2^{2m} & \rightarrow & (1/(3m+3))_2^{2m} & \rightarrow & (1/(3m+4))_2^{2m} & \rightarrow & \cdots \\ \downarrow & & \downarrow & & \downarrow & & \\ (1/(3m+4))_3^{2m} & \rightarrow & (1/(3m+5))_3^{2m} & \rightarrow & (1/(3m+6))_3^{2m} & \rightarrow & \cdots \\ \vdots & & \vdots & & \vdots & & \end{array}$$

**Proof.** The proof is similar to that of Lemma 1, and thus is omitted.(Q.E.D.)

From now on, the dynamical ordering in Lemma 3 will be called the dynamical ordering on the  $m$ -th floor.  $m$  is the number of turning-back or equivalently turning-forward points in an orbit. Consequently, the dynamical ordering on the  $m$ -th floor is that for orbits with  $2m$  turning points. Using Lemma 2, we can construct the dynamical ordering between the orderings on adjacent two floors. We introduce three dimensional notation  $(i, j, m)$ , and specify the position of NBOs, for example,  $(1/(3m))_1^{2m}$  at  $(1, 1, m)$ ,  $(1/(3m+1))_1^{2m}$  at  $(1, 2, m)$  and  $(1/(3m+2))_2^{2m}$  at  $(2, 1, m)$ . As a result, we have theorem 1 on the three dimensional dynamical ordering for NBOs.

**Theorem 1.** The following dynamical orderings for NBOs hold.

$$(i, j, m) \rightarrow (i, j + 1, m), \quad (23)$$

$$(i, j, m) \rightarrow (i + 1, j, m), \quad (24)$$

$$(i, j, m) \rightarrow (i, j, m + 1), \quad (25)$$

where  $i, j, m \geq 1$  and an NBO of  $(i, j, m)$  element has a rotation number  $\nu = 1/(2i + j + 3(m - 1))$ .

## 2.4 Existence of BOs and their dynamical ordering

According to a theorem by Boyland,<sup>6),7)</sup> if a  $1/n$ -NBO exists, then there is a rotation band defined by  $[0/1, 1/(n - 1)]$ , and there exists a BO with a rotation number in the rotation band. Combining our results and this theorem, we have Proposition 2.

**Proposition 2.** A  $1/n$ -NBO ( $n \geq 3$ ) implies a  $1/(n - 1)$ -BO.

Let us denote a  $1/q$ -Birkhoff periodic orbit ( $q \geq 2$ ) by  $(1/q)_B$ . There exist  $q$  points in the region  $[0, \theta_l) \cup (\theta_r, 1) \in \mathbf{S}^1$ . Using this fact and the condition [3] of  $f$ , we can determine the dynamical order relation of them.

**Proposition 3.** The following dynamical ordering for BOs holds.

$$(1/2)_B \rightarrow (1/3)_B \rightarrow (1/4)_B \rightarrow \cdots$$

**Proof.** We prove the relation  $(1/2)_B \rightarrow (1/3)_B$ . The others are similarly proved and then the proof is omitted. The primitive tight mapping allowing  $(1/2)_B$  is displayed in Fig. 6 where two fixed points are located in  $I_2$  due to [3] and this interval is divided into  $I_2^l$  and  $I_2^r$ . The oriented graph is obtained in Fig. 7. The existence of loop  $I_2^l \succ I_2^l$  depends on the parameters and thus this loop is omitted. There exist two cycles not containing turning points  $I_1 \succ I_2^l \succ I_2^r \succ SI_1$  and  $I_1 \succ I_2^r \succ I_2^l \succ SI_1$ . This implies  $(1/3)_B$ . (Q.E.D.)

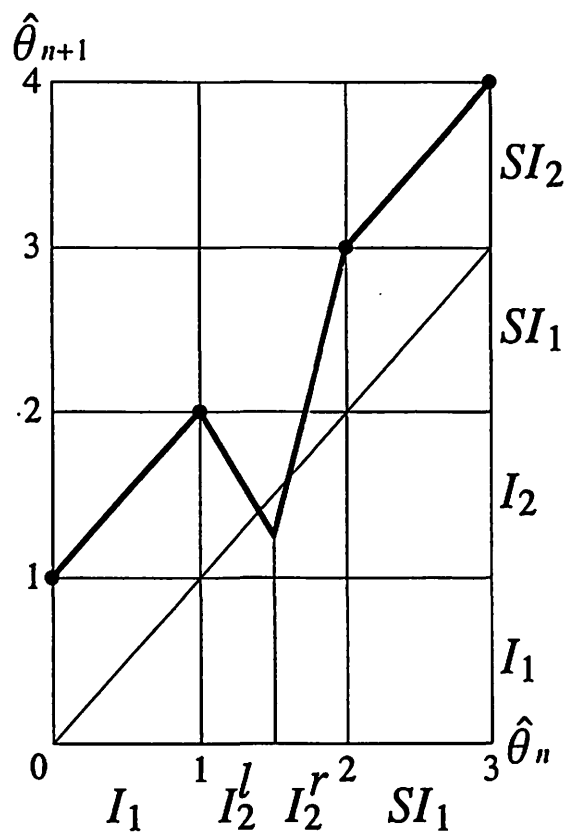


Fig. 6. Primitive tight mapping allowing  $(1/2)_B$ .

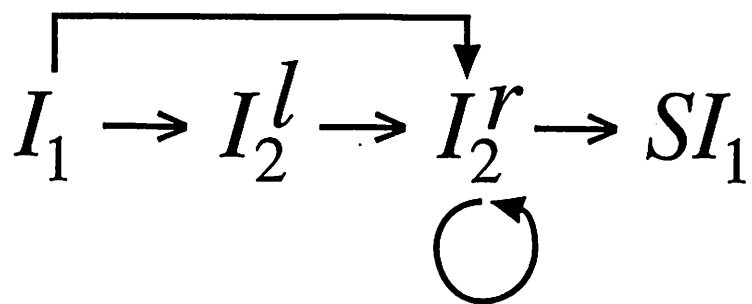


Fig. 7. The oriented graph of intervals constructed by using Fig. 6.

### 3 Braid and topological entropy

#### 3.1 Braid

We construct braids from periodic orbits by using information on the order of orbital points.<sup>6-8,12)</sup> From now on, we use NBOs in the first floor  $(i, j, 1)$ . The periodic points are located in the circle. Thus we connect  $\theta_n$  and  $\theta_{n+1}$  by an arrow. Examples are displayed in Fig. 6 where symbol  $i$  stands for  $\theta_i$ . In Fig. 8(a), an arrow from 1 to 2 intersects that from  $k$  to  $k+1$ . In the braid, a string from 1 to 2 does not intersect that from  $k$  to  $k+1$ . This displays a braid for BO. Figs. (b) and (c) correspond to braids for NBOs. In each braid, two strings intersect each other. If the orbit (fast orbit) goes over the slow orbit or the backward orbit, the string of fast orbit passes behind the string of slow orbit or backward orbit. Using this rule, two braids of Figs. (b) and (c) are constructed.

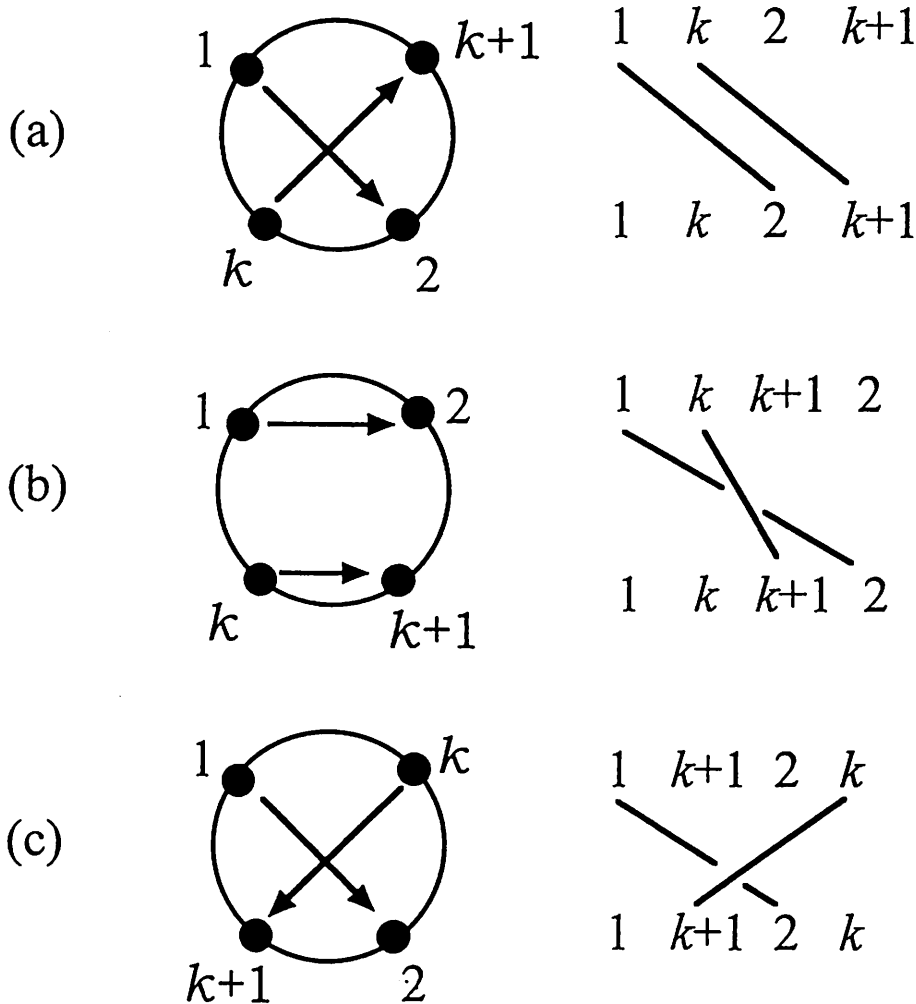


Fig. 8 Fig.(a) shows a part of BO and its braid, and Figs. (b) and (c) display the intersection of braids due to non-Birkhoffness.

We show two braids for 1/5-NBOs expressed by the generator of braid.<sup>11)</sup>

$$\beta(1, 3, 1) = \sigma_2^{-1} \sigma_1^{-1} \zeta_5 = \sigma_1^{-2} \zeta_5, \quad (26)$$

$$\beta(2, 1, 1) = \sigma_2^{-1} \sigma_3^{-1} \sigma_1^{-1} \sigma_2^{-1} \zeta_5 = \sigma_1^{-1} \sigma_2^{-1} \sigma_2^{-1} \sigma_1^{-1} \zeta_5 \quad (27)$$

where  $\zeta_5 = \sigma_4 \cdots \sigma_1$ , and Reidemeister and Markov moves<sup>9)</sup> are operated to derive the second expression. The difference of braids comes from that of  $k_B$ . We have the braid for an NBO of  $(i, j, 1)$  element in Lemma 1.

$$\beta(i, j, 1) = \zeta_{i+1}^{-1} \rho_{i+1}^{-1} \zeta_{2i+j} \quad (28)$$

where  $\rho_i = \sigma_1 \cdots \sigma_{i-1}$  and  $\zeta_i = \sigma_{i-1} \cdots \sigma_1$ . The part  $\zeta_{2i+j}$  is the braid of 1/5-Birkhoff orbit and the part  $\zeta_{i+1}^{-1} \rho_{i+1}^{-1}$  represents the non-Birkhoffness.

### 3.2 Topological entropy

We show the procedure to estimate the lower bound of topological entropy by using NBOs in Lemma 1. First, we construct the Burau matrix representation<sup>8),9)</sup> corresponding to the braid of an NBO. Next, we calculate the eigenvalues of Burau matrix. The maximum ( $\lambda_{max}$ ) of the absolute values of eigenvalues gives the lower bound of topological entropy,<sup>10),11)</sup> expressed by  $h = \ln \lambda_{max}$ .

Numerical results for topological entropy are shown in Table I. The maximum value is  $\ln(\sqrt{5} + 3)/2 = 0.962 \cdots$  estimated by using  $(1/3)_1^2$ .<sup>8)</sup> The entropy  $h(1, j, 1)$  is not a strictly decreasing function of  $j$ , but it accumulates at  $\ln 2$  in the limit  $j \rightarrow \infty$ . This fact implies that the entropy is larger than  $\ln 2$  for finite  $j$ . Finally it is noted that we can not determine the forcing relation of NBOs by using the topological entropy estimated in Table I.

**Table I:** Topological entropy  $h(i, j, 1)$  calculated by using the program in Ref. 12).

	$j = 1$	2	3	4	5	6	7
$i = 1$	0.962	0.776	0.767	0.713	0.714	0.694	0.698
2	0.652	0.575	0.558	0.530	0.512	0.508	0.499
3	0.562	0.499	0.491	0.462	0.460	0.445	0.446
4	0.465	0.422	0.416	0.398	0.389	0.382	0.373
5	0.413	0.379	0.375	0.355	0.354	0.342	0.343
6	0.364	0.338	0.334	0.321	0.315	0.310	0.304
7	0.332	0.310	0.307	0.293	0.292	0.282	0.283

## 4 Remarks

Suppose that  $f$  has one bifurcation parameter  $a$ , and assume the existence of a critical value  $a_c$  such that the mixing of the local and global motions exists at  $a > a_c$ . The transition from a local state to a mixed state is called the *crisis*. Let  $a_c(1/q|_{k_a}^2)$  with  $q = 2i + j$  be a critical value at which an NBO of  $(1/q)_{k_a}^2$  appears due to the tangent bifurcation. For BOs, the critical values  $a_c(1/q|_B)$  ( $q \geq 2$ ) are also defined.

In the limit  $i \rightarrow \infty$ ,  $\theta_1$  tends to  $\theta_r$  from the left side. In the limit  $j \rightarrow \infty$ ,  $\theta_{i+2}$  tends to  $\theta_r$  from the right side. The converged situation is that of crisis. We have the relation:

$$\lim_{i \rightarrow \infty} a_c(1/q|_i^2) = a_c \text{ for fixed } j, \quad (29)$$

$$\lim_{j \rightarrow \infty} a_c(1/q|_i^2) = a_c \text{ for fixed } i, \quad (30)$$

$$\lim_{q \rightarrow \infty} a_c(1/q|_B) = a_c^* \quad (31)$$



where a critical value  $a_c^*$  is the value for which  $f(\theta_{max}) = \theta_r$  holds. The topological entropy is larger than  $\ln 2$  at  $a > a_c$  since the limiting value is  $\ln 2$  for  $i = 1$  and  $j \rightarrow \infty$ , and Eqs. (29) and (30) hold.

We have used only the continuity of mapping function to prove Lemmas. Then Theorem 1 holds for a climbing sine-mapping (CSM) defined by

$$f(\theta) = \theta + \frac{K}{2\pi} \sin 2\pi\theta + \Omega, \quad (32)$$

where  $K > 0$  and  $\Omega \geq 0$ . This mapping satisfies the conditions [1]-[3] where  $\theta_c = 1/2$ . In the case that  $\Omega$  is fixed, we can regard  $K$  as a bifurcation parameter  $a$  mentioned above. Thus Eqs. (29)-(31) hold for CSM. There exists the parameter region satisfying  $a_c = a_c^*$ . Using CSM, we can draw the bifurcation diagram and confirm periodic windows corresponding to NBOs in Theorem 1 and to BOs in Proposition 3. However the structure of windows after the crisis is beyond all imagination.

The structure of dynamical ordering in Lemma 1 is similar to those derived in the standard mapping,<sup>12)</sup> the standard-like mappings,<sup>13)</sup> and the forced oscillator.<sup>14)</sup> The dynamical ordering similar to Lemma 1 may hold in the systems possessing the mixed state.

## References

- 1) C. Robinson, *Dynamical Systems* (CRC Press, 1999).
- 2) L. Alsedà, J. Llibre and M. Misiurewicz, *Combinatorial Dynamics and Entropy in Dimension One* (World Scientific, 1993).
- 3) W. de Melo and S. van Strien, *One-Dimensional Dynamics*, (Springer-Verlag, 1993).
- 4) T. Horita, H. Hata and H. Mori, Prog. Theor. Phys. **84** (1990), 558.
- 5) A. Katok, Ergod. Theor. & Dynam. Sys. **2** (1982), 185; G. R. Hall, Ergod. Theor. & Dynam. Sys. **4** (1984), 585.
- 6) P. Boyland, Contemp. Math. **81** (1988), 119.
- 7) P. Boyland, Topology and its Appl. **58** (1994), 223.
- 8) T. Matsuoka, in *Dynamical System 1* (World Scientific, 1986), p. 58; Contemp. Math. **152** (1993), 229. See also Bussei Kenkyu(Kyoto) **67** (1996), 1.
- 9) S. Moran, *The Mathematical Theory of Knots and Braids* (North-Holland, 1983).
- 10) D. Fried, in *Geometric Dynamics*. ed. J. Palis Jr. Lecture Notes in Mathematics **1007** (Springer-Verlag, 1983). p. 261.
- 11) B. Kolev, C. R. Acad. Sci. Paris, **309**, Ser. I (1989), 835.
- 12) Y. Yamaguchi and K. Tanikawa (Submitted to Prog. Theor. Phys.).
- 13) K. Tanikawa and Y. Yamaguchi, *Chaos* **12** (2002), 33.
- 14) Y. Yamaguchi and K. Tanikawa, Prog. Theor. Phys. **106** (2001), 1097.

# Relaxation in Hamiltonian systems with long-range interactions

YAMAGUCHI Y. Yoshiyuki\*

Department of Applied Mathematics and Physics,  
Kyoto University, 606-8501, Kyoto, Japan

## Abstract

Relaxation process of kinetic energy to canonical temperature is investigated through anomalous diffusion and local Lyapunov exponent in Hamiltonian systems with long-range interactions, in which a second order phase transition occurs. We find, near and below the critical energy, that (i) anomalous diffusion occurs even in equilibrium states, and (ii) the value of local Lyapunov exponent goes to the value of Lyapunov exponent at the earlier time than the relaxation time of the kinetic energy.

## 1 はじめに

長距離相互作用をしている系は、臨界現象などの協同現象を起こし非常に興味深い。もしこれらの系が拡張性を持つ、つまり自由度に比例してエネルギーが増えるような系であれば、平衡状態は統計力学によって知ることができる。一方で、非平衡状態から平衡状態への緩和過程など、系の時間発展を追うためには正準方程式を数値的に積分するという方法がある。ここでは、系のダイナミクスを理解するため、後者によって緩和過程を調べる。

長距離相互作用ハミルトン系における緩和過程は、例えば重力シートモデル

$$H(q, p) = \frac{1}{2} \sum_{i=1}^N p_i^2 + (2\pi G m^2) \sum_{i < j} |x_i - x_j| \quad (1)$$

---

\*yyama@i.kyoto-u.ac.jp

において詳しく研究されており、緩和は何段階かに分かれて行われることが報告されている [TGK96]。また、平均場 XY モデル (Hamiltonian Mean-Field model)

$$H(q, p) = \frac{1}{2} \sum_{i=1}^N p_i^2 + \frac{1}{2N} \sum_{i,j=1}^N [1 - \cos(q_i - q_j)] \quad (2)$$

においては、自由度と温度の関数としての緩和時間が示されたり [AR95]、異常拡散との関係が議論されたりしている [LRR99b]。後者においては、運動エネルギーがカノニカル温度に向けて緩和している間は異常拡散が起こるが、緩和してしまうと異常拡散は起こらないとしている。しかし、重力シートモデルにおいて緩和が何段階かに分かれていることを考えると、温度の緩和と異常拡散がそれほど密接に関係しているかどうかは確認すべき事項であると思われる。そこで本報告の目的の一つは、温度が平衡に達した後にも異常拡散が見られるかどうかを調べることである。

一方、軌道不安定性解析、いわゆるリアプノフ解析によって緩和過程を調べる研究もある。大域結合している各粒子が平面上の卵パックのような形をしたポテンシャルを持っている系

$$H(q, p) = \frac{1}{2} \sum_{i=1}^N (p_{i,x}^2 + p_{i,y}^2) + \frac{1}{2N} \sum_{i,j}^N [3 - \cos(x_i - x_j) - \cos(y_i - y_j) - \cos(x_i - x_j) \cos(y_i - y_j)] \quad (3)$$

において、局所リアプノフ数の時間発展と温度ならびにオーダーパラメータの時間発展の様子が調べられている [TA99]。温度やオーダーパラメータが緩和を開始する時刻に軌道不安定性が増大し始め、緩和が終了すると不安定性も一定値に落ち着くことが示された。しかし、系 (3) は一次相転移系であるため、二次相転移系である系 (2) で同じことが観測されるかどうかは自明ではない。これを調べるのが第二の目的である。

## 2 モデルと初期条件

考えるモデルは、系 (2) である。ここで、 $\mathbf{m}_j = (\cos q_j, \sin q_j)$  を導入すると、

$$H(p, q) = \frac{1}{2} \sum_{j=1}^N p_j^2 + \frac{1}{2N} \sum_{i,j=1}^N [1 - \mathbf{m}_i \cdot \mathbf{m}_j]$$

と書けるため、位相  $q_j$  を持つ単位長さの回転子が内積相互作用をしている系と見ることができる。ポテンシャルを  $N$  で割っているのは拡張性を保つためである [AT98]。

オーダーパラメータとして、つぎの量を導入する。

$$\mathbf{M} = \frac{1}{N} \sum_{i=1}^N \mathbf{m}_i = (M \cos \phi, M \sin \phi) \quad (4)$$

系 (2) は大域結合系であるため、運動方程式は平均場  $M$ 、すなわち  $M$  と  $\phi$  のみを用いて

$$\frac{d^2 q_i}{dt^2} = -M \sin(q_i - \phi)$$

と書ける。

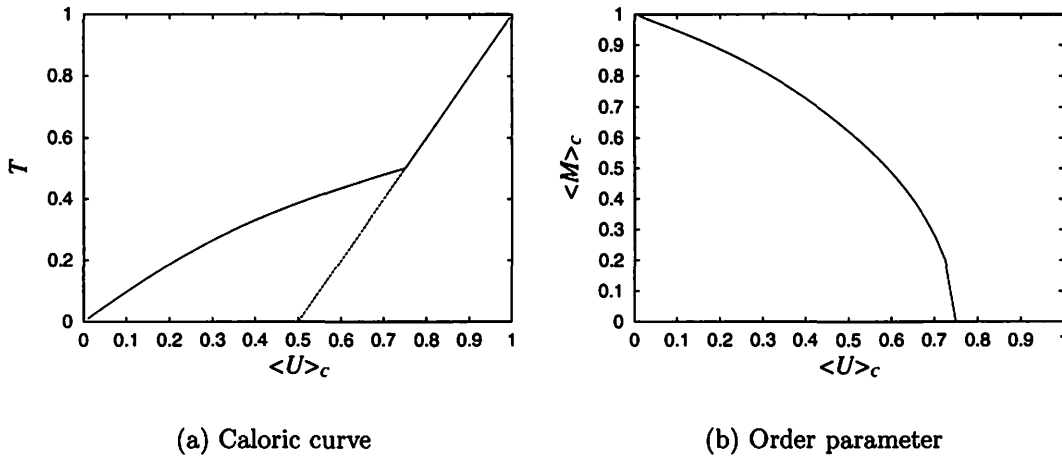


図 1: Fig.1. (a) Caloric curve (temperature  $T$  vs. canonical average of energy  $\langle U \rangle_C$ ). The solid curve is obtained from both of canonical and microcanonical ensembles, and we find a jump in the specific heat  $C_V = (\partial T / \partial \langle U \rangle_C)^{-1}$  at  $\langle U \rangle_C = 0.75$ . The dashed line is another stationary branch obtained from microcanonical ensemble, but it is not stable since the corresponding entropy is not maximal. (b) The modulus of order parameter  $\langle M \rangle_C$  vs. energy  $\langle U \rangle_C$ .

初期条件は以下のようにして設定する。系 (2) は拡張性を持つため統計力学を適用することができる [LRR99a]。カノニカル統計によって得たカロリック曲線、つまり温度をエネルギーの関数として書いた曲線より、 $\langle U \rangle_C = \langle E \rangle_C / N = 3/4$  において二次相転移を起こすことがわかる (図 1(a))。ここに、 $\langle \cdot \rangle_C$  はカノニカル統計による平均を表す。一方でミクロカノニカル統計を用いると、エントロピーの停留点は極大と極小の二種類あることがわかる。安定なブランチであるエントロピー極大条件からはカノニカル統計と同じカロリック曲線が得られるが、メタ安定なブランチであるエントロピー極小条件からはそれとは別の  $M = 0$  となるブランチ  $\langle T \rangle_{MC} = 2U - 1$  が得られる [AHR02]。記号  $\langle \cdot \rangle_{MC}$  はミクロカノニカル統計

による平均を表す。そこで数値計算の初期条件としてはこのメタ安定なブランチを選択し、安定なブランチへの緩和過程を観察することとする。

カノニカル平衡状態においては、 $\langle K \rangle_C$  を運動エネルギーのカノニカル平均値とすると  $T = 2\langle K \rangle_C / N$  なる関係が成り立つため、緩和の様子は運動エネルギー (の  $2/N$  倍である)  $2K(t)/N$  の時間変化を通して観察する。

また、異常拡散については、位相  $q_i$  の平均分散

$$\sigma_q^2(t; \tau) = \frac{1}{N} \sum_{i=1}^N (q_i(t) - q_i(\tau))^2 \quad (5)$$

を調べる。ここで、分散を調べ始める時刻を  $t = 0$  と限定せず  $\tau$  としていることに注意されたい。通常拡散では分散が  $t$  に比例して大きくなるが、異常拡散では  $t^\alpha (\alpha \neq 1)$  に比例して大きくなる。 $\alpha > 1$  のことを速い拡散、 $\alpha < 1$  のことを遅い拡散などと言う。

以下の計算では、自由度を  $N = 1000$  とし、エネルギーを臨界エネルギー  $U = E/N = 0.75$  よりわずかに小さい  $U = 0.69$  とする。数値計算のアルゴリズムは、4 次のシンプレクティックインテグレーター [Yos93] を用い、時間刻み幅は  $\Delta t = 0.2$  と設定した。このときエネルギーの誤差  $\Delta E$  は  $|\Delta E/E| < 5 \times 10^{-7}$  となる。

### 3 運動エネルギーの緩和と異常拡散

本節では、運動エネルギー  $2K(t)/N$  の緩和と位相の異常拡散について述べる。なお、これら2つの量の関係を議論している文献 [LRR99b] においては、 $2K(t)/N$  のかわりに

$$\frac{1}{t} \int_0^t \frac{2K(s)}{N} ds$$

が用いられているが、ここでは異常拡散と相空間内の状態との関連を重視して、時間平均ではなく各時刻での運動エネルギーを観察することとする。

図 2(a) には運動エネルギーがカノニカル温度に緩和する様子が、図 2(b) には位相の分散  $\sigma_q^2(t; 0)$  の時間変化がそれぞれ示されている。運動エネルギーは、およそ  $t_{\text{relax}} = 6 \times 10^4$  でカノニカル温度に緩和しており、一方位相の分散は  $t_{\text{cross}} = 2 \times 10^5$  で異常拡散から通常拡散へ移行しており、おおざっぱには温度の緩和と異常拡散とが同じ時間スケールで起こっているように見える。これら二つの現象が密接に関係しているかどうかを調べるため、 $\sigma_q^2(t; \tau)$  において分散を計算し始める時刻  $\tau$  を 0 でない次の 2 つの時刻に設定してみる。一つは運動エネルギーの緩和が終了した時刻  $t_{\text{relax}}$  で、もう一つは異常拡散から通常拡散へ移行した時刻  $t_{\text{cross}}$  よりも先の時刻  $10^6$  である。

それぞれの結果を図 3 に示した。これら 2 つの時刻  $\tau$  に対しても、 $\sigma_q^2(t; \tau)$  は異常拡散から通常拡散への移行を示しており、異常拡散がかならずしも緩和過程においてのみ現れるわけではないことがわかる。

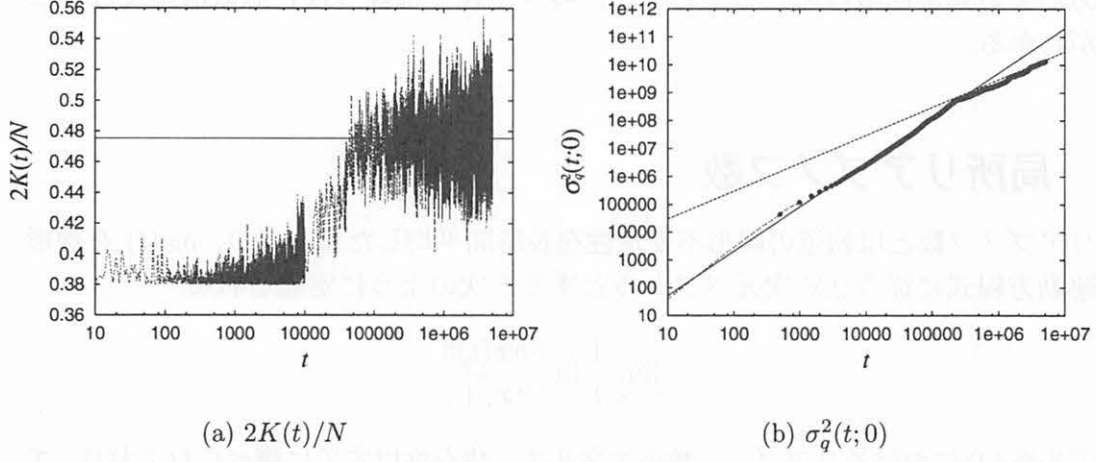


Fig.2. (a) Time series of  $2K(t)/N$ . The horizontal line represents the value of canonical temperature. (b) Time series of  $\sigma_q^2(t;0)$ . We find the crossover time  $t_{cross}$  around  $2 \times 10^5$ . The variance  $\sigma_q^2(t;0)$  is approximated as  $t^{1.6}$  and  $t$  before and after  $t_{cross}$ , respectively.  $N = 1000$  and  $U = 0.69$ .

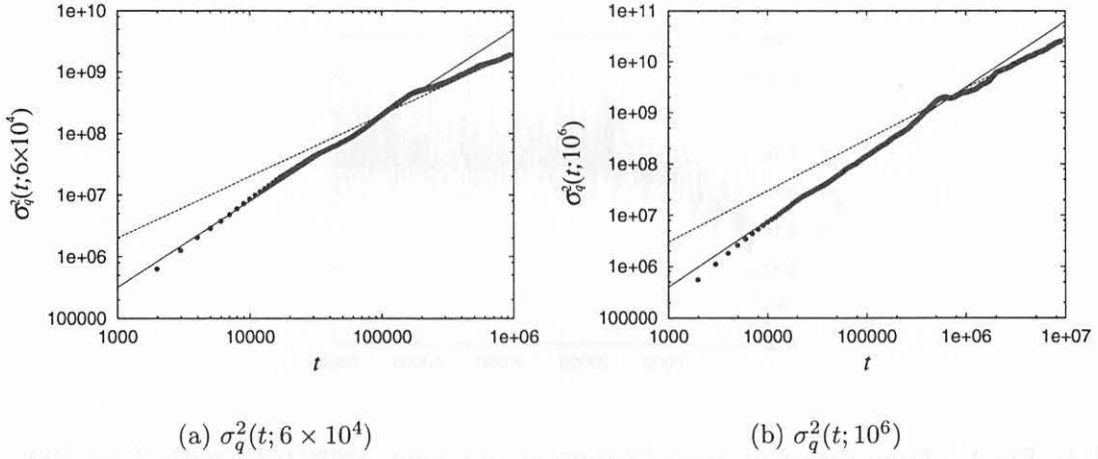


Fig.3. Time series of variance of phases  $\sigma_q^2(t;\tau)$ . (a)  $\tau = 6 \times 10^4$ . The two lines represent  $t^{1.4}$  and  $t$ . (b)  $\tau = 10^6$ . The two lines represent  $t^{1.3}$  and  $t$ .  $N = 1000$  and  $U = 0.69$ . The same initial condition with Fig.2.

なお、 $\tau = 0, 6 \times 10^4, 10^6$  のそれぞれに対する異常拡散の指数  $\alpha$  は、それぞれおおよそ 1.6, 1.4, 1.3 となる。つまり、運動エネルギーがカノニカル温度に達したあとでも異常拡散は起こっているが、時間が経てば経つほど拡散は遅くなることがわかる。

## 4 局所リアプノフ数

リアプノフ数とは軌道の線形不安定性を長時間平均した量であり、 $\delta \mathbf{x}(t)$  を線形化運動方程式に従う  $2N$  次元ベクトルとすると次のように定義される:

$$\lambda = \lim_{t \rightarrow \infty} \frac{1}{t} \ln \frac{\|\delta \mathbf{x}(t)\|}{\|\delta \mathbf{x}(0)\|}.$$

モデル系 (2) におけるリアプノフ数のエネルギー依存性はすでに調べられており、エネルギーが臨界エネルギー  $U_c = 3/4$  に向かって大きくなるにつれてリアプノフ数も大きくなり、熱力学極限 ( $N \rightarrow \infty$ ) においては  $U > U_c$  で 0 となる [Yam96, Fir98]。ここでは、軌道不安定性と緩和過程の関係を調べるため、局所リアプノフ数と言われる次の量を観測する:

$$\lambda^{local}(t; T) = \frac{1}{T} \ln \frac{\|\delta \mathbf{x}(t)\|}{\|\delta \mathbf{x}(t-T)\|}. \quad (6)$$

これは、時間間隔  $T$  を設け、時間  $T$  での軌道不安定性を時系列として表した量である。

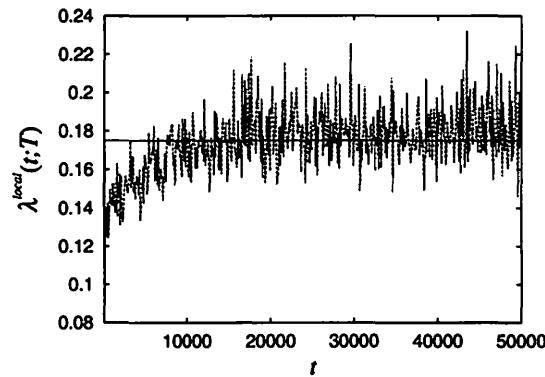


図 4: Fig.4. Time series of local Lyapunov exponent  $\lambda^{local}(t; T)$  with  $T = 100$ . It converges around  $t = 10^4$ , which is much earlier than the first passage time of  $2K(t)/N$  to the canonical temperature.  $N = 1000$  and  $U = 0.69$ .

図 4 には、図 2 と同じ初期条件のもとでの局所リアプノフ数の時系列を示してある。これによると、局所リアプノフ数は、運動エネルギーが初めてカノニカル

温度に達する時刻 (およそ  $t = 6 \times 10^4$ ) よりもずいぶん早く  $t = 10^4$  あたりに収束しており、緩和の終了と軌道不安定性の収束に関連がないことが見られる。これは、一次相転移を起こす系 (3) との大きな違いである [TA99]。

## 5 まとめと課題

長距離相互作用のモデルのひとつである、平均場相互作用系において運動エネルギーの緩和過程と異常拡散、局所リアプノフ数との関係を調べた。その結果、次のことがわかった。(i) 運動エネルギーの緩和過程において異常拡散が見られる。しかし、異常拡散は運動エネルギーの緩和過程が終了したあとでも見られる。(ii) 異常拡散を観測し始める時刻を後ろにずらしていくと、異常拡散のベキ指数は徐々に小さくなっていく。(iii) 一次相転移系では緩和の終了時刻と局所リアプノフ数の収束時刻がほぼ同一となるが、二次相転移系では後者の方が前者よりも早い。これらの結果から、二次相転移系においては緩和過程を異常拡散や局所リアプノフ数といった概念で捉えることは難しいことがわかる。

課題としては、重力シートモデルにおいて見出された何段階かの緩和過程をこのモデル系においても見出すことにより、今みている緩和がどの段階にあたるのかを同定しなければならない。また、結果 (ii) をより精緻に調べることで、定性的ではなく定量的に緩和過程を特徴付ける可能性について考察することが挙げられる。

## Acknowledgements

I thank to Alessandro Toricini and Stefano Ruffo for valuable discussions and comments.

## 参考文献

- [AHR02] Mickaël Antoni, Haye Hinrichsen, and Stefano Ruffo. “On the micro-canonical solution of a system of fully coupled particles”. *Chaos, Solitons and Fractals*, 13:393–99, 2002.
- [AR95] M. Antoni and S. Ruffo. “Clustering and relaxation in Hamiltonian long-range dynamics”. *Phys. Rev. E*, 52(3):2361–74, 1995.
- [AT98] Celia Anteneodo and Constantino Tsallis. “Breakdown of exponential sensitivity to initial conditions: Role of the range of interactions”. *Phys. Rev. Lett.*, 80(24):5313–6, 1998.



- [Fir98] Marie-Christine Firpo. “Analytic estimation of the Lyapunov exponent in a mean-field model undergoing a phase transition”. *Phys. Rev. E*, 57:6599–603, 1998.
- [LRR99a] Vito Latora, Andrea Rapisarda, and Stefano Ruffo. “Chaos and statistical mechanics in the Hamiltonian mean field model”. *Physica D*, 131:38–54, 1999.
- [LRR99b] Vito Latora, Andrea Rapisarda, and Stefano Ruffo. “Superdiffusion and out-of-equilibrium chaotic dynamics with many degrees of freedoms”. *Phys. Rev. Lett.*, 83(11):2104–7, 1999.
- [TA99] Alessandro Torcini and Michaël Antoni. “Equilibrium and dynamical properties of two-dimensional N-body systems with long-range attractive interactions”. *Phys. Rev. E*, 59:2746–63, 1999.
- [TGK96] Toshio Tsuchiya, Naoteru Gouda, and Tetsuro Konishi. “Relaxation processes in one-dimensional self-gravitating many-body systems”. *Phys. Rev. E*, 53:2210–6, 1996.
- [Yam96] Y. Y. Yamaguchi. “Slow relaxation at critical point of second order phase transition in a highly chaotic hamiltonian system”. *Prog. Theor. Phys.*, 95:717–31, 1996.
- [Yos93] H. Yoshida. “Recent progress in the theory and application of symplectic integrators”. *Celestial Mechanics and Dynamical Astronomy*, 56:27–43, 1993.

# A MODIFIED HERMITE INTEGRATOR FOR PLANETARY DYNAMICS

Eiichiro Kokubo

*Division of Theoretical Astrophysics,  
National Astronomical Observatory,  
Osawa, Mitaka, Tokyo, 181-8588, Japan  
E-mail: kokubo@th.nao.ac.jp*

and

Junichiro Makino

*Department of Astronomy, Faculty of Science, University of Tokyo,  
Hongo, Bunkyo-Ku, Tokyo, 113-0033, Japan*

## ABSTRACT

We describe modified time-symmetric Hermite integrators specialized for long-term integration of planetary orbits. Our time-symmetric integrators have no secular errors in the semimajor axis and the eccentricity for the integration of two-body Kepler problems as usual time-symmetric and symplectic integrators. Usual time-symmetric or symplectic integrators, however, show a secular drift in the argument of pericenter. Our new family of integrators has one free parameter, which we can adjust to eliminate the error in the argument of pericenter without breaking the time-symmetry or changing the order of the integrator. The value of the free parameter for which the error is eliminated shows very weak dependence on the size of the timestep and the eccentricity. We analytically show that the leading term of the error vanishes for a unique value of the parameter. We describe the second- and the fourth-order schemes. Extension to higher order is straightforward.

# 1 INTRODUCTION

Time-symmetric (e.g., Quinlan and Tremaine 1990, Calvo and Sanz-Serna 1994) and symplectic (e.g., Kinoshita, Yoshida, and Nakai 1991, Saha and Tremaine 1992) integrators have recently been used for the study of the long-term stability of planetary orbits. Both integrators have a desirable property that they have no secular errors in the semimajor axis  $a$  and the eccentricity  $e$  when Keplerian orbit is integrated with a constant timestep, while widely used high-order multi-step integrators, such as the Störmer-Cowell integrator, show secular errors. On the other hand, there are linear errors in the argument of pericenter  $\omega$  and the time of pericenter passage  $T$  for the time-symmetric and the symplectic integrators, while the high-order multi-step integrators show quadratic errors in the time of pericenter passage. For the long-term orbital stability problem, in particular for resonance related problems, an integrator that has no secular errors in either action variables  $(a, e)$  or angle variables  $(\omega, T)$  would be extremely useful.

In the field of the structure engineering, numerical integration methods for the dynamic vibration equation (second-order linear differential equation) has been investigated in detail (e.g., Wood 1990). The Newmark method (Newmark 1959) is the most popular integrator for the dynamic vibration equation. It is a family of second-order one-step integrators that has two parameters, which includes the leapfrog (Verlet) scheme and the trapezoidal formula as special cases. It is known that by choosing appropriate parameters, we can remove the error in the phase of the vibration.

In the present paper, we apply this concept of the Newmark method to higher order time-symmetric Hermite integrators. The 4th-order Hermite integrator (Makino and Aarseth 1992, Kokubo, Yoshinaga, and Makino 1998, hereafter referred to as KYM98) is widely used for the relatively long-term calculation of planetary systems from its simple yet accurate algorithm (e.g, Alexander and Agnor 1998, Yoshinaga, Kokubo, and Makino 1999, Iwasaki et al. 2001, Kominami and Ida 2002). KYM98 presented the time-symmetric Her-

mite integrator that shows no secular errors in  $(a, e)$  when two-body Kepler problems are integrated and demonstrated that their integrator is effective in planetary  $N$ -body problems. Our purpose here is to remove or reduce the errors of the angle orbital elements in the time-symmetric Hermite scheme without losing the time-symmetry. We focus on the error of the argument of pericenter in this paper.

The time-symmetric integrators require constant timesteps to avoid secular errors in  $(a, e)$ . In collisional  $N$ -body problems, however, the individual (variable) timestep scheme (Aarseth 1985) is indispensable to resolve close encounters of particles accurately and economically. KYM98 showed that in planetary  $N$ -body systems such as planetesimal and protoplanet systems, particles share the same constant timestep except for the relatively rare cases of close encounters even though the hierarchical individual timestep (Makino 1991) is adopted. This almost constant timesteps over particles and time are the reason the time-symmetric Hermite integrator greatly reduces the integration error of the planetary  $N$ -body systems. In case of close encounters, the orbits are integrated accurately with small timesteps in the individual timestep scheme. Thus, their time-symmetric integrator with the individual timestep is effective in planetary  $N$ -body problems. For the same reason, our new integrators are also effective not only in two-body Kepler problems but also in planetary  $N$ -body problems.

In section 2, we describe the Newmark method and demonstrate how to optimize the integrator for the Kepler problem. In section 3, the fourth-order scheme based on the time-symmetric Hermite scheme (KYM98) is presented and its behavior is investigated. Section 4 is devoted for summary and discussion.

## 2 THE TIME-SYMMETRIC NEWMARK METHOD

We describe the time-symmetric Newmark method and apply it to planetary orbital calculation. The Newmark method (Newmark 1959) takes the form

$$\mathbf{x}_1 = \mathbf{x}_0 + \mathbf{v}_0\Delta t + \frac{1}{2}\mathbf{a}_0\Delta t^2 + \beta\dot{\mathbf{a}}_0\Delta t^3, \quad (1)$$

$$\mathbf{v}_1 = \mathbf{v}_0 + \mathbf{a}_0\Delta t + \gamma\dot{\mathbf{a}}_0\Delta t^2, \quad (2)$$

where  $\mathbf{x}$ ,  $\mathbf{v}$ , and  $\mathbf{a}$  are the position, the velocity, and the acceleration, respectively. The subscripts 0 and 1 indicate time,  $\Delta t$  is the timestep, and  $\beta$  and  $\gamma$  are the Newmark parameters. The time derivative of the acceleration  $\dot{\mathbf{a}}_0$  is given by

$$\dot{\mathbf{a}}_0 = \frac{\mathbf{a}_1 - \mathbf{a}_0}{\Delta t}. \quad (3)$$

In order to keep the time-symmetry,  $\gamma$  should be  $1/2$ . The time-symmetric Newmark scheme takes the implicit form

$$\mathbf{x}_1 = \mathbf{x}_0 + \frac{1}{2}(\mathbf{v}_1 + \mathbf{v}_0)\Delta t + \left(\beta - \frac{1}{4}\right)(\mathbf{a}_1 - \mathbf{a}_0)\Delta t^2, \quad (4)$$

$$\mathbf{v}_1 = \mathbf{v}_0 + \frac{1}{2}(\mathbf{a}_1 + \mathbf{a}_0)\Delta t. \quad (5)$$

The  $\beta = 0$  scheme corresponds to the leapfrog scheme and  $\beta = 1/4$  the trapezoidal formula.

The accuracy of the Newmark method is second-order as seen in Eq. (2). This means that the last term of the r.h.s. of Eq. (1) does not affect the order of accuracy. In other words,  $\beta$  is a free parameter. We can vary the value of  $\beta$  without changing the order of accuracy or losing the time-symmetry. Skeel, Zhang, and Schlick (1997) analyzed the same algorithm but in a slightly different representation. They found that this family of schemes can be regarded symplectic when associated with an adequate variable transformation.

Figure 1 shows the error in the osculating orbital elements for the two-body Kepler problem solved by the time-symmetric Newmark scheme with  $\beta = 0, 1/12, 1/6$  for 10 orbital periods ( $20\pi$  time units). The initial orbital elements are  $a = 1$ ,  $e = 0.1$ ,  $\omega = \pi$ ,

and  $T = \pi$ . The orbit is integrated with a constant timestep  $\Delta t = 2^{-5}$ , in other words, about 200 integration steps per orbit. There are no secular errors for  $a$  and  $e$  but periodic changes for all the values of  $\beta$ . In time-symmetric integrators with a constant timestep, the periodic errors in  $a$  and  $e$  cancel out in an orbit.

For  $\omega$ , the result for  $\beta = 1/12$  shows no secular error for this time scale, while those for  $\beta = 0$  and  $1/6$  show linear drifts. Figure 2 shows the result of time integration same as Fig. 1 but for  $10^6$  orbital periods with  $\beta = 1/12$ . In Fig. 2,  $\Delta\omega$  is visible. However, the error is far smaller than those by the  $\beta = 0$  and  $1/6$  schemes. For  $\beta = 0$  and  $1/6$ ,  $|\Delta\omega|$  per orbit is of the order of  $10^{-3}$ , while for  $\beta = 1/12$ , it is of the order of  $10^{-7}$ .

This “magic number”,  $\beta = 1/12$ , can be derived analytically through the evaluation of the leading term of the error. The leading terms of the local truncation errors of the time-symmetric Newmark scheme are given by

$$\Delta x = \left(\beta - \frac{1}{6}\right) \dot{a}_0 \Delta t^3, \quad (6)$$

$$\Delta v = \frac{1}{12} a_0^{(2)} \Delta t^3. \quad (7)$$

As the eccentricity has no secular error, it is sufficient to evaluate the error of one component of the eccentricity vector. Here we take the  $y$ -component of the eccentricity vector  $e_y = e \sin \omega$  which is given by

$$e_y = -v_x(xv_y - yv_x) - \frac{y}{r}, \quad (8)$$

where  $r = (x^2 + y^2)^{1/2}$ . The leading term of the error of  $e_y$  is given by

$$\Delta e_y = \frac{\partial e_y}{\partial x} \Delta x + \frac{\partial e_y}{\partial y} \Delta y + \frac{\partial e_y}{\partial v_x} \Delta v_x + \frac{\partial e_y}{\partial v_y} \Delta v_y. \quad (9)$$

The symmetry of the Kepler orbit requires that  $\Delta\omega$  should cancel out in half an orbital period. This leading error of  $\omega$  vanishes if  $\beta$  satisfies the following linear equation:

$$\int_0^{\frac{T_K}{2}} \Delta e_y dt = \Delta t^3 \int_0^{\frac{T_K}{2}} \left[ \left(\beta - \frac{1}{6}\right) \left( \frac{\partial e_y}{\partial x} \dot{a}_x + \frac{\partial e_y}{\partial y} \dot{a}_y \right) + \frac{1}{12} \left( \frac{\partial e_y}{\partial v_x} a_x^{(2)} + \frac{\partial e_y}{\partial v_y} a_y^{(2)} \right) \right] dt = 0, \quad (10)$$

where  $T_K$  is the Kepler period. We see in Eq. (10) that  $\beta$  is not a function of  $\Delta t$ . Eq. (10) leads to

$$\beta = \frac{1}{6} - \frac{1}{12} \frac{\int_0^{\frac{T_K}{2}} \left( \frac{\partial e_y}{\partial v_x} a_x^{(2)} + \frac{\partial e_y}{\partial v_y} a_y^{(2)} \right) dt}{\int_0^{\frac{T_K}{2}} \left( \frac{\partial e_y}{\partial x} \dot{a}_x + \frac{\partial e_y}{\partial y} \dot{a}_y \right) dt} \quad (11)$$

Substituting Kepler solutions for Eq. (11) and integrating for half an orbit (for details, see Appendix A), we have

$$\beta = \frac{1}{12}. \quad (12)$$

For  $\beta = 1/12$ , only the leading term of  $\Delta\omega$  vanishes. Therefore, for a finite step size, the next term in the error should become visible. Figure 3 shows the dependence of  $|\Delta\omega|$  per orbit on  $\Delta t$ . The initial conditions are the same as those for Fig. 1. It is clearly shown that  $|\Delta\omega|$  for  $\beta = 1/12$  is  $O(\Delta t^4)$ , while  $|\Delta\omega|$  for other values of  $\beta$  is  $O(\Delta t^2)$ . This means that not only the  $O(\Delta t^3)$  term of the local truncation error but also the  $O(\Delta t^4)$  term vanishes for  $\beta = 1/12$ . The global error is one order lower than the local truncation error since the number of steps for a period is inversely proportional to  $\Delta t$ . This behavior is quite natural. For any symmetric scheme, the coefficient of the local error is exactly zero for even orders (Kinoshita 1968, Cano and Sanz-Serna 1997). Thus, if the  $O(\Delta t^3)$  term vanishes, the next term is  $O(\Delta t^5)$ .

We have shown that with an appropriate value of  $\beta$  we can reduce the secular error for  $\omega$  drastically. On the other hand, there are linear errors for  $T$  for all the  $\beta$  schemes as opposed the quadratic errors of non-time-symmetric schemes.

The second-order schemes are often insufficient in simulations with a wide range of length and time scales because of its low accuracy. We investigate the property of the fourth-order generalized scheme in the next section.

### 3 FOURTH-ORDER INTEGRATOR

The time-symmetric Newmark method can be extended to higher order integrators. As the simplest example, we apply it to the implicit Hermite scheme, a fourth-order time-symmetric scheme (KYM98).

The Hermite integrator is based on the Taylor series up to the order of the third time derivative of the acceleration  $\mathbf{a}_0^{(3)}$ , given by

$$\mathbf{x}_1 = \mathbf{x}_0 + \mathbf{v}_0 \Delta t + \frac{\mathbf{a}_0}{2} \Delta t^2 + \frac{\dot{\mathbf{a}}_0}{6} \Delta t^3 + \frac{\mathbf{a}_0^{(2)}}{24} \Delta t^4 + \alpha \frac{\mathbf{a}_0^{(3)}}{120} \Delta t^5, \quad (13)$$

$$\mathbf{v}_1 = \mathbf{v}_0 + \mathbf{a}_0 \Delta t + \frac{\dot{\mathbf{a}}_0}{2} \Delta t^2 + \frac{\mathbf{a}_0^{(2)}}{6} \Delta t^3 + \frac{\mathbf{a}_0^{(3)}}{24} \Delta t^4, \quad (14)$$

where  $\mathbf{a}_0^{(2)}$  and  $\mathbf{a}_0^{(3)}$  are obtained by the 3rd-order Hermite interpolation constructed from  $\mathbf{a}$  and  $\dot{\mathbf{a}}$  at time  $t_0$  and  $t_1$  as

$$\mathbf{a}_0^{(2)} = \frac{-6(\mathbf{a}_0 - \mathbf{a}_1) - \Delta t(4\dot{\mathbf{a}}_0 + 2\dot{\mathbf{a}}_1)}{\Delta t^2}, \quad (15)$$

$$\mathbf{a}_0^{(3)} = \frac{12(\mathbf{a}_0 - \mathbf{a}_1) + 6\Delta t(\dot{\mathbf{a}}_0 + \dot{\mathbf{a}}_1)}{\Delta t^3}, \quad (16)$$

and  $\alpha$  is a new parameter introduced to control integration errors. The role of  $\alpha$  is the same as that of  $\beta$  of the Newmark method. The order of the accuracy of the scheme is determined by the  $O(\Delta t^5)$  term of the velocity. We can, therefore, change the weight of the  $O(\Delta t^5)$  term of the position without changing the order of the accuracy.

From Eqs. (13) through (16), the Hermite integrator can be rewritten in an implicit form as

$$\mathbf{x}_1 = \mathbf{x}_0 + \frac{1}{2}(\mathbf{v}_1 + \mathbf{v}_0)\Delta t - \frac{\alpha}{10}(\mathbf{a}_1 - \mathbf{a}_0)\Delta t^2 + \frac{6\alpha - 5}{120}(\dot{\mathbf{a}}_1 + \dot{\mathbf{a}}_0)\Delta t^3, \quad (17)$$

$$\mathbf{v}_1 = \mathbf{v}_0 + \frac{1}{2}(\mathbf{a}_1 + \mathbf{a}_0)\Delta t - \frac{1}{12}(\dot{\mathbf{a}}_1 - \dot{\mathbf{a}}_0)\Delta t^2. \quad (18)$$

It is clear in this formula that the Hermite integrator is time-symmetric, in other words, the physical values with subscripts 0 and 1 are used symmetrically. Note that the  $\alpha = 1$



scheme corresponds to the Hermite scheme of Makino and Aarseth (1992) and  $\alpha = 5/6$  that of Hut *et al.* (1995).

We determine  $\alpha$  in the same way as the time-symmetric Newmark scheme. The leading terms of the local truncation errors of the implicit Hermite scheme are given by

$$\Delta x = \frac{\alpha - 1}{120} a_0^{(3)} \Delta t^5, \quad (19)$$

$$\Delta v = -\frac{1}{720} a_0^{(4)} \Delta t^5. \quad (20)$$

The leading term of the error in  $e_y$  is  $O(\Delta t^5)$  that is given by

$$\int_0^{\frac{T_K}{2}} \Delta e_y dt = \frac{\Delta t^5}{120} \int_0^{\frac{T_K}{2}} \left[ (\alpha - 1) \left( \frac{\partial e_y}{\partial x} a_x^{(3)} + \frac{\partial e_y}{\partial y} a_y^{(3)} \right) - \frac{1}{6} \left( \frac{\partial e_y}{\partial v_x} a_x^{(4)} + \frac{\partial e_y}{\partial v_y} a_y^{(4)} \right) \right] dt. \quad (21)$$

Therefore, the leading term vanishes when  $\alpha$  satisfies (for details, see Appendix B)

$$\alpha = 1 + \frac{\frac{1}{6} \int_0^{\frac{T_K}{2}} \left( \frac{\partial e_y}{\partial v_x} a_x^{(4)} + \frac{\partial e_y}{\partial v_y} a_y^{(4)} \right) dt}{\int_0^{\frac{T_K}{2}} \left( \frac{\partial e_y}{\partial x} a_x^{(3)} + \frac{\partial e_y}{\partial y} a_y^{(3)} \right) dt} = \frac{7}{6}. \quad (22)$$

As  $\beta$  for the time-symmetric Newmark method,  $\alpha$  is independent of  $e$  and  $\Delta t$ .

Figure 4 shows the error in the osculating orbital elements against time for 10 orbital periods for the  $\alpha = 5/6, 1, 7/6$  schemes. Time-symmetry is realized by applying 5 iteration of corrections, namely P(EC)<sup>5</sup> scheme. The initial conditions and the timestep are the same as those for Fig. 1. No schemes have secular errors in  $a$  and  $e$ . The scheme with  $\alpha = 7/6$  has no secular error in  $\omega$ , either. In this time scale there seems to be no secular error in  $T$  for the case of  $\alpha = 7/6$ . Figure 5 is the same as Fig. 2 but for the  $\alpha = 7/6$  scheme. In this time scale, we can see the small linear error in  $\omega$ . This is due to higher order errors of  $\omega$  that do not cancel out in half an orbit with  $\alpha = 7/6$  because  $\alpha$  is determined so that the leading error term of  $\omega$  vanishes.

The dependence of  $\Delta\omega$  per orbit on  $\Delta t$  is shown in Fig. 6 for  $\alpha = 5/6, 1, 7/6$ . It is clear that  $|\Delta\omega|$  for  $\alpha = 7/6$  is  $O(\Delta t^6)$ , while  $|\Delta\omega|$  for other  $\alpha$  is  $O(\Delta t^4)$ . This means that not

only the  $O(\Delta t^5)$  term of the local truncation error but also the  $O(\Delta t^6)$  term cancels out in an orbital period as does the  $O(\Delta t^4)$  term of the time-symmetric Newmark scheme. We also show the same plots for the 4th-order symplectic integrator (e.g., Kinoshita, Yoshida, and Nakai 1991) for comparison. The error  $|\Delta\omega|$  for the symplectic integrator is  $10^{4-5}$  times larger than that for the Hermite scheme with  $\alpha = 7/6$ . Figure 7 shows  $|\Delta\omega|$  per orbit against  $e$ . The errors monotonically increase with  $e$ . Since  $|\Delta\omega|$  increases linearly with  $t$  as shown in Fig. 5, we can estimate  $\Delta\omega$  after any orbital periods based on Fig. 7.

The fact that our scheme with  $\alpha = 7/6$  gives very high accuracy for pure Kepler problems does not necessarily guarantee that it gives good results for systems with more than one planet. In order to test our scheme, we integrate a planar Jupiter-Saturn like system for 300,000 years. The masses and the initial orbital elements of the two planets are  $M_1 = 2 \times 10^{30}\text{g}$ ,  $a_1 = 5\text{AU}$ ,  $e_1 = 0.05$ ,  $\omega_1 = \pi/2$ ,  $T_1 = 0$  and  $M_2 = 5 \times 10^{29}\text{g}$ ,  $a_2 = 10\text{AU}$ ,  $e_2 = 0.05$ ,  $\omega_2 = 3\pi/2$ ,  $T_2 = 0$  (the sun is fixed at the coordinate origin). The evolution of the eccentricity vector  $(e \cos \omega, e \sin \omega)$  averaged over 500 years is plotted in Fig. 8. As the reference, the result obtained by the standard  $\alpha = 1$  scheme with the timestep  $\Delta t = 0.125$  (562 steps per  $T_K$  for planet 1) is plotted. The results for the implicit Hermite scheme with  $\alpha = 1$  and  $7/6$  and the 4th-order symplectic integrator with  $\Delta t = 1$  are plotted. Though the timesteps are different by a factor of 8, the reference result and the result for  $\alpha = 7/6$  are indistinguishable. Their difference in  $\omega$  after 300,000 years is about 0.00041. On the other hand, the difference in  $\omega$  between the reference result and the result for  $\alpha = 1$  is about 0.15. The drift in  $\omega$  is faster than that of the reference result by about 3%, which is due to the integration error. It is shown that the good property of the  $\alpha = 7/6$  scheme is retained even when it is used for this kind of three-body problem. The result for the symplectic integrator shows a completely different pattern. In particular, the eccentricity vector rotates in the opposite direction of the results for the implicit Hermite scheme.

## 4 SUMMARY AND DISCUSSION

We developed new modified time-symmetric Hermite integrators specialized for planetary dynamics. The new time-symmetric integrators have no secular errors in the semimajor axis and the eccentricity. The new time-symmetric schemes have one free parameter,  $\alpha$ , which we can adjust to reduce the error of the argument of pericenter drastically without breaking the time-symmetry and changing the order of accuracy. The free parameter  $\alpha$  is the coefficient of the highest order term of the position that is one order higher than that of the velocity. In the case of the second- or the fourth-order scheme, the error order of  $\omega$  becomes two-order higher for a unique value of  $\alpha$ . These values of  $\alpha$  are independent of orbital elements or the size of timestep. We presented the second- and the fourth-order schemes. It is, however, straightforward to apply our method to higher order time-symmetric integrators.

It is possible to adjust  $\alpha$  to reduce the error in the time of pericenter passage. However, in this case, the optimal values depend on the eccentricity. It is also possible to choose  $\alpha$  so that the secular error in  $\omega$  becomes exactly zero for a given timestep and eccentricity. This value is slightly different from the value for which the leading term of the error vanishes and depends on the eccentricity and the size of the timestep.

Our method is not limited to the Kepler problem. It is possible to apply the method to any periodic systems. We can optimize integrators for a given problem easily by choosing an adequate value of one free parameter so that one additional constant of motion is conserved well.

We can also use the time-symmetric variable timestep scheme (Hut *et al.* 1995) with our method. The merit of our method remains with the time-symmetric variable timestep. However, with the time-symmetric variable timestep, the optimal value of  $\alpha$  depends on both the eccentricity and the way to determine the timestep.

We showed that the accuracy of the 4th-order Hermite integrator with  $\alpha = 7/6$  is higher

than the 4th-order symplectic integrator of Kinoshita, Nakai, and Yoshida (1991). The mixed-variable symplectic integrator (Kinoshita, Nakai, Yoshida 1991) or Wisdom-Holman map (Wisdom and Holman 1991), however, has the accuracy much higher than our 4th-order Hermite integrator when orbits are nearly Keplerian. We agree that for the very-long-term integration of relatively stable planetary systems such as the solar system, the mixed variable method is suitable. For the moderately-long-term integration of relatively unstable planetary systems such as protoplanet systems, the Hermite integrator has been used by many authors (e.g., Alexander and Agnor 1998, Yoshinaga, Kokubo, and Makino 1999). We can not use the above symplectic integrators for this collisional (unstable)  $N$ -body systems where close encounters of bodies are important, because they do not allow the individually variable timestep that is necessary to resolve close encounters accurately and economically.

We emphasize that one of the important merits of our method is that this scheme can be used for planetary  $N$ -body simulation. As discussed in KYM98, in planetary  $N$ -body systems, most particles are on nearly circular orbits because the gravitational interaction among particles are weak compared with the gravity of the central body except for close encounters of particles. Those orbits are integrated by constant timesteps even though we allow variable timesteps. In this case, by adopting  $\alpha$  scheme, we can integrate individual orbits more precisely than by usual time-symmetric scheme, which is important for resonance related problems. The application of our method easily improves the accuracy of the calculation such as the stability of protoplanet systems (e.g., Yoshinaga, Kokubo, and Makino 1998) and the accretion of protoplanets (e.g., Alexander and Agnor 1998). The implementation of our method is simple and easy. One has only to introduce  $\alpha$  in time-symmetric integrators and set an adequate value.

## A DERIVATION OF $\beta$

In the orbital reference system, the Kepler solution takes the form,

$$x = a(\cos u - e), \quad (23)$$

$$y = a\sqrt{1 - e^2} \sin u, \quad (24)$$

$$v_x = -\frac{na^2}{r} \sin u, \quad (25)$$

$$v_y = \frac{na^2\sqrt{1 - e^2}}{r} \cos u, \quad (26)$$

where  $u$  is the eccentric anomaly and  $n$  is the mean motion. The time derivative of  $u$  is given by

$$\frac{du}{dt} = \frac{n}{1 - e \cos u}. \quad (27)$$

The acceleration by the central body is

$$\mathbf{a} = -\mu \frac{\mathbf{r}}{r^3}, \quad (28)$$

where  $\mu$  is constant. The first and the second time derivatives of the acceleration are given by

$$\dot{\mathbf{a}} = -\mu \frac{\mathbf{v}}{r^3} + \mu \frac{3(\mathbf{r} \cdot \mathbf{v})\mathbf{r}}{r^5}, \quad (29)$$

$$\mathbf{a}^{(2)} = -\mu \frac{\mathbf{a}}{r^3} + \mu \frac{3(v^2 + \mathbf{r} \cdot \mathbf{a})\mathbf{r} + 6(\mathbf{r} \cdot \mathbf{v})\mathbf{v}}{r^5} - \mu \frac{15(\mathbf{r} \cdot \mathbf{v})^2\mathbf{r}}{r^7}. \quad (30)$$

The derivatives of the  $y$ -component of the eccentricity vector are given by

$$\frac{\partial e_y}{\partial x} = -v_x v_y + \frac{xy}{r^3}, \quad (31)$$

$$\frac{\partial e_y}{\partial y} = v_x^2 - \frac{x^2}{r^3}, \quad (32)$$

$$\frac{\partial e_y}{\partial v_x} = -xv_y + 2yv_x, \quad (33)$$

$$\frac{\partial e_y}{\partial v_y} = -xv_x. \quad (34)$$

Substituting the Kepler solutions into Eq. (11) with Eqs. (29) through (34) and transforming the integral variable from  $t$  to  $u$  with Eq. (27), we obtain

$$\beta = \frac{1}{6} - \frac{1}{12} \frac{\int_0^\pi \frac{-3e \cos^2 u + \cos u + 2e}{(1 - e \cos u)^5} du}{\int_0^\pi \frac{(12e^3 - 14e) \cos^2 u + (5e^2 - 1) \cos u - 15e^3 + 13e}{(1 - e \cos u)^6} du} = \frac{1}{12}. \quad (35)$$

## B DERIVATION OF $\alpha$

The third and the fourth time derivatives of the acceleration are given by

$$\begin{aligned} \mathbf{a}^{(3)} = & -\mu \frac{\dot{\mathbf{a}}}{r^3} + \mu \frac{3(3\mathbf{v} \cdot \mathbf{a} + \mathbf{r} \cdot \dot{\mathbf{a}})\mathbf{r} + 9(v^2 + \mathbf{r} \cdot \mathbf{a})\mathbf{v} + 9(\mathbf{r} \cdot \mathbf{v})\mathbf{a}}{r^5} \\ & - \mu \frac{45(\mathbf{r} \cdot \mathbf{v})(v^2 + \mathbf{r} \cdot \mathbf{a})\mathbf{r} + 45(\mathbf{r} \cdot \mathbf{v})^2\mathbf{v} + \mu \frac{105(\mathbf{r} \cdot \mathbf{v})^3\mathbf{r}}{r^9}}{r^7}, \end{aligned} \quad (36)$$

$$\begin{aligned} \mathbf{a}^{(4)} = & -\mu \frac{\mathbf{a}^{(2)}}{r^3} + \mu \frac{3(3a^2 + 4\mathbf{v} \cdot \dot{\mathbf{a}} + \mathbf{r} \cdot \mathbf{a}^{(2)})\mathbf{r} + 12(3\mathbf{v} \cdot \mathbf{a} + \mathbf{r} \cdot \dot{\mathbf{a}})\mathbf{v} + 18(v^2 + \mathbf{r} \cdot \mathbf{a})\mathbf{a} + 12(\mathbf{r} \cdot \mathbf{v})\dot{\mathbf{a}}}{r^5} \\ & - \mu \frac{45(v^2 + \mathbf{r} \cdot \mathbf{a})^2\mathbf{r} + 60(\mathbf{r} \cdot \mathbf{v})(3\mathbf{v} \cdot \mathbf{a} + \mathbf{r} \cdot \dot{\mathbf{a}})\mathbf{r} + 180(\mathbf{r} \cdot \mathbf{v})(v^2 + \mathbf{r} \cdot \mathbf{a})\mathbf{v} + 90(\mathbf{r} \cdot \mathbf{v})^2\mathbf{a}}{r^7} \\ & + \mu \frac{630(\mathbf{r} \cdot \mathbf{v})^2(v^2 + \mathbf{r} \cdot \mathbf{a})\mathbf{r} + 420(\mathbf{r} \cdot \mathbf{v})^3\mathbf{v}}{r^9} - \mu \frac{945(\mathbf{r} \cdot \mathbf{v})^4\mathbf{r}}{r^{11}}. \end{aligned} \quad (37)$$

Substituting the above derivatives to Eq. (22) and in the same way as the derivation of  $\beta$ , we obtain  $\alpha = 7/6$ .

## REFERENCES

- Aarseth S. J., 1985, in Brackhill J. U., Cohen B. I., eds, Multiple Time Scales (Academic Press, New York), p.377
- Alexander S. G., Agnor C. B., 1998, *Icarus*, 132, 113
- Calvo M. P., Sanz-Serna J. M., 1994, preprint
- Cano B., Sanz-Serna J. M., 1997, *SIAM J. Num. Anal.*, 34, 1391

- Hut P., Makino J., McMillan S., 1995, *ApJL*, 443, 93
- Iwasaki K., Tanaka H., Nakazawa K., Emori H., 2001, *PASJ*, 53, 321
- Kinoshita H., 1968, *PASJ*, 20, 1
- Kinoshita H., Yoshida H., Nakai H., 1991, *Celest. Mech. Dyn. Astron.*, 50, 59
- Kokubo E., Yoshinaga K., Makino J., 1998, *MNRAS*, 297, 1067
- Kominami J., Ida S., 2002, *Icarus*, 157, 43
- Makino J., 1991, *PASJ*, 43, 859
- Makino J., Aarseth S. J., 1992, *PASJ*, 44, 141
- Newmark N. M., 1959, in *Proceedings of ASCE, Journal of Engineering Mechanics (EM3)*, 85, 67
- Quinlan G. D., Tremaine S., 1990, *AJ*, 100, 1694
- Saha P., Tremaine S., 1992, *AJ*, 104, 4
- Skeel R. D., Zhang G., Schlick T., 1997, *SIAM J. Sci. Comput.*, 18, 203
- Wisdom J., Holman M., 1991, *AJ*, 102, 1528
- Wood W. L., 1990, *Practical Time-stepping Schemes* (Clarendon, Oxford)
- Yoshinaga K., Kokubo E., Makino J., 1998, *Icarus*, 139, 328

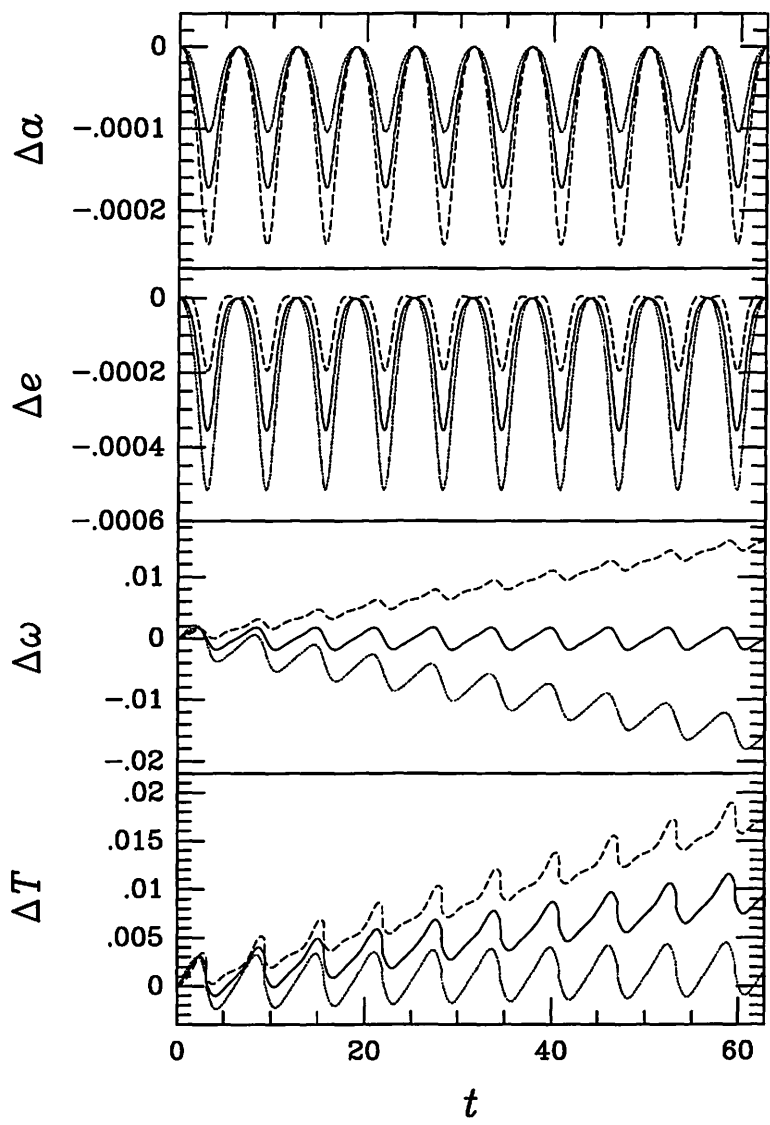


Fig. 1: The errors of orbital elements  $a$ ,  $e$ ,  $\omega$ , and  $T$  from top to bottom for the time-symmetric Newmark scheme with  $\beta = 0$  (dotted curve),  $\beta = 1/12$  (solid curve), and  $\beta = 1/6$  (dashed curve) against time for 10 orbital periods. The Kepler period is  $2\pi$  and the timestep is  $\Delta t = 2^{-5}$ .



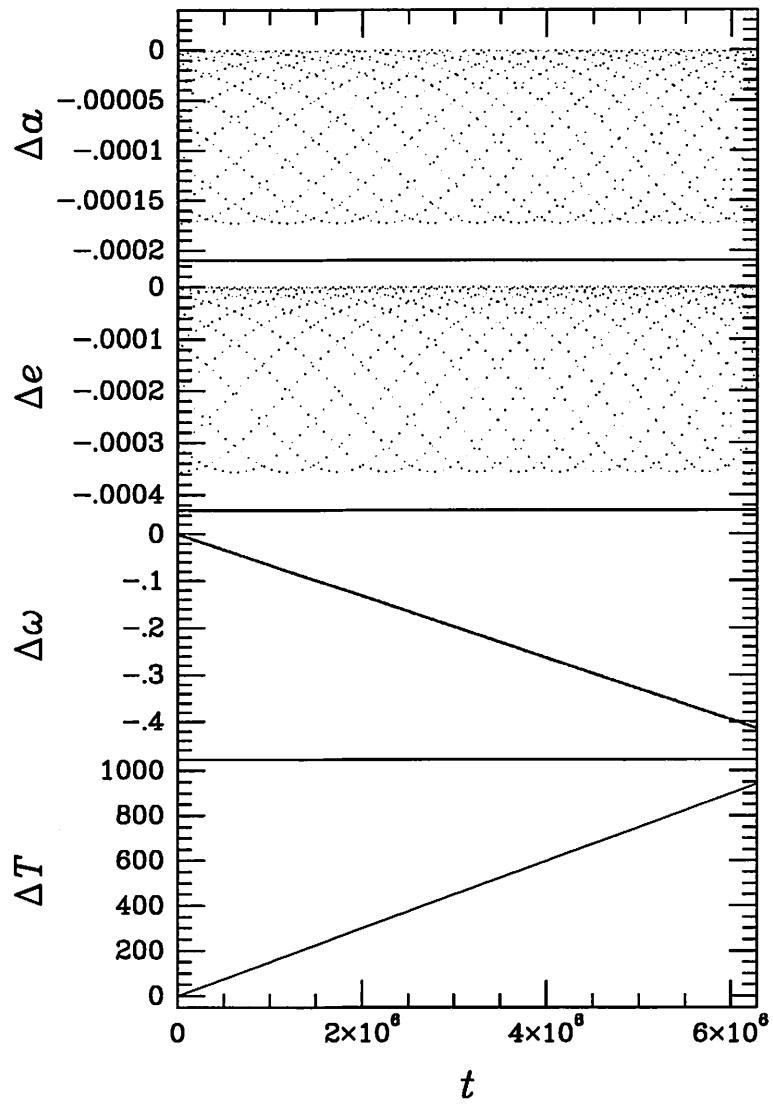


Fig. 2: The same as Fig. 1 but for  $\beta = 1/12$  for  $10^6$  orbital periods.

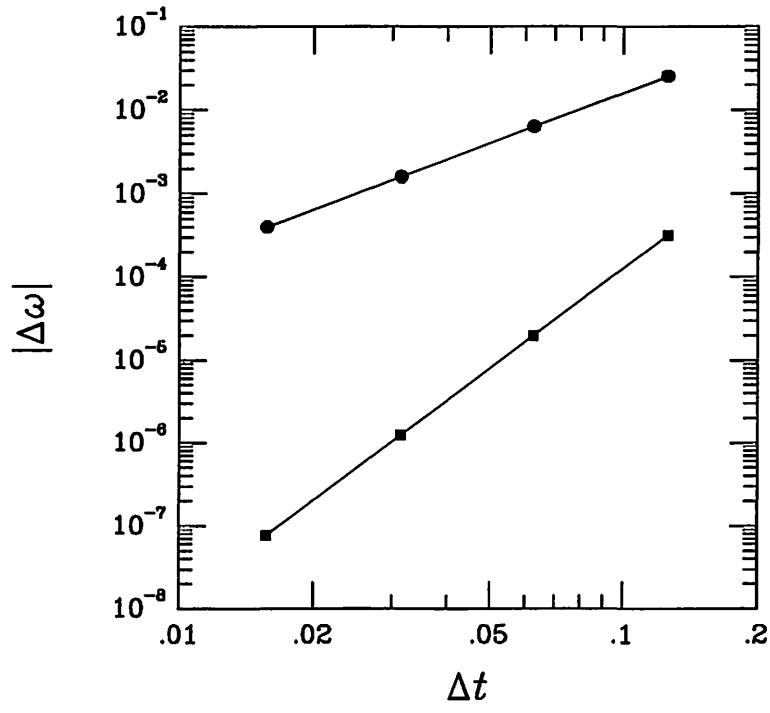


Fig. 3: The error in the argument of pericenter  $\Delta\omega$  after an orbital period is plotted against the timestep for  $\beta = 0$  (triangles),  $1/12$  (squares), and  $1/6$  (circles). The plots of  $\beta = 0$  and  $1$  overlap.

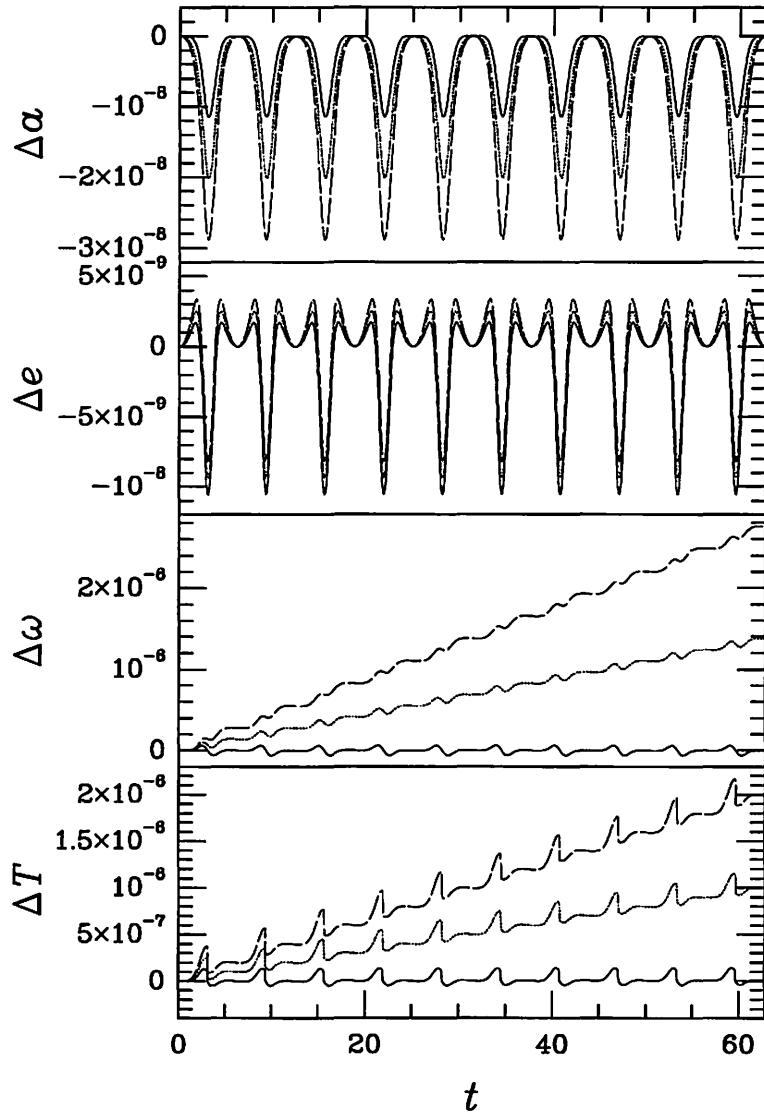


Fig. 4: The same as Fig. 1 but for the errors of the implicit Hermite scheme with  $\alpha = 5/6$  (dashed curve),  $\alpha = 1$  (dotted curve), and  $\alpha = 7/6$  (solid curve).

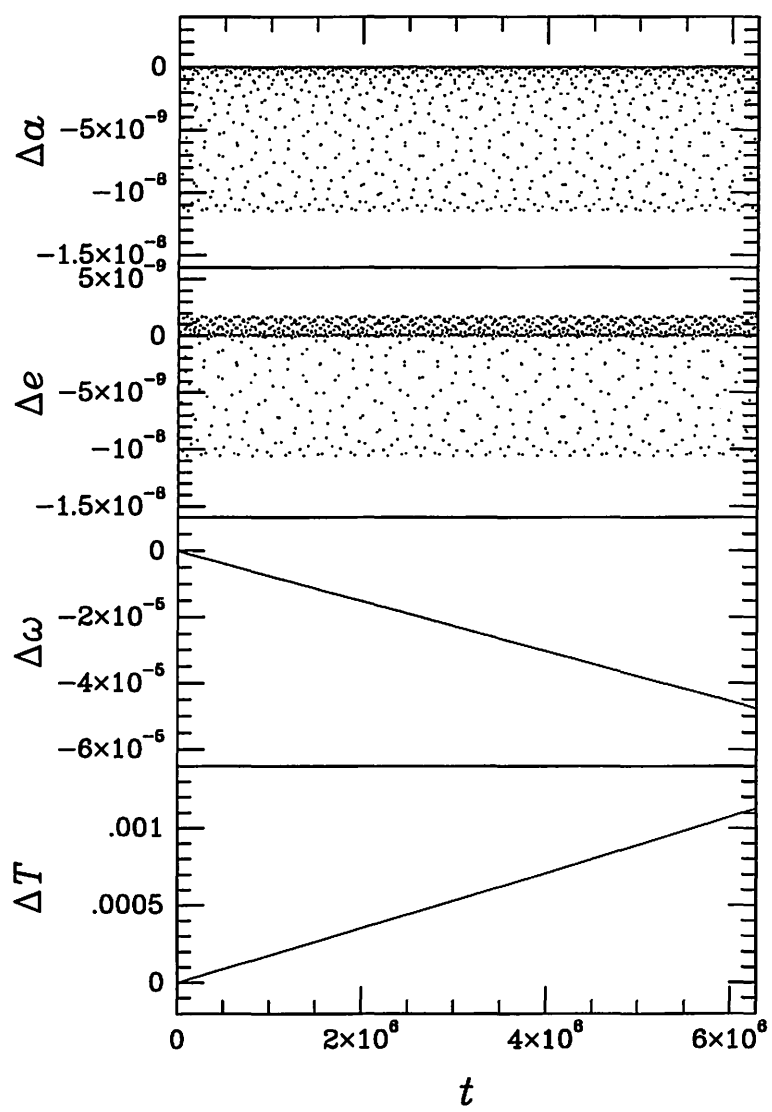


Fig. 5: The same as Fig. 2 but for the Hermite scheme with  $\alpha = 7/6$ .

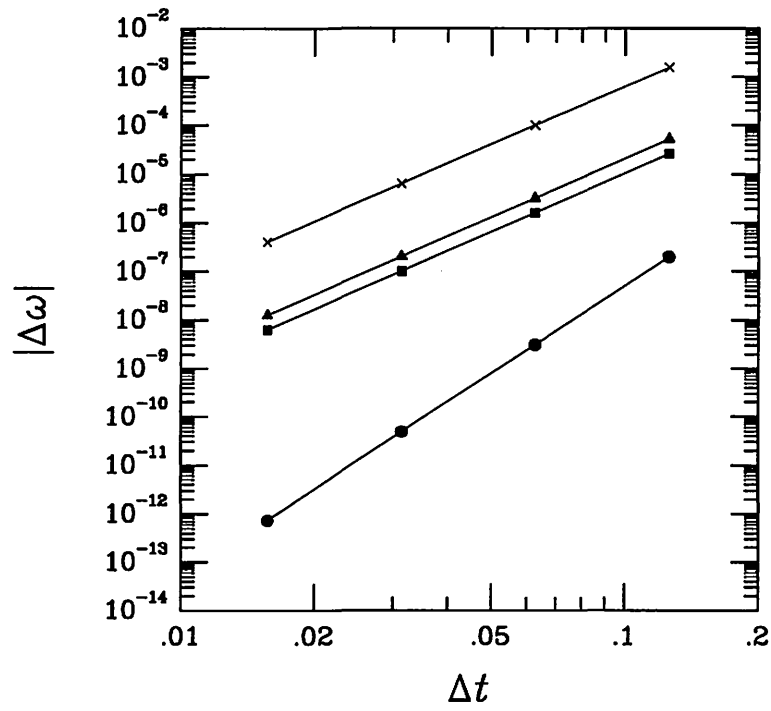


Fig. 6: The same as Fig. 3 but for the implicit Hermite scheme with  $\alpha = 5/6$  (triangles),  $\alpha = 1$  (squares), and  $\alpha = 7/6$  (circles). The result for the 4th-order symplectic integrator is also plotted (crosses).

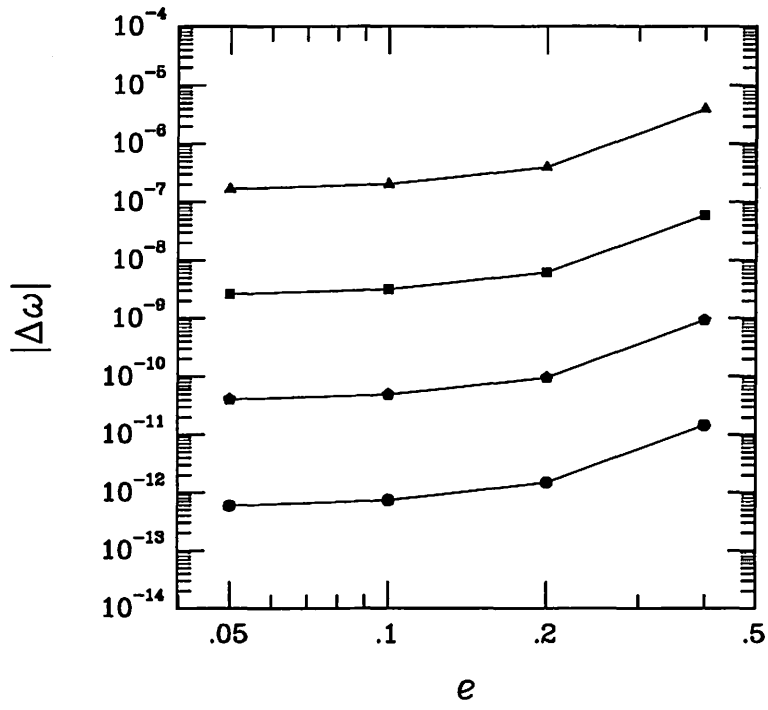


Fig. 7: The error in the argument of pericenter  $\Delta\omega$  after one orbital period is plotted against the eccentricity for the  $\alpha = 7/6$  scheme. The timesteps are  $\Delta t = 2\pi/50$  (triangles),  $2\pi/100$  (squares),  $2\pi/200$  (squares), and  $2\pi/400$  (circles).

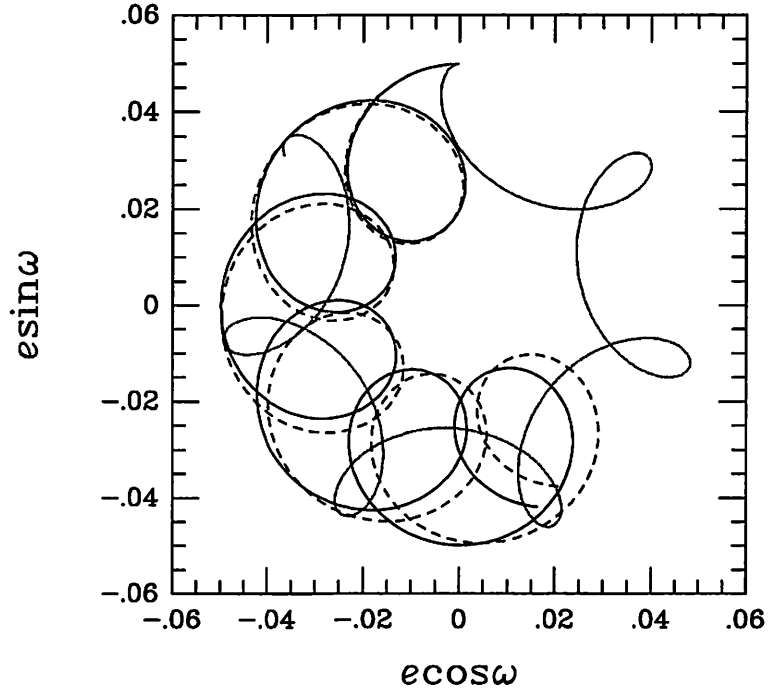


Fig. 8: Time evolution of the eccentricity vector of planet 1 for 300,000 years. The eccentricity vector is averaged over 500 years. A three-body system is integrated by the  $\alpha = 1$  (dashed curve) and the  $\alpha = 7/6$  (solid curve) schemes with the timestep  $\Delta t = 1$ . As the reference, the result obtained by the  $\alpha = 1$  scheme with  $\Delta t = 0.125$  (solid curve) is shown. The result for the 4th-order symplectic integrator with  $\Delta t = 1$  is also plotted (dotted curve).

# Secular numerical error in $H = T(p) + V(q)$ symplectic integrator: simple analysis for error reduction

Takashi Ito and Kiyotaka Tanikawa

*National Astronomical Observatory, Mitaka, Tokyo 181-8588, Japan*

tito@cc.nao.ac.jp

A simple error analysis for  $H = T(p) + V(q)$  symplectic integrator (i.e. not mixed-variable type) is presented. The truncation error in a second-order integrator is analytically analyzed up to first-order approximation for the two-body problem using a canonical perturbation theory. In  $H = T(p) + V(q)$  type integrators, we cannot employ sophisticated techniques such as warm start or symplectic corrector for the reduction of secular numerical error. But we have confirmed that so-called “iterative start,” where we repeat many short-term numerical integrations while gradually changing initial orbital configuration and searching a point with minimum numerical error, may reduce the secular numerical error in angle variables under certain conditions. According to our numerical integrations on two kinds of three-body planetary systems (weakly and strongly perturbed), simple  $H = T(p) + V(q)$  symplectic integrators are still useful when employed together with the iterative start. To obtain a simple interpretation how the errors are reduced (or not reduced in most systems), we take a nonlinear pendulum system with one degree of freedom for an example, and illustrate that the reduction of the numerical error in  $H = T(p) + V(q)$  symplectic integrator occurs when the potential energy of the system is not an “isochrone” one — when the fundamental frequency of the system depends on initial amplitude of oscillation.

## 1. Introduction

In dynamical studies of solar and extrasolar planetary objects, analytical complexity of perturbation techniques and development of fast computers has led us to the investigation by numerical methods. One of the promising ways for long-term numerical integrations is symplectic integrator designed specifically to maintain the Hamiltonian structure of equations of motion (Yoshida, 1990b; Gladman *et al.*, 1991; Kinoshita *et al.*, 1991; Kinoshita and Nakai, 1992; Yoshida, 1993; Sanz-Serna and Calvo, 1994). One of the typical types of the symplectic integrators splits the Hamiltonian  $H$  into two integrable parts as  $H = T(p) + V(q)$  where  $T$ ,  $V$ ,  $q$ , and  $p$  are kinetic energy, potential energy, canonical coordinate and conjugate momentum, respectively. On the other hand, so called Wisdom-Holman map (also called “mixed variable symplectic integrator,” and hereafter we call it “WH map”) by Wisdom and Holman (1991, 1992) can be more accurate by a factor of the ratio of planetary to central mass than the general-purpose symplectic integrators of  $H = T + V$  type. The principle behind the WH map is to split the Hamiltonian into an unperturbed Kepler part and a perturbation part as  $H = H_{\text{kep}}(L) + H_{\text{int}}(l)$ , where  $L$  and  $l$  denote canonical variables for Keplerian motion symbolically. In each step of the integration, the system is first moved forward in time according to Kepler motion  $H_{\text{kep}}$ , and then a kick in momentum is applied which is derived from the perturbation part of the Hamiltonian,  $H_{\text{int}}$ . Not only the Keplerian part, but also the interaction part is analytical since the perturbation Hamiltonian is basically a function of only relative Cartesian coordinates. The coordinate transformation between  $L$  in  $H_{\text{kep}}$  and  $l$  in  $H_{\text{int}}$  is efficiently encapsulated using the Gauß’s  $f$ - and  $g$ -functions (cf. Danby, 1992).



Also, there are many peripheral techniques for the WH map for the purpose of reducing its numerical error especially in angle variables. They work mostly owing to the smallness of the perturbed part of Hamiltonian ( $H_{\text{int}}$ ) than Keplerian part ( $H_{\text{kep}}$ ). Saha and Tremaine (1992) have devised a special start-up procedure to reduce the truncation error of angle variables, called “warm start,” utilizing one of the characteristics of Hamiltonian system — existence of adiabatic invariant. They have also invented a symplectic scheme with individual stepsizes (though not variable stepsizes), dividing the whole Hamiltonian into each planet’s Keplerian and perturbation parts (Saha and Tremaine, 1994). Wisdom et al. (1996) found a canonical transformation of variables that eliminates error Hamiltonian and greatly improves the accuracy of symplectic integration. This transformation is expressed in terms of a Lie operator that must be applied before each step, and an inverse transformation at the end of each step. This operation is called “symplectic correction.” Some notes on the dependence of the numerical error on initial starting conditions in the WH-type symplectic integrators is also mentioned in Michel and Valsecchi (1996).

There are also many other variants of and modifications to the WH map for the applications in dynamical astronomy. Those kinds of enhancement in symplectic integrators, especially of the WH map, now enable us to perform very long-term numerical integrations ten to hundred times faster than before. With the WH map, timescale of numerical integrations of solar system planetary orbits have reached the age of the solar system, i.e. 4.5 Gyr (Ito *et al.*, 1996; Duncan and Lissauer, 1998; Ito and Tanikawa, 2002).

Now, let us be back at the general-purpose symplectic integrators,  $H = T(p) + V(q)$  type. For problems proxy to the Keplerian motion, the general-purpose method is less efficient than the WH map is. However, the general-purpose method has their literal advantage, i.e. generality: they can be adapted to general dynamical problems which are far from integrable and whose zeroth order approximate solutions are not known. We can think of many of such far-integrable problems, as we discover more and more extrasolar planetary systems, since many of the extrasolar planetary orbital configurations so far discovered are significantly unlike ours (Boss, 1996; Marcy *et al.*, 2000; Marcy and Butler, 2000). Typical ones are the planetary systems in or around binaries. In such systems where the ratio of interaction Hamiltonian  $H_{\text{int}}$  and Kepler Hamiltonian  $H_{\text{kep}}$  is generally not sufficiently small, we cannot exploit the near-integrability of the system which the WH map requires. Also, it is generally not easy to apply the WH map to situations with a lot of close encounters among particles. Though several variants of the WH map are now proposed to handle such collisional systems (Levison and Duncan, 1994; Mikkola, 1997; Lee *et al.*, 1997; Duncan *et al.*, 1998; Chambers, 1999; Mikkola and Tanikawa, 1999), such symplectic schemes might be highly complicated and lose computational efficiency.

Standing on the above viewpoints, we present in this paper a simple error analysis on the  $T(p) + V(q)$  type symplectic integrator (hereafter we call it “TV method” in contrast to the “WH map”). Since many researches have been already done so far on characteristics of the TV symplectic method, this paper may be in a sense an expository one. Hence, to make the way of the error analysis as transparent and general as possible, we take a few very simple dynamical systems as examples: two-body problem, perturbed three-body problems, harmonic oscillator with low degrees of freedom, and a nonlinear pendulum with one degree of freedom. Most of them are nearly integrable, or even analytical solutions are already known. In the former half

of this manuscript, our approach is somewhat close to Kinoshita et al. (1991)’s one. Thus we have felt it advisable to give more details than would otherwise be necessary. This is also in keeping with the view of this paper as an expository one.

In Section 2., we present a brief review of symplectic integrators, especially of the TV method. In Section 3., we will discuss the error Hamiltonian for a first- and second-order TV symplectic methods for the planar two-body problem. Based on the result obtained in this section, we demonstrate to calculate some analytical expressions of numerical symplectic solutions using a canonical perturbation theory in Section 4.. Particularly in the subsection 4.7, we confirm the dependence of numerical longitudinal error on initial orbital configuration. In certain configurations we can significantly reduce the longitudinal error arising from the symplectic integrator; in other words, the iterative start works. While in most configurations, we cannot. In Section 5., we argue on the way of error reduction in angle variables by the iterative start in the TV method. Demonstrations by some numerical experiments in two kinds of three-body systems are described: One is the orbital motion of a “massive” asteroid perturbed by Jupiter, and the other is the orbital motion of an extrasolar planet orbiting around a binary system, MACHO-97-BLG-41. We also mention slightly the “warm start” and its relationship to the topic in this manuscript. Finally in Section 6., we try to illustrate how the errors are reduced or not reduced in various dynamical systems. To explain this qualitatively, we have taken a few systems with low degrees of freedom as examples. We have so far found that we can possibly reduce the numerical error in TV symplectic method considerably by the iterative start when the potential energy  $V$  of the system is not “isochrone” — when fundamental frequency of the system depends on initial amplitude of oscillation.

## 2. Symplectic integrator

First we present a brief review of the generic type of symplectic integrator.

According to Yoshida (1993), explicit symplectic integrators can be reformulated by the Lie algebra (Neri, 1987). We rewrite the Hamilton equations

$$\frac{dq}{dt} = \frac{\partial H}{\partial p}, \quad \frac{dp}{dt} = -\frac{\partial H}{\partial q}, \quad (1)$$

in the form as

$$\frac{dz}{dt} = \{z, H(z)\}, \quad (2)$$

where  $z = q$  or  $p$ , and the braces  $\{, \}$  stands for the Poisson bracket. When we introduce a differential operator  $D_G$  by

$$D_GF \equiv \{F, G\}, \quad (3)$$

then (2) is rewritten as

$$\frac{dz}{dt} = D_H z, \quad (4)$$

so the formal solution, or the exact time evolution of  $z(t)$  from  $t = 0$  to  $t = \tau$  is given by

$$z(\tau) = \left[ e^{\tau D_H} \right] z(0). \quad (5)$$

For a Hamiltonian of the form

$$H = T(p) + V(q), \quad (6)$$

$D_H = D_T + D_V$  and we have a formal solution

$$z(\tau) = \left[ e^{\tau(A+B)} \right] z(0), \quad (7)$$

where  $A \equiv D_T$  and  $B \equiv D_V$ . Operators  $A, B$  are non-commuting in general.

Kinetic energy  $T(p)$  and potential energy  $V(q)$  are individually integrable, so we can get the exact solutions

$$z_A(\tau) = \left[ e^{\tau A} \right] z_A(0), \quad (8)$$

$$z_B(\tau) = \left[ e^{\tau B} \right] z_B(0). \quad (9)$$

Here we should remark that the time evolution of  $z$  under  $e^{\tau A}$  or  $e^{\tau B}$  keeps the symplecticity of the system. This fact is one of the most essential cores of symplectic integration theory. For example, let us take (8) as an example and see how the symplecticity is kept. The symplectic map (8) is a kind of contact canonical transformation under the Hamiltonian  $T(p)$ . Writing down the canonical equation of motion concerning  $T$ , we have

$$\frac{dq}{dt} = \frac{\partial T(p)}{\partial p}, \quad \frac{dp}{dt} = -\frac{\partial T(p)}{\partial q}. \quad (10)$$

Since  $T(p)$  does not contain  $q$ , we get

$$\frac{dp}{dt} = 0, \quad (11)$$

$$\therefore p = \text{constant}, \quad (12)$$

hence we know that  $T(p)$  is also a constant. This leads us to

$$\frac{dq}{dt} = \frac{\partial T(p)}{\partial p}, \quad (13)$$

is a function of  $p$  and also a constant.

$$\therefore q = Ct + q_0, \quad (14)$$

where  $C$  and  $q_0$  are certain constants. (14) means that a particle in phase-space  $(q, p)$  moves linearly with time having a constant velocity. Such an equi-velocity linear motion in phase-space obviously preserves any volume in phase-space. Then the contact transformation (8) preserves the symplecticity of the system whatever value  $\tau$  has. We can apply the same discussion on the contact transformation (9), leading to the conclusion that (9) preserves the symplecticity of the system. Since a product of two canonical transformations is found to be canonical, a map  $e^{\tau A} e^{\tau B}$  also preserves the symplecticity. This argument is applicable to other systems with any degree of freedom. This fact ensures us the area (or volume) preservation property of symplectic schemes. However, note that this character does not directly lead to the conservation of total energy and total angular momentum of the system in symplectic integration.

Now, what we need is the solution under  $H = T(p) + V(q)$ . But since the operators  $A$  and  $B$  are not commutable, we have to find a product which approximates  $e^{\tau(A+B)}$  to an appropriate order.

There is a formula which exactly answers to our question: Baker-Campbell-Hausdorff (BCH) formula (Dragt and Finn, 1976; Varadarajan, 1974) about a product of two exponential functions of non-commuting operators  $X$  and  $Y$ ; when we write the product as

$$e^X e^Y = e^Z, \quad (15)$$

$Z$  turns out to be as follows according to the BCH formula

$$Z = X + Y + \frac{1}{2}[X, Y] + \frac{1}{12}([X, [X, Y]] + [Y, [Y, X]]) + \frac{1}{24}[X, [Y, [Y, X]]] + \dots, \quad (16)$$

where  $[X, Y] \equiv XY - YX$ . For a first-order symplectic integrator, we can apply the BCH formula to  $e^{\tau(A+B)}$  as

$$e^{\tau D_T} e^{\tau D_V} = e^{\tau D_{\tilde{H}_{1st}}}, \quad (17)$$

and obtain

$$\tilde{H}_{1st} = T + V + \frac{\tau}{2}\{V, T\} + \frac{\tau^2}{12}(\{\{T, V\}, V\} + \{\{V, T\}, T\}) + O(\tau^3). \quad (18)$$

For a second-order symplectic integrator, we obtain

$$e^{\frac{\tau}{2} D_T} e^{\tau D_V} e^{\frac{\tau}{2} D_T} = e^{\tau D_{\tilde{H}_{2nd}}}, \quad (19)$$

where

$$\tilde{H}_{2nd} = T + V + \tau^2 \left( \frac{1}{12} \{\{T, V\}, V\} - \frac{1}{24} \{\{V, T\}, T\} \right) + O(\tau^4). \quad (20)$$

Similarly in general, for an  $n$ -th order symplectic integrator, we find the Hamiltonian  $\tilde{H}_n$  as

$$\tilde{H}_n = H + H_{\text{err}} + O(\tau^{n+1}), \quad (21)$$

where  $H = T(p) + V(q)$  and  $H_{\text{err}} = O(\tau^n)$ . We call hereafter  $\tilde{H}$  the Hamiltonian of a surrogate system. We notice that the error of the total energy ( $\tilde{H} - H$ ) remains of the order of  $H_{\text{err}}$ , i.e.  $\tau^n$ .  $H_{\text{err}}$  is a set of terms which consist of  $n$ -fold Poisson brackets and called error Hamiltonian. Note that rigorous convergence of the series (18)(20)(21) is not guaranteed for general nonlinear systems.

### 3. Error Hamiltonian of the two-body problem

The purpose of this section is to express the error Hamiltonian  $H_{\text{err}}$  in a second-order symplectic integrator (20) as a function of Kepler orbital elements. We need the error Hamiltonian to estimate numerical error by symplectic integration in later sections. We take a planar two-body problem (masses  $m_0$  and  $m_1$  with  $\mu = G(m_0 + m_1)$ ) as an example whose dynamical characteristics is very well known.

The Hamiltonian for a two-body problem is written in the heliocentric coordinates  $(\mathbf{q}, \mathbf{v})$

$$H = \tilde{m} \left( \frac{v^2}{2} - \frac{\mu}{r} \right), \quad (22)$$

where

$$\tilde{m} = \frac{m_0 m_1}{m_0 + m_1}, \quad v = |\mathbf{v}|, \quad r = |\mathbf{q}|,$$

with a set of canonical variables  $(\mathbf{q}, \mathbf{p})$  and  $\mathbf{p} = \tilde{m}\mathbf{v}$ . Without loss of generality, we can consider the factor  $\tilde{m}$  in the right-hand side of (22), the reduced mass of the two-body system, as unity. This is possible because there are three units to be determined for a two-body system to dynamically work: mass, length, and time. The determination of  $\mu$ , semimajor axis  $a$  and the reduced mass  $\tilde{m}$  corresponds to the determination of these three units we use. Thus the Hamiltonian of the two-body problem (22) is reduced to that of a system where a infinitesimally small mass particle orbits around a central mass (say, the Sun) whose mass is  $m_0 + m_1$  as

$$H = \frac{v^2}{2} - \frac{\mu}{r}, \quad (23)$$

with canonical variables  $(\mathbf{q}, \mathbf{v})$ . See Appendix A for more details.

Since the kinetic energy  $T(\mathbf{v}) = v^2/2$  is a function only of canonical momentum  $\mathbf{v}$ , and the potential energy  $V(\mathbf{q}) = -\mu/r$  is a function only of canonical coordinate  $\mathbf{q}$ , it is clear

$$\frac{\partial T(\mathbf{v})}{\partial q_i} = 0, \quad \frac{\partial V(\mathbf{q})}{\partial v_i} = 0, \quad (i = 1, 2) \quad (24)$$

Hence the actual expression of the second-order error Hamiltonian up to  $O(\tau^2)$  approximation becomes from (20)

$$\begin{aligned} \frac{H_{\text{err}}}{\tau^2} &= \frac{1}{12} \{ \{T(\mathbf{v}), V(\mathbf{q})\}, V(\mathbf{q}) \} - \frac{1}{24} \{ \{V(\mathbf{q}), T(\mathbf{v})\}, T(\mathbf{v}) \} \\ &= \frac{1}{12} \left[ \left( \frac{\partial V}{\partial q_1} \right)^2 \frac{\partial^2 T}{\partial v_1^2} + 2 \frac{\partial V}{\partial q_1} \frac{\partial V}{\partial q_2} \frac{\partial^2 T}{\partial v_1 \partial v_2} + \left( \frac{\partial V}{\partial q_2} \right)^2 \frac{\partial^2 T}{\partial v_2^2} \right] \\ &\quad - \frac{1}{24} \left[ \left( \frac{\partial T}{\partial v_1} \right)^2 \frac{\partial^2 V}{\partial q_1^2} + 2 \frac{\partial T}{\partial v_1} \frac{\partial T}{\partial v_2} \frac{\partial^2 V}{\partial q_1 \partial q_2} + \left( \frac{\partial T}{\partial v_2} \right)^2 \frac{\partial^2 V}{\partial q_2^2} \right]. \end{aligned} \quad (25)$$

We need partial derivatives of the kinetic energy

$$T(\mathbf{v}) = \frac{v_1^2 + v_2^2}{2}, \quad (26)$$

and the potential energy

$$V(\mathbf{q}) = -\frac{\mu}{r} = -\mu (q_1^2 + q_2^2)^{-\frac{1}{2}}, \quad (27)$$

in (25). They become as follows:

$$\frac{\partial T}{\partial v_1} = v_1, \quad \frac{\partial T}{\partial v_2} = v_2, \quad (28)$$

$$\left( \frac{\partial T}{\partial v_1} \right)^2 = v_1^2, \quad \left( \frac{\partial T}{\partial v_2} \right)^2 = v_2^2, \quad (29)$$

$$\frac{\partial^2 T}{\partial v_1^2} = \frac{\partial^2 T}{\partial v_2^2} = 1, \quad (30)$$

$$\frac{\partial^2 T}{\partial v_1 \partial v_2} = 0, \quad (31)$$

$$\frac{\partial V}{\partial q_1} = \mu q_1 (q_1^2 + q_2^2)^{-\frac{3}{2}} = \frac{\mu q_1}{r^3}, \quad \frac{\partial V}{\partial q_2} = \mu q_2 (q_1^2 + q_2^2)^{-\frac{3}{2}} = \frac{\mu q_2}{r^3}, \quad (32)$$

$$\left(\frac{\partial V}{\partial q_1}\right)^2 = \mu^2 q_1^2 (q_1^2 + q_2^2)^{-3} = \frac{\mu^2 q_1^2}{r^6}, \quad \left(\frac{\partial V}{\partial q_2}\right)^2 = \mu^2 q_2^2 (q_1^2 + q_2^2)^{-3} = \frac{\mu^2 q_2^2}{r^6}, \quad (33)$$

$$\frac{\partial^2 V}{\partial q_1^2} = \mu \left( \frac{1}{r^3} - \frac{3q_1^2}{r^5} \right) = \frac{\mu}{r^5} (-2q_1^2 + q_2^2), \quad (34)$$

$$\frac{\partial^2 V}{\partial q_2^2} = \mu \left( \frac{1}{r^3} - \frac{3q_2^2}{r^5} \right) = \frac{\mu}{r^5} (-2q_2^2 + q_1^2). \quad (35)$$

Substituting (28)(29)(30)(31)(32)(33)(34)(34) into (25), we get

$$\begin{aligned} \frac{H_{\text{err}}}{r^2} &= \frac{1}{12} \left( \frac{\mu^2 q_1^2}{r^6} + \frac{\mu^2 q_2^2}{r^6} \right) - \frac{1}{24} \left[ v_1^2 \frac{\mu}{r^5} (-2q_1^2 + q_2^2) + 2v_1 v_2 \left( -\frac{3\mu q_1 q_2}{r^5} \right) + v_2^2 \frac{\mu}{r^5} (-2q_2^2 + q_1^2) \right] \\ &= \frac{1}{12} \frac{\mu^2}{r^4} - \frac{1}{24} \frac{\mu}{r^5} (-2v_1^2 q_1^2 + v_1^2 q_2^2 - 6v_1 v_2 q_1 q_2 + v_2^2 q_1^2 - 2v_2^2 q_2^2). \end{aligned} \quad (36)$$

Now, expressing the angular momentum integral  $h$  as

$$\begin{aligned} h^2 &= |\mathbf{q} \times \mathbf{v}|^2 \\ &= (q_1 v_2 - q_2 v_1)^2 \\ &= q_1^2 v_2^2 - 2q_1 q_2 v_1 v_2 + q_2^2 v_1^2. \end{aligned} \quad (37)$$

We know that  $h$  can be also expressed by Kepler orbital elements as

$$h = \sqrt{\mu a(1 - e^2)}. \quad (38)$$

Similarly, the energy integral can be expressed as

$$\frac{v^2}{2} - \frac{\mu}{r} = -\frac{\mu^2}{2h^2} (1 - e^2) \quad (39)$$

$$\therefore v^2 = \frac{2\mu}{r} - \frac{\mu^2}{h^2} (1 - e^2), \quad (40)$$

Using (37), (38) and (40), we can rewrite the quantity in the parentheses in the second term of the right-hand side of (36) as

$$\begin{aligned} &-2v_1^2 q_1^2 + v_1^2 q_2^2 - 6v_1 v_2 q_1 q_2 + v_2^2 q_1^2 - 2v_2^2 q_2^2 \\ &= 3(q_1^2 v_2^2 - 2q_1 q_2 v_1 v_2 + q_2^2 v_1^2) - 2(q_1^2 v_2^2 + q_2^2 v_1^2 + v_1^2 q_1^2 + v_2^2 q_2^2) \\ &= 3h^2 - 2(q_1^2 + q_2^2)(v_1^2 + v_2^2) \\ &= 3h^2 - 2r^2 v^2 \\ &= 3\mu a(1 - e^2) - 2r^2 \cdot 2 \left( \frac{\mu}{r} - \frac{\mu^2(1 - e^2)}{2\mu a(1 - e^2)} \right) \\ &= 3\mu a(1 - e^2) - 4\mu r + \frac{2\mu r^2}{a}. \end{aligned} \quad (41)$$

From (36) and (41), the final form of the error Hamiltonian expressed by the Kepler orbital elements becomes

$$\begin{aligned} H_{\text{err}} &= \frac{\tau^2 \mu^2}{12 r^4} - \frac{\tau^2 \mu}{24 r^5} \left( 3\mu a (1 - e^2) - 4\mu r + \frac{2\mu r^2}{a} \right) \\ &= \frac{\tau^2 \mu^2}{24} \left( \frac{6}{r^4} - \frac{3a(1 - e^2)}{r^5} - \frac{2}{ar^3} \right). \end{aligned} \quad (42)$$

The unperturbed or Keplerian part of Hamiltonian is  $-\frac{\mu^2}{2L^2}$ , so the surrogate Hamiltonian  $\tilde{H}$  in the second-order symplectic integrator ends up with

$$\begin{aligned} \tilde{H} &= H + H_{\text{err}} + O(\tau^4) \\ &= -\frac{\mu^2}{2L^2} + \frac{\tau^2 \mu^2}{24} \left( \frac{6}{r^4} - \frac{3a(1 - e^2)}{r^5} - \frac{2}{ar^3} \right) + O(\tau^4). \end{aligned} \quad (43)$$

Next, let us calculate the secular (i.e. time-averaged) value of the error Hamiltonian (42). To do this, time-averaged values of  $\frac{1}{r^3}$ ,  $\frac{1}{r^4}$ ,  $\frac{1}{r^5}$  are necessary. Using the relationship

$$\frac{dl}{df} = \frac{r^2}{a^2 \sqrt{1 - e^2}}, \quad (44)$$

they can be obtained as follows:

$$\begin{aligned} \left\langle \frac{a^3}{r^3} \right\rangle &= \frac{1}{2\pi} \int_0^{2\pi} \frac{a^3}{r^3} dl \\ &= (1 - e^2)^{-\frac{3}{2}}, \end{aligned} \quad (45)$$

$$\begin{aligned} \left\langle \frac{a^4}{r^4} \right\rangle &= \frac{1}{2\pi} \int_0^{2\pi} \frac{a^4}{r^4} dl \\ &= (1 - e^2)^{-\frac{5}{2}} \left( 1 + \frac{e^2}{2} \right), \end{aligned} \quad (46)$$

$$\begin{aligned} \left\langle \frac{a^5}{r^5} \right\rangle &= \frac{1}{2\pi} \int_0^{2\pi} \frac{a^5}{r^5} dl \\ &= (1 - e^2)^{-\frac{7}{2}} \left( 1 + \frac{3e^2}{2} \right). \end{aligned} \quad (47)$$

Substituting (45)(46)(47) into (42), the secular part of the error Hamiltonian becomes

$$\begin{aligned} \langle H_{\text{err}} \rangle &= \left\langle \frac{\mu^2 \tau^2}{24} \left( \frac{6}{r^4} - \frac{3a(1 - e^2)}{r^5} - \frac{2}{ar^3} \right) \right\rangle \\ &= \frac{\mu^2 \tau^2}{24a^4 \eta^5} \left( 1 + \frac{e^2}{2} \right), \end{aligned} \quad (48)$$

where

$$\eta \equiv \sqrt{1 - e^2}. \quad (49)$$

Now we can calculate  $H_{\text{err}}$  for the first-order symplectic integrator as well. Similar to (20), the error Hamiltonian for the first-order symplectic integrator becomes according to the BCH formula as

$$H_{\text{err},1\text{st}} = \frac{\tau}{2} \{V, T\} + \frac{\tau^2}{12} (\{\{T, V\}, V\} + \{\{V, T\}, T\}) + O(\tau^3). \quad (50)$$

For the first term in right-hand side of (50), we get

$$\begin{aligned} \{V(\mathbf{q}), T(\mathbf{p})\} &= \left( \frac{\partial V}{\partial q_1} \frac{\partial T}{\partial p_1} - \frac{\partial V}{\partial p_1} \frac{\partial T}{\partial q_1} \right) + \left( \frac{\partial V}{\partial q_2} \frac{\partial T}{\partial p_2} - \frac{\partial V}{\partial p_2} \frac{\partial T}{\partial q_2} \right) \\ &= \frac{\partial V}{\partial q_1} \frac{\partial T}{\partial p_1} + \frac{\partial V}{\partial q_2} \frac{\partial T}{\partial p_2} \\ &= \frac{\mu q_1}{r^3} \cdot v_1 + \frac{\mu q_2}{r^3} \cdot v_2 \\ &= \frac{\mu}{r^3} (\mathbf{r} \cdot \mathbf{v}). \end{aligned} \quad (51)$$

Each component of the velocity vector  $\mathbf{v}$  is expressed by the Kepler orbital elements on orbital plane as (Danby, 1992)

$$v_1 = -\frac{an}{\eta} \sin f, \quad v_2 = \frac{an}{\eta} (\cos f + e), \quad (52)$$

with

$$q_1 = r \cos f, \quad q_2 = r \sin f, \quad (53)$$

which leads to

$$\begin{aligned} \mathbf{r} \cdot \mathbf{v} &= q_1 v_1 + q_2 v_2 \\ &= r \cos f \left( -\frac{an}{\eta} \sin f \right) + r \sin f \left( \frac{an}{\eta} (\cos f + e) \right) \\ &= -\frac{anr}{2\eta} \sin 2f + \frac{anr}{2\eta} \sin 2f + \frac{earn}{\eta} \sin f \\ &= \frac{earn}{\eta} \sin f. \end{aligned} \quad (54)$$

We already knew the components of the second term in right-hand side of (50) by (42) as

$$\{\{T, V\}, V\} = \frac{\mu^2}{r^4}, \quad (55)$$

and

$$\{\{V, T\}, T\} = \frac{\mu}{r^5} \left( 3\mu a (1 - e^2) - 4\mu r + \frac{2\mu r^2}{a} \right). \quad (56)$$

Adding all the relevant terms, we get the final form of the error Hamiltonian up to  $O(\tau^2)$  as

$$\begin{aligned} H_{\text{err},1\text{st}} &= \frac{\tau}{2} \{V, T\} + \frac{\tau^2}{12} (\{\{T, V\}, V\} + \{\{V, T\}, T\}) \\ &= \frac{\tau}{2} \frac{\mu}{r^3} \frac{earn}{\eta} \sin f + \frac{\tau^2}{12} \left( \frac{\mu^2}{r^4} + \frac{\mu}{r^5} \left( 3\mu a (1 - e^2) - 4\mu r + \frac{2\mu r^2}{a} \right) \right) \\ &= \frac{\tau \mu e a r n}{2\eta r} \sin f + \frac{\tau^2 \mu^2}{12} \left( -\frac{3}{r^4} + \frac{3a(1 - e^2)}{r^5} + \frac{2}{ar^3} \right). \end{aligned} \quad (57)$$



We can calculate the secular part of  $H_{\text{err},1\text{st}}$  by averaging (57) over the period of mean anomaly. First, as for the term of  $O(\tau)$ ,

$$\begin{aligned}
\left\langle \frac{\tau \mu_{\text{ean}}}{2\eta r} \sin f \right\rangle &= \frac{1}{2\pi} \frac{\tau \mu_{\text{ean}}}{2\eta} \int_0^{2\pi} \frac{\sin f}{r} dl \\
&= \frac{1}{2\pi} \frac{\tau \mu_{\text{ean}}}{2\eta} \int_0^{2\pi} \frac{\sin f}{r} \frac{r^2}{a^2 \eta} df \\
&= \frac{1}{2\pi} \frac{\tau \mu_{\text{ean}}}{2\eta} \frac{a(1-e^2)}{a^2 \eta} \int_0^{2\pi} \frac{\sin f}{1+e \cos f} df \\
&= 0.
\end{aligned} \tag{58}$$

As for the terms of  $O(\tau^2)$ , we can calculate them as in the same way with (48). Hence

$$\begin{aligned}
\langle H_{\text{err},1\text{st}} \rangle &= \left\langle \frac{\tau \mu_{\text{ean}}}{2\eta r} \sin f + \frac{\tau^2 \mu^2}{12} \left( -\frac{3}{r^4} + \frac{3a(1-e^2)}{r^5} + \frac{2}{ar^3} \right) \right\rangle \\
&= \frac{\mu^2 \tau^2}{12} \frac{(1-e^2)^{-\frac{5}{2}}}{a^4} \left[ -3 \left( 1 + \frac{e^2}{2} \right) + 3 \left( 1 + \frac{3e^2}{2} \right) + 2(1-e^2) \right] \\
&= \frac{\mu^2 \tau^2}{12} \frac{(1-e^2)^{-\frac{5}{2}}}{a^4} (2+e^2) \\
&= \frac{\mu^2 \tau^2}{12a^4 \eta^5} (2+e^2).
\end{aligned} \tag{59}$$

Therefore we can proceed the same discussion using the secular error Hamiltonian for the first-order symplectic integrator (59) as well as that for the second-order symplectic integrator (48). In the following discussion we focus on the second-order symplectic integrator, so-called “leap frog,” using the error Hamiltonian (48).

## 4. Analytic solution by a canonical perturbation theory

Based on the result obtained in Section 3., we demonstrate to calculate analytical expressions of symplectic numerical solutions using a canonical perturbation theory. We can apply the treatment in this section to symplectic integrators of any order and to perturbation theory to any order, though it immediately leads to a terrible increase of relevant terms.

### 4.1 Canonical perturbation theory by the Lie transformation

Hori (1966, 1967) has developed a perturbation method with unspecified canonical variables utilizing the Lie transformation, and has presented several sample problems by his method (Hori, 1970; Hori, 1971). Hori’s method is characterized for its explicitness of canonical variables; the bothering inversion of old and new variables after obtaining analytical solution is no longer necessary. This is one of the major differences of Hori’s new method from traditional canonical perturbation theories such as by Delaunay or von Zeipel (cf. von Zeipel 1916, Shniad 1970, Yuasa 1971. Consult textbooks by Boccaletti and Pucacco (1998) or Lichtenberg and Lieberman (1992) for their general introduction). Let us briefly summarize the Hori’s perturbation method before applying it to our problem in the next subsection.

Let  $\xi, \eta$  be a set of  $2n$  canonical variables and  $f(\xi, \eta), S(\xi, \eta)$  be arbitrary functions of  $\xi, \eta$ . Differential operators  $D_s^n (n = 0, 1, 2, \dots)$  are defined as

$$D_s^0 f = f, \quad (60)$$

$$D_s^1 f = \{f, S\}, \quad (61)$$

$$D_s^n f = D_s^{n-1}(D_s^1 f). \quad (n \geq 2) \quad (62)$$

Then, the following theorem is due to Lie (1888): A set of  $2n$  variables  $x, y$  defined by the equation

$$f(x, y) = \sum_{n=0}^{\infty} \frac{\epsilon^n}{n!} D_s^n f(\xi, \eta), \quad (63)$$

is canonical if the series in the right-hand side of (63) converges.  $\epsilon$  is a small constant independent of  $\xi$  and  $\eta$ .

Let us consider a nearly integrable Hamiltonian system which is described by two-dimensional Delaunay variables  $(L, G, l, g)$  as

$$H(L, G, l, g) = H_0(L) + H_1(L, G, l, g), \quad (64)$$

where  $H_0(L)$  is the integrable part and  $H_1(L, G, l, g)$  is the perturbation part. Then, let us apply the Hori's perturbation method to the Hamiltonian system (64) to obtain the solution of the system. The general policy to apply canonical transformation here is to remove all angles and to make the system be integrable, such as

$$H^*(L^*, G^*) = H_0^*(L^*) + H_1^*(L^*, G^*), \quad (65)$$

where a superscript  $*$  symbolically denotes that the variables (or functions) have been canonically transformed.

As for the zeroth order Hamiltonian, the function form is the same before and after the canonical transformation:

$$H_0^*(L^*) = H_0(L^*). \quad (66)$$

Next we introduce a parameter  $t^*$  which satisfies the following relationship and is removed later on as

$$\frac{dL^*}{dt^*} = -\frac{\partial H_0^*}{\partial l^*}, \quad \frac{dl^*}{dt^*} = \frac{\partial H_0^*}{\partial L^*}, \quad (67)$$

$$\frac{dG^*}{dt^*} = -\frac{\partial H_0^*}{\partial g^*}, \quad \frac{dg^*}{dt^*} = \frac{\partial H_0^*}{\partial G^*}. \quad (68)$$

It is clearly seen that the system  $H_0^*$  is integrable since  $H_0^*$  is a function of only  $L^*$ . Then we get the following solution with constants of integration  $C_1, C_2, C_3$  and  $C_4$  as

$$L^* = C_1, \quad g^* = C_2, \quad G^* = C_3, \quad (69)$$

$$\frac{dl^*}{dt^*} = \frac{\partial H_0^*}{\partial L^*} = \text{constant}, \quad (70)$$

$$\therefore l^* = \text{constant} \times t^* + C_4. \quad (71)$$

As you can see,  $t^*$  is a time-like variable which describes the evolution of the non-perturbed (or integrable) system,  $H_0^*$ .

To the first-order,  $H_1^*$  becomes the  $t^*$ -averaged part of  $H_1$  as

$$H_1^* = \langle H_1(L^*, G^*, l^*, g^*) \rangle = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T H_1(L^*, G^*, l^*, g^*) dt^*, \quad (72)$$

or, if  $H_1$  is a periodic function of  $t^*$ , then

$$H_1^* = \langle H_1(L^*, G^*, l^*, g^*) \rangle = \frac{1}{T_p} \int_0^{T_p} H_1(L^*, G^*, l^*, g^*) dt^*, \quad (73)$$

where  $T_p$  is the period. Fortunately, the error Hamiltonians  $H_{\text{err}}$  we consider here is nearly periodic in most cases of planetary dynamics, so (73) is convenient instead of (72).

In the actual two-body problem, time  $t$  is related only to the mean anomaly  $l$ . Hence  $dt^*$  can be transformed into  $dl^*/n^*$  in (73) as

$$\begin{aligned} H_1^* = \langle H_1(L^*, G^*, l^*, g^*) \rangle &= \frac{1}{2\pi n^*} \int_0^{2\pi} H_1(L^*, G^*, l^*, g^*) dl^* \\ &= H_1^*(L^*, G^*, -, g^*), \end{aligned} \quad (74)$$

where the sign “-” in  $H_1^*$  denotes the absence of a variable  $l^*$  by elimination.

Thus the canonically transformed Hamiltonian  $H^*$  in (65) up to the first-order finally becomes

$$H^*(L^*, G^*, -, g^*) = H_0^*(L^*) + H_1^*(L^*, G^*, -, g^*), \quad (75)$$

using the first-order generating function  $S_1$  to transform  $H$  into  $H^*$

$$S_1(L^*, G^*, l^*, g^*) = \int (H_1(L^*, G^*, l^*, g^*) - H_1^*(L^*, G^*, -, g^*)) dt^*. \quad (76)$$

Higher-order solutions can be obtained by similar ways.

## 4.2 Solution for $L$

Now let us apply the Hori’s method to symplectic integrators. Here we consider that the integrable part of Hamiltonian  $H_0$  in (64) corresponds to  $H = -\frac{\mu^2}{2L^2}$  in (43), and the perturbed part of Hamiltonian  $H_1$  in (64) corresponds to  $H_{\text{err}}$  in (43). Henceforward we use the notation  $H_0$  as the integrable part and  $H_1$  as the perturbed part of Hamiltonian in this section. In summary,

$$\begin{cases} \tilde{H} = H + H_{\text{err}} & \cdots (43) \\ \uparrow & \uparrow \quad \uparrow \\ H = H_0 + H_1 & \cdots (64) \end{cases} \quad (77)$$

Using the Lie transformation (63), we obtain final solutions for  $L, G, l, g$ . As for  $L$  up to the first-order,

$$\begin{aligned} L &= L^* + \{L^*, S_1\} \\ &= L^* + \left( \frac{\partial L^*}{\partial l^*} \frac{\partial S_1}{\partial L^*} - \frac{\partial L^*}{\partial L^*} \frac{\partial S_1}{\partial l^*} + \frac{\partial L^*}{\partial g^*} \frac{\partial S_1}{\partial G^*} - \frac{\partial L^*}{\partial G^*} \frac{\partial S_1}{\partial g^*} \right) \end{aligned}$$

$$\begin{aligned}
&= L^* - \frac{\partial S_1}{\partial l^*} \\
&= L^* - \frac{\partial}{\partial l^*} \int (H_1 - H_1^*) dt^* \\
&= L^* - \frac{1}{n^*} \frac{d}{dt^*} \int (H_1 - H_1^*) dt^* \\
&= L^* - \frac{1}{n^*} (H_1 - H_1^*), \tag{78}
\end{aligned}$$

where  $\frac{\partial}{\partial l^*}$  is replaced by  $\frac{1}{n^*} \frac{d}{dt^*}$  since  $t^*$  affects on  $H_1$  only through  $l^*$ . Note that since the second term of the right-hand side of the analytic solution (78) is a small quantity of first-order, we can replace  $n^*$  by  $n_0$  here.

In (78),  $L^*$  and  $n^*$  are constants which should be determined by their initial conditions (or observation values). Representing the initial conditions by subscript 0 as  $L_0$  and  $l_0$ , we get

$$L_0 = L^* - \frac{1}{n^*} (H_{1,t=0} - H_1^*), \tag{79}$$

$$\therefore L^* = L_0 + \frac{1}{n^*} (H_{1,t=0} - H_1^*), \tag{80}$$

where  $H_{1,t=0}$  is the initial value of  $H_1$  when  $t = 0$ .

Now we can compare the analytic solution of  $L$  (80) with a solution by numerical symplectic integration. Substituting  $L^*$  of (80) into (78), we have plotted the time variation of the analytic solution of  $L$  in Figure 1 together with a solution by numerical symplectic integration using the second-order symplectic integrator. We have chosen the value of stepsize  $\tau$  as  $1/100$  of the orbital period  $T$ , i.e.  $\tau/T = 0.01$ . Initial conditions of the two-body system are listed in Table 1. The analytical solution by (78) and the numerical solution coincide very well within the first-order approximation. A higher order analytical solution will further reduce the difference between these two solutions indicated in the lower panel of Figure 1.

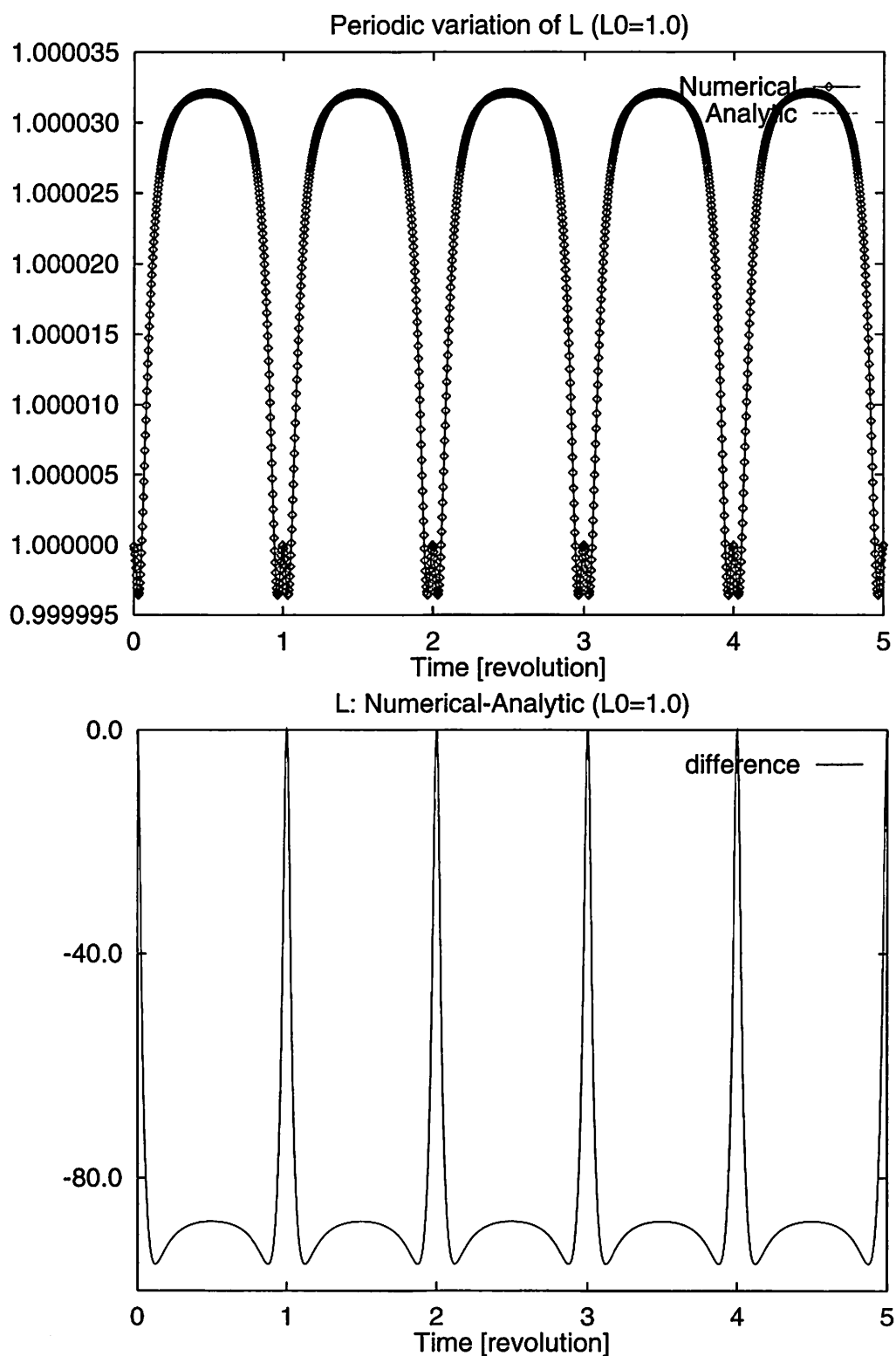
semimajor axis	$a$	1.0
eccentricity	$e$	0.5
argument of pericenter (degrees)	$\omega$	20.0
initial mean anomaly	$l_0$	0.0
mass coefficient	$\mu$	1.0

**Table 1.** Initial conditions for the two-body system used in this section.

Incidentally, from (78) and (80) we obtain

$$\begin{aligned}
L &= L_0 + \frac{1}{n^*} (H_{1,t=0} - H_1^*) - \frac{1}{n^*} (H_1 - H_1^*) \\
&= L_0 + \frac{1}{n^*} (H_{1,t=0} - H_1). \tag{81}
\end{aligned}$$

This means that the secular error of the action  $L$  can be removed up to the first-order by an appropriate selection of the initial value of  $H_1$  so that  $\langle H_{1,t=0} - H_1 \rangle = H_{1,t=0} - H_1^* = 0$ . It is the essential idea of the “iterative start” in Saha and Tremaine (1992) intending to reduce the secular numerical error in the angle  $l$ . We will discuss this fact in later sections.



**Figure 1.** (Upper) analytic and numerical solutions of the Delaunay element  $L$  in the system described in Table 1. The squares denote numerical solution by the second-order symplectic integrator and the lines denote solution by the first-order perturbation theory. (Lower) the difference of the two solutions (numerical – analytical) magnified by  $10^{10}$ .

### 4.3 Solution for $G$

Similar to  $L$ , solution for  $G$  can be obtained up to the first-order as

$$\begin{aligned}
G &= G^* + \{G^*, S_1\} \\
&= G^* + \left( \frac{\partial G^*}{\partial l^*} \frac{\partial S_1}{\partial L^*} - \frac{\partial G^*}{\partial L^*} \frac{\partial S_1}{\partial l^*} + \frac{\partial G^*}{\partial g^*} \frac{\partial S_1}{\partial G^*} - \frac{\partial G^*}{\partial G^*} \frac{\partial S_1}{\partial g^*} \right) \\
&= G^* - \frac{\partial S_1}{\partial g^*} \\
&= G^* - \frac{\partial}{\partial g^*} \int (H_1 - H_1^*) dt^*.
\end{aligned} \tag{82}$$

However, since  $H_1$  does not contain  $g^*$  at all, it becomes

$$\frac{\partial}{\partial g^*} \int (H_1 - H_1^*) dt^* = 0, \tag{83}$$

which means

$$G = G^* = G_0 = \text{constant}. \tag{84}$$

Hence there are no secular nor periodic numerical errors in  $G$  in the second-order symplectic integrator considered here. Actually, it is proved that any type of explicit symplectic integrator rigorously preserves the total angular momentum of system within the range of round-off errors (Yoshida, 1990a; Gladman *et al.*, 1991). In Figure 2, we have plotted relative error of the angular momentum of the two-body system,  $G/G_0 - 1$ , by the symplectic numerical integration. We can see the relative error of  $G$  is very close to the order of the round-off of the computation system,  $O(10^{-16})$ .

### 4.4 Solution for $l$

Same as  $L$  and  $G$ ,

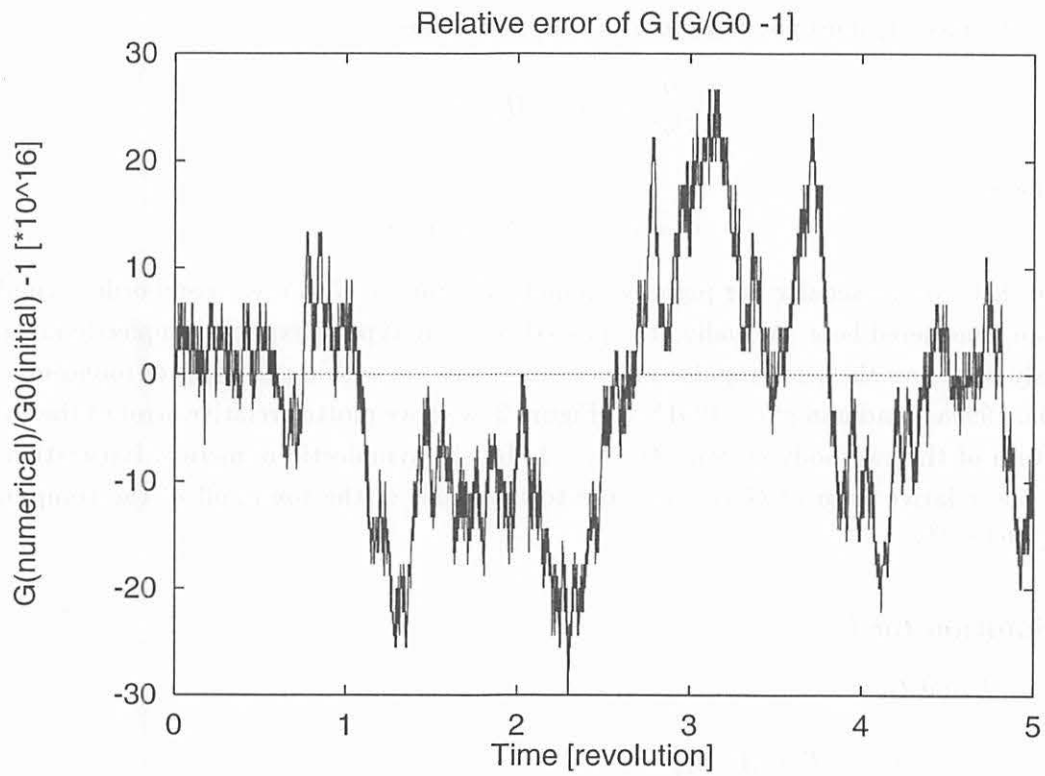
$$\begin{aligned}
l &= l^* + \{l^*, S_1\} \\
&= l^* + \left( \frac{\partial l^*}{\partial l^*} \frac{\partial S_1}{\partial L^*} - \frac{\partial l^*}{\partial L^*} \frac{\partial S_1}{\partial l^*} + \frac{\partial l^*}{\partial g^*} \frac{\partial S_1}{\partial G^*} - \frac{\partial l^*}{\partial G^*} \frac{\partial S_1}{\partial g^*} \right) \\
&= l^* + \frac{\partial S_1}{\partial L^*} \\
&= l^* + \frac{\partial}{\partial L^*} \int (H_1 - H_1^*) dt^*.
\end{aligned} \tag{85}$$

Now the secular error of  $l$  is caused by  $l^*$ , and the periodic error of  $l$  is caused by  $\frac{\partial S_1}{\partial L^*}$ . We show the specific derivation for each of them in the next sections.

#### 4.4.1 Secular error of $l$

The canonical equation of motion on  $l^*$  using the new Hamiltonian (75) is

$$\frac{dl^*}{dt} = \frac{\partial H^*}{\partial L^*}$$



**Figure 2.** The relative error of total angular momentum,  $G/G_0 - 1$ . The quantization around  $10^{-16}$  denotes that the error is due only to round-off because the machine-epsilon of double precision is  $2.2204460492503131 \times 10^{-16}$  in our system (HP-UX 11).

$$\begin{aligned}
&= \frac{\partial}{\partial L^*} \left( -\frac{\mu^2}{2L^{*2}} + \frac{\tau^2 \mu^2}{24a^{*4}\eta^{*5}} \left( 1 + \frac{e^{*2}}{2} \right) \right) \\
&= \frac{\mu^2}{L^{*3}} + \frac{\tau^2 \mu^2}{24} \frac{\partial}{\partial L^*} \left[ a^{*-4} (1 - e^{*2})^{-\frac{2}{5}} \left( 1 + \frac{e^{*2}}{2} \right) \right] \\
&= \frac{\mu^2}{L^{*3}} - \frac{\tau^2 \mu n^*}{12a^{*3}\eta^{*5}} \left( 1 + \frac{5e^{*2}}{4} \right). \quad \left( \because \frac{\mu}{L^*} = \frac{n^{*2}a^{*3}}{n^*a^{*2}} = n^*a^* \right) \quad (86)
\end{aligned}$$

Now we substitute  $L^*$  in (80) into the first term of the right-hand side of (86),

$$\begin{aligned}
\frac{\mu^2}{L^{*3}} &= \frac{\mu^2}{\left( L_0 + \frac{1}{n^*}(H_{1,t=0} - H_1^*) \right)^3} \\
&\simeq \frac{\mu^2}{L_0^3} - \frac{3\mu^2}{L_0^4 n^*} (H_{1,t=0} - H_1^*), \quad (87)
\end{aligned}$$

up to the leading order term of  $H_1/H_0$ .

Therefore (86) becomes

$$\begin{aligned}
\frac{dl^*}{dt} &= \frac{\mu^2}{L_0^3} - \frac{3\mu^2}{L_0^4 n^*} (H_{1,t=0} - H_1^*) - \frac{\mu \tau^2 n^*}{12a^{*3}\eta^{*5}} \left( 1 + \frac{5e^{*2}}{4} \right), \\
&= \frac{\mu^2}{L_0^3} + \tau^2 n^* \left( -\frac{3a^*}{\mu} H_{1,t=0} + \frac{\mu}{24a^{*3}\eta^{*3}} \right). \quad (88)
\end{aligned}$$

$$\therefore l^* = \frac{\mu^2}{L_0^3} t + \tau^2 n^* \left( -\frac{3a^*}{\mu} H_{1,t=0} + \frac{\mu}{24a^{*3}\eta^{*3}} \right) t + l_0. \quad (89)$$

The second term in the right-hand side of (89) represents the secular error of  $l$  up to the first-order of the perturbation theory.

#### 4.4.2 Periodic error of $l$

The periodic error of  $l$  is more complex to calculate. From (85),

$$\begin{aligned}
\frac{\partial S_1}{\partial L^*} &= \frac{\partial}{\partial L^*} \int (H_1 - H_1^*) dt^* \\
&= \frac{\partial}{\partial L^*} \int \left[ \frac{\tau^2 \mu^2}{24} \left( \frac{6}{r^{*4}} - \frac{3a^* (1 - e^{*2})}{r^{*5}} - \frac{2}{a^* r^{*3}} \right) - \frac{\tau^2 \mu^2}{24} \frac{1 + \frac{e^{*2}}{2}}{a^{*4} \eta^{*5}} \right] dt^*. \quad (90)
\end{aligned}$$

It is clear that we have to perform following two calculations successively:

1. Indefinite integration of

$$\int \frac{dt^*}{r^{*n}} = \frac{1}{n^*} \int \frac{dl^*}{r^{*n}}, \quad (n = 3, 4, 5)$$

2. Partial differentiation of the indefinite integrals by  $L^*$ .

Since all the periodic terms are of the first-order, the variables with superscript  $*$  can be replaced by those without  $*$ . We neglect most of the superscripts  $*$  in the following discussion for simplicity.



**Indefinite integrals of  $1/r^{*n}$**  We know relationships between  $r$ ,  $f$ , and  $l$  as

$$r = \frac{a(1 - e^2)}{1 + e \cos f}, \quad \frac{dl}{df} = \frac{r^2}{a^2 \eta}, \quad (91)$$

and the relationships of cosines

$$\cos^2 f = \frac{1}{2}(1 + \cos 2f), \quad \cos^3 f = \frac{1}{4}(\cos 3f + 3 \cos f). \quad (92)$$

Using above equations,

$$\begin{aligned} \int \frac{dl}{r^3} &= \int \frac{1}{r^3} \frac{r^2}{a^2 \eta} df \\ &= \frac{1}{a^3 \eta^3} (f + e \sin f) + \text{constant}, \end{aligned} \quad (93)$$

$$\begin{aligned} \int \frac{dl}{r^4} &= \int \frac{1}{r^4} \frac{r^2}{a^2 \eta} df \\ &= \frac{1}{a^4 \eta^5} \left[ \left(1 + \frac{e^2}{2}\right) f + 2e \sin f + \frac{e^2}{4} \sin 2f \right] + \text{constant}, \end{aligned} \quad (94)$$

$$\begin{aligned} \int \frac{dl}{r^5} &= \int \frac{1}{r^5} \frac{r^2}{a^2 \eta} df \\ &= \frac{1}{a^5 \eta^7} \left[ \left(1 + \frac{3e^2}{2}\right) f + \left(3e + \frac{3e^3}{4}\right) \sin f + \frac{3e^2}{4} \sin 2f + \frac{e^3}{12} \sin 3f \right] + \text{constant}. \end{aligned} \quad (95)$$

Substituting (93)(94)(95) into (86),

$$\begin{aligned} \int \left( \frac{6}{r^4} - \frac{3a(1 - e^2)}{r^5} - \frac{2}{ar^3} \right) dl &= 6 \int \frac{dl}{r^4} - 3a(1 - e^2) \int \frac{dl}{r^5} - \frac{2}{a} \int \frac{dl}{r^3} \\ &= \frac{1}{a^4 \eta^5} \left[ \left(1 + \frac{e^2}{2}\right) f + \left(e - \frac{e^3}{4}\right) \sin f - \frac{3e^2}{4} \sin 2f - \frac{e^3}{4} \sin 3f \right]. \end{aligned} \quad (96)$$

Hence, the first-order generating function  $S_1$  becomes by (76) with the superscript \* as

$$\begin{aligned} S_1 &= \int (H_1(L^*, G^*, l^*, g^*) - H_1^*) dt^* \\ &= \int \left[ \frac{\tau^2 \mu^2}{24} \left( \frac{6}{r^{*4}} - \frac{3a^*(1 - e^{*2})}{r^{*5}} - \frac{2}{a^* r^{*3}} \right) - \frac{\tau^2 \mu^2}{24} \frac{1 + \frac{e^{*2}}{2}}{a^{*4} \eta^{*5}} \right] dt^* \\ &= \frac{1}{n^*} \int \left[ \frac{\tau^2 \mu^2}{24} \left( \frac{6}{r^{*4}} - \frac{3a^*(1 - e^{*2})}{r^{*5}} - \frac{2}{a^* r^{*3}} \right) - \frac{\tau^2 \mu^2}{24} \frac{1 + \frac{e^{*2}}{2}}{a^{*4} \eta^{*5}} \right] dl^* \\ &= \frac{\tau^2 \mu^2}{24 n^* a^{*4} \eta^{*5}} \left[ \left(1 + \frac{e^{*2}}{2}\right) f^* + \left(e^* - \frac{e^{*3}}{4}\right) \sin f^* - \frac{3e^{*2}}{4} \sin 2f^* - \frac{e^{*3}}{4} \sin 3f^* - \left(1 + \frac{e^{*2}}{2}\right) l^* \right]. \end{aligned} \quad (97)$$

**Partial derivatives by  $L^*$**  Next we have to calculate the partial derivative  $\frac{\partial S_1}{\partial L^*}$ . All the necessary partial derivatives are given in Appendix C. As for the coefficient part in (97), it becomes (neglecting superscript \*)

$$\frac{\tau^2 \mu^2}{24n^* a^{*4} \eta^{*5}} = \frac{\mu^4}{(\mu a (1 - e^2))^{\frac{5}{2}}} = \frac{\mu^4}{G^5}, \quad (98)$$

Hence

$$\frac{\partial}{\partial L} \frac{\mu^2}{na^4 \eta^5} = \mu^4 \frac{\partial}{\partial L} G^{-5} = 0. \quad (99)$$

Similarly, periodic terms of  $f$  and  $l$  can be differentiated using the relationship

$$\frac{\partial f}{\partial L} = \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f, \quad \frac{\partial l}{\partial L} = 0, \quad (100)$$

as

$$\begin{aligned} \frac{\partial}{\partial L} \left( 1 + \frac{e^2}{2} \right) f &= f \frac{\partial}{\partial L} \left( 1 + \frac{e^2}{2} \right) + \left( 1 + \frac{e^2}{2} \right) \frac{\partial f}{\partial L} \\ &= \frac{G^2}{L^3} f + \frac{G^2}{eL^3} \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f, \end{aligned} \quad (101)$$

$$\begin{aligned} \frac{\partial}{\partial L} \left( e - \frac{e^3}{4} \sin f \right) &= \sin f \frac{\partial}{\partial L} \left( e - \frac{e^3}{4} \right) + \left( e - \frac{e^3}{4} \right) \cos f \frac{\partial f}{\partial L} \\ &= \frac{G^2}{eL^3} \left[ \left( 1 - \frac{3e^2}{4} \right) \sin f + \left( e - \frac{e^3}{4} \right) \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f \cos f \right], \end{aligned} \quad (102)$$

$$\begin{aligned} \frac{\partial}{\partial L} (e^2 \sin 2f) &= \sin 2f \frac{\partial}{\partial L} e^2 + e^2 \cdot 2 \cos 2f \frac{\partial f}{\partial L} \\ &= \frac{G^2}{eL^3} \cdot 2e \left[ \sin 2f + e \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \cos 2f \sin f \right], \end{aligned} \quad (103)$$

$$\begin{aligned} \frac{\partial}{\partial L} (e^3 \sin 3f) &= \sin 3f \frac{\partial e^3}{\partial L} + e^3 \cdot 3 \cos 3f \frac{\partial f}{\partial L} \\ &= \frac{G^2}{eL^3} \cdot 3e^2 \left[ \sin 3f + e \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \cos 3f \sin f \right], \end{aligned} \quad (104)$$

$$\begin{aligned} \frac{\partial}{\partial L} \left( 1 + \frac{e^2}{2} \right) l &= l \frac{\partial}{\partial L} \left( 1 + \frac{e^2}{2} \right) \\ &= \frac{G^2}{L^3} l. \end{aligned} \quad (105)$$

Adapting (99)(101)(102)(103)(104) (105) for (97), we get the partial derivative of  $S_1$  by  $L$  (or  $L^*$ ) as

$$\begin{aligned}
\frac{\partial S_1}{\partial L} &= \frac{\tau^2 \mu^2}{24na^4 \eta^5} \left[ \frac{G^2}{L^3} f + \frac{G^2}{eL^3} \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f \right. \\
&\quad + \frac{G^2}{eL^3} \left\{ \left( 1 - \frac{3e^2}{4} \right) \sin f + \left( e - \frac{e^3}{4} \right) \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f \cos f \right\} \\
&\quad - \frac{3}{4} \frac{G^2}{eL^3} \cdot 2e \left\{ \sin 2f + e \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f \cos 2f \right\} \\
&\quad \left. - \frac{1}{4} \frac{G^2}{eL^3} \cdot 3e^2 \left\{ \sin 3f + e \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f \cos 3f \right\} - \frac{G^2}{L^3} l \right] \\
&= \frac{\tau^2 \mu^2}{24na^4 \eta^5} \left[ \frac{G^2}{L^3} (f - l) + \frac{G^2}{eL^3} \left\{ \left( 1 - \frac{3e^2}{4} \right) \sin f + \left( e - \frac{e^3}{4} \right) \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f \right. \right. \\
&\quad + \left( 1 + \frac{e^2}{2} \right) \left( \frac{a}{r} + \frac{L^2}{G^2} \right) - \frac{3}{4} \cdot 2e \left( \sin 2f + e \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \cos 2f \right) \\
&\quad \left. \left. - \frac{1}{4} \cdot 3e^2 \left( \sin 3f + e \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \cos 3f \right) \right\} \right]. \tag{106}
\end{aligned}$$

Therefore, from (85) and (89), the final solution for  $l$  up to the first-order perturbation in this theory becomes as follows:

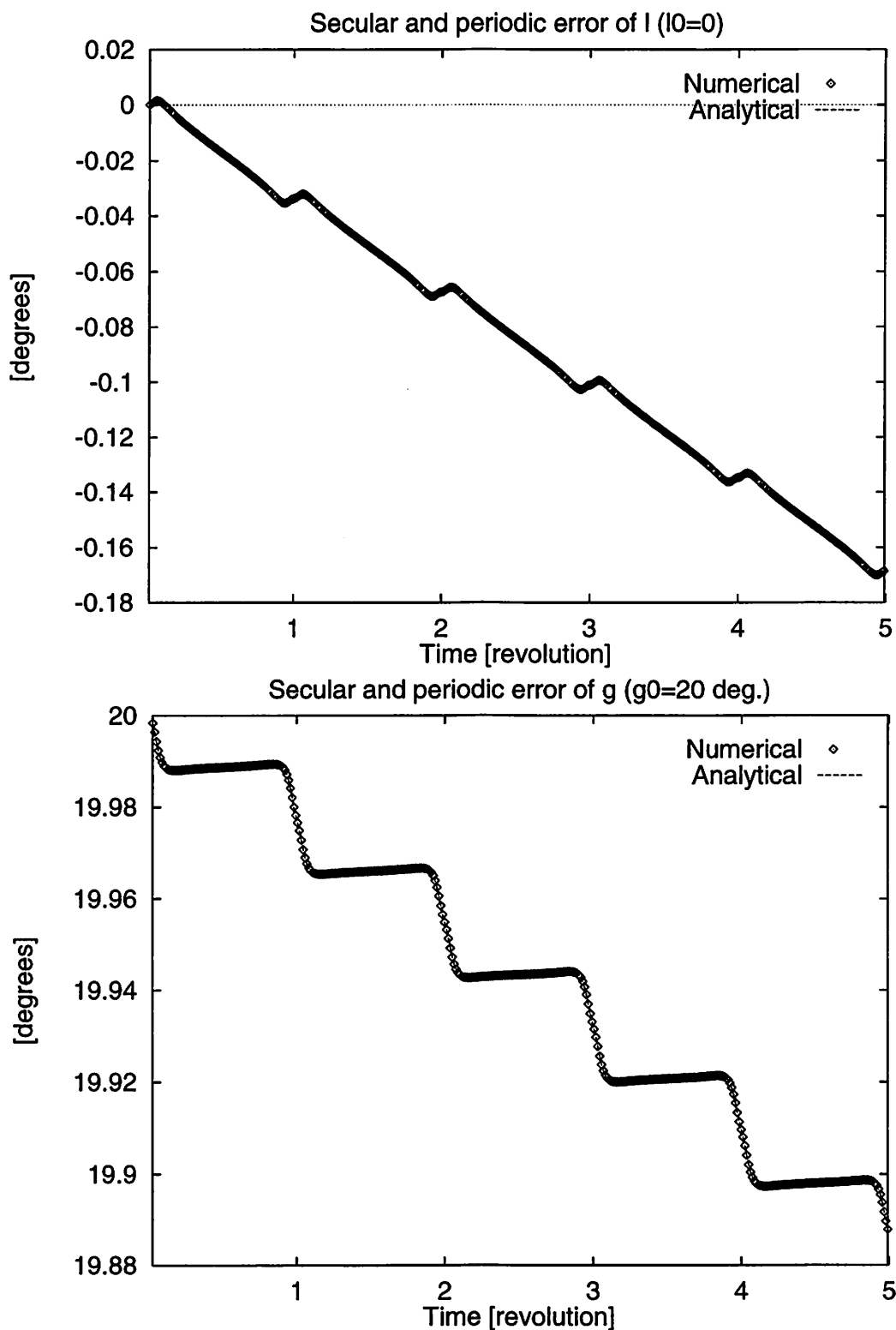
$$\begin{aligned}
l &= l^* + \{l^*, S_1\} \\
&= l^* + \frac{\partial S_1}{\partial L^*} \\
&= l_0 + \frac{\mu^2}{L_0^3} t + \tau^2 n^* \left( -\frac{3a^*}{\mu} H_{1,t=0} + \frac{\mu}{24a^{*3} \eta^{*3}} \right) t \\
&\quad + \frac{\tau^2 \mu^2}{24na^4 \eta^5} \left[ \frac{G^2}{L^3} (f - l) + \frac{G^2}{eL^3} \left\{ \left( 1 - \frac{3e^2}{4} \right) \sin f + \left( e - \frac{e^3}{4} \right) \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f \right. \right. \\
&\quad + \left( 1 + \frac{e^2}{2} \right) \left( \frac{a}{r} + \frac{L^2}{G^2} \right) - \frac{3}{4} \cdot 2e \left( \sin 2f + e \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \cos 2f \right) \\
&\quad \left. \left. - \frac{1}{4} \cdot 3e^2 \left( \sin 3f + e \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \cos 3f \right) \right\} \right]. \tag{107}
\end{aligned}$$

The errors of  $l$  are plotted in upper panels of Figures 3, 4, and 5. The upper panel of Figures 3 shows the secular and periodic errors of  $l$  by the numerical symplectic integration and the analytical perturbation theory, compared with the exact solution of the Keplerian motion. The upper panel of Figures 4 shows only the periodic errors of  $l$ . The upper panel of Figures 5 shows the difference of the periodic errors of  $l$  by the numerical integration and analytical perturbation theory. Higher-order analytical solution will reduce the difference between these two.

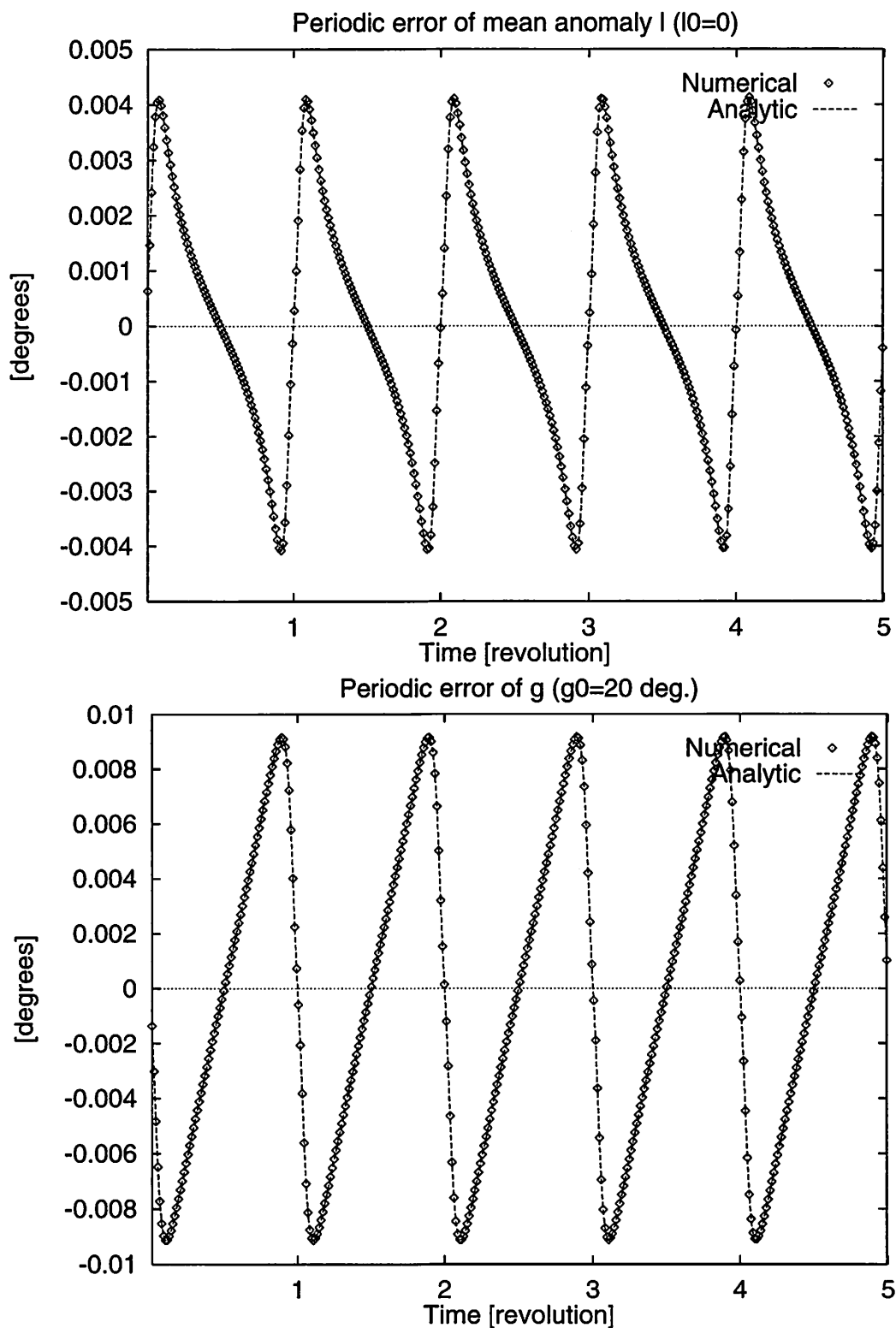
#### 4.5 Solution for $g$

Same as  $l$ ,

$$g = g^* + \{g^*, S_1\}$$



**Figure 3.** The numerical and analytical solution of the secular + periodic errors in  $l$  (upper) and  $g$  (lower), compared with the exact solution of the Keplerian motion. The squares denote the numerical solution by the second-order symplectic integrator, the lines denote the solution by the first-order perturbation theory.



**Figure 4.** The numerical and analytical solution of the periodic errors in  $l$  (upper) and  $g$  (lower), subtracting the secular error shown in Figure 3. The squares denote the numerical solution by the second-order symplectic integrator, the lines denote the solution by the first-order perturbation theory.

$$\begin{aligned}
&= g^* + \left( \frac{\partial g^*}{\partial l^*} \frac{\partial S_1}{\partial L^*} - \frac{\partial g^*}{\partial L^*} \frac{\partial S_1}{\partial l^*} + \frac{\partial g^*}{\partial g^*} \frac{\partial S_1}{\partial G^*} - \frac{\partial g^*}{\partial G^*} \frac{\partial S_1}{\partial g^*} \right) \\
&= g^* + \frac{\partial S_1}{\partial G^*} \\
&= g^* + \frac{\partial}{\partial G^*} \int (H_1 - H_1^*) dt^*. \tag{108}
\end{aligned}$$

Now we know that the secular error of  $g$  is caused from  $g^*$ , and the periodic error of  $g$  is caused from  $\frac{\partial S_1}{\partial G^*}$ .

#### 4.5.1 Secular error of $g$

The canonical equation of motion using  $g^*$  using the new Hamiltonian (75) becomes

$$\begin{aligned}
\frac{dg^*}{dt} &= \frac{\partial H^*}{\partial G^*} \\
&= \frac{\partial}{\partial G^*} \left( -\frac{\mu^2}{L^{*2}} + \frac{\tau^2 \mu^2}{24 a^{*4} \eta^{*5}} \left( 1 + \frac{e^{*2}}{2} \right) \right) \\
&= -\frac{\tau^2 \mu}{4 a^{*3} \eta^{*6}} \left( 1 + \frac{e^{*2}}{4} \right) n. \tag{109}
\end{aligned}$$

$$\therefore g^* = -\frac{\tau^2 \mu}{4 a^{*3} \eta^{*6}} \left( 1 + \frac{e^{*2}}{4} \right) n^* t + g_0. \tag{110}$$

The second term in the right-hand side of (110) represents the the secular error of  $g$  up to the first-order of the perturbation theory.  $g_0$  is the initial value of  $g$  when  $t = 0$ .

#### 4.5.2 Periodic error of $g$

The periodic error of  $g$  can be obtained as the same way as  $l$ . From (108),

$$\begin{aligned}
\frac{dS_1}{dG^*} &= \frac{\partial}{\partial G^*} \int (H_1 - H_1^*) dt^* \\
&= \frac{\partial}{\partial G^*} \int \left[ \frac{\tau^2 \mu^2}{24} \left( \frac{6}{r^{*4}} - \frac{3a^* (1 - e^{*2})}{r^{*5}} - \frac{2}{a^* r^{*3}} \right) - \frac{\tau^2 \mu^2}{24} \frac{1 + \frac{e^{*2}}{2}}{a^{*4} \eta^{*5}} \right] dt^*. \tag{111}
\end{aligned}$$

Same as  $l$ , it is clear that we have to perform the following two calculations successively:

1. Indefinite integration of

$$\int \frac{dt^*}{r^{*n}} = \frac{1}{n^*} \int \frac{dl^*}{r^{*n}}, \quad (n = 3, 4, 5)$$

2. Partial differentiation of the indefinite integrals by  $G^*$ .

The first task has been already done. The second task is given as follows:

$$\begin{aligned}
\frac{\partial S_1}{\partial G^*} &= \frac{\partial}{\partial G^*} \int (H_1 - H_1^*) dt^* \\
&= \frac{\partial}{\partial G^*} \frac{1}{n^*} \int (H_1 - H_1^*) dl^* \\
&= \frac{\partial}{\partial G^*} \left[ \frac{\tau^2 \mu^2}{24 n^* a^{*4} \eta^{*5}} \mathcal{F}(f^*, l^*) \right], \tag{112}
\end{aligned}$$

where  $\mathcal{F}(f^*, l^*)$  denotes the periodic function of  $f^*$  and  $l^*$  described in the integrand of (97) as

$$\mathcal{F}(f^*, l^*) = \left(1 + \frac{e^{*2}}{2}\right) f^* + \left(e^* - \frac{e^{*3}}{4}\right) \sin f^* - \frac{3e^{*2}}{4} \sin 2f^* - \frac{e^{*3}}{4} \sin 3f^* - \left(1 + \frac{e^{*2}}{2}\right) l^*. \quad (113)$$

Henceforward, the variables with superscript  $*$  are replaced by those without  $*$  since all the periodic terms are of the first-order of perturbation. We neglect  $*$  in the following discussion for simplicity.

As for the coefficient part in (112),

$$\frac{\mu^2}{na^4\eta^5} = \frac{a^{\frac{3}{2}}}{\mu^{\frac{1}{2}}} \frac{\mu^2}{a^4\eta^5} = \mu^4 \left[ \mu a (1 - e^2) \right] = \frac{\mu^4}{G^5}. \quad (114)$$

Therefore

$$\begin{aligned} \frac{\partial S_1}{\partial G} &= \frac{\partial}{\partial G} \int (H_1 - H_1^*) dt^* \\ &= \frac{\tau\mu^4}{24} \left[ \frac{\partial}{\partial G} \left( \frac{\mu}{G^5} \right) \mathcal{F}(f, l) + \frac{1}{G^5} \frac{\partial}{\partial G} \mathcal{F}(f, l) \right]. \end{aligned} \quad (115)$$

It is possible to calculate the partial derivatives of  $\mathcal{F}(f^*, l^*)$  in the same way as in  $l$ :

$$\begin{aligned} \frac{\partial}{\partial G} \left( 1 + \frac{e^2}{2} \right) f &= e \frac{\partial e}{\partial G} f + \left( 1 + \frac{e^2}{2} \right) \frac{\partial f}{\partial G} \\ &= e \left( -\frac{G}{eL^2} \right) f + \frac{\partial e}{\partial G} \left( 1 + \frac{e^2}{2} \right) \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f, \end{aligned} \quad (116)$$

$$\frac{\partial}{\partial G} \left( 1 + \frac{e^2}{2} \right) l = e \frac{\partial e}{\partial G} l = e \left( -\frac{G}{eL^2} \right) l, \quad (117)$$

$$\begin{aligned} \frac{\partial}{\partial G} \left( e - \frac{e^3}{4} \right) \sin f &= \left( 1 - \frac{3e^2}{4} \right) \frac{\partial e}{\partial G} \sin f + \left( e - \frac{e^3}{4} \right) \cos f \frac{\partial f}{\partial G} \\ &= \frac{\partial e}{\partial G} \left[ \left( e - \frac{3e^2}{4} \right) \sin f + \left( e - \frac{e^3}{4} \right) \cos f \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f \right], \end{aligned} \quad (118)$$

$$\begin{aligned} \frac{\partial}{\partial G} (e^2 \sin 2f) &= 2e \frac{\partial e}{\partial G} \sin 2f + 2e^2 \cos 2f \frac{\partial f}{\partial G} \\ &= \frac{\partial e}{\partial G} \left[ 2e \sin 2f + 2e^2 \cos 2f \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f \right], \end{aligned} \quad (119)$$

$$\begin{aligned} \frac{\partial}{\partial G} (e^3 \sin 3f) &= 3e^2 \frac{\partial e}{\partial G} \sin 3f + 3e^3 \cos 3f \frac{\partial f}{\partial G} \\ &= \frac{\partial e}{\partial G} \left[ 3e^2 \sin 3f + 3e^3 \cos 3f \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f \right]. \end{aligned} \quad (120)$$

Substituting (116)(117)(118)(119)(120) into (112), periodic errors of  $g$  becomes

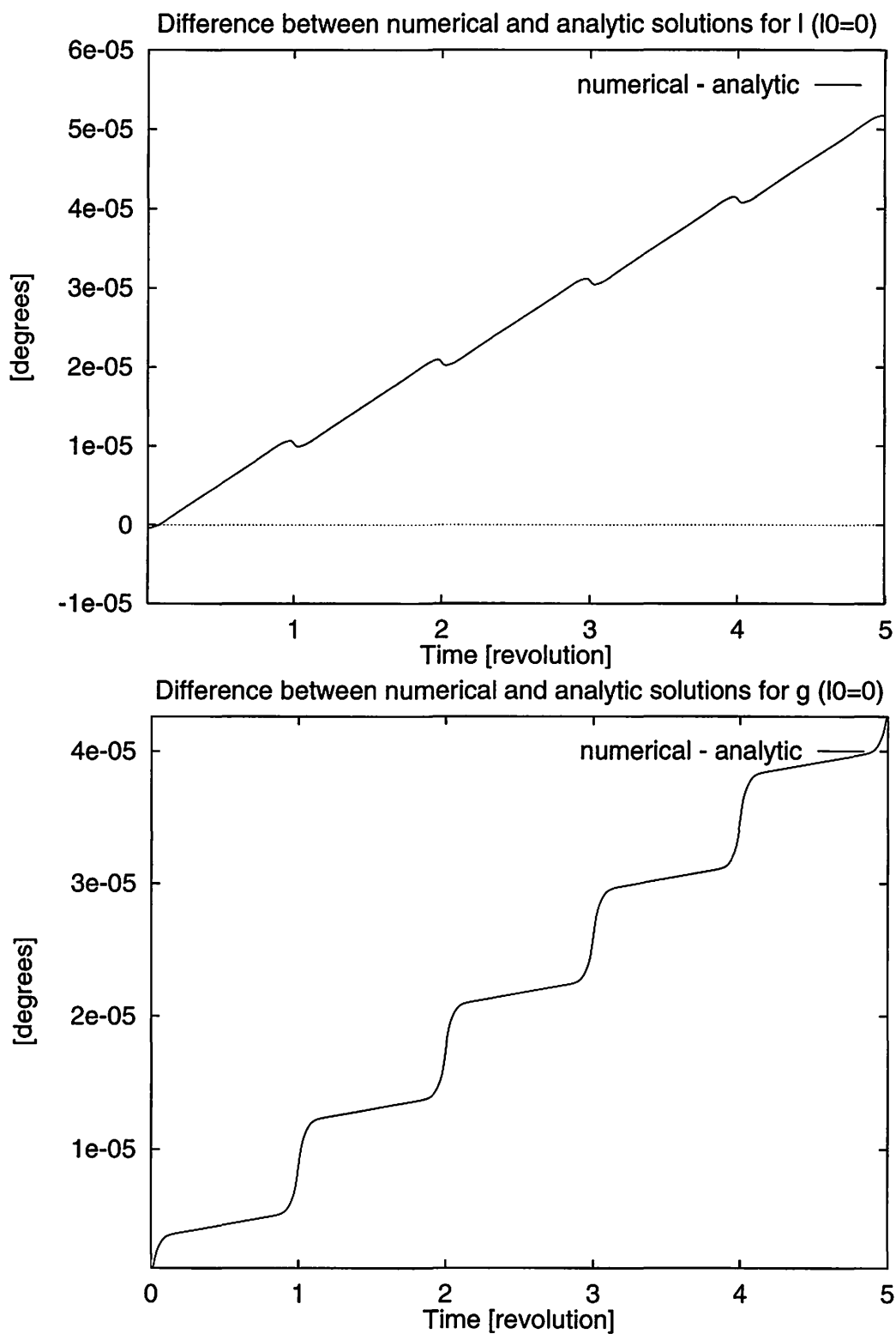
$$\begin{aligned}
\frac{\partial S_1}{\partial G^*} &= \frac{\partial}{\partial G^*} \int (H_1 - H_1^*) dt^* \\
&= \frac{\tau^2 \mu^4}{24} \left[ -\frac{5}{G^6} \left\{ \left(1 + \frac{e^2}{2}\right) (f - l) + \left(e - \frac{e^3}{4}\right) \sin f - \frac{3e^2}{4} \sin 2f - \frac{e^3}{4} \sin f \right\} \right. \\
&\quad + \frac{1}{G^5} \left\{ -\frac{G}{L^2} (f - l) - \frac{G}{eL^2} \left( \left(1 + \frac{e^2}{2}\right) \left(\frac{a}{r} + \frac{L^2}{G^2}\right) \right. \right. \\
&\quad \left. \left. + \left(1 - \frac{3e^2}{4}\right) \sin f + \left(e - \frac{e^3}{4}\right) \cos f \left(\frac{a}{r} + \frac{L^2}{G^2}\right) \right. \right. \\
&\quad \left. \left. - \frac{3}{2} e \left( \sin 2f + e \cos 2f \left(\frac{a}{r} + \frac{L^2}{G^2}\right) \right) \right. \right. \\
&\quad \left. \left. \left. - \frac{3}{4} e^2 \left( \sin 3f + e \cos 3f \left(\frac{a}{r} + \frac{L^2}{G^2}\right) \right) \right) \right\} \right]. \tag{121}
\end{aligned}$$

Therefore, from (108) and (121), final solution for  $g$  up to the first-order perturbation can be obtained as follows:

$$\begin{aligned}
g &= g^* + \{g^*, S_1\} \\
&= g^* + \frac{\partial S_1}{\partial G^*} \\
&= g_0 - \frac{\tau^2 \mu}{4a^{*3} \eta^{*6}} \left(1 + \frac{e^{*2}}{4}\right) n^* t \\
&\quad + \frac{\tau^2 \mu^4}{24} \left[ -\frac{5}{G^6} \left\{ \left(1 + \frac{e^2}{2}\right) (f - l) + \left(e - \frac{e^3}{4}\right) \sin f - \frac{3e^2}{4} \sin 2f - \frac{e^3}{4} \sin f \right\} \right. \\
&\quad + \frac{1}{G^5} \left\{ -\frac{G}{L^2} (f - l) - \frac{G}{eL^2} \left( \left(1 + \frac{e^2}{2}\right) \left(\frac{a}{r} + \frac{L^2}{G^2}\right) \right. \right. \\
&\quad \left. \left. + \left(1 - \frac{3e^2}{4}\right) \sin f + \left(e - \frac{e^3}{4}\right) \cos f \left(\frac{a}{r} + \frac{L^2}{G^2}\right) \right. \right. \\
&\quad \left. \left. - \frac{3}{2} e \left( \sin 2f + e \cos 2f \left(\frac{a}{r} + \frac{L^2}{G^2}\right) \right) \right. \right. \\
&\quad \left. \left. \left. - \frac{3}{4} e^2 \left( \sin 3f + e \cos 3f \left(\frac{a}{r} + \frac{L^2}{G^2}\right) \right) \right) \right\} \right]. \tag{122}
\end{aligned}$$

The solution for  $g$  is plotted in the lower panels of Figures 3, 4, and 5. The lower panel of Figures 3 shows the secular and periodic errors of  $g$  by the numerical symplectic integration and analytical perturbation theory, compared with the exact solution of the Keplerian motion. The lower panel of Figures 4 shows only the periodic errors of  $g$ . The lower panel of Figures 5 shows the difference of the periodic errors of  $g$  by the numerical integration and the analytical perturbation theory. Higher-order analytical solution will reduce the difference between these two.





**Figure 5.** The differences of the periodic errors of  $l$  (upper) and  $g$  (lower) obtained by the numerical and analytical methods, which are equivalent to the differences in two data (squares and lines) in Figure 3.

#### 4.6 Another interpretation of the numerical error source

The principle source of the numerical error by symplectic integration is that of mean anomaly  $l$ , which grows linearly in time. According to Kinoshita et al. (1991), we can derive the source of the secular numerical error of  $l$  (89) as follows: The mean anomaly  $\tilde{l}$  of the surrogate system is dominated by the surrogate Hamiltonian  $\tilde{H}$ , and the mean anomaly  $l$  of the real system is dominated by the real Hamiltonian  $H$ . The equations of motion which  $\tilde{l}$  and  $l$  follow except for periodic parts would be

$$\frac{d\tilde{l}}{dt} = \frac{\partial \tilde{H}}{\partial L}, \quad (123)$$

$$\frac{dl}{dt} = \frac{\partial H}{\partial L}, \quad (124)$$

respectively. Subtracting (124) from (123), we get

$$\begin{aligned} \frac{d}{dt}(\tilde{l} - l) &= \frac{\partial \tilde{H}}{\partial L} - \frac{\partial H}{\partial L} \\ &= \frac{\partial H_{\text{err}}}{\partial L} \\ &= -\frac{\tau^2 \mu n}{12a^3 \eta^5} \left(1 + \frac{5e^2}{4}\right), \end{aligned} \quad (125)$$

which is equal to Eq. (19) in Kinoshita et al. (1991), and also coincides with the second term of the right-hand side of our (86) except for the superscript \*.

In addition to (125), there is another source of the secular truncation error in the mean anomaly  $l$  due to the constant part of the truncation error the total energy,  $E$ . Since the surrogate Hamiltonian  $\tilde{H}$  is strictly preserved by symplectic integration, we have

$$\tilde{H} = H(q_0, p_0) + H_{\text{err}}(q_0, p_0) = H(q, p) + H_{\text{err}}(q, p), \quad (126)$$

to  $O(\tau^2)$  approximation.  $(q_0, p_0)$  are initial values, and  $(q, p)$  are the approximate solutions obtained by the symplectic integration. From (126), the truncation error of the total energy (i.e. secular part of Hamiltonian) becomes

$$\Delta E = H(q, p) - H(q_0, p_0) = H_{\text{err}}(q_0, p_0) - H_{\text{err}}(q, p). \quad (127)$$

Since  $H_{\text{err}}(q_0, p_0)$  is fixed, constant part of  $\Delta E$  is given by

$$\begin{aligned} \Delta E_c \equiv \langle \Delta E \rangle &= \langle H_{\text{err}}(q_0, p_0) - H_{\text{err}}(q, p) \rangle \\ &= H_{\text{err}}(q_0, p_0) - \langle H_{\text{err}}(q, p) \rangle \\ &= H_{\text{err}}(q_0, p_0) - \frac{\tau^2 \mu^2}{24a^4 \eta^5} \left(1 + \frac{e^2}{2}\right). \end{aligned} \quad (128)$$

In general, we can derive the constant bias in semimajor axis ( $\Delta a$ ) due to the constant part of orbital energy offset ( $\Delta E_c$ ) as follows:

$$E = -\frac{\mu}{2a}, \quad (129)$$

$$\therefore a = -\frac{\mu}{2E}, \quad (130)$$

$$\therefore \Delta a = \frac{\mu}{2E^2} \Delta E_c. \quad (131)$$

Hereafter we use  $\Delta E_c$  instead of  $\Delta E$  in order to remark explicitly it is constant.

From the Kepler's third law,  $n^2 a^3 = \mu$  is fixed in the gravitational two-body problem. Then the secular error of  $l$  due to  $\Delta a$  can be obtained by taking a variation of the Kepler's third law:

$$2\Delta n \cdot a^3 + n^2 \cdot 3a\Delta a = 0, \quad (132)$$

$$\therefore 2\Delta n a + 3n\Delta a = 0. \quad (133)$$

If we are to express the variation of the time derivative of  $l$  as

$$\frac{d}{dt} (\tilde{l} - l) = \Delta \dot{l} = \Delta n, \quad (134)$$

the secular error of  $l$  becomes from (133)

$$\Delta \dot{l} = \Delta n = -\frac{3n\Delta a}{2a}. \quad (135)$$

Hence the additional secular truncation error of the mean anomaly  $l$  due to  $\Delta E_c$  is from (128) and (131)

$$\begin{aligned} \Delta \dot{l} &= -\frac{3n\Delta a}{2a} \\ &= -\frac{3n}{2a} \cdot \frac{\mu}{2E^2} \cdot \Delta E_c \\ &= -\frac{3n\mu}{4a} \left(-\frac{2a}{\mu}\right)^2 \left[ H_{\text{err}}(q_0, p_0) - \frac{\tau^2 \mu^2}{24a^4 \eta^5} \left(1 + \frac{e^2}{2}\right) \right] \\ &= -\frac{3an}{\mu} H_{\text{err}}(q_0, p_0) + \frac{\tau^2 \mu n}{8a^3 \eta^5} \left(1 + \frac{e^2}{2}\right). \end{aligned} \quad (136)$$

Now we have the total secular truncation errors for the mean anomaly  $l$  by adding (125) and (136) as

$$\begin{aligned} \Delta \dot{l} &= -\frac{\tau^2 \mu n}{12a^3 \eta^5} \left(1 + \frac{5e^2}{4}\right) + \left[ -\frac{3an}{\mu} H_{\text{err}}(q_0, p_0) + \frac{\tau^2 \mu n}{8a^3 \eta^5} \left(1 + \frac{e^2}{2}\right) \right] \\ &= -\frac{3an}{\mu} H_{\text{err}}(q_0, p_0) + \frac{\tau^2 \mu n}{24a^3 \eta^3}, \end{aligned} \quad (137)$$

which is equal to the second term of the right-hand side of our (89).

However, this derivation in Kinoshita et al. (1991) is somewhat confusing in spite of its correct solution in (137). Especially, the appearance of the second source of the secular error  $\Delta \dot{l}$  (136) seems to be too abrupt. This is caused by the confusion of  $L$  and  $L^*$  in (125). After the operation of time-averaging (by a certain canonical transformation),  $L^*$ , instead of  $L$ , must be used to describe the canonical equation of motion; (125) should be derived correctly from the canonically transformed equations of motion such as

$$\frac{dl^*}{dt} = -\frac{\partial H^*}{\partial L^*}. \quad (138)$$

Moreover, the averaged value of  $L$  is not equal to  $L^*$ ;  $L^*$  has a constant bias to  $L$ , which is the true reason of the “another source of the secular truncation error in the mean anomaly due to the constant part of the truncation error in the energy” in Kinoshita et al. (1991). The detailed and exact form of  $L^*$  is presented in (80) in the previous sections.

#### 4.7 Dependence on initial configuration

As you can see in the equation (107), the secular numerical error in  $l$  arises from the coefficient of  $t$ , namely

$$\delta n_{\text{sec}} \equiv \tau^2 n \left( -\frac{3a}{\mu} H_{1,t=0} + \frac{\mu}{24a^3\eta^3} \right). \quad (139)$$

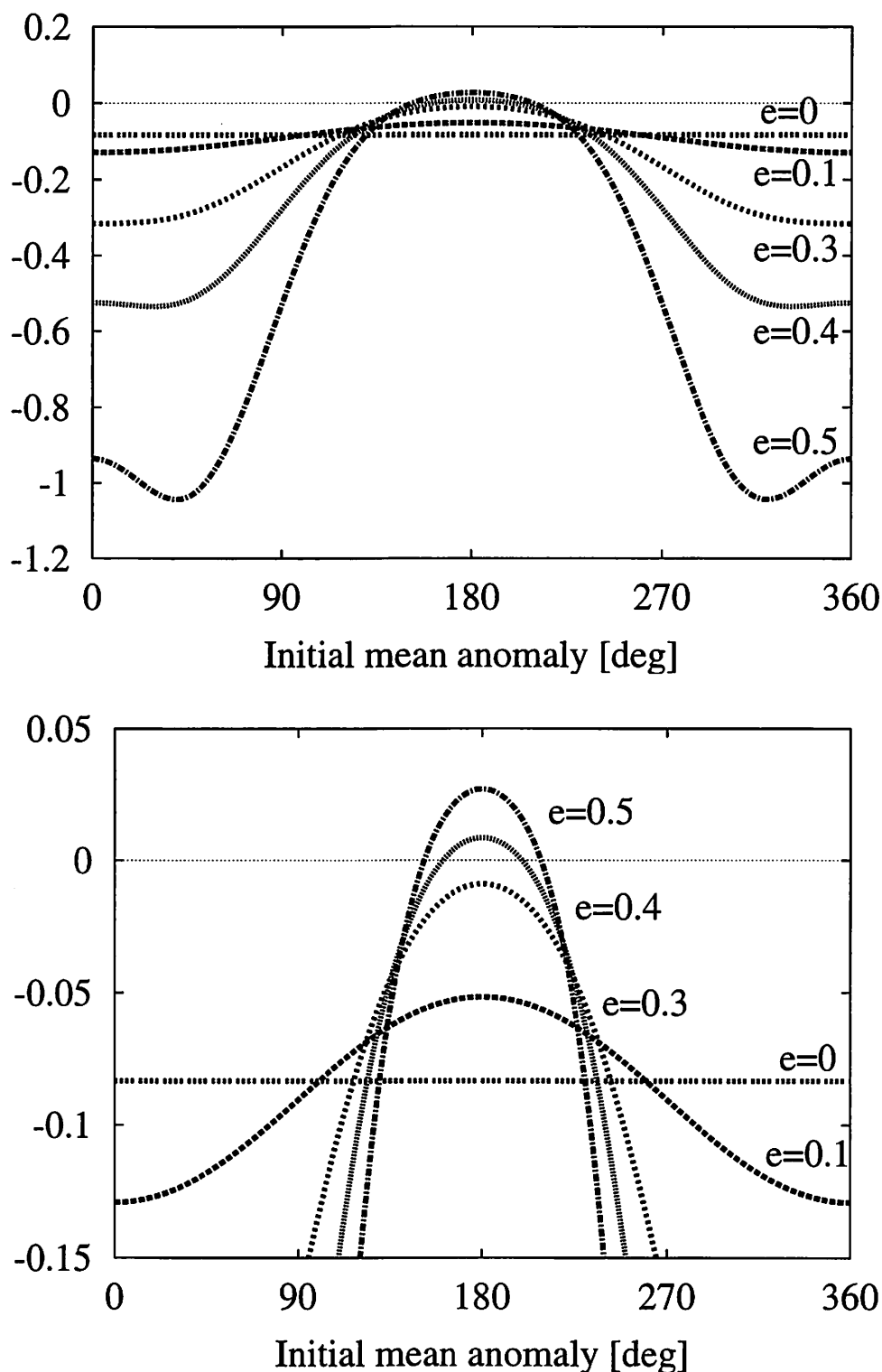
Hereafter we call the coefficient (139)  $\delta n_{\text{sec}}$ . Note that in the expression of  $\delta n_{\text{sec}}$  in (139) we neglected all of the superscripts  $*$ .

Not only  $\delta n_{\text{sec}}$  has a dependence on the initial longitude in  $H_{1,t=0}$  as discussed in the previous sections, but this has a dependence on the initial orbital shape,  $e$ . The effect of  $a$  is only to scale the unit of time, so we can neglect it from the discussion of numerical error here. We plot this  $\delta n_{\text{sec}}$ 's dependence on initial eccentricity as well as initial starting longitude in the two-body problem in Figure 6. We can anticipate from (139) that there are certain initial mean anomalies ( $l_0$ ) which make  $\delta n_{\text{sec}}$  very small, possible zero. Actually in some cases of higher eccentricities,  $\delta n_{\text{sec}}$  becomes zero at certain values of initial mean anomaly in Figure 6. However, generally  $\delta n_{\text{sec}}$  does not become zero whatever we change the initial mean anomaly.

In Figure 7, we exaggeratedly illustrate the trajectories of  $(l, L)$  of the two-body system described in this section. “Exact” denotes the exact solution which goes from  $(1, 0)$  and comes back at  $(1, 0)$  again. “Numerical” denotes the symplectic numerical solution which goes from a different point from  $(1, 0)$  and does not come back at the starting point. “Synthetic” denotes the analytical secular solution obtained by the perturbation theory which is close to the time-average of the numerical solution. As we see, the exact and the synthetic solutions are far from coincidence when the initial mean anomaly  $l_0 = 0$  (left panel), while they coincide pretty well when the initial mean anomaly  $l_0 = 180^\circ$  (right panel). This result corresponds to the result shown in Figure 6 in terms of the numerical error of the symplectic integrator.

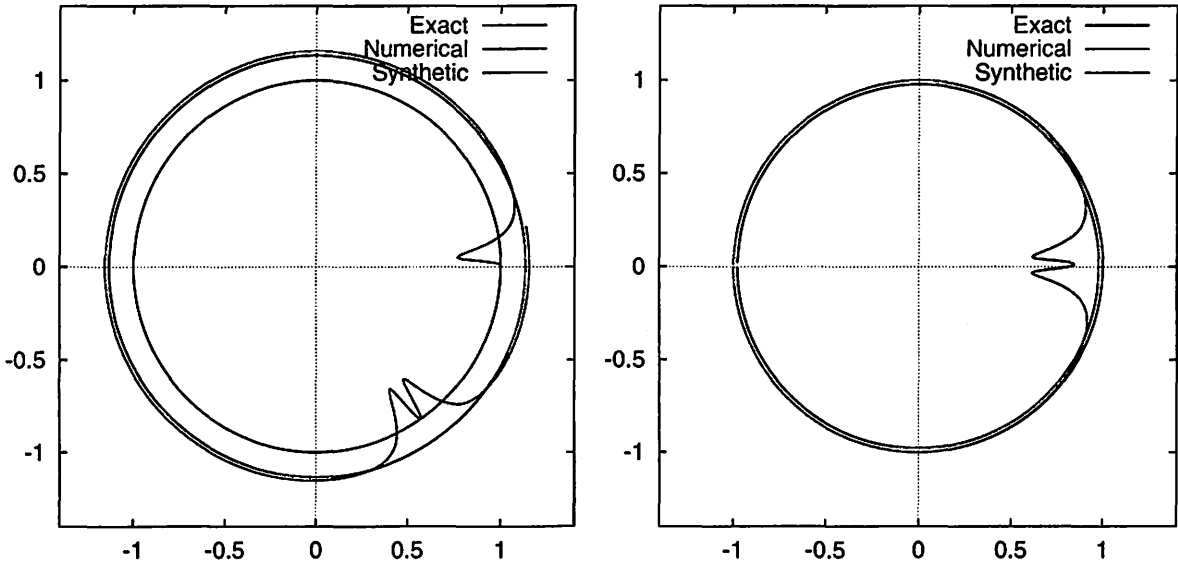
### 5. Reduction of errors by the iterative start

In the examples of Kepler problem in the previous sections, we could find an approximate analytical solution of symplectic numerical error by a perturbation theory. However in general many-body systems, it is quite difficult, or virtually impossible, to obtain an analytical form of the numerical error; hence we cannot know which initial mean anomaly would reduce the numerical error of symplectic integrator in analytical way. Thus we have to depend on a numerical way to look for the initial conditions which reduce numerical errors. This kind of numerical method we use when we start our integration has been originally proposed by Saha & Tremaine (1992) under the name of “iterative start.” Our method is essentially the same with their iterative start, but a bit different in its way of implementation. We mention our method and results in two kinds of three-body dynamical system: Sun–Jupiter and a fictitious



**Figure 6.** (Upper) the secular error coefficient of mean anomaly in the two-body problem (139) as a function of initial mean anomaly  $l_0$ . Five curves show the results when  $e = 0, 0.1, 0.3, 0.4$ , and  $0.5$ . (Lower) an enlarged panel of the central part of the upper one.

middle-sized body in the asteroidal belt, and a planet orbiting around a binary system named MACHO-97-BLG-41 which has been discovered by a gravitational microlensing event.



**Figure 7.** Exaggerated illustration of the trajectories of  $(l, L)$  of the two-body system (see Table 1). (Left) when  $l_0 = 0$ . (Right) when  $l_0 = 180^\circ$ . The line denoted “Exact” shows the exact solution which goes from  $(1, 0)$  and comes back at  $(1, 0)$  again. The line “Numerical” denotes the symplectic numerical solution which goes from a different point from  $(1, 0)$  and does not come back at the starting point. The line “Synthetic” means the analytical secular solution obtained by the perturbation theory which is close to the time-average of the numerical solution. We have intentionally exaggerated the deviation of the numerical and synthetic solutions from the exact ones in order to bring out the difference.

## 5.1 Perturbed motion of a middle-sized planet

### 5.1.1 Settings of numerical experiments

First we consider a weakly perturbed three-body system, Sun–Jupiter and a fictitious middle-sized body in the asteroidal belt (see Figure 8). The fictitious middle-sized body has a finite mass of  $1/10 M_{\text{Jupiter}}$ , hence the problem is not a restricted one. The initial orbital elements of the middle-sized body are similar to those of Ceres: when  $e = 0.1$  and  $e = 0.4$ ,  $a = 2.6\text{AU}$ . When  $e = 0.6$ ,  $a = 2.2\text{AU}$  so that we avoid its close encounters with the outer planet. The values of other orbital elements than  $a$  or  $e$  are the same as those of Ceres. The mass and initial orbital elements of the outer massive planet in this system is just the same as those of Jupiter. These initial orbital elements of the bodies are basically taken from the Development Ephemeris of JPL, DE245 (Standish, 1990).

A principle way to execute the “iterative start” in this section is as follows:

1. Given a nominal set of initial orbital elements, perform an integration with a very high

accuracy covering a shorter timespan than the main integration.

2. Choose several initial conditions of all relevant bodies from the results of the accurate integration.
3. Perform several short-term integrations with a normal accuracy using the initial conditions selected above.
4. Calculate the numerical differences between the accurate integration and each of the short-term integration.
5. Select an initial condition which produces the least numerical error as the set of starting orbital elements of the main integration.
6. Perform the main integration using the initial condition selected above.

Note that the integration periods of the accurate integration and the short-term integrations are much shorter than the period of the main integration. For example, when the period of the main integration is  $1 \times 10^8$  years, we would take a  $10^4$ -year for the accurate and the short-term integration periods. The initial conditions for the short-term integrations are chosen while the Jupiter-like planet orbits the Sun once (about twelve years). Interval among each initial condition is  $\sim 1^\circ$  in the Jupiter-like planet's longitude, or about ten days in time. We have illustrated the situation in Figure 9.

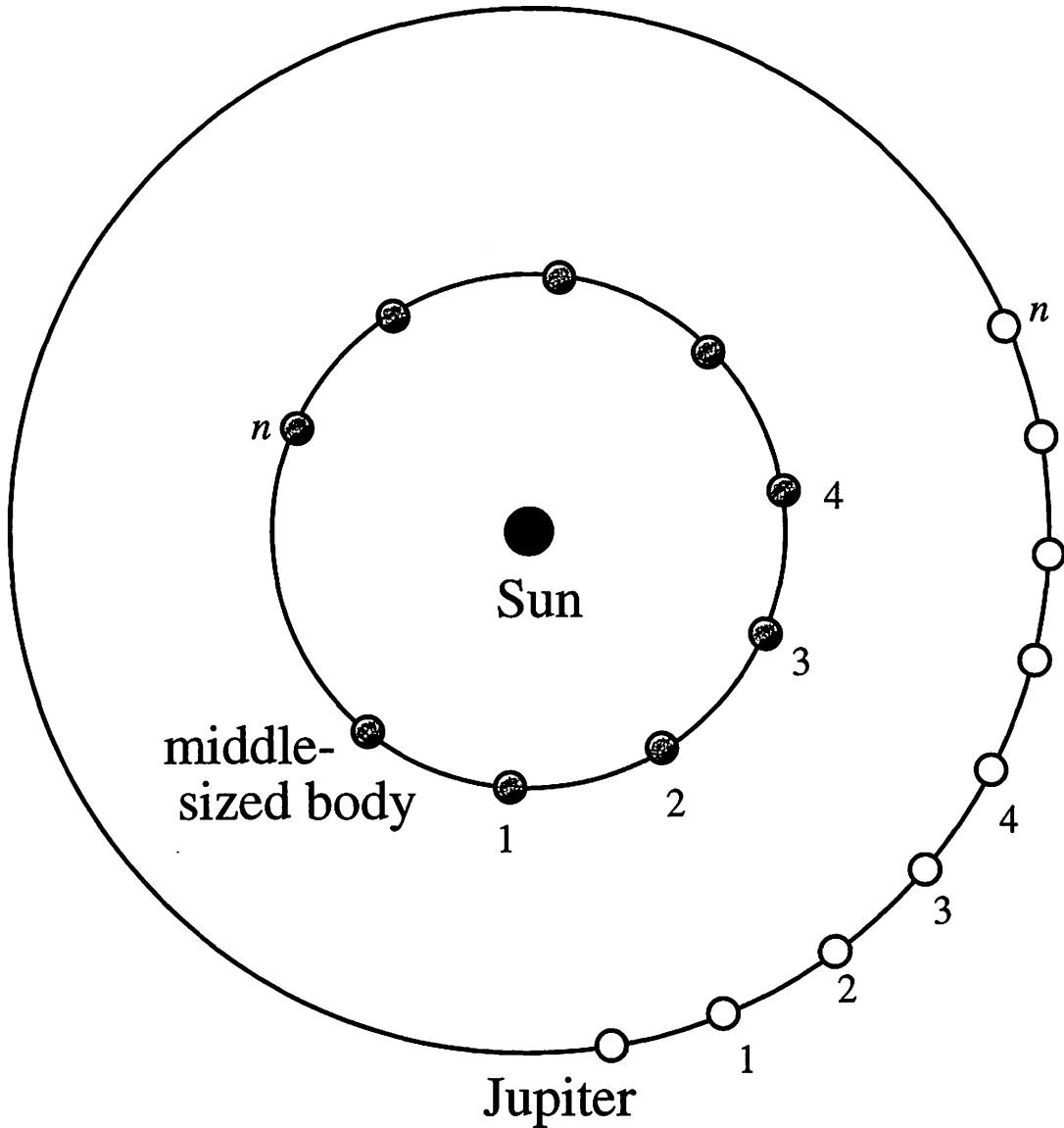
When perturbation to the Kepler motion is very small, we may be able to implement the "iterative start" on the system in a simpler way as follows:

1. Given a nominal set of initial orbital elements, we can fix the Keplerian osculating orbital orbits of each body.
2. Let each body move on its osculating orbit by a small interval (see Figure 8).
3. Name each position as  $1, 2, 3, \dots, n$ . At each position, perform two sets of numerical integrations of a short period: one is a very accurate integration, and the other's accuracy is the same as that of the main integration.
4. Compare the two sets of integrations and calculate their numerical difference at each position from 1 to  $n$ .
5. Repeat the above comparison until we reach a certain point,  $n$ .
6. Select an initial condition which produces the least numerical difference as the starting orbital elements of the main integration.
7. Perform the main integration using the set of initial condition selected above.

In the discussions below, we choose the latter procedure. Note again that the procedure is valid only for a slightly disturbed system such as the Sun–Jupiter–a middle-sized planet. We cannot apply this simplified procedure to a significantly perturbed system like the planetary

system around a binary which we will discuss later. This is due to a stronger perturbation on the planetary orbit from the short-term orbital motion of binary.

We fix the period of the accurate and the short-term integrations as  $2 \times 10^4$  years. We take the time interval of each set of initial condition as  $P_{\text{Jupiter}}/360$  where  $P_{\text{Jupiter}}$  is the orbital period of Jupiter. We choose 360 sets of initial conditions for comparison, meanwhile Jupiter rotates around the Sun once, and the middle-sized planet does twice or more.

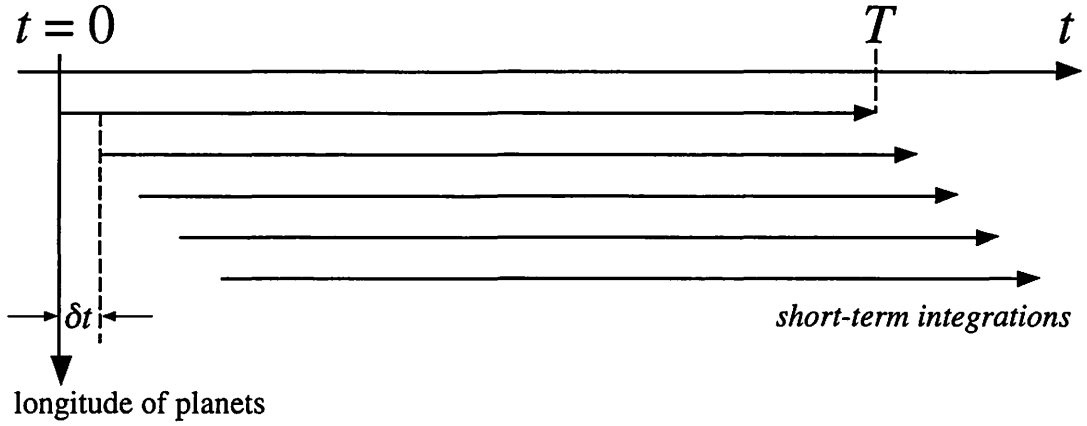


**Figure 8.** A schematic illustration of the Sun–Jupiter–a middle-sized body system. Each short-term integration starts at the numbered position from  $1, \dots, n$ . See also Figure 9.

### 5.1.2 Results of the numerical experiments

We have performed numerical experiments for the three-body planetary system using the second-order explicit symplectic integrator described in Section 3.. As for canonical variables





**Figure 9.** A schematic illustration of our way of implementation of the “iterative start.” Each short-term integration starts at a different time (or a different orbital position) on a same dynamical trajectory. See also Figure 8.

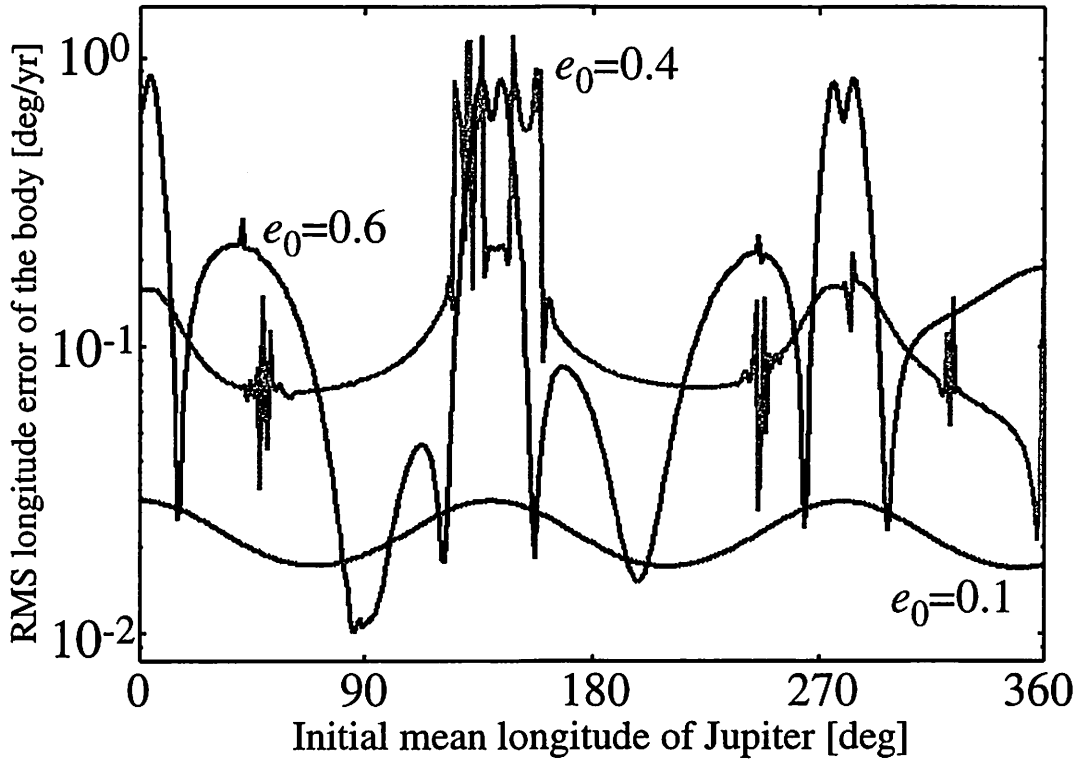
used in the scheme, we have adopted the Jacobi coordinate. The shorter integrations are done with a stepsize of 4 days. As the standard numerical integration with a high accuracy, we have performed an integration over  $2 \times 10^4$ -year period with a stepsize of  $0.0625 = 1/16$  days. We consider this standard integration much more accurate than the shorter integrations, and calculate the longitudinal difference of the middle-sized planet between the standard integration and the shorter integrations. We have chosen 360 sets of initial orbital conditions while Jupiter revolves once around the Sun from its initial position. We have tested three sets of numerical integrations, changing initial eccentricity of the middle-sized planet  $e_0$ ;  $e_0 = 0.1, 0.4$ , and  $0.6$ . For  $e_0 = 0.6$  set, the initial semimajor axis of the middle-sized planet is set as  $a_0 = 2.2\text{AU}$  so that we avoid its close encounter with Jupiter. In other sets,  $a_0 = 2.2\text{AU}$  which is similar to the semimajor axis of Ceres.

The root-mean-square (RMS) of the longitudinal error of the middle-sized planet per year is shown in Figure 10. Here the horizontal axis is denoted as the initial mean longitude of Jupiter, but note that the initial mean longitude of the middle-sized planet changes accordingly. If we fix the mean longitude of one planet and change that of another, it ends up with integration of the orbital motion in many different dynamical systems. This is not what we mean to study in this manuscript.

In Figure 10, we notice three interesting characters. First, root mean square (RMS) of the longitudinal error of the middle-sized planet differs a lot, depending on initial starting point on the dynamical system. Second, the rate of error reduction is much larger when the initial eccentricity of the middle-sized planet ( $e_0$ ) is large. When  $e_0 = 0.6$ , maximum difference of the RMS is nearly two orders of magnitude. Third, there are some spike-like features on the RMS curves, especially when  $e_0 = 0.4$ .

As for the first point, we can understand it in analogy with the similar analysis in the two-body problem discussed before. Since in Figure 10 there is no zero-axis in ordinate because the figure draws the RMS of the planet’s longitudinal error, we drawn a simple average of the longitudinal error when  $e_0 = 0.6$  in Figure 11. We clearly see that under certain values of

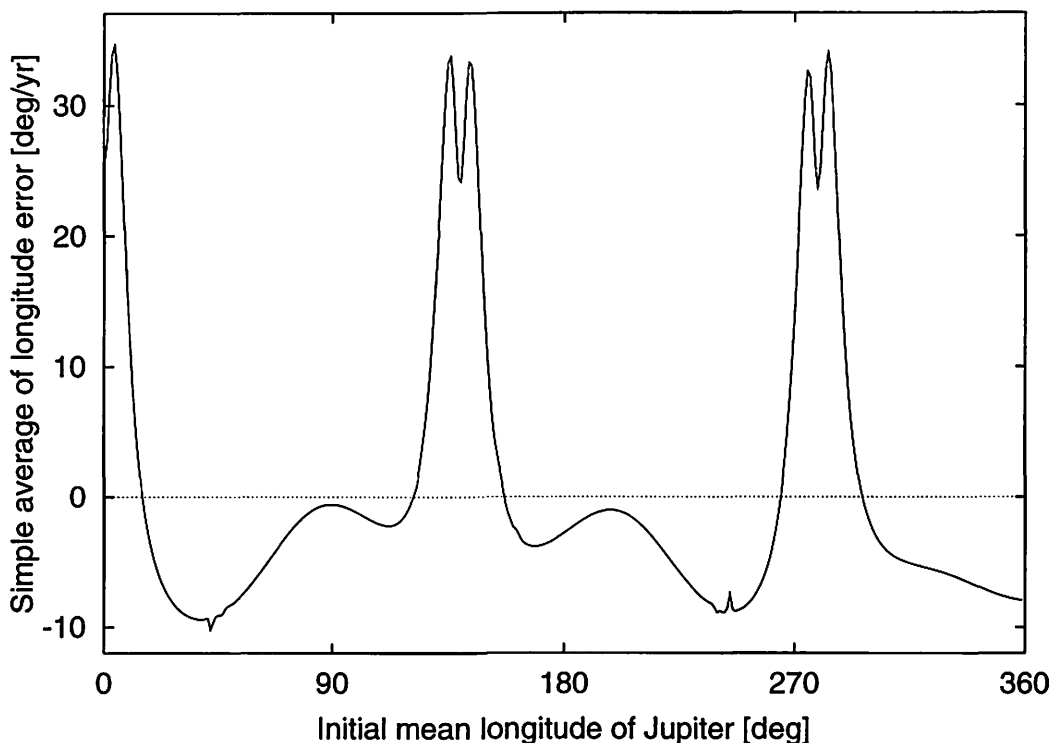
the initial conditions, the longitudinal error crosses zero axis. This means when we start the integrations from such initial conditions that lies on the zero-axis in Figure 11, we can reduce the longitudinal error of the planet to a large extent. In contrast when we choose a bad initial condition, the longitudinal error of the planet increases terribly, which degrades the accuracy of numerical integration very much.



**Figure 10.** The RMS numerical error of the middle-sized planet (deg/year) as a function of Jupiter’s initial mean longitude.  $e_0$  is the initial eccentricity of the the middle-sized planet.

As for the second point, we see the same trend as in the two-body problem discussed below (see Figure 6). We chose some of the typical numerical results of time-series and showed them in Figures 12 ( $e_0 = 0.1$ ), 13 ( $e_0 = 0.4$ ), and 14 ( $e_0 = 0.6$ ). When the initial eccentricity of the middle-sized planet is not so large as  $e_0 = 0.1$ , the degrees of the error reduction by the iterative start is not prominent. However as  $e_0$  grows, the degree of the error reduction becomes larger, even up to two orders of magnitude (the lower panel in Figure 14). We could find five initial conditions where the longitudinal error of the middle-sized planet becomes very small (or possible zero) when  $e_0 = 0.6$  as in Figure 11. But there is only one condition when  $e_0 = 0.4$  (near  $l_{J0} = 360^\circ$ ). When  $e_0 = 0.1$ , we could not find any of such appropriate initial conditions.

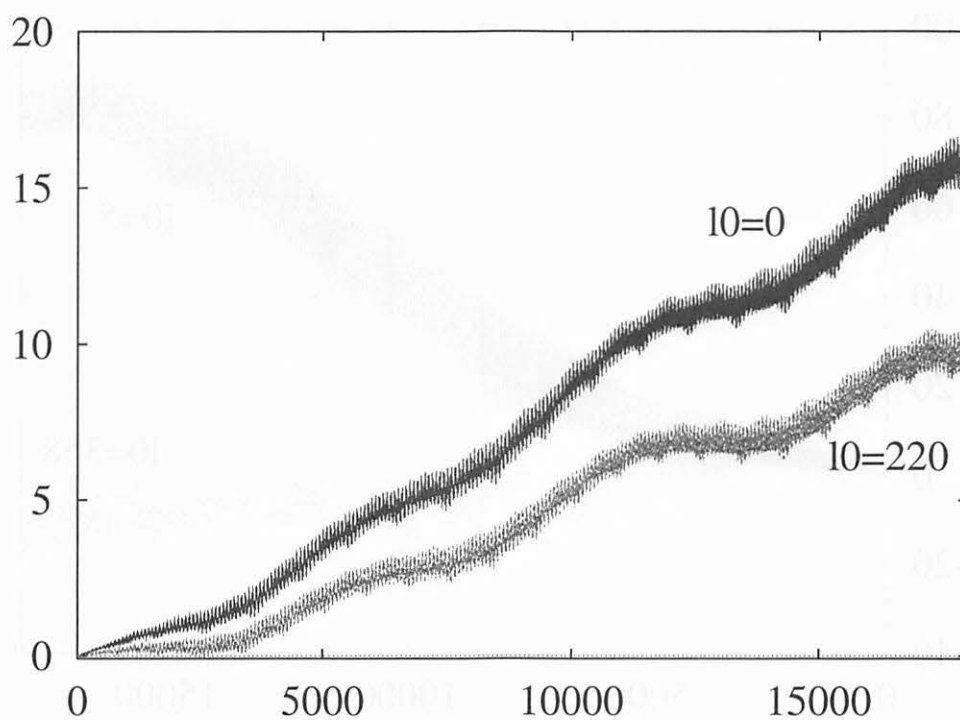
At present we do not have a definite answer why these numerical errors can be significantly reduced much when the initial eccentricity is low in the two-body and weakly perturbed three-body systems such as discussed above. It may be related to a kind of geometry in phase-space of the dynamical system. A simple guess goes on like this: when the eccentricities of bodies are large, geometry of the trajectory in phase-space is somewhat distorted or warped. This distortion could be common in both the real system dominated by Hamiltonian  $H$  and the



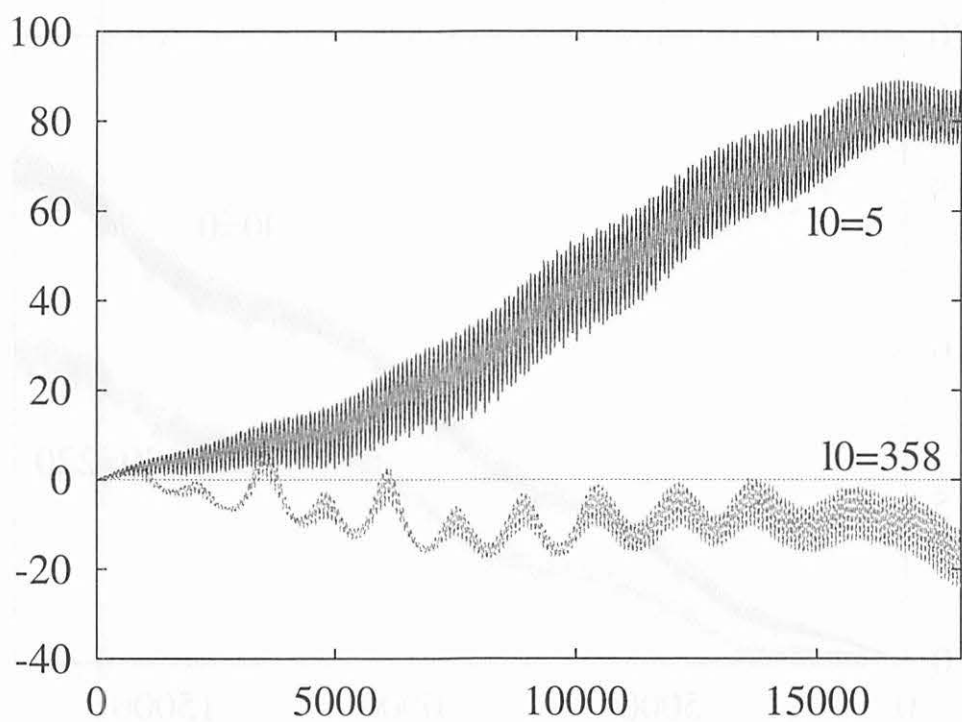
**Figure 11.** The averaged (simple sum) numerical error of the middle-sized planet (deg/year) as a function of Jupiter’s initial mean longitude when  $e_0 = 0.6$  in Figure 10.

surrogate system dominated by surrogate Hamiltonian  $\tilde{H}$  in (21). The initial conditions by which we can significantly reduce numerical error may be intersection points (possibly lines or plains if dimension of the phase-space is large) of the two distorted trajectories. On the other hand when the initial eccentricities of the bodies are small, the trajectory may be very smooth, and the possibility that the two trajectories intersect with each other may become lower, which leads to the non-existence of the initial conditions that reduce the longitudinal error to nearly zero in our numerical experiments. This discussion is still a simple guess. We have to seek a definitive answer confirming the structure of the phase-space in detail.

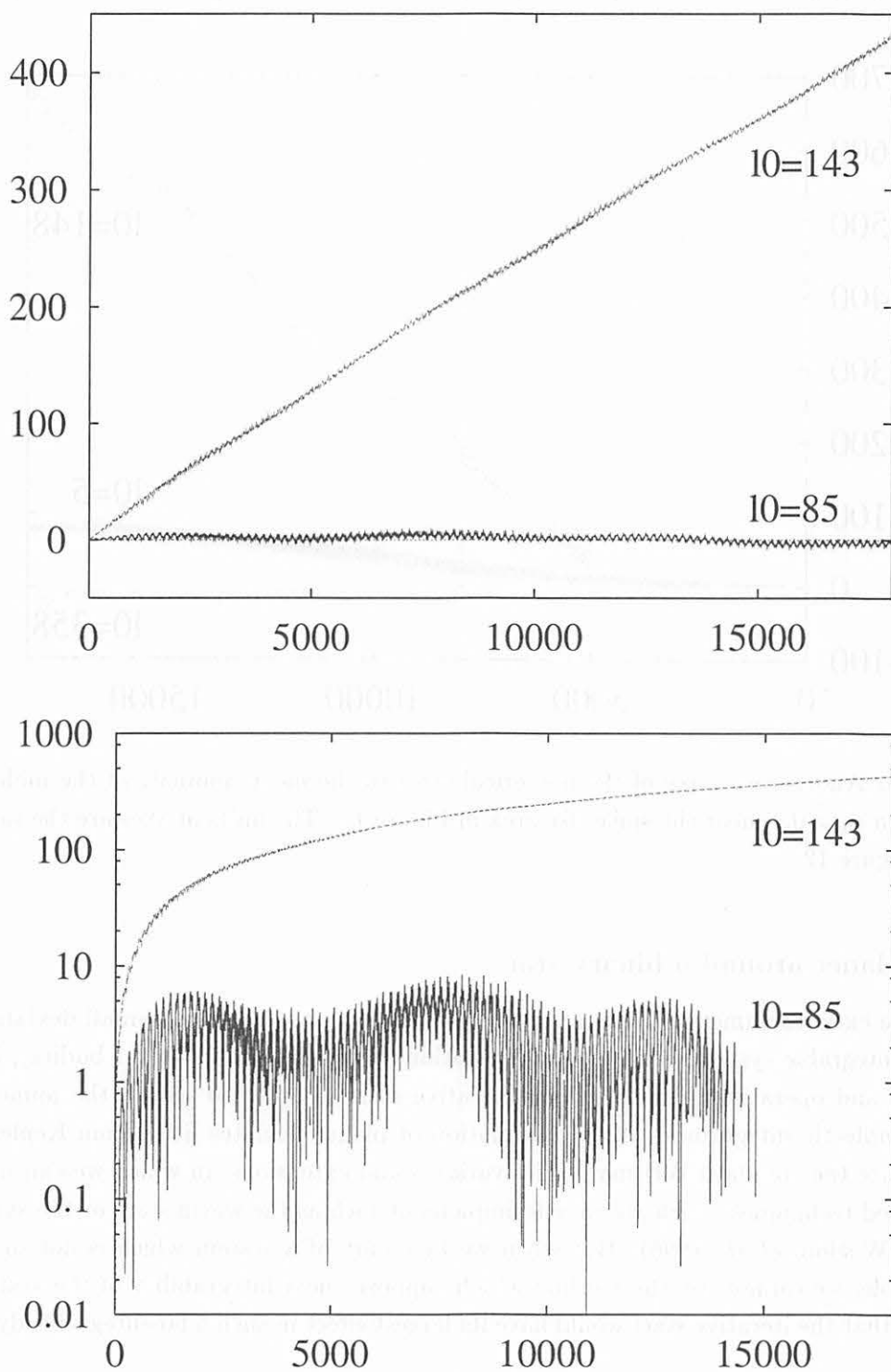
As for the reason of the third point of the spike-like feature, it is not obvious to explain. We took an example result of such spike-like features and showed what was going on there (Figure 15). The figure can be a counterpart of Figure 13 which shows an example time-series when  $e_0 = 0.4$  with  $l_{J0} = 358^\circ$  and  $l_{J0} = 5^\circ$ . When  $l_{J0} = 148^\circ$  that is just on the spike-like area, at the beginning the longitudinal error increases similarly to the results when  $l_{J0} = 5^\circ$ . However when  $t > 5000$  years, the slope of the curve of  $l_{J0} = 148^\circ$  suddenly increases, and the error increases rapidly afterwards. This nonlinear behavior of the numerical error may be similar to the stepsize resonance phenomena reported in WH-type symplectic integrator (Wisdom and Holman, 1992; Rauch and Holman, 1999) or symmetric multistep methods (Quinlan and Tremaine, 1990; Fukushima, 1998; Fukushima, 1999). But it would be not easy nor straightforward to understand why these spike-like features occur only in  $e_0 = 0.4$  systems, and why the spikes can work as not only to decrease the errors but also to increase



**Figure 12.** An example of the numerical error in the mean anomaly of the middle sized planet when  $e_0 = 0.1$ .  $l_0$  denotes the initial mean longitude of Jupiter. The unit of the vertical axis is deg/year, and the unit of the horizontal axis is year.

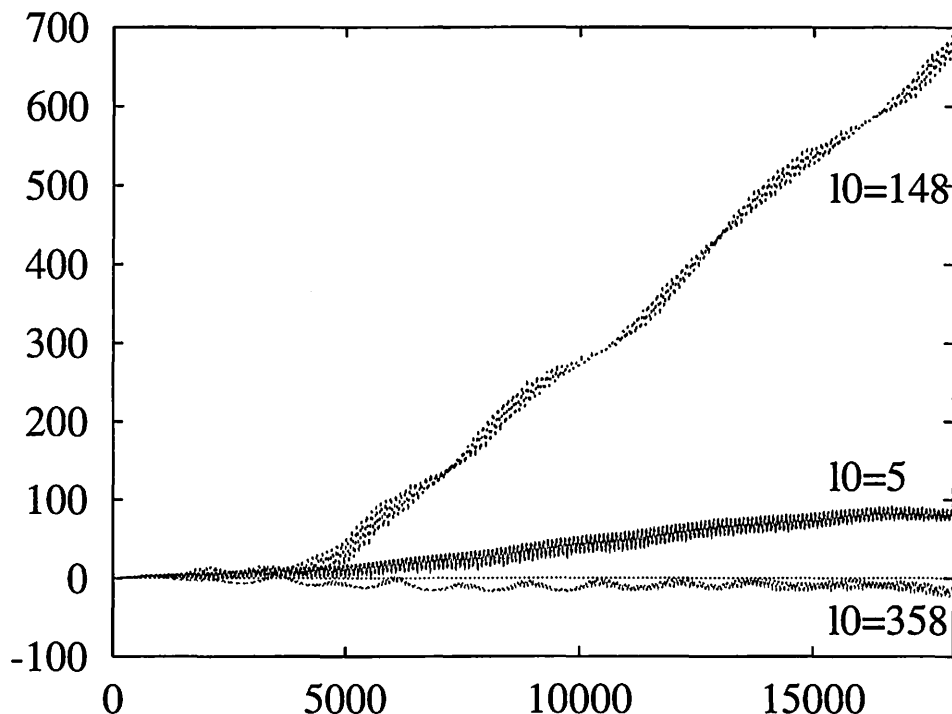


**Figure 13.** An example of the numerical error in the mean anomaly of the middle sized planet when  $e_0 = 0.4$ .  $l_0$  denotes the initial mean longitude of Jupiter. The units of axes are the same with those of Figure 12.



**Figure 14.** An example of the numerical error in the mean anomaly of the middle sized planet when  $e_0 = 0.6$ . The lower panel is a logarithmic version of the upper one. The units of axes are the same with those of Figure 12.

them (cf. Figure 10). Detailed inspection of what is going on in these spike-like regions and of its dependence on various orbital parameters are necessary to definitely answer our question.



**Figure 15.** Another example of the numerical error in the mean anomaly of the middle sized planet when  $e_0 = 0.4$ , near the spike-like area in Figure 10. The units of axes are the same with those of Figure 12.

## 5.2 A planet around a binary star

When we execute numerical integrations in a dynamical system with a small deviation from a certain integrable system (such as Kepler motion or free rotation of rigid bodies), it is not so efficient and operative to resort to the iterative start in order to reduce the numerical errors in symplectic integrators. When the motion of planet deviates little from Keplerian, we should utilize the standard WH map or its variants and extensions, in which we can use many sophisticated techniques which are easy to implement such as the warm start or the symplectic corrector (Wisdom *et al.*, 1996). But when we take care of a system which is not so close to be integrable, we cannot use the method which supposes near-integrability of the system. We anticipate that the iterative start would have its largest effect in such a far-integrable dynamical system.

From this viewpoint, we already have a good example of such non-integrable dynamical systems: An extrasolar planet orbiting around a binary star system, MACHO-97-BLG-41 (Bennett *et al.*, 1999; Albrow *et al.*, 2000). The planetary system around MACHO-97-BLG-41 was discovered by a gravitational microlensing event. Another planetary system (single planet + single star) was also found around MACHO-96-BLG-35 (Rhie *et al.*, 2000), which is expected

to be a kind of solar system kin comprising a low-mass terrestrial planet and a solar-type star.

As for the planetary system around MACHO-97-BLG-41, the lens system is expected to consist of a planet of about three Jupiter masses orbiting a binary stellar system comprising a late-K dwarf star and an M dwarf star. The stars are separated by  $1 \sim 2\text{AU}$  (nominally  $\sim 1.6\text{AU}$  in Bennett et al. paper), and the planet is orbiting around them at a distance of about several astronomical unit (nominally  $7\text{AU}$  in Bennett et al. paper). Since binary stars are expected to be much more common in the universe than single stars, it is likely that we find many more of this type of extrasolar planetary systems in the future.

One of the demerits of the extrasolar planet detection by microlensing events is that the accuracy of orbital determination is not so high. This is because it is quite hard and generally impossible for us to re-observe a planetary system which has been found by a microlensing event. We have to determine the orbital elements of planets through a set of observational data covering a very short range. Thus orbital elements of extrasolar planets found by microlensing events should contain large errors. We list the possible range of dynamical parameters of MACHO-97-BLG-41 planetary system which are taken from Bennett et al. (1999) in Table 2.

Distance to the lens	$6.3^{+0.6}_{-1.3} \text{ kpc}$
Total mass of the lens	$0.8 \pm 0.4 M_{\odot}$
Mass of the primary star	$M_1 = 0.6 \pm 0.3 M_{\odot}$
Mass of the secondary star	$M_2 = 0.16 \pm 0.08 M_{\odot}$
Mass of the possible planet	$M_3 = 0.033 \pm 0.017 M_{\odot}$ ( $= 3.5 \pm 1.8 M_J$ )
Separation between two stars	$1.5^{+0.1}_{-0.3} \text{ AU}$
Distance between planet and the center of mass of the lens	$5.7^{+0.6}_{-1.1} \text{ AU}$

**Table 2.** Masses and orbital parameters of MACHO-97-BLG-41 planetary system taken from Bennett et al. (1999).  $M_{\odot}$  is the Sun's mass, and  $M_J$  is the Jupiter's mass.

Among the rather uncertain orbital elements of the planetary system, we have chosen a set which is shown in Figure 16: the mass of the primary star  $M_1 = 0.6 M_{\odot}$ , the mass of the secondary star  $M_2 = 0.16 M_{\odot}$ , the mass of the planet  $M_3 = 0.028 M_{\odot} = 3 M_J$ , the separation between two stars is  $1.6\text{AU}$ , and the semimajor axis of the planet in terms of the barycentric frame of the lens system is  $6\text{AU}$ . Since our present integrations for this system are still preliminary, we have fixed the initial eccentricities  $e$ , longitudes of ascending nodes  $\Omega$ , inclinations  $I$ , arguments of perihelion  $\varpi$ , and mean anomalies  $l$  of the secondary star and of the planet as

$$\begin{aligned}
e_{\text{planet},0} &= 0, & e_{\text{secondary},0} &= 0.1, \\
I_{\text{planet},0} &= I_{\text{secondary},0} = 0.1 \text{ degrees}, \\
\varpi_{\text{planet},0} &= \varpi_{\text{secondary},0} = 0, \\
\Omega_{\text{planet},0} &= \Omega_{\text{secondary},0} = 0, \\
l_{\text{planet},0} &= 0, & l_{\text{secondary},0} &= 180 \text{ degrees}.
\end{aligned}$$



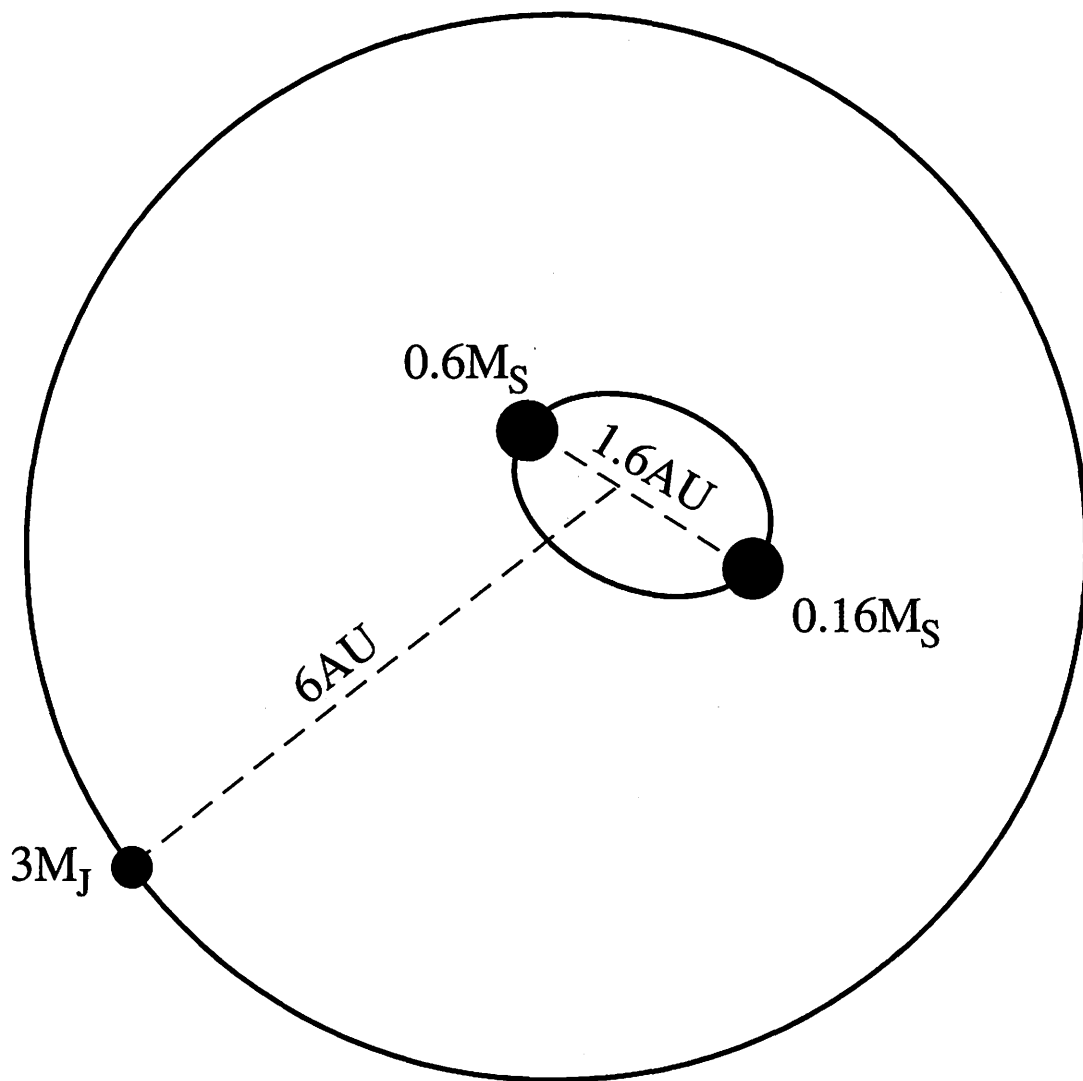
We have selected the non-zero values for inclinations  $I$  so that the system composes a three-dimensional point mass system. By choosing these orbital elements, the ratio of  $\epsilon H_{\text{int}}$  and  $H_{\text{kep}}$  in (141) becomes  $\sim 0.3$ , which is much larger than the perturbed Keplerian motions as described in the previous sections. Thus we think it is worth applying the iterative start when numerically integrating this system.

Several researches have been already done on the dynamical stability of planetary motion in or around binary stars (Wiegert and Holman, 1997; Holman *et al.*, 1997; Mazeh *et al.*, 1997; Holman and Wiegert, 1999). Based on these previous researches, Moriwaki (2001) has investigated on the stability and instability of the MACHO-97-BLG-41 planetary system using a high-order symplectic integrator. Also, Moriwaki and Nakagawa (2002) have performed long-term numerical integrations of the planetary motion with various initial conditions of binary eccentricities, planetary semimajor axis, planetary eccentricity, and longitude of planetary perihelion. The aim of their research was to see which kind of initial configuration produces stable orbits over the whole timespan of their integrations ( $10^6$  binary periods). Their numerical results show that the upper limit of the binary eccentricity in this system is 0.4 when the planet starts from a circular orbit. When the initial eccentricity of the planet becomes large, the stability of the planetary motion is deteriorated. Hence their numerical integration gives us a hint to deduce the upper limit of the initial planetary eccentricity ( $\sim 0.3$  in Moriwaki and Nakagawa's estimate).

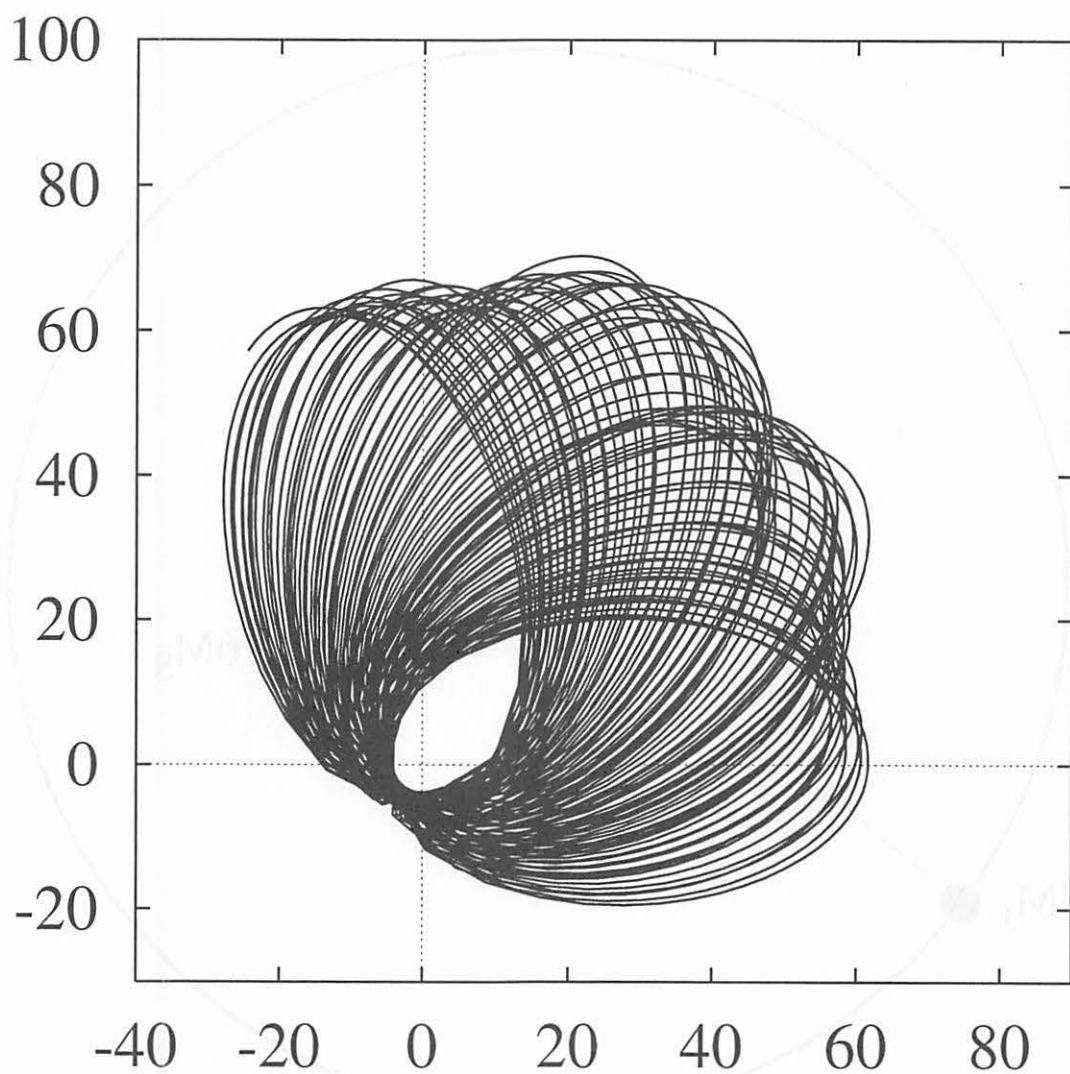
The specific way to implement the iterative start on the planetary system around MACHO-97-BLG-41 is just the same as before:

1. Given a nominal set of initial orbital elements of the planet and the binary stars, we perform an integration with a very high accuracy covering a shorter span than main integrations. We fix the period of this accurate integration as  $2 \times 10^4$  years, similar to the integrations described in the previous section.
2. Choose several initial conditions among the results of the accurate integration for all relevant bodies. We determine the interval of initial conditions for each shorter-term integration as  $P_{\text{planet}}/360$  where  $P_{\text{planet}}$  is the orbital period of the planet around the center of mass of the system. Until the planet gets back to its original longitude, the secondary star rotates around the primary several times.
3. Perform short-term integrations with a normal accuracy using the initial conditions obtained above. We fix the period of this short-term integrations as  $2 \times 10^4$  years. This period should be nearly equal to that of the accurate integration so that we can compare them afterwards.
4. Calculate numerical difference between the accurate and each of the short-term integrations.

In Figure 17 we show the result of the accurate numerical integration of the planetary system around MACHO-97-BLG-41. More specifically speaking, the “accurate” calculation means an numerical integration using a fourth-order generic (TV type) symplectic integrator with a stepsize of 0.0625 days. The planetary orbit rapidly precesses due to the strong perturbation from the binary stars inside. In contrast, we show several examples of the short-term numerical



**Figure 16.** A schematic illustration of the planetary system around MACHO-97-BLG-41 binary. The orbital elements are taken from Table 2.  $M_S$  denotes the Sun's mass ( $= M_\odot$ ).



**Figure 17.** The result of our accurate numerical integration of the planetary system around MACHO-97-BLG-41 binary stars. Eccentricity of the binary stars is 0.1. Only the planetary orbit is drawn here, omitting orbit of the binary stars. The origin of the coordinate is fixed on the center of the mass of the system. The unit of axes is AU.

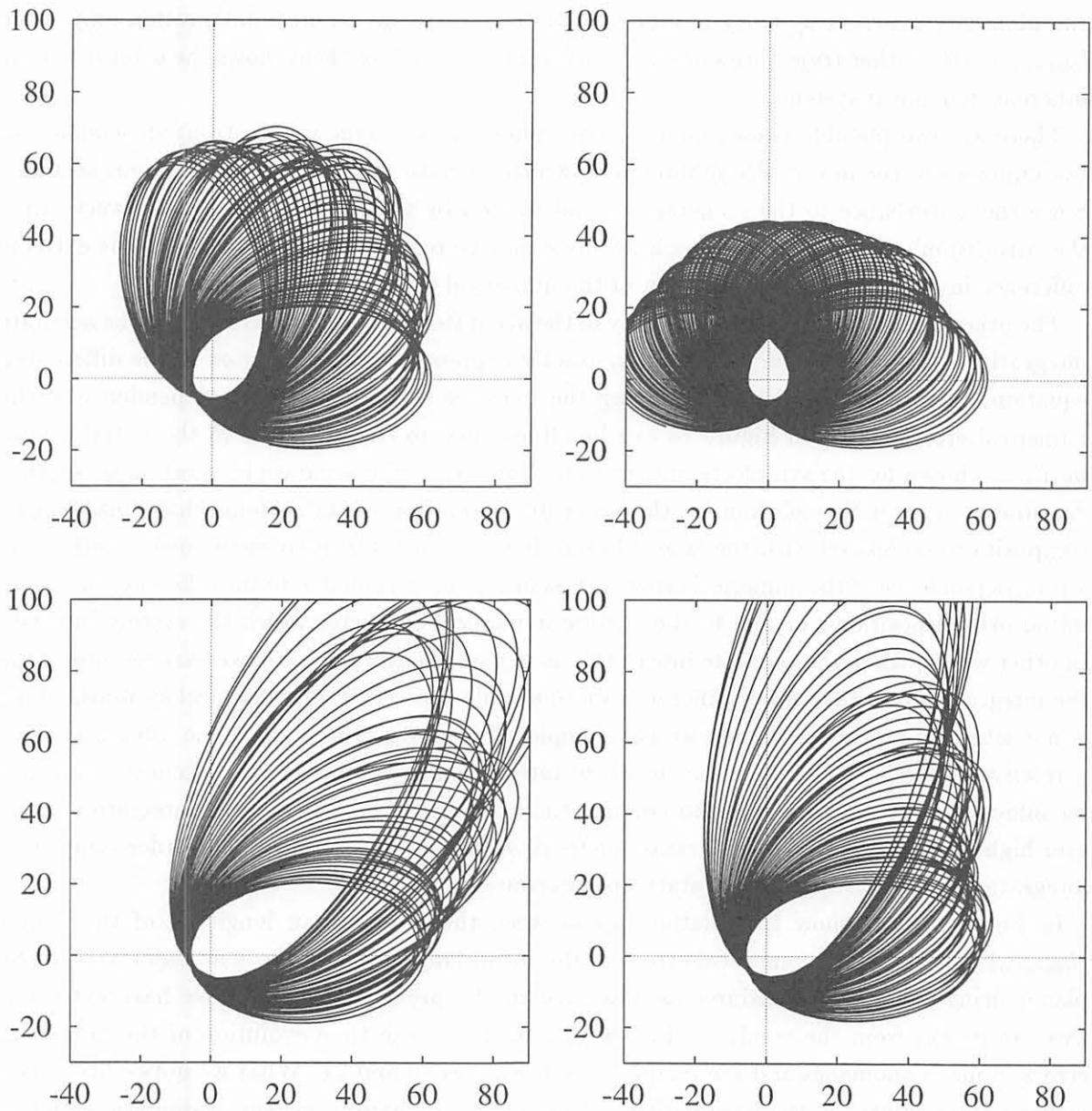
integrations with a moderate accuracy (using the second-order symplectic integrator with a stepsize of 1.0 days) in Figure 18. Although our numerical integrations are still preliminary, we can easily find some very interesting results. Among the short-term numerical integrations, the planetary orbital trajectory is very similar to that of the accurate integration only when  $l_{\text{planet},0} = 16^\circ$ . Other trajectories are quite different, as if each of them shows the orbit in totally different dynamical systems.

There are two possible causes for the above phenomenon. One is the strong dependence of the numerical error in symplectic integrator which we have discussed in the previous sections. Since the disturbance to the planetary orbital motion by the inner binary stars is very large, the “distortion” of the trajectory in phase-space may be remarkable, which leads to the extreme difference in enhancement or reduction of the numerical error shown in Figure 18.

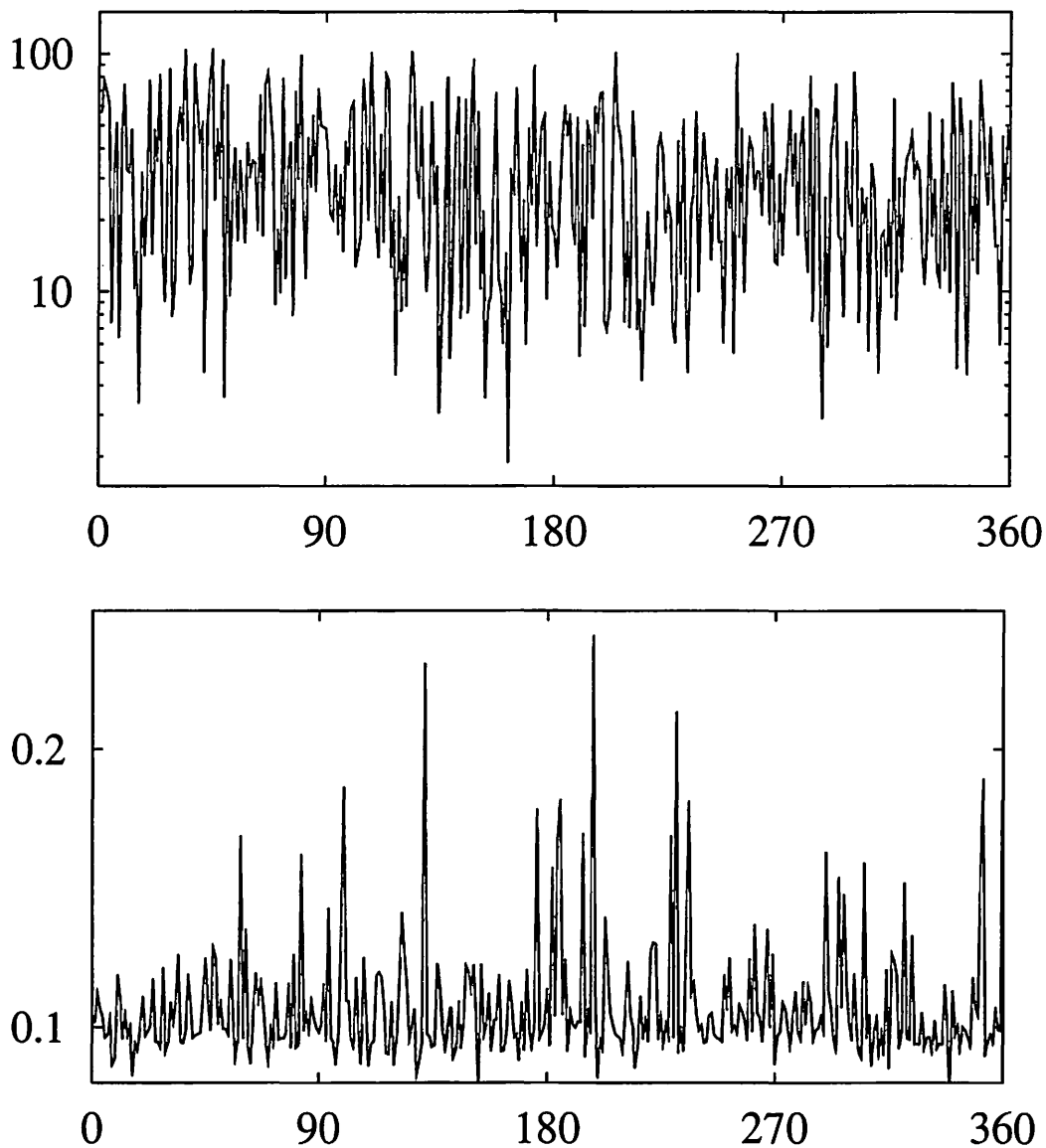
The other reason concerns the reliability of the accurate numerical integration. If the accurate integration is literally “accurate”; that is, exactly expresses the true solution of the differential equation, there is no problem when using the iterative start. The strong dependence of the numerical error such as in Figure 18 can be all ascribed to the difference in the initial orbital positions chosen for the symplectic integration. However, if the accurate integration is not that “accurate”: i.e. if the solution by the accurate integration contains somewhat “inaccurate” compositions compared with the true solution, it is not easy for us to distinguish whether the strong dependence of the numerical error such as in Figure 18 is all due to the difference of chosen initial orbital positions, or due to the chaotic dynamical character which the system involves. In other words, when the accurate integration is not sufficiently accurate, we may be comparing the integration results of many different (and mutually less relevant) dynamical systems, which is not what we meant. Although we have employed a fourth-order symplectic integrator with a relatively smaller stepsize in our accurate integration in Figure 17, the accuracy may not be sufficient. We are now going to confirm the accuracy of our “accurate” integration using ever higher and more precise integration method, such as a sixth- or an eighth-order symplectic integrator with some appropriate start-up procedures.

In Figure 19, we show the relationship between the initial mean longitude of the planet ( $l_{\text{planet},0}$ ) and the RMS numerical errors of the mean longitude and the semimajor axis of the planet using the same procedures as described in the previous sections. We have extracted three examples from the results in Figure 19 and drawn the time evolution of the numerical errors in mean anomalies and semimajor axes in Figures 20 and 21. What we notice first when viewing these figures is the irregularities of lines; we can hardly see any systematic trend in the graphs, unlike when we have taken care of a slightly perturbed three-body system in the previous section. This is we think probably due to the larger perturbation to the planetary orbit by the inner binary stars. We have to check whether this irregular trend is real or not using more accurate numerical integrations.

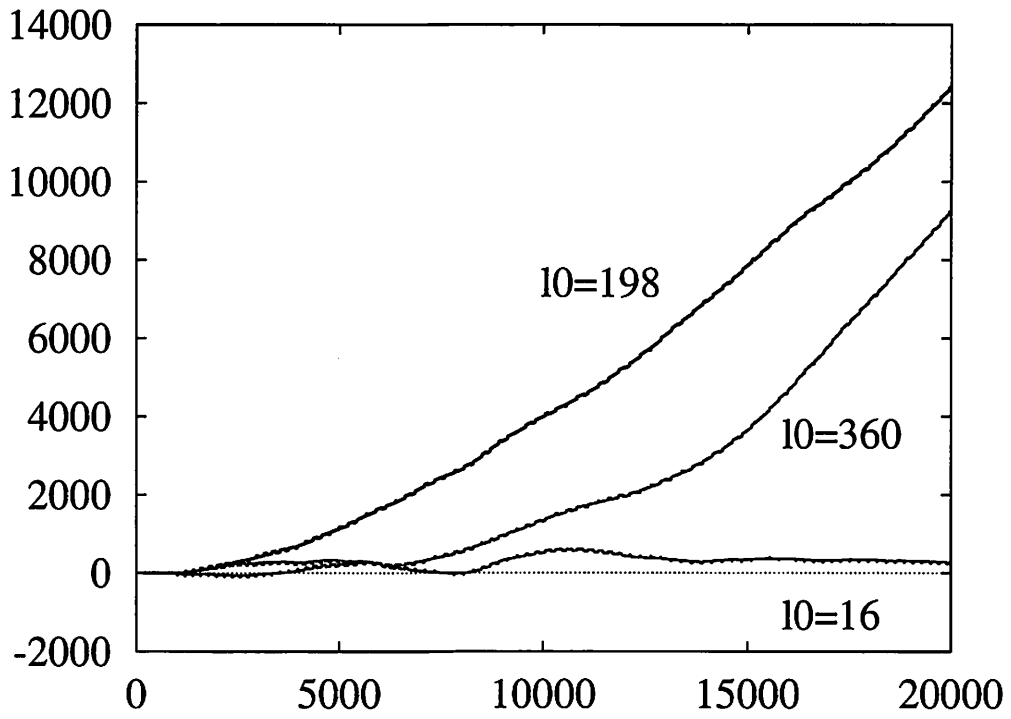
Another difficulty when we adopt the iterative start on a strongly perturbed dynamical system such as the planetary system around MACHO-97-BLG-41 binary stars is related to its non-integrability. A weakly perturbed system such as what we have described in the previous sections is also a non-integrable system, but we call it “nearly integrable.” In nearly-integrable systems, especially those which are close to the superposition of the two-body system, we know that many of angle variables degenerate. Hence there is a considerable difference in the timescale of



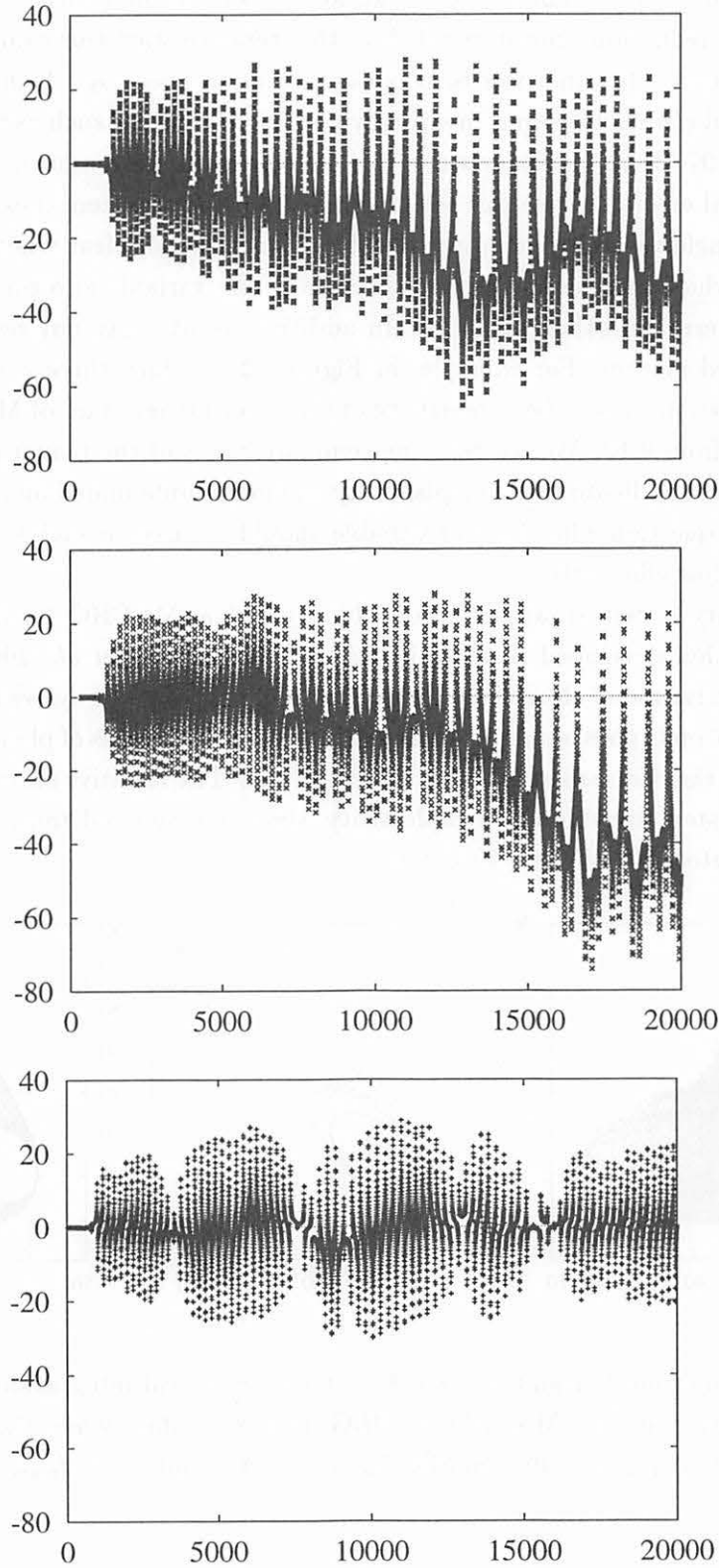
**Figure 18.** Some example results of the short-term numerical integrations ( $2 \times 10^4$  years) of the planetary motion around MACHO-97-BLG-41 binary stars when the eccentricity of the binary stars is 0.1. Upper left:  $l_{\text{planet},0} = 16^\circ$ , upper right:  $l_{\text{planet},0} = 34^\circ$ , lower left:  $l_{\text{planet},0} = 198^\circ$ , and lower right:  $l_{\text{planet},0} = 360^\circ$ . Only the planetary orbit is drawn here, omitting the orbit of the binary stars. The origin of the coordinate is fixed on the center of mass of the system. The unit of axes is AU.



**Figure 19.** The relationship between the initial mean longitude of the planet ( $l_{\text{planet},0}$ ; degree) and the RMS numerical errors of the mean longitude (upper; degree) and semimajor axis (lower; AU) of the planetary orbit in the MACHO-97-BLG-41 system. Each length of numerical data used for the comparison between the accurate and the short-term integrations is  $2 \times 10^4$  years. The origin of the orbital elements is the center of mass of the system.



**Figure 20.** The numerical difference of the planetary mean anomaly in the MACHO-97-BLG-41 system when  $l_{\text{planet},0} = 198^\circ$ ,  $l_{\text{planet},0} = 360^\circ$ , and  $l_{\text{planet},0} = 16^\circ$ . The unit of the vertical axis is degree, and the unit of the horizontal axis is year. The origin of the orbital elements is the center of mass of the system.

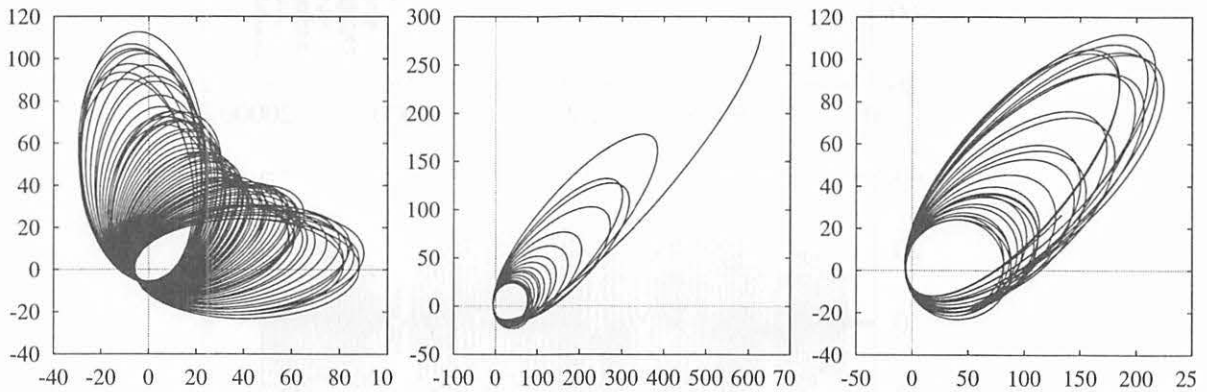


**Figure 21.** The numerical difference of the planetary semimajor axis in the MACHO-97-BLG-41 system: (upper) when  $l_{\text{planet},0} = 198^\circ$ , (middle) when  $l_{\text{planet},0} = 360^\circ$ , (lower) when  $l_{\text{planet},0} = 16^\circ$ . The unit of the vertical axis is AU, and the unit of the horizontal axis is year. The origin of the orbital elements is the center of mass of the system.



each variable in the system. Then it is enough for us to focus on a variable having the shortest timescale in error reduction procedures such as the iterative start (for example, mean anomaly in the Kepler motion). In other words, it is easy for us to guess by which variables we should evaluate numerical errors. However, in strongly perturbed systems such as the planetary system around MACHO-97-BLG-41 binary stars, it is not so easy to determine a variable by which we evaluate numerical errors. As we can see in Figure 18, such a system dose not degenerate any longer, with all angle variables changing quickly. Then, it is not clear whether or not an initial orbital position which reduces the numerical error of one variable also produces the minimum of the numerical error of other variables. In addition, orbits may not be bounded in such a strongly perturbed system: For example, in Figure 22 we show three example results of the “accurate” integrations when the eccentricity of the inner binary stars of MACHO-97-BLG-41 increases to 0.15 from 0.10. We see that the semimajor axis of the planet changes rapidly and secularly, which may indicate that the planetary motion is unbounded and unstable. We think it is still an open question what kind of variable should be used in such systems to reduce the numerical error most efficiently.

As for a planetary system around or inside a binary such as MACHO-97-BLG-41, a dedicated symplectic algorithm developed recently is available (Chambers *et al.*, 2002; Quintana *et al.*, 2002). However, Chambers’ algorithm still exploits the fact that the system is nearly integrable and bounded. If planets goes out or into binary stars, or if the actions of planets change secularly, it is not sure that the dedicated algorithm works as well. The iterative start works possibly well even in such a system keeping the symplecticity, since it is derived from the general-purpose symplectic integrator,  $H = T(p) + V(q)$  type.



**Figure 22.** Three example results of the short-term numerical integrations ( $2 \times 10^4$  years) of the planetary motion around MACHO-97-BLG-41 binary stars when the eccentricity of the binary is 0.15. Left:  $l_{\text{planet},0} = 40^\circ$ , middle:  $l_{\text{planet},0} = 96^\circ$ , and right:  $l_{\text{planet},0} = 221^\circ$ . The unit of axes is AU.

### 5.3 Relationship to the “warm start”

It is meaningful to consider whether there is any relationship between the items in this manuscript and a special start-up procedure called “warm start” (Saha and Tremaine, 1992;

Saha and Tremaine, 1994).

The discussion below follows Saha & Tremaine (1992). As we have seen in the previous sections, Hamiltonian  $\tilde{H}$  for a surrogate system dominating the WH-type symplectic map can be written as follows.

$$\tilde{H} = H + H_{\text{err}}, \quad (140)$$

$$H = H_{\text{kep}} + \epsilon H_{\text{int}}, \quad (141)$$

$$H_{\text{err}} = \frac{\epsilon \tau^2}{24} \{ \{H_{\text{kep}}, H_{\text{int}}\}, H_{\text{kep}} \} + O(\epsilon^2 \tau^2). \quad (142)$$

We define the actions and the angles in a real system as  $\mathbf{J} = (J_1, J_2, J_3)$  and  $\boldsymbol{\theta} = (\theta_1, \theta_2, \theta_3)$ . Also, the Kepler Hamiltonian  $H_{\text{kep}}$  is described by Delaunay variables as usual  $\mathbf{L} = (L, G, H)$  and  $\mathbf{l} = (l, g, h)$ . Although the Delaunay variables are dedicated to the Kepler motion, the discussion here can be extended to any kind of nearly integrable problems.

We know that the difference between  $H$  and  $H_{\text{kep}}$  is only  $O(\epsilon)$  where  $\epsilon$  is the order of magnitude of the perturbation we consider now. Then we can write as

$$\begin{aligned} J_1 &= L + O(\epsilon), & \theta_1 &= l + O(\epsilon), \\ J_2 &= G + O(\epsilon), & \theta_2 &= g + O(\epsilon), \\ J_3 &= H + O(\epsilon), & \theta_3 &= h + O(\epsilon). \end{aligned} \quad (143)$$

Since  $H_{\text{kep}}$  is a function of  $L$  only (i.e.  $H_{\text{kep}} = -\sum_i \frac{\mu_i^2}{2L_i^2}$ ), all the dependencies of  $H_{\text{kep}}$  on  $J_2, J_3, \theta_1, \theta_2, \theta_3$  are restricted to  $O(\epsilon)$ . Then we get

$$\frac{\partial H_{\text{kep}}}{\partial J_2} = O(\epsilon), \quad \frac{\partial H_{\text{kep}}}{\partial J_3} = O(\epsilon), \quad (144)$$

$$\frac{\partial H_{\text{kep}}}{\partial \theta_1} = O(\epsilon), \quad \frac{\partial H_{\text{kep}}}{\partial \theta_2} = O(\epsilon), \quad \frac{\partial H_{\text{kep}}}{\partial \theta_3} = O(\epsilon). \quad (145)$$

Substituting (144) and (145) into (142) and evaluating the Poisson bracket  $\{, \}$  by canonical variables of the real system  $(\mathbf{J}, \boldsymbol{\theta})$ , we can rewrite the error Hamiltonian  $H_{\text{err}}$ . Then we immediately find that only a term including  $\frac{\partial H_{\text{kep}}}{\partial J_1}$  remains up to  $O(\epsilon)$  approximation since the error Hamiltonian  $H_{\text{err}}$  originally has a factor  $\epsilon$  (142).

$$H_{\text{err}} = -\frac{\epsilon \tau^2}{24} \left( \frac{\partial H_{\text{kep}}}{\partial J_1} \right)^2 \left( \frac{\partial^2 H_{\text{int}}}{\partial \theta_1^2} \right) + O(\epsilon^2 \tau^2). \quad (146)$$

Here we are supposed to concern only bounded motions such as a stable planetary dynamics or a rotational motion of planet. Then,  $\frac{\partial^2 H_{\text{int}}}{\partial \theta_1^2}$  in (146) consists of only periodic terms: if there is a constant term in  $\frac{\partial^2 H_{\text{int}}}{\partial \theta_1^2}$ , e.g.  $C_1$  as

$$\frac{\partial^2 H_{\text{int}}}{\partial \theta_1^2} = C_1, \quad (147)$$

which leads to

$$\frac{\partial H_{\text{int}}}{\partial \theta_1} = C_1 \theta_1 + C_2, \quad (148)$$

where  $C_2$  is another constant of integration. Let us define  $J_{1,H_{\text{int}}}$  as an action due to  $H_{\text{int}}$ . Then, from (148),

$$\frac{dJ_{1,H_{\text{int}}}}{dt} = -\frac{\partial H_{\text{int}}}{\partial \theta_1} = -C_1\theta_1 - C_2, \quad (149)$$

$$\therefore J_{1,H_{\text{int}}} = -C_2t + f(\theta_1; t), \quad (150)$$

where  $f(\theta_1; t)$  is a function of  $\theta_1$  and  $t$ . Equation (150) means that there appears a secular motion in the action  $J_{1,H_{\text{int}}}$  if there is a constant term in  $\frac{\partial^2 H_{\text{int}}}{\partial \theta_1^2}$  in (146), which leads to a possible secular collapse of the system. Thus  $\frac{\partial^2 H_{\text{int}}}{\partial \theta_1^2}$  in (146) must consist of only periodic terms. The fact that  $\frac{\partial^2 H_{\text{int}}}{\partial \theta_1^2}$  is expressed only by periodic terms indicates that there are no “raw” angles  $\theta$  in the disturbing function which describes planetary perturbation. All of the angles appear as the form of periodic functions such as  $\sin \theta$  or  $\cos \theta$ .

From (146), we can decompose the error Hamiltonian  $H_{\text{err}}$  into the superposition of Fourier components as

$$H_{\text{err}} = \epsilon\tau^2 \sum_{\mathbf{m}} X_{\mathbf{m}}(\mathbf{J}) e^{i\mathbf{m} \cdot \boldsymbol{\theta}} + O(\epsilon^2\tau^2), \quad (151)$$

where  $\mathbf{m} = (m_1, m_2, m_3)$  is a integer vector and each  $m_j$  takes any value from  $-\infty$  to  $+\infty$ , and  $X_{\mathbf{m}}$  is the coefficients of Fourier transformation. Since  $H_{\text{err}}$  has no secular term other than  $O(\epsilon^2\tau^2)$ , and since only  $\theta_1$  is directly related to time,  $X_{\mathbf{m}}$  should be zero when  $m_1 = 0$ . Hence the time-averaged (i.e. secular) value of  $H_{\text{err}}$  should be

$$\langle H_{\text{err}} \rangle = O(\epsilon^2\tau^2). \quad (152)$$

The fact described by (152) plays an essential role in the principle of the warm start.

If we define the canonical frequency of the real system  $H$  as

$$\omega \equiv \frac{\partial H}{\partial \mathbf{J}}, \quad (153)$$

the formal solution for the action would be

$$\tilde{\mathbf{J}} = \mathbf{J} - \epsilon\tau^2 \sum_{\mathbf{m}} \frac{X_{\mathbf{m}}\mathbf{m}}{\omega \cdot \mathbf{m}} e^{i\mathbf{m} \cdot \boldsymbol{\theta}}, \quad (154)$$

and the canonical frequency of the surrogate system  $\tilde{H}$  becomes

$$\begin{aligned} \tilde{\omega}(\tilde{\mathbf{J}}) \equiv \frac{\partial \tilde{H}(\tilde{\mathbf{J}})}{\partial \tilde{\mathbf{J}}} &= \frac{\partial}{\partial \tilde{\mathbf{J}}} \left( H(\tilde{\mathbf{J}}) + \langle H_{\text{err}}(\tilde{\mathbf{J}}) \rangle \right) \\ &= \frac{\partial H(\tilde{\mathbf{J}})}{\partial \tilde{\mathbf{J}}} + O(\epsilon^2\tau^2) \\ &= \omega(\tilde{\mathbf{J}}) + O(\epsilon^2\tau^2). \end{aligned} \quad (155)$$

The fundamental idea of the warm start described in Saha & Tremaine (1992) is to make the difference between  $\tilde{\mathbf{J}}$  and  $\mathbf{J}$  as small as possible at the start of symplectic integration, and keep the difference small all through the integration using the adiabatic invariant character of Hamiltonian systems.

Thus we can understand that the warm start is effective only when the error Hamiltonian consists only of periodic terms as in (151), and when its averaged value becomes far smaller

(as in (152)). In this case, the warm start is more effective and easier to implement than the iterative start. However, when the error Hamiltonian does not consist only of periodic terms such as in (48), the warm start is no longer valid. We have to resort to other general method, such as the iterative start.

## 6. Interpretation and other analytical examples

As discussed above, we can reduce the numerical error of  $H = T(p) + V(q)$  type symplectic integrator by choosing certain starting conditions. But mostly, the numerical error does not decrease significantly (cf. when  $e_0 = 0.1$  in Figure 10). What is the difference between systems where we can and cannot reduce the numerical error of the symplectic integrator?

To answer this question, we have applied a canonical perturbation theory to several simple dynamical problems. We consider a symplectic integrator dominated by a surrogate Hamiltonian  $\tilde{H} = H + H_{\text{err}}$  as a kind of nearly-integrable, disturbed Hamiltonian problem. If the Hamiltonian of the real system  $H$  is integrable (e.g. the Kepler motion or the harmonic oscillator) and  $H_{\text{err}}$  is sufficiently small, we can obtain approximate solution of the system, i.e. numerical solution by symplectic integrator including numerical errors.

Suppose the surrogate Hamiltonian takes the following form:

$$\tilde{H}(\theta, J) = H(J) + H_{\text{err}}(\theta, J), \quad (156)$$

where  $\theta$  and  $J$  are the angle and action variables of the real system,  $H$ .  $J$  is an integral (or constant) in  $H$ , hence  $H$  includes no angle,  $\theta$ .  $J$  would be no longer a constant (nor action) in the surrogate system  $\tilde{H}$ , but  $J$  and  $\theta$  still serve as canonical variables in the surrogate system as  $\tilde{H}(\theta, J)$ .

Here we suppose that the degree of freedom of the system is one for simplicity. Extension to problems with larger degrees of freedom is in principle possible, though it generally requires formidable algebra. Since we focus ourselves on a first-order solution here, using a traditional theory (von Zeipel, 1916) or a standard theory exploiting Lie series (Hori, 1966; Deprit, 1969) makes no difference.

Averaging  $H_{\text{err}}$  by the angle  $\theta$ , we obtain the new Hamiltonian  $\tilde{H}^*$  in general as

$$\tilde{H}^*(J^*) = H^*(J^*) + \langle H_{\text{err}}(\theta^*, J^*) \rangle \quad (157)$$

$$= H^*(J^*) + H_{\text{err}}^*(J^*). \quad (158)$$

Then we get the canonical equations of motion for the new system as

$$\frac{dJ^*}{dt} = -\frac{\partial \tilde{H}^*(J^*)}{\partial \theta^*}, \quad (159)$$

$$\begin{aligned} \frac{d\theta^*}{dt} &= \frac{\partial \tilde{H}^*}{\partial J^*} \\ &= \frac{\partial H^*(J^*)}{\partial J^*} + \frac{\partial H_{\text{err}}^*(J^*)}{\partial J^*}, \end{aligned} \quad (160)$$

but since  $\tilde{H}^*$  includes only  $J^*$  and does not include the angle variable  $\theta^*$ ,  $J^*$  is a constant (or an integral) and  $\theta^*$  exhibits a equi-velocity linear motion in phase-space as

$$J^* = \text{constant}, \quad \theta^* = \omega t + \text{constant}, \quad (161)$$

where  $\omega$  is the canonical frequency of the new system defined as

$$\omega(J^*) \equiv \frac{\partial \tilde{H}^*(J^*)}{\partial J^*}. \quad (162)$$

In the following subsections, we take several simple Hamiltonian systems as examples, and see how the solution (i.e. the numerical solution by the symplectic integrator) would be by the canonical perturbation theory.

## 6.1 Harmonic oscillator

Let us begin with the simplest system, the harmonic oscillator with one degree of freedom. The Hamiltonian of the system is

$$H = T(p) + V(q) = \frac{p^2}{2m} + \frac{m\omega_0^2}{2}q^2, \quad (163)$$

where  $q$  and  $p$  are canonical coordinate and momentum.  $\omega_0$  is a constant which describes the oscillation frequency of the system, and  $m$  corresponds to the mass of the oscillator. The angle/action variables  $(\theta, J)$  of the harmonic oscillator are obtained, for example, through the Poincaré transformation as

$$q = \sqrt{\frac{2J}{m\omega_0}} \sin \theta, \quad p = \sqrt{2m\omega_0 J} \cos \theta. \quad (164)$$

The Hamiltonian (163) becomes now

$$H = \omega_0 J. \quad (165)$$

Thus we found the system integrable.

### 6.1.1 Numerical errors of the first-order formula

According to the BCH formula, the error Hamiltonian of the first-order symplectic integrator becomes as follows:

$$H_{\text{err}} = \frac{\tau}{2} \{V, T\} + \frac{\tau^2}{12} (\{ \{T, V\}, V \} + \{ \{V, T\}, T \}) + O(\tau^3), \quad (166)$$

where  $\tau$  is the stepsize of integration. Each Poisson bracket becomes

$$\begin{aligned} \{V, T\} &= \frac{\partial T}{\partial q} \frac{\partial V}{\partial p} - \frac{\partial T}{\partial p} \frac{\partial V}{\partial q} \\ &= -\omega_0^2 qp \\ &= -\omega_0^2 J \sin 2\theta, \end{aligned} \quad (167)$$

$$\{T, V\} = \omega_0^2 qp = \omega_0^2 J \sin 2\theta, \quad (168)$$

$$\begin{aligned} \{\{T, V\}, V\} &= \left\{ \omega_0^2 qp, \frac{m\omega_0^2}{2} q^2 \right\} \\ &= -m\omega_0^4 q^2 \\ &= -2\omega_0^3 J \sin^2 \theta, \end{aligned} \quad (169)$$

$$\begin{aligned} \{\{V, T\}, T\} &= \left\{ -\omega_0^2 qp, \frac{p^2}{2m} \right\} \\ &= -\frac{\omega_0^2 p^2}{m} \\ &= -2\omega_0^3 J \cos^2 \theta. \end{aligned} \quad (170)$$

Hence

$$\begin{aligned} H_{\text{err}} &= \frac{\tau}{2} \{V, T\} + \frac{\tau^2}{12} (\{\{T, V\}, V\} + \{\{V, T\}, T\}) + O(\tau^3) \\ &= -\frac{\tau}{2} \omega_0^2 J \sin 2\theta + \frac{\tau^2}{12} (-2\omega_0^3 J \sin^2 \theta - 2\omega_0^3 J \cos^2 \theta) + O(\tau^3) \\ &= -\frac{\tau}{2} \omega_0^2 J \sin 2\theta - \frac{\tau^2}{12} \omega_0^3 J + O(\tau^3). \end{aligned} \quad (171)$$

Hereafter we write the Hamiltonian  $\tilde{H} = H + H_{\text{err}}$  for simplicity as

$$H(\theta, J) = H_0(\theta, J) + H_1(\theta, J), \quad (172)$$

according to the way in the previous section (see (77)).

Then we get

$$H(\theta, J) = \omega_0 J - \frac{\tau}{2} \omega_0^2 J \sin 2\theta - \frac{\tau^2}{12} \omega_0^3 J + O(\tau^3). \quad (173)$$

We transform the original Hamiltonian  $H(\theta, J)$  into an integrable form  $H^*(J^*)$  through a canonical transformation. If we assume that the new Hamiltonian  $H^*$  would be expanded into the form

$$H^*(J^*) = H_0^*(J^*) + H_1^*(J^*) + H_2^*(J^*) + \cdots, \quad (174)$$

and the generating function  $S$  of the transformation could be the form

$$S(\theta^*, J^*) = S_0(\theta^*, J^*) + S_1(\theta^*, J^*) + S_2(\theta^*, J^*) + \cdots, \quad (175)$$

where the order of magnitude of  $H_{i+1}/H_i$  and  $S_{i+1}/S_i$  ( $i = 0, 1, \dots$ ) equals to that of the perturbation, i.e.  $O(\tau)$ .

Omitting all the algebra on the way, we would obtain

$$H_0^*(J^*) = H_0(J^*) = \omega_0 J^*, \quad (176)$$

hence the zeroth-order Hamiltonian is uniquely determined. As for the first-order Hamiltonian,

$$H_1^*(J^*) = \{H_0(\theta^*, J^*), S_0(\theta^*, J^*)\} + H_1(\theta^*, J^*). \quad (177)$$

Here we have to determine both  $H_1^*(J^*)$  and  $S_0(\theta^*, J^*)$  simultaneously. Thus we request  $H_1^*(J^*)$  not to contain any angle variables as

$$H_1^*(J^*) = \langle H_1(\theta^*, J^*) \rangle_{\theta^*}, \quad (178)$$

and let  $S_0$  include the rest of periodic terms as

$$\{H_0(\theta^*, J^*), S_1(\theta^*, J^*)\} = \langle H_1(\theta^*, J^*) \rangle_{\theta^*} - H_1(\theta^*, J^*). \quad (179)$$

Now let us remember the canonical equations of motion for unperturbed part of the Hamiltonian in the transformed system

$$\frac{d\theta^*}{dt^*} = \frac{\partial H_0}{\partial J^*}, \quad \frac{dJ^*}{dt^*} = -\frac{\partial H_0}{\partial \theta^*}, \quad (180)$$

where  $t^*$  is the “time” for this system. Using (180), the Poisson bracket in (177) becomes

$$\begin{aligned} \{H_0(\theta^*, J^*), S_0(\theta^*, J^*)\}_{(\theta^*, J^*)} &= \frac{\partial H_0}{\partial \theta^*} \frac{\partial S_0}{\partial J^*} - \frac{\partial H_0}{\partial J^*} \frac{\partial S_0}{\partial \theta^*} \\ &= -\left( \frac{\partial S_0}{\partial J^*} \frac{dJ^*}{dt^*} + \frac{\partial S_0}{\partial \theta^*} \frac{d\theta^*}{dt^*} \right) \\ &= -\frac{dS_0}{dt^*}. \end{aligned} \quad (181)$$

Thus (179) becomes

$$\frac{dS_0}{dt^*} = -\{H_0(\theta^*, J^*), S_0(\theta^*, J^*)\} = H_1(\theta^*, J^*) - \langle H_1(\theta^*, J^*) \rangle_{\theta^*}, \quad (182)$$

$$\therefore S_0 = \int (H_1 - H_1^*). \quad (183)$$

Substituting the specific form of  $H_1$  (173) into (178), we obtain

$$\begin{aligned} H_1^*(J^*) &= \left\langle -\frac{\tau}{2} \omega_0^2 J^* \sin 2\theta^* - \frac{\tau^2}{12} \omega_0^3 J^* \right\rangle_{\theta^*} + O(\tau^3) \\ &= -\frac{\tau^2}{12} \omega_0^3 J^* + O(\tau^3). \end{aligned} \quad (184)$$

The final Hamiltonian up to the first-order approximation is thus

$$H^*(J^*) = H_0^*(J^*) + H_1^*(J^*) = \omega_0 J^* - \frac{\tau^2}{12} \omega_0^3 J^* + O(\tau^3). \quad (185)$$

The canonical equations of motion for  $(\theta^*, J^*)$  system are

$$\frac{d\theta^*}{dt} = \frac{\partial H^*(J^*)}{\partial J^*} = \omega_0 - \frac{\tau^2}{12} \omega_0^3, \quad (186)$$

and

$$\frac{dJ^*}{dt} = -\frac{\partial H^*(J^*)}{\partial \theta^*} = 0. \quad (187)$$

Thus  $J^*$  is a constant, and the canonical frequency of the new system  $\omega(J^*)$  becomes

$$\omega(J^*) \equiv \frac{\partial H^*(J^*)}{\partial J^*} = \dot{\theta}^* = \omega_0 - \frac{\tau^2}{12} \omega_0^3, \quad (188)$$

which is also a constant. Since the relationship between the new and old variable has the form

$$\theta = \theta^* + \text{periodic terms}, \quad J = J^* + \text{periodic terms}, \quad (189)$$

the secular numerical error in the angle  $\theta$  is produced from the second term in the right-hand side of (188),  $-\frac{\tau^2}{12}\omega_0^3$ , which is a constant. In this sense, we cannot manage to reduce the secular numerical error of the angle  $\theta$  when we solve the motion of the harmonic oscillator by the first-order symplectic integrator; the error does not depend on initial starting conditions.

### 6.1.2 Numerical errors of the second-order formula

The error Hamiltonian of the second-order symplectic integrator can be obtained by the BCH formula up to  $O(\tau^2)$  as

$$H_{\text{err}} = \frac{\tau^2}{12} \left( \{ \{T, V\}, V \} - \frac{1}{2} \{ \{V, T\}, T \} \right). \quad (190)$$

We have already calculated the Poisson brackets  $\{ \{T, V\}, V \}$  and  $\{ \{V, T\}, T \}$ , so

$$\begin{aligned} \frac{12}{\tau^2} H_{\text{err}} &= \{ \{T, V\}, V \} - \frac{1}{2} \{ \{V, T\}, T \} \\ &= m\omega_0^4 q^2 - \frac{\omega_0^2 p^2}{m} \\ &= m\omega_0^4 \left( \frac{2J}{m\omega_0} \right) \sin^2 \theta - \frac{\omega_0^2}{2m} (2m\omega_0 J) \cos^2 \theta \\ &= 2\omega_0^3 J \sin^2 \theta - \omega_0^3 J \cos^2 \theta \\ &= \frac{\omega_0^3 J}{2} (1 - 3 \cos 2\theta). \end{aligned} \quad (191)$$

Averaging this error Hamiltonian, we obtain

$$H_1^*(J^*) = \langle H_{\text{err}}(\theta^*, J^*) \rangle = \frac{\tau^2 \omega_0^3 J^*}{24}, \quad (192)$$

which leads us to the final Hamiltonian as

$$H^*(J^*) = H_0^*(J^*) + H_1^*(J^*) = \omega_0 J^* + \frac{\tau^2}{24} \omega_0^3 J^* + O(\tau^4). \quad (193)$$

The canonical equations of motion for  $(\theta^*, J^*)$  system are

$$\frac{d\theta^*}{dt} = \frac{\partial H^*(J^*)}{\partial J^*} = \omega_0 + \frac{\tau^2}{24} \omega_0^3 + O(\tau^4), \quad (194)$$

and

$$\frac{dJ^*}{dt} = -\frac{\partial H^*(J^*)}{\partial \theta^*} = 0. \quad (195)$$

Thus  $J^*$  is again a constant, and the canonical frequency of the new system becomes

$$\omega(J^*) \equiv \frac{\partial H^*(J^*)}{\partial J^*} = \dot{\theta}^* = \omega_0 + \frac{\tau^2}{24} \omega_0^3 \quad (196)$$

which is also a constant. Since the relationship between the new and old variable has the form

$$\theta = \theta^* + \text{periodic terms}, \quad J = J^* + \text{periodic terms}, \quad (197)$$

the secular numerical error in the angle  $\theta$  is produced from the second term in the right-hand side of (196),  $\frac{\tau^2}{24}\omega_0^3$ , which is a constant. Thus we know that we cannot reduce the secular numerical error of the angle  $\theta$  arising from the second-order symplectic integrator.



## 6.2 Nonlinear pendulum

As seen before, the harmonic oscillator produces no secular numerical error in the first- and second-order symplectic integrator. We guess this is ascribed to the isochrone character of the potential of the harmonic oscillator—eigenfrequency of the system does not depend on the initial amplitude of oscillation. The eigenfrequency of the harmonic oscillator does not contain action variables;  $\omega_0$  is a pure constant, not a function of action variable such as  $\omega_0(J)$ . The isochrone characteristic is typical of the harmonic oscillator. In contrast, the Keplerian potential,  $-\mu/r$ , is not isochrone. The eigenfrequency of the Keplerian motion (i.e. mean motion  $n$ ) depends on the initial amplitude of oscillation (i.e. semimajor axis  $a$ ) as  $n = \sqrt{\mu/a^3}$ . We already knew that the Keplerian motion produces the secular numerical error in angle variables when we use a first- and second-order symplectic integrator. In this subsection, we check whether a nonlinear and non-isochrone pendulum causes any secular numerical error in symplectic integrators.

### 6.2.1 Angle and action variables

We begin with the following Hamiltonian

$$H(q, p) = \frac{p^2}{2} + \frac{\omega_0^2 q^2}{2} - \frac{\epsilon q^4}{4!} + \frac{\epsilon^2 q^6}{6!} + \dots \quad (198)$$

which is originally derived from the Hamiltonian of a pendulum having the form of  $H = \frac{p^2}{2} + b \cos q$ , but with a slight modification (Boccaletti and Pucacco, 1998). We recognize  $\epsilon$  as a small and constant parameter. The unperturbed part of the Hamiltonian (198) is

$$H_0(q, p) = \frac{p^2}{2} + \frac{\omega_0^2 q^2}{2}, \quad (199)$$

which is identical to the Hamiltonian of the harmonic oscillator. Since we already knew the relationship between  $(q, p)$  and the action-angle variables  $(\theta, J)$  of the harmonic oscillator as (164), we can rewrite the Hamiltonian (199) as

$$H_0(J) = \omega_0 J, \quad (200)$$

which leads to the new form of the whole Hamiltonian as

$$H(\theta, J) = \omega_0 J - \frac{\epsilon J^2}{16} + \frac{\epsilon^2 J^3}{256\omega_0} + \dots \quad (201)$$

Now we start to obtain the action ( $J^*$ ) for the Hamiltonian system (201) which allows us to write

$$\begin{aligned} H(\theta, J) &= H^*(J^*) \\ &= H_0^*(J^*) + \epsilon H_1^*(J^*) + \epsilon^2 H_2^*(J^*) + \dots, \end{aligned} \quad (202)$$

through the canonical transformation of Hori (1966). Note that we implicitly assume that the Hamiltonian  $H(\theta, J)$  nor  $H^*(J^*)$  does not contain time explicitly,

Let the generating function  $S$  of the transformation  $(\theta, J) \rightarrow (\theta^*, J^*)$  as

$$S(\theta^*, J^*) = S_1(\theta^*, J^*) + \epsilon S_2(\theta^*, J^*) + \epsilon^2 S_3(\theta^*, J^*) + \dots \quad (203)$$

We apply the Lie's expansion theorem to  $H(\theta, J)$ , and get

$$H(\theta, J) = H(\theta^*, J^*) + \epsilon \{H(\theta^*, J^*), S(\theta^*, J^*)\} + \frac{\epsilon^2}{2} \{\{H(\theta^*, J^*), S(\theta^*, J^*)\}, S(\theta^*, J^*)\} + \dots \quad (204)$$

Substituting (201) and (203) into (204), we get

$$\begin{aligned} H(\theta, J) &= H_0(\theta^*, J^*) + \epsilon H_1(\theta^*, J^*) + \epsilon^2 H_2(\theta^*, J^*) + \dots \\ &\quad + \epsilon \{H_0 + \epsilon H_1 + \epsilon^2 H_2 + \dots, S_1 + \epsilon S_2 + \epsilon^2 S_3 + \dots\} \\ &\quad + \frac{\epsilon^2}{2} \{\{H_0 + \epsilon H_1 + \epsilon^2 H_2 + \dots, S_1 + \epsilon S_2 + \epsilon^2 S_3 + \dots\}, S_1 + \epsilon S_2 + \epsilon^2 S_3 + \dots\} + \dots \\ &= H_0(\theta^*, J^*) + \epsilon H_1(\theta^*, J^*) + \epsilon^2 H_2(\theta^*, J^*) + \dots \\ &\quad + \epsilon \{H_0(\theta^*, J^*), S_1(\theta^*, J^*)\} + \epsilon^2 (\{H_0(\theta^*, J^*), S_2(\theta^*, J^*)\} + \{H_1(\theta^*, J^*), S_1(\theta^*, J^*)\}) \\ &\quad + \frac{\epsilon^2}{2} \{\{H_0(\theta^*, J^*), S_1(\theta^*, J^*)\}, S_1(\theta^*, J^*)\} + O(\epsilon^3). \end{aligned} \quad (205)$$

Equating the two Hamiltonians  $H(\theta, J)$  and  $H^*(J^*)$ , comparing the terms of each order of  $\epsilon$  in (202) and (205), we obtain the result for  $O(\epsilon^0)$  as

$$H_0^*(J^*) = H_0(J^*) = \omega_0 J^*. \quad (206)$$

For  $O(\epsilon^1)$ , we get

$$H_1^*(J^*) = \{H_0(\theta^*, J^*), H_1(\theta^*, J^*)\} + H_1(\theta^*, J^*). \quad (207)$$

Let us consider the canonical equations of motion of the unperturbed system  $H_0^*(J^*)$ , introducing a time-like variable  $t^*$  as

$$\frac{d\theta^*}{dt^*} = \frac{\partial H_0^*(J^*)}{\partial J^*}, \quad \frac{dJ^*}{dt^*} = -\frac{\partial H_0^*(J^*)}{\partial \theta^*}, \quad (208)$$

$$\therefore J^* = \text{constant}, \quad \theta^* = \omega_0 t^* + \theta_0^*, \quad (209)$$

which leads to

$$\begin{aligned} \{H_0^*(J^*), S_1(\theta^*, J^*)\} &= \frac{\partial H_0^*}{\partial \theta^*} \frac{\partial S_1}{\partial J^*} - \frac{\partial H_0^*}{\partial J^*} \frac{\partial S_1}{\partial \theta^*} \\ &= -\left(\frac{\partial S_1}{\partial J^*} \frac{dJ^*}{dt^*} + \frac{\partial S_1}{\partial \theta^*} \frac{d\theta^*}{dt^*}\right) \\ &= -\frac{dS_1}{dt^*}. \end{aligned} \quad (210)$$

As a rule of usual perturbation theory, we request  $H_1^*(J^*)$  not to contain  $t^*$  as

$$\begin{aligned} H_1^*(J^*) &= \langle H_1(\theta^*, J^*) \rangle_{t^*} \\ &= \frac{1}{T} \int_0^T \left[ -\frac{J^{*2}}{6} \sin^4(\omega_0 t^* + \theta_0^*) \right] dt^* \\ &= \frac{\omega}{2\pi} \int_0^{2\pi} \left[ -\frac{J^{*2}}{6} \sin^4 \theta^* \right] \frac{dt^*}{d\theta^*} d\theta^* \\ &= \frac{\omega_0}{2\pi} \int_0^{2\pi} \left[ -\frac{J^{*2}}{48} \sin^4 \theta^* \right] d\theta^* \\ &= -\frac{J^{*2}}{16}. \end{aligned} \quad (211)$$

The leading term of the generating function  $S_1$  becomes

$$\begin{aligned}
S_1(\theta^*, J^*) &= \int (H_1^* - \langle H_1^* \rangle) dt^* \\
&= \int \frac{dS_1}{dt^*} t^* \\
&= \int \frac{J^{*2}}{12} \left( \cos 2\theta^* - \frac{1}{4} \cos 4\theta^* \right) \frac{d\theta^*}{\omega_0} \\
&= \frac{J^*}{24\omega_0} \left( \sin 2\theta^* - \frac{1}{8} \sin 4\theta^* \right), \tag{212}
\end{aligned}$$

neglecting a constant of integration.

For  $O(\epsilon^2)$  terms, we get

$$\begin{aligned}
H_2^*(J^*) &= \{H_0(\theta^*, J^*), S_2(\theta^*, J^*)\} + \{H_1(\theta^*, J^*), S_1(\theta^*, J^*)\} \\
&+ \frac{1}{2} \{ \{H_0(\theta^*, J^*), S_1(\theta^*, J^*)\}, S_1(\theta^*, J^*) \} + H_2(\theta^*, J^*), \tag{213}
\end{aligned}$$

and we request again

$$H_2^*(J^*) = \langle H_2(\theta^*, J^*) \rangle_{t^*}, \tag{214}$$

and

$$\{H_0^*(J^*), S_2(\theta^*, J^*)\} = -\frac{dS_2}{dt^*}. \tag{215}$$

Thus we can calculate ever higher-order components of  $H_i^*$  and  $S_i$  in similar ways. The new Hamiltonian and canonical frequency of the system up to  $O(\epsilon)$  now become

$$H^*(J^*) = H_0^*(J^*) + \epsilon H_1^*(J^*) = \omega_0 J^* - \epsilon \frac{J^{*2}}{16}. \tag{216}$$

and

$$\omega(J^*) \equiv \frac{\partial H^*(J^*)}{\partial J^*} = \omega_0 - \epsilon \frac{J^*}{8}. \tag{217}$$

The canonical equations of motion written in the new variables are

$$\begin{aligned}
\frac{d\theta^*}{dt} &= \frac{\partial H^*(J^*)}{\partial J^*} \\
&= \omega(J^*) = \text{constant}, \tag{218}
\end{aligned}$$

$$\frac{dJ^*}{dt} = -\frac{\partial H^*(J^*)}{\partial \theta^*} = 0, \tag{219}$$

$$\therefore J^* = \text{constant}. \tag{220}$$

The relationship between old  $(\theta, J)$  and new  $(\theta^*, J^*)$  variables is now explicit using the Poisson bracket operator  $D_S \equiv \{, S\}$  as

$$\begin{aligned}
\theta &= e^{\epsilon D_S} \theta^* \\
&= \theta^* + \epsilon \frac{\partial}{\partial J^*} S(\theta^*, J^*) + \frac{\epsilon^2}{2} \{ \{ \theta^*, S(\theta^*, J^*) \}, S(\theta^*, J^*) \} + \dots, \tag{221}
\end{aligned}$$

$$\begin{aligned}
J &= e^{\epsilon D_S} J^* \\
&= J^* - \epsilon \frac{\partial}{\partial \theta^*} S(\theta^*, J^*) + \frac{\epsilon^2}{2} \{ \{ J^*, S(\theta^*, J^*) \}, S(\theta^*, J^*) \} + \dots. \tag{222}
\end{aligned}$$

Substituting the specific form of  $S_1$  (212) into (221) and (222), we get

$$\theta = \theta^* + \epsilon \frac{J^*}{12\omega_0} \left( \sin 2\theta^* - \frac{1}{8} \sin 4\theta^* \right) + O(\epsilon^2), \quad (223)$$

$$J = J^* - \epsilon \frac{J^*}{24\omega_0} \left( \cos 2\theta^* - \frac{1}{2} \cos 4\theta^* \right) + O(\epsilon^2). \quad (224)$$

Since  $(\theta, J)$  are the angle and action variables in the Hamiltonian system  $H_0$  of the harmonic oscillator, the relationship between  $(q, p)$  and  $(\theta, J)$  has been described in (164). Now in the perturbed system (198),  $(\theta, J)$  are no longer angle/action variables, but are still canonical variables described in (223) and (224). Hence we obtain the final solution  $(q, p)$  for the system (198) by substituting  $\theta$  and  $J$  of (223) and (224) into the  $(q, p)$ -( $\theta, J$ ) relationship of the system  $H_0$ , i.e. (164). Thus we get

$$q = \sqrt{\frac{2J^*}{\omega_0} \left[ 1 - \frac{\epsilon J^*}{12\omega_0} \left( \cos 2\theta^* - \frac{1}{4} \cos 4\theta^* \right) \right]} \sin \left[ \theta^* + \frac{\epsilon J^*}{12\omega_0} \left( \sin 2\theta^* - \frac{1}{8} \sin 4\theta^* \right) \right] + O(\epsilon^2), \quad (225)$$

$$p = \sqrt{2J\omega_0 \left[ 1 - \frac{\epsilon J}{12\omega_0} \left( \cos 2\theta - \frac{1}{4} \cos 4\theta \right) \right]} \cos \left[ \theta + \frac{\epsilon J}{12\omega_0} \left( \sin 2\theta - \frac{1}{8} \sin 4\theta \right) \right] + O(\epsilon^2). \quad (226)$$

Hereafter we consider  $(\theta, J)$  as  $(\theta^*, J^*)$ , neglecting the superscript  $*$  attached. Hence we get

$$q = \sqrt{\frac{2J}{\omega_0} \left[ 1 - \frac{\epsilon J}{12\omega_0} \left( \cos 2\theta - \frac{1}{4} \cos 4\theta \right) \right]} \sin \left[ \theta + \frac{\epsilon J}{12\omega_0} \left( \sin 2\theta - \frac{1}{8} \sin 4\theta \right) \right] + O(\epsilon^2), \quad (227)$$

$$p = \sqrt{2J\omega_0 \left[ 1 - \frac{\epsilon J}{12\omega_0} \left( \cos 2\theta - \frac{1}{4} \cos 4\theta \right) \right]} \cos \left[ \theta + \frac{\epsilon J}{12\omega_0} \left( \sin 2\theta - \frac{1}{8} \sin 4\theta \right) \right] + O(\epsilon^2). \quad (228)$$

### 6.2.2 Error Hamiltonian in symplectic integrator

Let us consider a situation where we numerically calculate the solutions of a system dominated by the Hamiltonian (198) up to first-order, namely

$$H(q, p) = \frac{p^2}{2} + \frac{\omega_0^2 q^2}{2} - \frac{\epsilon q^4}{24}. \quad (229)$$

When we write the Hamiltonian as  $H = T(p) + V(q)$ ,

$$T(p) = \frac{p^2}{2}, \quad V(q) = \frac{\omega_0^2 q^2}{2} - \frac{\epsilon q^4}{24}. \quad (230)$$

The error Hamiltonian for the first-order symplectic integrator thus becomes

$$\begin{aligned} H_{\text{err,1st}} &= \frac{\tau}{2} \{V(q), T(p)\} \\ &= \frac{\tau}{2} \frac{\partial V}{\partial q} \frac{\partial T}{\partial p} \\ &= \frac{\tau}{2} \left( \omega_0^2 q - \frac{\epsilon}{6} q^3 \right) p. \end{aligned} \quad (231)$$

By (227),

$$\begin{aligned}
 q^3 &= \sqrt{\frac{8J^3}{\omega_0^3}} \sin^3 \theta + O(\epsilon) \\
 &= \frac{1}{4} \sqrt{\frac{8J^3}{\omega_0^3}} (3 \sin \theta - \sin 3\theta) + O(\epsilon),
 \end{aligned} \tag{232}$$

which leads to

$$\begin{aligned}
 \omega_0^2 q - \frac{\epsilon}{6} q^3 &= \sqrt{2J\omega_0^3} \left[ \sin \theta + \frac{\epsilon J}{24\omega_0} \left( \frac{5}{2} \sin \theta + \frac{1}{4} \sin 3\theta \right) \right] - \frac{\epsilon}{6} \frac{1}{4} \sqrt{\frac{8J^3}{\omega_0^3}} (3 \sin \theta - \sin 3\theta) \\
 &= \sin \theta \left[ \sqrt{2J\omega_0^3} \left( 1 + \frac{\epsilon J}{24\omega_0} \frac{5}{2} \right) - \frac{\epsilon}{6} \frac{3}{4} \sqrt{\frac{8J^3}{\omega_0^3}} \right] + \sin 3\theta \left[ \sqrt{2J\omega_0^3} \frac{\epsilon J}{24\omega_0} \frac{1}{4} + \frac{\epsilon}{6} \frac{1}{4} \sqrt{\frac{8J^3}{\omega_0^3}} \right] \\
 &= \sin \theta \left[ \sqrt{2J\omega_0^3} + \epsilon \left( \frac{5}{48} \sqrt{2J^3\omega_0} - \frac{1}{8} \sqrt{8J^3\omega_0} \right) \right] + \epsilon \sin 3\theta \left( \frac{\sqrt{2J^3\omega_0}}{96} + \frac{1}{24} \sqrt{\frac{8J^3}{\omega_0^2}} \right),
 \end{aligned} \tag{233}$$

up to  $O(\epsilon)$ . Thus the error Hamiltonian in (231) becomes

$$\begin{aligned}
 \frac{2}{\tau} H_{\text{err,1st}} &= \left( \omega_0^2 q - \frac{\epsilon}{6} q^3 \right) p \\
 &= \left[ \sin \theta \left( \sqrt{2J\omega_0^3} + \epsilon \left( \frac{5}{48} \sqrt{2J^3\omega_0} - \frac{1}{8} \sqrt{8J^3\omega_0} \right) \right) + \epsilon \sin 3\theta \left( \frac{\sqrt{2J^3\omega_0}}{96} + \frac{1}{24} \sqrt{\frac{8J^3}{\omega_0^2}} \right) \right] \\
 &\quad \times \sqrt{2J\omega_0} \left[ \cos \theta + \frac{\epsilon J}{24\omega_0} \left( \frac{5}{2} \sin \theta + \frac{1}{4} \sin 3\theta \right) \right] \\
 &= \sqrt{2J\omega_0} \left\{ \sin \theta \cos \theta \left[ \sqrt{2J\omega_0^3} + \epsilon \left( \frac{5}{48} \sqrt{2J^3\omega_0} - \frac{1}{8} \sqrt{8J^3\omega_0} \right) \right] \right. \\
 &\quad + \epsilon \left[ \left( \frac{\sqrt{2J^3\omega_0}}{96} + \frac{1}{24} \sqrt{\frac{8J^3}{\omega_0^2}} \right) \sin 3\theta \cos \theta \right. \\
 &\quad \left. \left. + \frac{J}{24\omega_0} \sqrt{2J\omega_0^3} \sin \theta \left( \frac{5}{2} \sin \theta + \frac{1}{4} \sin 3\theta \right) \right] \right\} \\
 &= \sqrt{2J\omega_0} \left\{ \frac{1}{2} \sin 2\theta \left[ \sqrt{2J\omega_0^3} + \epsilon \left( \frac{5}{48} \sqrt{2J^3\omega_0} - \frac{1}{8} \sqrt{8J^3\omega_0} \right) \right] \right. \\
 &\quad + \epsilon \left[ \left( \frac{\sqrt{2J^3\omega_0}}{96} + \frac{1}{24} \sqrt{\frac{8J^3}{\omega_0^2}} \right) \frac{1}{2} (\sin 4\theta + \sin 2\theta) \right. \\
 &\quad \left. \left. + \frac{J}{24\omega_0} \sqrt{2J\omega_0^3} \left( \frac{5}{4} - \frac{5}{4} \cos 2\theta - \frac{1}{8} \cos 4\theta + \frac{1}{8} \cos 2\theta \right) \right] \right\},
 \end{aligned} \tag{234}$$

up to  $O(\epsilon)$ . Hence if we average  $H_{\text{err,1st}}$  over  $\theta$ , then we get

$$\begin{aligned}
 \langle H_{\text{err,1st}} \rangle &= \epsilon \frac{5}{4} \frac{\tau}{2} \sqrt{2J\omega_0} \frac{J}{24\omega_0} \sqrt{2J\omega_0^3} \\
 &= \epsilon \tau \frac{5}{96} \omega_0 J^2.
 \end{aligned} \tag{235}$$

In order to estimate the numerical error of the first-order symplectic integrator when adopting to the nonlinear pendulum system, we now apply the Hori's perturbation theory to a Hamiltonian

$$H = \frac{p^2}{2} + \frac{\omega_0^2 q^2}{2} - \frac{\epsilon q^4}{24} + H_{\text{err,1st}}. \quad (236)$$

We rewrite the above Hamiltonian as

$$H = H_0(J) + H_1(\theta, J), \quad (237)$$

where

$$H_0 = \frac{p^2}{2} + \frac{\omega_0^2 q^2}{2} - \frac{\epsilon q^4}{24} = \omega_0 J - \epsilon \frac{J^2}{16}, \quad (238)$$

from (216), and

$$H_1(\theta, J) = H_{\text{err,1st}}(\theta, J), \quad (239)$$

from (234). Since we already knew the averaged perturbation Hamiltonian as (235), applying the perturbation theory gives us a new transformed Hamiltonian as

$$H^*(J^*) = \omega_0 J^* - \epsilon \frac{J^{*2}}{16} + \epsilon \tau \frac{5}{96} \omega_0 J^{*2}, \quad (240)$$

and a new canonical frequency as

$$\omega(J^*) \equiv \frac{\partial H^*(J^*)}{\partial J^*} = \omega_0 - \frac{\epsilon}{8} J^* - \epsilon \tau \frac{5}{192} \omega_0 J^*. \quad (241)$$

Now that the canonical frequency  $\omega$  includes the action  $J^*$  as  $\omega(J^*)$ , we may be able to make it approach the real value,  $\omega_0 - \frac{\epsilon}{8} J_0$  where  $J_0$  is the initial (or "observed") value of the action in real system.

The canonical equation of motion for  $J^*$  becomes as

$$\frac{dJ^*}{dt} = -\frac{\partial H^*(J^*)}{\partial \theta^*}, \quad (242)$$

from which we know that  $J^*$  is a constant. According to the relationship between old and new variables of (222), we get

$$\begin{aligned} J &= e^{\epsilon D_S} J^* \\ &= J^* - \epsilon \frac{\partial}{\partial \theta^*} S(\theta^*, J^*) + \frac{\epsilon^2}{2} \{ \{ J^*, S(\theta^*, J^*) \}, S(\theta^*, J^*) \} \\ &= J^* - \frac{\partial}{\partial \theta^*} \int (H_1 - H_1^*) dt^* \\ &= J^* - \frac{1}{\omega_0} (H_1 - H_1^*). \end{aligned} \quad (243)$$

Suppose  $J = J_0$  when  $t = 0$ , then

$$J_0 = J^* - \frac{1}{\omega_0} (H_{1,t=0} - H_1^*), \quad (244)$$

since  $J^*$  and  $H_1^*$  are constants, so

$$\therefore J^* = J_0 + \frac{1}{\omega_0} (H_{1,t=0} - H_1^*). \quad (245)$$

As for  $\theta^*$ , the canonical equation of motion becomes

$$\begin{aligned}
\frac{d\theta^*}{dt} &= \frac{\partial H^*(J^*)}{\partial J^*} \\
&= \omega_0 - \frac{\epsilon}{8} J^* + \epsilon \tau \frac{5\omega_0}{96} J^* \\
&= \omega_0 - \frac{\epsilon}{8} \left( J_0 + \frac{1}{\omega_0} (H_{1,t=0} - H_1^*) \right) + \epsilon \tau \frac{5\omega_0}{96} \left( J_0 + \frac{1}{\omega_0} (H_{1,t=0} - H_1^*) \right) \\
&= \omega_0 - \frac{\epsilon}{8} J_0 + \left( \epsilon \tau \frac{5\omega_0}{96} - \frac{\epsilon}{8} \right) \frac{H_{1,t=0} - H_1^*}{\omega_0} + \epsilon \tau \frac{5\omega_0}{96} J_0.
\end{aligned} \tag{246}$$

In (246), the first two terms in the right-hand side ( $\omega_0 - \frac{\epsilon}{8} J_0$ ) denote the canonical frequency of the unperturbed Hamiltonian system, (238). Other terms denote numerical secular error of the angle variable,  $\theta$ . What is most important here is that the secular (or constant) numerical error

$$\left( \epsilon \tau \frac{5\omega_0}{96} - \frac{\epsilon}{8} \right) \frac{H_{1,t=0} - H_1^*}{\omega_0} + \epsilon \tau \frac{5\omega_0}{96} J_0, \tag{247}$$

may depend strongly on initial values of  $J$ ,  $\theta$ , and parameters  $\epsilon$ ,  $\tau$ , and  $\omega_0$ . Although we do not demonstrate here, certain combinations of these parameters may reduce the secular error of (247) nearly equal to zero. Other combinations may terribly increase the secular numerical error. This kind of error reduction happens only when the canonical frequency of a Hamiltonian system depends on its initial oscillation amplitude; in other words, when the potential of the Hamiltonian is not isochrone, and dependent on  $J$ . The reduction or non-reduction of the secular numerical error in planetary longitudes seen in the previous sections is thus qualitatively understood to some extent.

## Appendix A. Jacobi coordinate

We request a Hamiltonian used in the WH map to have the following properties:

1. As for the Keplerian part, it should have the same form as the Hamiltonian of the two-body problem,  $\frac{p^2}{2m} - \frac{\mu}{r}$ , or the sum of this form.
2. As for the interaction part, it should be described only by the relative distance, such as  $V(r)$ .
3. The magnitude of the interaction part should be much less than that of the Keplerian part ( $H_{\text{kep}} \gg H_{\text{int}}$ ).

However, simple heliocentric or barycentric coordinate does not satisfy the above requests. For example, writing the Hamiltonian using the barycentric coordinate ends up with

$$\begin{aligned}
 H &= \sum_{i=0}^N \frac{p_i^2}{2m_i} + \left( - \sum_{i=1}^N \frac{Gm_0m_i}{|\mathbf{r}_i - \mathbf{r}_0|} + \sum_{i=1}^N \frac{Gm_0m_i}{|\mathbf{r}_i - \mathbf{r}_0|} \right) - \sum_{i=0}^N \sum_{j=i+1}^N \frac{Gm_im_j}{|\mathbf{r}_i - \mathbf{r}_j|} \\
 &= \frac{p_0^2}{2m_0} + \sum_{i=1}^N \left( \frac{p_i^2}{2m_i} - \frac{Gm_0m_i}{|\mathbf{r}_i - \mathbf{r}_0|} \right) + \sum_{i=1}^N \frac{Gm_0m_i}{|\mathbf{r}_i - \mathbf{r}_0|} - \left( \sum_{j=1}^N \frac{Gm_0m_j}{|\mathbf{r}_0 - \mathbf{r}_j|} + \sum_{i=1}^N \sum_{j=i+1}^N \frac{Gm_im_j}{|\mathbf{r}_i - \mathbf{r}_j|} \right) \\
 &= \frac{p_0^2}{2m_0} + \sum_{i=1}^N \left( \frac{p_i^2}{2m_i} - \frac{Gm_0m_i}{|\mathbf{r}_i - \mathbf{r}_0|} \right) - \sum_{i=1}^N \sum_{j=i+1}^N \frac{Gm_im_j}{|\mathbf{r}_i - \mathbf{r}_j|}, \tag{248}
 \end{aligned}$$

where we cannot classify the kinetic energy of the Sun ( $p_0^2/2m_0$ ) into  $H_{\text{kep}}$  nor  $H_{\text{int}}$ . One of the canonical variables which suits our request is the Jacobi coordinates (Plummer, 1960). The Jacobi coordinates  $\tilde{\mathbf{r}}_i$  are defined as

$$\tilde{\mathbf{r}}_i = \mathbf{r}_i - \frac{1}{\sigma_{i-1}} \sum_{j=1}^{i-1} m_j \mathbf{r}_j, \tag{249}$$

where

$$\sigma_i = \sigma_{i-1} + m_i, \quad \sigma_0 = m_0 = M_\odot, \quad \sigma_{-1} \equiv 0, \tag{250}$$

$$\tilde{m}_i = \frac{\sigma_{i-1}}{\sigma_i} m_i, \tag{251}$$

$$\tilde{\mu}_i = \frac{\sigma_i}{\sigma_{i-1}} G, \tag{252}$$

Canonical momenta which are conjugate to  $\tilde{\mathbf{r}}_i$  are

$$\tilde{\mathbf{p}}_i = \tilde{m}_i \tilde{\mathbf{v}}_i, \tag{253}$$

and the velocities are

$$\tilde{\mathbf{v}}_i = \frac{d\tilde{\mathbf{r}}_i}{dt}. \tag{254}$$

In this manuscript, we have utilized the Jacobi coordinate system in our symplectic numerical integration of the  $H = T(p) + V(q)$  type. Thus the discussion below has a certain sense to describe our method of numerical integration in detail.



The advantage of transforming to the Jacobi variables is that in the barycentric frame, the kinetic energy of the  $N + 1$ -body system becomes

$$\sum_{i=1}^N \frac{\tilde{p}_i^2}{2\tilde{m}_i}, \quad (255)$$

without terms of the central mass as follows (Plummer, 1960).

When we represent the position of the barycenter of the particle of 1 to  $i$  as  $\mathbf{R}_i$ ,

$$\sigma_i \mathbf{R}_i = m_0 \mathbf{r}_0 + m_1 \mathbf{r}_1 + m_2 \mathbf{r}_2 + m_{i-1} \mathbf{r}_{i-1} + m_i \mathbf{r}_i, \quad (256)$$

$$\sigma_{i-1} \mathbf{R}_{i-1} = m_0 \mathbf{r}_0 + m_1 \mathbf{r}_1 + m_2 \mathbf{r}_2 + m_{i-1} \mathbf{r}_{i-1}. \quad (257)$$

Subtracting (257) from (256),

$$\sigma_i \mathbf{X}_i - \sigma_{i-1} \mathbf{X}_{i-1} = m_i \mathbf{r}_i = (\sigma_i - \sigma_{i-1}) \mathbf{r}_i, \quad (258)$$

which is due to the definition of  $\sigma_i$  (250). By the definition of the Jacobi coordinates  $\tilde{\mathbf{r}}_i$ ,

$$\tilde{\mathbf{r}}_i = \mathbf{r}_i - \mathbf{R}_{i-1}, \quad (259)$$

$$\therefore \mathbf{r}_i = \tilde{\mathbf{r}}_i + \mathbf{R}_{i-1}. \quad (260)$$

Substituting (260) into (258),

$$\sigma_i \mathbf{R}_i - \sigma_{i-1} \mathbf{R}_{i-1} = (\sigma_i - \sigma_{i-1})(\tilde{\mathbf{r}}_i + \mathbf{R}_{i-1}), \quad (261)$$

$$\therefore \sigma_i (\mathbf{R}_i - \mathbf{R}_{i-1}) = (\sigma_i - \sigma_{i-1}) \tilde{\mathbf{r}}_i. \quad (262)$$

Hereafter we concentrate on the  $x$ -components of the vectors  $\mathbf{r}_i$ ,  $\tilde{\mathbf{r}}_i$ ,  $\mathbf{R}_i$ :  $x_i$ ,  $\tilde{x}_i$ , and  $X_i$ , respectively. Taking the square of (258) and (262),

$$(\sigma_i - \sigma_{i-1})^2 x_i^2 = (\sigma_i X_i - \sigma_{i-1} X_{i-1})^2, \quad (263)$$

$$(\sigma_i - \sigma_{i-1})^2 \tilde{x}_i^2 = \sigma_i^2 (X_i - X_{i-1})^2. \quad (264)$$

Performing the operation (263) - (264)  $\times \frac{\sigma_{i-1}}{\sigma_i}$ , we get

$$\begin{aligned} & (\sigma_i - \sigma_{i-1})^2 \left( x_i^2 - \frac{\sigma_{i-1}}{\sigma_i} \tilde{x}_i^2 \right) \\ &= (\sigma_i X_i - \sigma_{i-1} X_{i-1})^2 - \frac{\sigma_{i-1}}{\sigma_i} (X_i - X_{i-1})^2 \\ &= \sigma_i^2 X_i^2 - 2\sigma_i \sigma_{i-1} X_i X_{i-1} + \sigma_{i-1}^2 X_{i-1}^2 - \sigma_{i-1} \sigma_i (X_i^2 - 2X_i X_{i-1} + X_{i-1}^2) \\ &= X_i^2 (\sigma_i^2 - \sigma_i \sigma_{i-1}) + X_{i-1}^2 (\sigma_{i-1}^2 - \sigma_i \sigma_{i-1}) \\ &= (\sigma_i - \sigma_{i-1}) (\sigma_i X_i^2 - \sigma_{i-1} X_{i-1}^2). \end{aligned} \quad (265)$$

In the case of finite-mass system,  $\sigma_i \neq \sigma_{i-1}$ , then

$$(\sigma_i - \sigma_{i-1}) \left( x_i^2 - \frac{\sigma_{i-1}}{\sigma_i} \tilde{x}_i^2 \right) = \sigma_i X_i^2 - \sigma_{i-1} X_{i-1}^2. \quad (266)$$

Addition of all the equations of the type (266) from  $i = 0$  to  $i = N$ ,

$$\begin{aligned}
& \sum_{i=0}^N (\sigma_i - \sigma_{i-1}) \left( r_i^2 - \frac{\sigma_{i-1}}{\sigma_i} \tilde{r}_i^2 \right) \\
&= (\sigma_0 X_0^2 - \sigma_{-1} X_{-1}^2) + (\sigma_1 X_1^2 - \sigma_0 X_0^2) + (\sigma_2 X_2^2 - \sigma_1 X_1^2) \\
&\quad + \cdots + (\sigma_{N-1} X_{N-1}^2 - \sigma_{N-2} X_{N-2}^2) + (\sigma_N X_N^2 - \sigma_{N-1} X_{N-1}^2) \\
&= \sigma_N X_N^2 - \sigma_0 X_0^2 \\
&= \sigma_N X_N^2. \quad (\because \sigma_{-1} = 0)
\end{aligned} \tag{267}$$

Since  $\sigma_i - \sigma_{i-1} = m_i$ ,

$$\begin{aligned}
\sum_{i=0}^N m_i x_i^2 &= \sum_{i=0}^N m_i \frac{\sigma_{i-1}}{\sigma_i} \tilde{x}_i^2 + \sigma_N X_N^2 \\
&= \sum_{i=1}^N m_i \frac{\sigma_{i-1}}{\sigma_i} \tilde{x}_i^2 + \sigma_N X_N^2 \quad (\because \sigma_{-1} = 0) \\
&= \sum_{i=1}^N \tilde{m}_i \tilde{x}_i^2 + \sigma_N X_N^2.
\end{aligned} \tag{268}$$

The relations between the coordinates have been written down for one kind only. But they are linear and the same for all three coordinates  $\mathbf{r}_i = (x_i, y_i, z_i)$ ,  $\tilde{\mathbf{r}}_i = (\tilde{x}_i, \tilde{y}_i, \tilde{z}_i)$ , and  $\mathbf{R}_i = (X_i, Y_i, Z_i)$  separately as

$$\sum_{i=0}^N m_i y_i^2 = \sum_{i=1}^N \tilde{m}_i \tilde{y}_i^2 + \sigma_N Y_N^2, \tag{269}$$

$$\sum_{i=0}^N m_i z_i^2 = \sum_{i=1}^N \tilde{m}_i \tilde{z}_i^2 + \sigma_N Z_N^2. \tag{270}$$

Above derivation can be applied also to the velocity components

$$\mathbf{v}_i = \frac{d\mathbf{r}}{dt} = (\dot{x}_i, \dot{y}_i, \dot{z}_i), \tag{271}$$

$$\tilde{\mathbf{v}}_i = \frac{d\tilde{\mathbf{r}}}{dt} = (\dot{\tilde{x}}_i, \dot{\tilde{y}}_i, \dot{\tilde{z}}_i), \tag{272}$$

$$\mathbf{V}_i = \frac{d\mathbf{R}}{dt} = (\dot{X}_i, \dot{Y}_i, \dot{Z}_i), \tag{273}$$

as

$$\sum_{i=0}^N m_i \dot{x}_i^2 = \sum_{i=1}^N \tilde{m}_i \dot{\tilde{x}}_i^2 + \sigma_N \dot{X}_N^2, \tag{274}$$

$$\sum_{i=0}^N m_i \dot{y}_i^2 = \sum_{i=1}^N \tilde{m}_i \dot{\tilde{y}}_i^2 + \sigma_N \dot{Y}_N^2, \tag{275}$$

$$\sum_{i=0}^N m_i \dot{z}_i^2 = \sum_{i=1}^N \tilde{m}_i \dot{\tilde{z}}_i^2 + \sigma_N \dot{Z}_N^2. \tag{276}$$

Adding (274), (275), (276), we can represent the kinetic energy of the system

$$\sum_{i=0}^N \frac{m_i}{2} (\dot{y}_i^2 + \dot{y}_i^2 + \dot{z}_i^2) = \sum_{i=1}^N \frac{\tilde{m}_i}{2} (\dot{x}_i^2 + \dot{y}_i^2 + \dot{z}_i^2) + \frac{\sigma_N}{2} (\dot{X}_N^2 + \dot{Y}_N^2 + \dot{Z}_N^2), \quad (277)$$

or using the momentum  $\mathbf{p}_i$  and  $\tilde{\mathbf{p}}_i$

$$\sum_{i=0}^N \frac{\mathbf{p}_i^2}{2m_i} = \sum_{i=1}^N \frac{\tilde{\mathbf{p}}_i^2}{2\tilde{m}_i} + \frac{\mathbf{p}_R^2}{2M_{\text{tot}}}, \quad (278)$$

where  $\mathbf{p}_R$  is the total momentum of the barycenter (total momentum of the whole system), and  $M_{\text{tot}}$  is the total mass of the system,  $M_{\text{tot}} = \sum_{i=0}^N m_i$ . Then, the general Hamiltonian for the  $N + 1$ -body (one central mass and  $N$  planets) becomes

$$H = \frac{\mathbf{p}_R^2}{2M_{\text{tot}}} + \sum_{i=1}^N \frac{\tilde{\mathbf{p}}_i^2}{2\tilde{m}_i} - \sum_{i=0}^N \sum_{j=i+1}^N \frac{k^2 m_i m_j}{|\mathbf{r}_i - \mathbf{r}_j|}. \quad (279)$$

By construction, the total momentum  $\mathbf{p}_R$  is an integral of the motion, which means that the center of mass moves as a free particle. Hence the center of mass contribution to the Hamiltonian  $\mathbf{p}_R^2/2M_{\text{tot}}$  will be omitted. Thus the problem of  $N + 1$  bodies is reduced to a problem of  $N$  fictitious bodies with mass  $\tilde{m}_i$ , and the total order of the differential equations of motion is reduced by 6. In view of (253) and (254), we can rewrite the full Hamiltonian as

$$\begin{aligned} H &= \sum_{i=1}^N \frac{\tilde{\mathbf{p}}_i^2}{2\tilde{m}_i} - \sum_{i=0}^N \sum_{j=i+1}^N \frac{k^2 m_i m_j}{|\mathbf{r}_i - \mathbf{r}_j|} \\ &= \sum_{i=1}^N \frac{\tilde{\mathbf{p}}_i^2}{2\tilde{m}_i} - \sum_{i=1}^N \frac{k^2 m_i m_0}{r_i} - \sum_{i=1}^N \sum_{j=i+1}^N \frac{k^2 m_i m_j}{|\mathbf{r}_i - \mathbf{r}_j|}. \end{aligned} \quad (280)$$

where  $r_j$  denote the heliocentric distance,  $|\mathbf{r}_i - \mathbf{r}_0|$ . Note that we have changed the index  $j$  to  $i$  for simplicity in the second sum of (280). Adding and subtracting the quantity

$$\sum_{i=1}^N \frac{k^2 m_i m_0}{\tilde{r}_i}, \quad (281)$$

into the righthand-side of (280), the Hamiltonian becomes

$$\begin{aligned} H &= \sum_{i=1}^N \frac{\tilde{\mathbf{p}}_i^2}{2\tilde{m}_i} + \left( - \sum_{i=1}^N \frac{k^2 m_i m_0}{\tilde{r}_i} + \sum_{i=1}^N \frac{k^2 m_i m_0}{\tilde{r}_i} \right) - \sum_{i=1}^N \frac{k^2 m_i}{r_i} - \sum_{i=1}^N \sum_{j=i+1}^N \frac{k^2 m_i m_j}{|\mathbf{r}_i - \mathbf{r}_j|} \\ &= \sum_{i=1}^N \left( \frac{\tilde{\mathbf{p}}_i^2}{2\tilde{m}_i} - \frac{k^2 m_i m_0}{\tilde{r}_i} \right) + \sum_{i=1}^N \left( \frac{k^2 m_i m_0}{\tilde{r}_i} - \frac{k^2 m_i m_0}{r_i} \right) - \sum_{i=1}^N \sum_{j=i+1}^N \frac{k^2 m_i m_j}{|\mathbf{r}_i - \mathbf{r}_j|} \\ &= \sum_{i=1}^N \left( \frac{\tilde{\mathbf{p}}_i^2}{2\tilde{m}_i} - \tilde{\mu}_i \frac{\tilde{m}_i m_0}{\tilde{r}_i} \right) + k^2 \sum_{i=1}^N m_i m_0 \left( \frac{1}{\tilde{r}_i} - \frac{1}{r_i} \right) - \sum_{i=1}^N \sum_{j=i+1}^N \frac{k^2 m_i m_j}{|\mathbf{r}_i - \mathbf{r}_j|}. \end{aligned} \quad (282)$$

The relationship

$$k^2 m_i = k^2 \frac{m_i}{\tilde{m}_i} \tilde{m}_i$$

$$\begin{aligned}
&= k^2 \frac{\sigma_i}{\sigma_{i-1}} \tilde{m}_i \\
&= k^2 \frac{\tilde{\mu}_i}{k^2} \tilde{m}_i \\
&= \tilde{\mu}_i \tilde{m}_i,
\end{aligned} \tag{283}$$

is used in the first sum in the righthand-side of (282).

Thus the Hamiltonian in (282) becomes a desirable form for the WH map as

$$H = H_{\text{kep}} + \epsilon H_{\text{int}}, \tag{284}$$

where

$$H_{\text{kep}} = \sum_{i=1}^N \left( \frac{\tilde{p}_i^2}{2\tilde{m}_i} - \tilde{\mu}_i \frac{\tilde{m}_i m_0}{\tilde{r}_i} \right), \tag{285}$$

$$\epsilon H_{\text{int}} = H_{\text{direct}} + H_{\text{indirect}}, \tag{286}$$

$$H_{\text{direct}} = - \sum_{i=1}^N \sum_{j=i+1}^N \frac{k^2 m_i m_j}{|\mathbf{r}_i - \mathbf{r}_j|}, \tag{287}$$

$$H_{\text{indirect}} = k^2 \sum_{i=1}^N m_i m_0 \left( \frac{1}{\tilde{r}_i} - \frac{1}{r_i} \right). \tag{288}$$

The magnitude of  $H_{\text{direct}}$  is  $O(m^2)$ . The magnitude of  $H_{\text{indirect}}$  is also  $O(m^2)$  because of the difference of close terms  $1/\tilde{r}_i - 1/r_i$ ,  $O(m)$ . Hence the magnitude of  $\epsilon H_{\text{int}}$  becomes  $O(m)$  times smaller than that of the Kepler Hamiltonian,  $H_{\text{kep}}$ .

Note that some numerical inaccuracies can arise from  $H_{\text{indirect}}$  in which a subtraction of two quantities at the same order is performed. Straightforward evaluation of these expressions can be avoided by certain reformulation as used in Encke's method (Battin, 1987).

We can also obtain expressions for the angular momentum by the Jacobi coordinates. See Plummer (1960) for detail.

## Two-body Hamiltonian in the Jacobi coordinates and energy integral

Consider the Hamiltonian of the two-body problem written in the Jacobi coordinates using  $\mu = k^2(m_0 + m_1)$ . From (285) when  $N = 1$ , we can easily obtain

$$H_{2\text{body}} = \frac{\tilde{p}_1^2}{2\tilde{m}_1} - \tilde{\mu}_1 \frac{\tilde{m}_1 m_0}{\tilde{r}_1}. \tag{289}$$

Using the following relationships

$$\tilde{m}_1 = \frac{\sigma_0}{\sigma_1} m_1 = \frac{m_0 m_1}{m_0 + m_1}, \tag{290}$$

$$\tilde{\mu}_1 = \frac{\sigma_1}{\sigma_0} k^2 = \frac{m_0 + m_1}{m_0 m_1} k^2, \tag{291}$$

$$\tilde{r}_1 = r_1, \quad \tilde{v}_1 = v_1, \quad \tilde{p}_1 = \tilde{m}_1 v_1, \tag{292}$$

we get

$$\begin{aligned}
H_{2\text{body}} &= \frac{m_0 + m_1}{2m_0m_1} \left( \frac{m_0m_1}{m_0 + m_1} v_1 \right)^2 - \frac{m_0 + m_1}{m_0} k^2 \frac{m_0m_1}{m_0 + m_1} \frac{m_0}{r_1} \\
&= \frac{1}{2} \frac{m_0m_1}{m_0 + m_1} v_1^2 - \frac{k^2 m_1 m_0}{r_1} \\
&= \frac{m_0m_1}{m_0 + m_1} \left( \frac{v_1^2}{2} - \frac{\mu}{r_1} \right) \\
&= \tilde{m}_1 \left( \frac{v_1^2}{2} - \frac{\mu}{r_1} \right), \tag{293}
\end{aligned}$$

with a set of canonical variables  $(\tilde{r}_1, \tilde{p}_1) = (r_1, \tilde{m}_1 v_1)$ .

The equation (293) indicates that the general two-body Hamiltonian  $H_{2\text{body}}$  is not identical to the usual “energy integral,”  $\frac{v^2}{2} - \frac{\mu}{r}$ , by a factor of the reduced mass,  $\tilde{m}_1 = \frac{m_0m_1}{m_0+m_1}$ . When we analyzed the error Hamiltonian  $H_{\text{err}}$  in the previous sections, we have taken that this reduced mass  $\tilde{m}_1 = \frac{m_0m_1}{m_0+m_1}$  as unity for simplicity, and treated

$$H_{2\text{body}} = \frac{v_1^2}{2} - \frac{\mu}{r_1}, \tag{294}$$

with canonical variables  $(r_1, v_1)$ ; i.e. usual heliocentric position and velocity. We can normalize the two-body Hamiltonian (293) into the form of (294) through a certain conversion of units.

## Appendix B. Canonical relative (DH) coordinates

Let the coordinates and velocities of planets viewed from barycentric frame as  $\rho_i$  and  $\dot{\rho}_i$ . The relationship between the coordinates and velocities based on heliocentric frame  $r_i$  and  $\dot{r}_i$  is

$$\rho_i = r_i - \frac{\sum_{j=1}^N m_j r_j}{M + \sum_{j=1}^N m_j}, \quad \dot{\rho}_i = \dot{r}_i - \frac{\sum_{j=1}^N m_j \dot{r}_j}{M + \sum_{j=1}^N m_j}, \tag{295}$$

where  $N$  is the number of planets and  $M$  is the central mass. If we denote the position of the central mass by  $\rho_c$ ,

$$\begin{aligned}
\rho_c &= r_c - \frac{\sum_{j=1}^N m_j r_j}{M + \sum_{j=1}^N m_j} \\
&= - \frac{\sum_{j=1}^N m_j r_j}{M + \sum_{j=1}^N m_j}, \tag{296}
\end{aligned}$$

since obviously  $r_c = \mathbf{0}$  in heliocentric coordinate.

Below, we consider separately the kinetic energy and potential energy of the system.

Total kinetic energy of planets  $T_p$  is

$$T_p = \frac{1}{2} \sum_{j=1}^N m_j \dot{\rho}_j^2$$

$$\begin{aligned}
&= \frac{1}{2} \sum_{j=1}^N m_j \left( \dot{\mathbf{r}}_j - \frac{\sum_{i=1}^N m_i \dot{\mathbf{r}}_i}{M + \sum_{i=1}^N m_i} \right)^2 \\
&= \frac{1}{2} \sum_{j=1}^N m_j \left[ \dot{\mathbf{r}}_j^2 - \frac{2}{M + \sum_{i=1}^N m_i} \left\{ \dot{\mathbf{r}}_j \cdot \sum_{i=1}^N m_i \dot{\mathbf{r}}_i \right\} + \frac{\left( \sum_{i=1}^N m_i \dot{\mathbf{r}}_i \right) \cdot \left( \sum_{i=1}^N m_i \dot{\mathbf{r}}_i \right)}{\left( M + \sum_{i=1}^N m_i \right)^2} \right].
\end{aligned} \tag{297}$$

The kinetic energy of the central mass  $T_c$  is

$$\begin{aligned}
T_c &= \frac{M}{2} \dot{\boldsymbol{\rho}}_c^2 \\
&= \frac{M}{2} \left( - \frac{\sum_{i=1}^N m_i \dot{\mathbf{r}}_i}{M + \sum_{i=1}^N m_i} \right)^2 \\
&= \frac{M}{2} \frac{\left( \sum_{j=1}^N m_j \dot{\mathbf{r}}_j \right) \cdot \left( \sum_{j=1}^N m_j \dot{\mathbf{r}}_j \right)}{\left( M + \sum_{j=1}^N m_j \right)^2}.
\end{aligned} \tag{298}$$

Hence the total kinetic energy of the system  $T$  becomes

$$\begin{aligned}
T &= T_p + T_c \\
&= \frac{1}{2} \sum_{j=1}^N \left[ m_j \dot{\mathbf{r}}_j^2 - m_j \left( \frac{2 \sum_{i=1}^N m_i \dot{\mathbf{r}}_i}{M + \sum_{i=1}^N m_i} \right) \cdot \dot{\mathbf{r}}_j \right. \\
&\quad \left. + \frac{\left( \sum_{i=1}^N m_i \dot{\mathbf{r}}_i \right)^2 + \left( \sum_{i=1}^N m_i \dot{\mathbf{y}}_i \right)^2 + \left( \sum_{i=1}^N m_i \dot{\mathbf{z}}_i \right)^2}{\left( M + \sum_{i=1}^N m_i \right)^2} m_j \right] \\
&\quad + \frac{M}{2} \frac{\left( \sum_{i=1}^N m_i \dot{\mathbf{r}}_i \right)^2 + \left( \sum_{i=1}^N m_i \dot{\mathbf{y}}_i \right)^2 + \left( \sum_{i=1}^N m_i \dot{\mathbf{z}}_i \right)^2}{\left( M + \sum_{i=1}^N m_i \right)^2} \\
&= \frac{1}{2} \sum_{j=1}^N m_j \dot{\mathbf{r}}_j^2 - \frac{1}{2} \sum_{j=1}^N m_j \frac{2 \sum_{i=1}^N m_i \dot{\mathbf{r}}_i}{M + \sum_{i=1}^N m_i} \cdot \dot{\mathbf{r}}_j \\
&\quad + \frac{M + \sum_{j=1}^N m_j}{\left( M + \sum_{i=1}^N m_i \right)^2} \cdot \frac{1}{2} \left\{ \left( \sum_{i=1}^N m_i \dot{\mathbf{r}}_i \right)^2 + \left( \sum_{i=1}^N m_i \dot{\mathbf{y}}_i \right)^2 + \left( \sum_{i=1}^N m_i \dot{\mathbf{z}}_i \right)^2 \right\}. \tag{299}
\end{aligned}$$

The second term in the right-hand side of the equation (299) becomes

$$\begin{aligned}
&-\frac{1}{2} \sum_{j=1}^N m_j \underbrace{\left( \frac{2 \sum_{i=1}^N m_i \dot{\mathbf{r}}_i}{M + \sum_{i=1}^N m_i} \right)}_{\text{independent of } j} \cdot \dot{\mathbf{r}}_j = -\frac{1}{2} \left( \frac{2 \sum_{i=1}^N m_i \dot{\mathbf{r}}_i}{M + \sum_{i=1}^N m_i} \right) \sum_{j=1}^N m_j \dot{\mathbf{r}}_j \\
&= -\frac{1}{M + \sum_{i=1}^N m_i} \left( \sum_{i=1}^N m_i \dot{\mathbf{r}}_i \right) \cdot \left( \sum_{j=1}^N m_j \dot{\mathbf{r}}_j \right) \\
&= -\frac{1}{M + \sum_{i=1}^N m_i} \left[ \left( \sum_{i=1}^N m_i \dot{\mathbf{r}}_i \right)^2 + \left( \sum_{i=1}^N m_i \dot{\mathbf{y}}_i \right)^2 + \left( \sum_{i=1}^N m_i \dot{\mathbf{z}}_i \right)^2 \right], \tag{300}
\end{aligned}$$

which ends up with the final form of the total kinetic energy  $T$  as

$$\begin{aligned}
T &= T_p + T_c \\
&= \frac{1}{2} \sum_{j=1}^N m_j \dot{r}_j^2 - \frac{\left(\sum_{i=1}^N m_i \dot{x}_i\right)^2 + \left(\sum_{i=1}^N m_i \dot{y}_i\right)^2 + \left(\sum_{i=1}^N m_i \dot{z}_i\right)^2}{M + \sum_{i=1}^N m_i} \\
&\quad + \frac{1}{2} \frac{\left(\sum_{i=1}^N m_i \dot{x}_i\right)^2 + \left(\sum_{i=1}^N m_i \dot{y}_i\right)^2 + \left(\sum_{i=1}^N m_i \dot{z}_i\right)^2}{M + \sum_{i=1}^N m_i} \\
&= \frac{1}{2} \sum_{j=1}^N m_j \dot{r}_j^2 - \frac{1}{2} \frac{\left(\sum_{i=1}^N m_i \dot{x}_i\right)^2 + \left(\sum_{i=1}^N m_i \dot{y}_i\right)^2 + \left(\sum_{i=1}^N m_i \dot{z}_i\right)^2}{M + \sum_{i=1}^N m_i}. \tag{301}
\end{aligned}$$

Thus we can express the total kinetic energy  $T$  only by the heliocentric velocities,  $\dot{\mathbf{r}}_i$ .

If we consider  $\mathbf{r}$  as a canonical coordinate, the canonical conjugate momenta  $\mathbf{p}$  to  $\mathbf{r}$  are derived from Lagrangian

$$L(\mathbf{r}, \dot{\mathbf{r}}) = T(\dot{\mathbf{r}}) - V(\mathbf{r}), \tag{302}$$

where  $T(\dot{\mathbf{r}})$  is kinetic energy and  $V(\mathbf{r})$  is potential energy.

If we define a temporary coordinate  $\mathbf{r}^*$  as

$$\mathbf{r}^* \equiv \frac{1}{M + \sum_{i=1}^N m_i} \sum_{j=1}^N m_j \mathbf{r}_j, \tag{303}$$

then the kinetic energy  $T(\dot{\mathbf{r}})$  is expressed as

$$T(\mathbf{r}) = \frac{1}{2} \sum_{j=1}^N m_j \mathbf{r}_j^2 - \frac{M + \sum_{i=1}^N m_i}{2} \dot{\mathbf{r}}^{*2}. \tag{304}$$

Hence the canonical momenta  $\mathbf{p}$  become

$$\begin{aligned}
\mathbf{p}_i &\equiv \frac{\partial L}{\partial \dot{\mathbf{r}}_i} \\
&= \frac{\partial}{\partial \dot{\mathbf{r}}_i} (T(\dot{\mathbf{r}}) - V(\mathbf{r}_i)) \\
&= \frac{\partial T}{\partial \dot{\mathbf{r}}_i} \\
&= m_i \dot{\mathbf{r}}_i - \frac{M + \sum_{j=1}^N m_j}{2} 2\dot{\mathbf{r}}^* \cdot \frac{\partial}{\partial \dot{\mathbf{r}}_i} \left( \frac{1}{M + \sum_{j=1}^N m_j} \sum_{j=1}^N m_j \dot{\mathbf{r}}_j \right) \\
&= m_i \dot{\mathbf{r}}_i - \left( M + \sum_{j=1}^N m_j \right) \dot{\mathbf{r}}^* \cdot \frac{1}{M + \sum_{j=1}^N m_j} m_i \\
&= m_i (\dot{\mathbf{r}}_i - \dot{\mathbf{r}}^*) \\
&= m_i \left( \dot{\mathbf{r}}_i - \frac{1}{M + \sum_{j=1}^N m_j} \sum_{j=1}^N m_j \dot{\mathbf{r}}_j \right). \tag{305}
\end{aligned}$$

The canonical momenta  $\mathbf{p}_i$  in (305) are equivalent to those of what are used in so-called “democratic heliocentric” or “mixed-center” coordinate, which are

$$\mathbf{P}_i \equiv \mathbf{p}_i^{\text{inert}} - \frac{m_i}{M_{\text{total}}} \sum_{j=0}^N \mathbf{p}_j^{\text{inert}}, \tag{306}$$

$$M_{\text{total}} \equiv M + \sum_{i=1}^N m_i, \quad (307)$$

where  $\mathbf{p}_i^{\text{inert}}$  denotes momenta reckoned from a certain fixed point in the inertial frame (for example, barycenter, i.e. essentially the same as  $m_i \boldsymbol{\rho}_i$  in (295)). We now show that the momenta in (305) and momenta in (306) are rigorously equivalent.

Since  $\mathbf{r}$  are the heliocentric coordinate,

$$\dot{\mathbf{r}}_i = \dot{\mathbf{q}}_i - \dot{\mathbf{q}}_0 = \frac{\mathbf{p}_i^{\text{inert}}}{m_i} - \frac{\mathbf{p}_0^{\text{inert}}}{m_0}, \quad (308)$$

where we define  $\mathbf{q}_i$  as the coordinate reckoned from certain fixed point in the inertial frame (this is equivalent to  $\boldsymbol{\rho}_i$  in (295)). Then,  $\mathbf{r}^*$  becomes by (303) as

$$\dot{\mathbf{r}}^* = \frac{1}{M_{\text{total}}} \sum_{j=1}^N m_j \dot{\mathbf{r}}_j = \frac{1}{M_{\text{total}}} \sum_{j=1}^N m_j \left( \frac{\mathbf{p}_j^{\text{inert}}}{m_j} - \frac{\mathbf{p}_0^{\text{inert}}}{m_0} \right). \quad (309)$$

This leads to the expression of  $\mathbf{p}_i$  in (305) as

$$\begin{aligned} \mathbf{p}_i &= m_i \dot{\mathbf{r}}_i - m_i \dot{\mathbf{r}}^* \\ &= m_i \left( \frac{\mathbf{p}_i^{\text{inert}}}{m_i} - \frac{\mathbf{p}_0^{\text{inert}}}{m_0} \right) - m_i \frac{1}{M_{\text{total}}} \sum_{j=1}^N m_j \left( \frac{\mathbf{p}_j^{\text{inert}}}{m_j} - \frac{\mathbf{p}_0^{\text{inert}}}{m_0} \right) \\ &= \mathbf{p}_i^{\text{inert}} - \frac{m_i}{m_0} \mathbf{p}_0^{\text{inert}} - \frac{m_i}{M_{\text{total}}} \sum_{j=1}^N m_j \left( \frac{\mathbf{p}_j^{\text{inert}}}{m_j} - \frac{m_j}{m_0} \mathbf{p}_0^{\text{inert}} \right) \\ &= \mathbf{p}_i^{\text{inert}} - \frac{m_i}{M_{\text{total}}} \sum_{j=1}^N \mathbf{p}_j^{\text{inert}} - \frac{m_i}{m_0} \mathbf{p}_0^{\text{inert}} + \frac{m_i}{M_{\text{total}}} \sum_{j=1}^N \frac{m_j}{m_0} \mathbf{p}_0^{\text{inert}} \\ &= \mathbf{p}_i^{\text{inert}} - \frac{m_i}{M_{\text{total}}} \sum_{j=1}^N \mathbf{p}_j^{\text{inert}} + \left( -\frac{m_i}{m_0} + \frac{m_i}{M_{\text{total}}} \sum_{j=1}^N \frac{m_j}{m_0} \right) \mathbf{p}_0^{\text{inert}} \\ &= \mathbf{p}_i^{\text{inert}} - \frac{m_i}{M_{\text{total}}} \sum_{j=1}^N \mathbf{p}_j^{\text{inert}} + \left[ -\frac{m_i}{m_0} \left\{ 1 - \frac{1}{M_{\text{total}}} (M_{\text{total}} - m_0) \right\} \right] \mathbf{p}_0^{\text{inert}} \\ &= \mathbf{p}_i^{\text{inert}} - \frac{m_i}{M_{\text{total}}} \sum_{j=1}^N \mathbf{p}_j^{\text{inert}} - \frac{m_i}{m_0} \left( 1 - \left( 1 - \frac{m_0}{M_{\text{total}}} \right) \right) \mathbf{p}_0^{\text{inert}} \\ &= \mathbf{p}_i^{\text{inert}} - \frac{m_i}{M_{\text{total}}} \sum_{j=1}^N \mathbf{p}_j^{\text{inert}} - \frac{m_i}{M_{\text{total}}} \mathbf{p}_0^{\text{inert}} \\ &= \mathbf{p}_i^{\text{inert}} - \frac{m_i}{M_{\text{total}}} \sum_{j=0}^N \mathbf{p}_j^{\text{inert}}, \end{aligned} \quad (310)$$

which is rigorously equivalent to the momenta used in “democratic heliocentric” (DH) coordinate,  $\mathbf{P}_i$  in (306). The idea of the democratic heliocentric was first advocated by Poincaré and have been described in detail in Charlier (1902) as a name of “*Canonische relative Coordinaten*.” Later, this coordinate was born again to be used in SyMBA in Duncan et al. (1998) as follows.



Below,  $Q_i$  and  $P_i$  are canonically conjugate each other.

$$Q_i = \begin{cases} q_i - q_0, & (i = 1, \dots, N) \\ \frac{1}{M_{\text{tot}}} \sum_{j=0}^N m_j q_j, & (i = 0) \end{cases} \quad (311)$$

$$P_i = \begin{cases} p_i - \frac{m_i}{M_{\text{tot}}} \sum_{j=0}^N p_j, & (i = 1, \dots, N) \\ \sum_{j=0}^N p_j, & (i = 0) \end{cases} \quad (312)$$

The Hamiltonian described by the DH coordinate is as follows:

$$H = H_{\text{kep}} + H_{\text{sun}} + H_{\text{int}}, \quad (313)$$

$$H_{\text{kep}} = \sum_{i=1}^N \left( \frac{|P_i|^2}{2m_i} - \frac{Gm_i m_0}{|Q_i|} \right), \quad (314)$$

$$H_{\text{sun}} = \frac{1}{2m_0} \left| \sum_{i=1}^N P_i \right|^2, \quad (315)$$

$$H_{\text{int}} = - \sum_{i=1}^{N-1} \sum_{j=i+1}^N \frac{Gm_i m_j}{|Q_i - Q_j|}. \quad (316)$$

Note that the Hamiltonian is now divided into three parts, not into two parts as in the generic WH map in (141). This is because of the existence of the Sun's kinetic energy  $H_{\text{sun}}$ . The amount of computation does not increase significantly by the additional procedures due to this new division, since the procedure which involves  $H_{\text{sun}}$  is summing up of linear terms of  $P$ .

## Appendix C. Partial derivatives of Kepler orbital elements

The partial derivatives of true anomaly  $f$  by the Delaunay variables  $L$  and  $G$  end up with

$$\frac{\partial f}{\partial L} = \frac{G^2}{eL^3} \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f, \quad (317)$$

$$\frac{\partial f}{\partial G} = -\frac{G}{eL^2} \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f, \quad (318)$$

from the relationship

$$\frac{\partial f}{\partial e} = \frac{L^2}{G^2} \sin f. \quad (319)$$

However, we are likely to obtain wrong solutions such as

$$\frac{\partial f}{\partial L} = \frac{G^2}{eL^3} \left( \frac{2a}{r} + \frac{L^2}{G^2} \right) \sin f, \quad (320)$$

or

$$\frac{\partial f}{\partial G} = -\frac{G}{eL^2} \left( \frac{2a}{r} + \frac{L^2}{G^2} \right) \sin f, \quad (321)$$

instead of the correct (317) or (318) due to the usual expressions of  $\frac{\partial f}{\partial e}$  in some of the standard celestial mechanics textbooks (e.g. p. 567 in Brouwer and Clemence (1961), p. 349 in Nagasawa (1983), and Eq. (2.104) at p. 39 in Kinoshita (1998)) as

$$\frac{\partial f}{\partial e} = \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f. \quad (322)$$

To avoid such confusion, we must derive  $\frac{\partial f}{\partial L}$  and  $\frac{\partial f}{\partial G}$  through the definition of differential transformation: using Jacobian matrices.  $\frac{\partial f}{\partial L}$  and  $\frac{\partial f}{\partial G}$  are obtained through a differential transformation of variables from Kepler orbital elements to Delaunay elements as

$$(da, de, d\omega, dI, d\Omega, df) \rightarrow (dL, dG, dH, dl, dg, dh) \quad (323)$$

However, the relationship between  $f$  and  $l$  is not explicit because of the existence of Kepler's equation  $u - e \sin u = l$ . We have to calculate the following three conversion matrices

$$(da, de, d\omega, dI, d\Omega, df) \rightarrow (da, de, d\omega, dI, d\Omega, du), \quad (324)$$

$$(da, de, d\omega, dI, d\Omega, du) \rightarrow (da, de, d\omega, dI, d\Omega, dl), \quad (325)$$

$$(da, de, d\omega, dI, d\Omega, dl) \rightarrow (dL, dG, dH, dl, dg, dh), \quad (326)$$

and multiply them to reach our final goal, (323).

As a set of independent variables to describe the Kepler orbital motion, we consider the following four sets:

- Using mean anomaly  $l$  as  $(a, e, \omega, I, \Omega, l)$
- Using true anomaly  $f$  as  $(a, e, \omega, I, \Omega, f)$
- Using eccentric anomaly  $u$  as  $(a, e, \omega, I, \Omega, u)$
- Delaunay canonical variables  $(L, G, H, l, g, h)$

Hence there are  $4P_2 = 12$  differential transformations among them:

$$\begin{aligned} (da, de, d\omega, dI, d\Omega, df) &\rightarrow (da, de, d\omega, dI, d\Omega, du) \\ (da, de, d\omega, dI, d\Omega, df) &\rightarrow (da, de, d\omega, dI, d\Omega, dl) \\ (da, de, d\omega, dI, d\Omega, df) &\rightarrow (dL, dG, dH, dl, dg, dh) \\ (da, de, d\omega, dI, d\Omega, du) &\rightarrow (da, de, d\omega, dI, d\Omega, dl) \\ (da, de, d\omega, dI, d\Omega, du) &\rightarrow (da, de, d\omega, dI, d\Omega, df) \\ (da, de, d\omega, dI, d\Omega, du) &\rightarrow (dL, dG, dH, dl, dg, dh) \\ (da, de, d\omega, dI, d\Omega, dl) &\rightarrow (da, de, d\omega, dI, d\Omega, df) \\ (da, de, d\omega, dI, d\Omega, dl) &\rightarrow (da, de, d\omega, dI, d\Omega, du) \\ (da, de, d\omega, dI, d\Omega, dl) &\rightarrow (dL, dG, dH, dl, dg, dh) \\ (dL, dG, dH, dl, dg, dh) &\rightarrow (da, de, d\omega, dI, d\Omega, dl) \\ (dL, dG, dH, dl, dg, dh) &\rightarrow (da, de, d\omega, dI, d\Omega, du) \\ (dL, dG, dH, dl, dg, dh) &\rightarrow (da, de, d\omega, dI, d\Omega, df) \end{aligned}$$

This section gives the forward transformations such as

$$\begin{aligned}
& (da, de, d\omega, dI, d\Omega, df) \\
& \quad \downarrow \\
& (da, de, d\omega, dI, d\Omega, du) \\
& \quad \downarrow \\
& (da, de, d\omega, dI, d\Omega, dl) \\
& \quad \downarrow \\
& (dL, dG, dH, dl, dg, dh).
\end{aligned}$$

The latter half can be done similarly, and omitted here.

Note that the lines and columns in the Jacobian matrices described here may be transposed from those in usual textbooks. In the following discussion,

$$r = a(1 - e \cos u) = \frac{a(1 - e^2)}{1 + e \cos f}, \quad (327)$$

$$\eta \equiv \sqrt{1 - e^2}. \quad (328)$$

We frequently consult the relationship between eccentric anomaly  $u$  and true anomaly  $f$

$$\sin u = \frac{\eta \sin f}{1 + e \cos f}, \quad \cos u = \frac{e + \cos f}{1 + e \cos f}, \quad (329)$$

$$\sin f = \frac{\eta \sin u}{1 - e \cos u}, \quad \cos f = \frac{\cos u - e}{1 - e \cos u}. \quad (330)$$

#### Appendix C. 1 $(da, de, d\omega, dI, d\Omega, df) \rightarrow (da, de, d\omega, dI, d\Omega, du)$

Since the Kepler orbital elements  $a, e, \omega, I, \Omega$  are independent with each other,

$$\begin{aligned}
\frac{\partial a}{\partial e} &= \frac{\partial a}{\partial \omega} = \frac{\partial a}{\partial I} = \frac{\partial a}{\partial \Omega} = 0, & \frac{\partial a}{\partial a} &= 1, \\
\frac{\partial e}{\partial a} &= \frac{\partial e}{\partial \omega} = \frac{\partial e}{\partial I} = \frac{\partial e}{\partial \Omega} = 0, & \frac{\partial e}{\partial e} &= 1, \\
\frac{\partial \omega}{\partial a} &= \frac{\partial \omega}{\partial e} = \frac{\partial \omega}{\partial I} = \frac{\partial \omega}{\partial \Omega} = 0, & \frac{\partial \omega}{\partial \omega} &= 1, \\
\frac{\partial I}{\partial a} &= \frac{\partial I}{\partial e} = \frac{\partial I}{\partial \omega} = \frac{\partial I}{\partial \Omega} = 0, & \frac{\partial I}{\partial I} &= 1, \\
\frac{\partial \Omega}{\partial a} &= \frac{\partial \Omega}{\partial e} = \frac{\partial \Omega}{\partial I} = \frac{\partial \Omega}{\partial \omega} = 0, & \frac{\partial \Omega}{\partial \Omega} &= 1.
\end{aligned} \quad (331)$$

Hence the differential transformation matrix in this case becomes

$$\begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ df \end{pmatrix} = \frac{\partial(a, e, \omega, I, \Omega, f)}{\partial(a, e, \omega, I, \Omega, u)} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ du \end{pmatrix}$$

$$\begin{aligned}
&= \begin{pmatrix} \frac{\partial a}{\partial e} & \frac{\partial a}{\partial e} & \frac{\partial a}{\partial \omega} & \frac{\partial a}{\partial I} & \frac{\partial a}{\partial \Omega} & \frac{\partial a}{\partial u} \\ \frac{\partial e}{\partial a} & \frac{\partial e}{\partial e} & \frac{\partial e}{\partial \omega} & \frac{\partial e}{\partial I} & \frac{\partial e}{\partial \Omega} & \frac{\partial e}{\partial u} \\ \frac{\partial \omega}{\partial a} & \frac{\partial \omega}{\partial e} & \frac{\partial \omega}{\partial \omega} & \frac{\partial \omega}{\partial I} & \frac{\partial \omega}{\partial \Omega} & \frac{\partial \omega}{\partial u} \\ \frac{\partial I}{\partial a} & \frac{\partial I}{\partial e} & \frac{\partial I}{\partial \omega} & \frac{\partial I}{\partial I} & \frac{\partial I}{\partial \Omega} & \frac{\partial I}{\partial u} \\ \frac{\partial \Omega}{\partial a} & \frac{\partial \Omega}{\partial e} & \frac{\partial \Omega}{\partial \omega} & \frac{\partial \Omega}{\partial I} & \frac{\partial \Omega}{\partial \Omega} & \frac{\partial \Omega}{\partial u} \\ \frac{\partial f}{\partial a} & \frac{\partial f}{\partial e} & \frac{\partial f}{\partial \omega} & \frac{\partial f}{\partial I} & \frac{\partial f}{\partial \Omega} & \frac{\partial f}{\partial u} \end{pmatrix} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ du \end{pmatrix} \\
&= \begin{pmatrix} 1 & & & 0 & 0 \\ & 1 & & & 0 \\ & & 1 & & 0 \\ & & & 1 & 0 \\ 0 & & & & 1 & 0 \\ 0 & \frac{\partial f}{\partial e} & 0 & 0 & 0 & \frac{\partial f}{\partial u} \end{pmatrix} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ du \end{pmatrix}. \tag{332}
\end{aligned}$$

Only  $\frac{\partial f}{\partial e}$  and  $\frac{\partial f}{\partial u}$  should be taken into account in all the components of (332).  $\frac{\partial f}{\partial e}$  can be obtained as

$$\begin{aligned}
\frac{\partial}{\partial e} \cos f &= -\sin f \frac{\partial f}{\partial e} \\
&= \frac{\partial}{\partial e} \left( (\cos u - e)(1 - e \cos u)^{-1} \right) \\
&= -1 \cdot (1 - e \cos u)^{-1} + (-e + \cos u) \cdot (-1) \cdot (1 - e \cos u)^{-2} (-\cos u) \\
&= \frac{-1}{1 - e \cos u} + \frac{-e + \cos u}{(1 - e \cos u)^2} \cos u \\
&= \frac{1}{(1 - e \cos u)^2} [-1 + e \cos u + (-e + \cos u) \cos u] \\
&= \frac{1}{(1 - e \cos u)^2} (\cos^2 u - 1) \\
&= \frac{-\sin^2 u}{(1 - e \cos u)^2} \\
&= -\left( \frac{\sin u}{1 - e \cos u} \right)^2 \\
&= -\left( \frac{\sin f}{\eta} \right)^2 \quad (\because (329)) \tag{333}
\end{aligned}$$

$$\therefore \frac{\partial f}{\partial e} = \left( \frac{1}{-\sin f} \right) \left( -\frac{\sin^2 f}{\eta^2} \right) = \frac{1}{\eta^2} \sin f = \frac{L^2}{G^2} \sin f \tag{334}$$

Similarly,  $\frac{\partial f}{\partial u}$  is obtained as

$$\begin{aligned}
\frac{\partial}{\partial u} \cos f &= -\sin f \frac{\partial f}{\partial u} \\
&= \frac{\partial}{\partial u} \left( (\cos u - e)(1 - e \cos u)^{-1} \right) \\
&= -\sin u (1 - e \cos u)^{-1} + (-e + \cos u) \cdot (-1) \cdot (1 - e \cos u)^{-2} e \sin u \\
&= \frac{-\sin u}{1 - e \cos u} - \frac{-e + \cos u}{(1 - e \cos u)^2} e \sin u
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{(1 - e \cos u)^2} [-\sin u(1 - e \cos u) - e \sin u(-e + \cos u)] \\
&= \frac{\sin u}{(1 - e \cos u)^2} (1 - e^2) \\
&= \frac{\eta^2 \sin^2 u}{(1 - e \cos u)^2} \left( -\frac{1}{\sin u} \right) \\
&= -\frac{\sin^2 f}{\sin u} \quad (\because (329)) \tag{335}
\end{aligned}$$

$$\therefore \frac{\partial f}{\partial u} = \left( -\frac{1}{\sin f} \right) \left( -\frac{\sin^2 f}{\sin u} \right) = \frac{\eta \sin u}{1 - e \cos u} \frac{1}{\sin u} = \frac{a\eta}{r}. \quad (\because (327)) \tag{336}$$

Therefore (332) becomes

$$\begin{aligned}
\begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ df \end{pmatrix} &= \frac{\partial(a, e, \omega, I, \Omega, f)}{\partial(a, e, \omega, I, \Omega, u)} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ du \end{pmatrix} \\
&= \begin{pmatrix} 1 & & & & 0 & 0 \\ & 1 & & & & 0 \\ & & 1 & & & 0 \\ & & & 1 & & 0 \\ 0 & & & & 1 & 0 \\ 0 & \frac{\partial f}{\partial e} & 0 & 0 & 0 & \frac{\partial f}{\partial u} \end{pmatrix} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ du \end{pmatrix} \\
&= \begin{pmatrix} 1 & & & & 0 & 0 \\ & 1 & & & & 0 \\ & & 1 & & & 0 \\ & & & 1 & & 0 \\ 0 & & & & 1 & 0 \\ 0 & \frac{L^2}{G^2} \sin f & 0 & 0 & 0 & \frac{a\eta}{r} \end{pmatrix} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ du \end{pmatrix}. \tag{337}
\end{aligned}$$

Note that at this point

$$\frac{\partial f}{\partial e} = \frac{L^2}{G^2} \sin f, \tag{338}$$

in the final result (337). This expression is due that we consider the true anomaly  $f$  as a function of  $(e, u)$ , not  $(e, l)$  as in (322).

## Appendix C. 2 $(da, de, d\omega, dI, d\Omega, du) \rightarrow (da, de, d\omega, dI, d\Omega, dl)$

When we consider the independence of the mean anomaly  $l$  from any other Kepler orbital elements as

$$\frac{\partial l}{\partial a} = \frac{\partial l}{\partial e} = \frac{\partial l}{\partial \omega} = \frac{\partial l}{\partial I} = \frac{\partial l}{\partial \Omega} = 0, \quad \frac{\partial l}{\partial l} = 1, \tag{339}$$

the differential transformation matrix in this case becomes

$$\begin{aligned}
 \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ du \end{pmatrix} &= \frac{\partial(a, e, \omega, I, \Omega, u)}{\partial(a, e, \omega, I, \Omega, l)} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ dl \end{pmatrix} \\
 &= \begin{pmatrix} \frac{\partial a}{\partial a} & \frac{\partial a}{\partial e} & \frac{\partial a}{\partial \omega} & \frac{\partial a}{\partial I} & \frac{\partial a}{\partial \Omega} & \frac{\partial a}{\partial l} \\ \frac{\partial e}{\partial a} & \frac{\partial e}{\partial e} & \frac{\partial e}{\partial \omega} & \frac{\partial e}{\partial I} & \frac{\partial e}{\partial \Omega} & \frac{\partial e}{\partial l} \\ \frac{\partial \omega}{\partial a} & \frac{\partial \omega}{\partial e} & \frac{\partial \omega}{\partial \omega} & \frac{\partial \omega}{\partial I} & \frac{\partial \omega}{\partial \Omega} & \frac{\partial \omega}{\partial l} \\ \frac{\partial I}{\partial a} & \frac{\partial I}{\partial e} & \frac{\partial I}{\partial \omega} & \frac{\partial I}{\partial I} & \frac{\partial I}{\partial \Omega} & \frac{\partial I}{\partial l} \\ \frac{\partial \Omega}{\partial a} & \frac{\partial \Omega}{\partial e} & \frac{\partial \Omega}{\partial \omega} & \frac{\partial \Omega}{\partial I} & \frac{\partial \Omega}{\partial \Omega} & \frac{\partial \Omega}{\partial l} \\ \frac{\partial u}{\partial a} & \frac{\partial u}{\partial e} & \frac{\partial u}{\partial \omega} & \frac{\partial u}{\partial I} & \frac{\partial u}{\partial \Omega} & \frac{\partial u}{\partial l} \end{pmatrix} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ dl \end{pmatrix} \\
 &= \begin{pmatrix} 1 & & & 0 & 0 \\ & 1 & & & 0 \\ & & 1 & & 0 \\ & & & 1 & 0 \\ 0 & & & & 1 & 0 \\ 0 & \frac{\partial u}{\partial e} & 0 & 0 & 0 & \frac{\partial u}{\partial l} \end{pmatrix} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ dl \end{pmatrix}. \tag{340}
 \end{aligned}$$

Only  $\frac{\partial u}{\partial e}$  and  $\frac{\partial u}{\partial l}$  should be taken into account in all the components of (340).  $\frac{\partial u}{\partial e}$  can be obtained from the partial derivative of the Kepler's equation

$$u - e \sin u = l, \tag{341}$$

by  $e$

$$\frac{\partial u}{\partial e} - \left( \sin u + e \cos u \frac{\partial u}{\partial e} \right) = \frac{\partial l}{\partial e} = 0. \tag{342}$$

$$\therefore (1 - e \cos u) \frac{\partial u}{\partial e} = \sin u \tag{343}$$

$$\therefore \frac{\partial u}{\partial e} = \frac{\sin u}{1 - e \cos u} = \frac{\sin f}{\eta}. \quad (\because (330)) \tag{344}$$

Similarly,  $\frac{\partial u}{\partial l}$  can be obtained from the partial derivative of the Kepler's equation by  $l$  as

$$\frac{\partial u}{\partial l} - e \cos \frac{\partial u}{\partial l} = \frac{\partial l}{\partial l} = 1, \tag{345}$$

$$\therefore (1 - e \cos u) \frac{\partial u}{\partial l} = 1, \tag{346}$$

$$\therefore \frac{\partial u}{\partial l} = \frac{1}{1 - e \cos u} = \frac{a}{r}. \tag{347}$$

Therefore the final form of the transformation matrix (340) becomes

$$\begin{aligned}
 \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ du \end{pmatrix} &= \frac{\partial(a, e, \omega, I, \Omega, u)}{\partial(a, e, \omega, I, \Omega, l)} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ dl \end{pmatrix} \\
 &= \begin{pmatrix} 1 & & & & 0 & 0 \\ & 1 & & & & 0 \\ & & 1 & & & 0 \\ & & & 1 & & 0 \\ 0 & & & & 1 & 0 \\ 0 & \frac{\partial u}{\partial e} & 0 & 0 & 0 & \frac{\partial u}{\partial I} \end{pmatrix} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ dl \end{pmatrix} \\
 &= \begin{pmatrix} 1 & & & & 0 & 0 \\ & 1 & & & & 0 \\ & & 1 & & & 0 \\ & & & 1 & & 0 \\ 0 & & & & 1 & 0 \\ 0 & \frac{\sin f}{\eta} & 0 & 0 & 0 & \frac{a}{r} \end{pmatrix} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ dl \end{pmatrix}. \tag{348}
 \end{aligned}$$

### Appendix C. 3 $(da, de, d\omega, dI, d\Omega, df) \rightarrow (da, de, d\omega, dI, d\Omega, dl)$

This transformation matrix is obtained as a product of the two matrices

$$\frac{\partial(a, e, \omega, I, \Omega, f)}{\partial(a, e, \omega, I, \Omega, u)}$$

and

$$\frac{\partial(a, e, \omega, I, \Omega, u)}{\partial(a, e, \omega, I, \Omega, l)}$$

as

$$\begin{aligned}
 \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ df \end{pmatrix} &= \frac{\partial(a, e, \omega, I, \Omega, f)}{\partial(a, e, \omega, I, \Omega, l)} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ dl \end{pmatrix} = \begin{pmatrix} \frac{\partial a}{\partial a} & \frac{\partial a}{\partial e} & \frac{\partial a}{\partial \omega} & \frac{\partial a}{\partial I} & \frac{\partial a}{\partial \Omega} & \frac{\partial a}{\partial f} \\ \frac{\partial e}{\partial a} & \frac{\partial e}{\partial e} & \frac{\partial e}{\partial \omega} & \frac{\partial e}{\partial I} & \frac{\partial e}{\partial \Omega} & \frac{\partial e}{\partial f} \\ \frac{\partial \omega}{\partial a} & \frac{\partial \omega}{\partial e} & \frac{\partial \omega}{\partial \omega} & \frac{\partial \omega}{\partial I} & \frac{\partial \omega}{\partial \Omega} & \frac{\partial \omega}{\partial f} \\ \frac{\partial I}{\partial a} & \frac{\partial I}{\partial e} & \frac{\partial I}{\partial \omega} & \frac{\partial I}{\partial I} & \frac{\partial I}{\partial \Omega} & \frac{\partial I}{\partial f} \\ \frac{\partial \Omega}{\partial a} & \frac{\partial \Omega}{\partial e} & \frac{\partial \Omega}{\partial \omega} & \frac{\partial \Omega}{\partial I} & \frac{\partial \Omega}{\partial \Omega} & \frac{\partial \Omega}{\partial f} \\ \frac{\partial f}{\partial a} & \frac{\partial f}{\partial e} & \frac{\partial f}{\partial \omega} & \frac{\partial f}{\partial I} & \frac{\partial f}{\partial \Omega} & \frac{\partial f}{\partial f} \end{pmatrix} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ dl \end{pmatrix} \\
 &= \frac{\partial(a, e, \omega, I, \Omega, f)}{\partial(a, e, \omega, I, \Omega, u)} \frac{\partial(a, e, \omega, I, \Omega, u)}{\partial(a, e, \omega, I, \Omega, l)} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ dl \end{pmatrix}
 \end{aligned}$$

$$\begin{aligned}
&= \begin{pmatrix} 1 & & & 0 & 0 \\ & 1 & & & 0 \\ & & 1 & & 0 \\ & & & 1 & 0 \\ 0 & \frac{\partial f}{\partial e} & 0 & 0 & 0 \\ 0 & \frac{\partial f}{\partial u} & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 & & & 0 & 0 \\ & 1 & & & 0 \\ & & 1 & & 0 \\ & & & 1 & 0 \\ 0 & \frac{\partial u}{\partial e} & 0 & 0 & 0 \\ 0 & \frac{\partial u}{\partial l} & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ dl \end{pmatrix} \\
&= \begin{pmatrix} 1 & & & 0 & 0 \\ & 1 & & & 0 \\ & & 1 & & 0 \\ & & & 1 & 0 \\ 0 & \frac{L^2}{G^2} \sin f & 0 & 0 & 0 \\ 0 & \frac{a\eta}{r} & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 & & & 0 & 0 \\ & 1 & & & 0 \\ & & 1 & & 0 \\ & & & 1 & 0 \\ 0 & \frac{\sin f}{\eta} & 0 & 0 & 0 \\ 0 & \frac{a}{r} & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ dl \end{pmatrix} \\
&= \begin{pmatrix} 1 & & & 0 & 0 \\ & 1 & & & 0 \\ & & 1 & & 0 \\ & & & 1 & 0 \\ 0 & \frac{L^2}{G^2} \sin f + \frac{a\eta}{r} \frac{\sin f}{\eta} & 0 & 0 & 0 \\ 0 & \frac{a\eta}{r} \frac{a}{r} & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ dl \end{pmatrix} \\
&= \begin{pmatrix} 1 & & & 0 & 0 \\ & 1 & & & 0 \\ & & 1 & & 0 \\ & & & 1 & 0 \\ 0 & \left(\frac{a}{r} + \frac{L^2}{G^2}\right) \sin f & 0 & 0 & 0 \\ 0 & \frac{a^2 \eta}{r^2} & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ dl \end{pmatrix}. \tag{349}
\end{aligned}$$

Note that at this point

$$\frac{\partial f}{\partial e} = \left( \frac{a}{r} + \frac{L^2}{G^2} \right) \sin f, \tag{350}$$

in the final result (349). This expression is due that we consider the true anomaly  $f$  as a function of  $(e, l)$ , which is different from the results in (337) where  $f$  is a function of  $(e, u)$ .

#### Appendix C. 4 $(da, de, d\omega, dI, d\Omega, dl) \rightarrow (dL, dG, dH, dl, dg, dh)$

This transformation matrix contains the differential transformation from Kepler orbital elements to Delaunay canonical variables as

$$\begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ dl \end{pmatrix} = \frac{\partial(a, e, \omega, I, \Omega, l)}{\partial(L, G, H, l, g, h)} \begin{pmatrix} dL \\ dG \\ dH \\ dl \\ dg \\ dh \end{pmatrix} \tag{351}$$



$$= \begin{pmatrix} \frac{\partial a}{\partial L} & \frac{\partial a}{\partial G} & \frac{\partial a}{\partial H} & \frac{\partial a}{\partial l} & \frac{\partial a}{\partial g} & \frac{\partial a}{\partial h} \\ \frac{\partial e}{\partial L} & \frac{\partial e}{\partial G} & \frac{\partial e}{\partial H} & \frac{\partial e}{\partial l} & \frac{\partial e}{\partial g} & \frac{\partial e}{\partial h} \\ \frac{\partial \omega}{\partial L} & \frac{\partial \omega}{\partial G} & \frac{\partial \omega}{\partial H} & \frac{\partial \omega}{\partial l} & \frac{\partial \omega}{\partial g} & \frac{\partial \omega}{\partial h} \\ \frac{\partial I}{\partial L} & \frac{\partial I}{\partial G} & \frac{\partial I}{\partial H} & \frac{\partial I}{\partial l} & \frac{\partial I}{\partial g} & \frac{\partial I}{\partial h} \\ \frac{\partial \Omega}{\partial L} & \frac{\partial \Omega}{\partial G} & \frac{\partial \Omega}{\partial H} & \frac{\partial \Omega}{\partial l} & \frac{\partial \Omega}{\partial g} & \frac{\partial \Omega}{\partial h} \\ \frac{\partial l}{\partial L} & \frac{\partial l}{\partial G} & \frac{\partial l}{\partial H} & \frac{\partial l}{\partial l} & \frac{\partial l}{\partial g} & \frac{\partial l}{\partial h} \end{pmatrix} \begin{pmatrix} dL \\ dG \\ dH \\ dl \\ dg \\ dh \end{pmatrix}. \quad (352)$$

#### Appendix C. 4.1 Partial derivatives of $a$

Representing  $a$  by  $L$  using its definition,

$$a = \frac{L^2}{\mu}. \quad (353)$$

We then know that  $a$  depends only on  $L$ . Hence

$$\frac{\partial a}{\partial L} = \frac{2L}{\mu}, \quad (354)$$

$$\frac{\partial a}{\partial G} = \frac{\partial a}{\partial H} = \frac{\partial a}{\partial l} = \frac{\partial a}{\partial g} = \frac{\partial a}{\partial h} = 0. \quad (355)$$

#### Appendix C. 4.2 Partial derivatives of $e$

Representing  $e$  by Delaunay elements using its definition

$$\frac{G^2}{L^2} = 1 - e^2, \quad (356)$$

$$\therefore e = \sqrt{1 - \frac{G^2}{L^2}}. \quad (357)$$

We then know that  $e$  depends only on  $L$  and  $G$ . Hence

$$\frac{\partial e}{\partial L} = \frac{\partial}{\partial L} \sqrt{1 - \frac{G^2}{L^2}} = \frac{1}{2} \left( 1 - \frac{G^2}{L^2} \right)^{-\frac{1}{2}} \frac{2G^2}{L^3} = \frac{G^2}{eL^3}, \quad (358)$$

$$\frac{\partial e}{\partial G} = \frac{\partial}{\partial G} \sqrt{1 - \frac{G^2}{L^2}} = \frac{1}{2} \left( 1 - \frac{G^2}{L^2} \right)^{-\frac{1}{2}} \left( -\frac{2G}{L^2} \right) = -\frac{G}{eL^2}, \quad (359)$$

$$\frac{\partial e}{\partial H} = \frac{\partial e}{\partial l} = \frac{\partial e}{\partial g} = \frac{\partial e}{\partial h} = 0. \quad (360)$$

#### Appendix C. 4.3 Partial derivatives of $\omega$

Since the argument of perihelion  $\omega$  is equal to  $g$  by its definition,  $\omega$  is independent from all the other Delaunay variables than  $g$  as

$$\frac{\partial \omega}{\partial L} = \frac{\partial \omega}{\partial G} = \frac{\partial \omega}{\partial H} = \frac{\partial \omega}{\partial l} = \frac{\partial \omega}{\partial h} = 0, \quad \frac{\partial \omega}{\partial g} = 1. \quad (361)$$

#### Appendix C. 4.4 Partial derivatives of $I$

Representing  $I$  by Delaunay elements by its definition

$$\cos I = \frac{H}{G}, \quad (362)$$

which means that  $I$  depends only on  $G$  and  $H$ , which leads to

$$\frac{\partial I}{\partial L} = \frac{\partial I}{\partial l} = \frac{\partial I}{\partial g} = \frac{\partial I}{\partial h} = 0. \quad (363)$$

Partial derivative of (362) by  $G$  gives

$$-\frac{H}{G^2} = -\sin I \frac{\partial I}{\partial G}, \quad (364)$$

$$\therefore \frac{\partial I}{\partial G} = \frac{H}{G^2 \sin I} = \frac{G \cos I}{G^2 \sin I} = \frac{1}{G \tan I}. \quad (365)$$

Similarly, the partial derivative of (362) by  $H$  becomes

$$\frac{1}{G} = -\sin I \frac{\partial I}{\partial H}, \quad (366)$$

$$\therefore \frac{\partial I}{\partial H} = -\frac{1}{G \sin I}. \quad (367)$$

#### Appendix C. 4.5 Partial derivatives of $\Omega$

Since the longitude of ascending node  $\Omega$  is equal to  $h$  by its definition,  $\Omega$  is independent from all the other Delaunay variables than  $h$  as

$$\frac{\partial \Omega}{\partial L} = \frac{\partial \Omega}{\partial G} = \frac{\partial \Omega}{\partial H} = \frac{\partial \Omega}{\partial l} = \frac{\partial \Omega}{\partial g} = 0, \quad \frac{\partial \omega}{\partial h} = 1. \quad (368)$$

#### Appendix C. 4.6 Partial derivatives of $l$

Since the mean anomaly  $l$  is identical to a Delaunay variable  $l$ ,

$$\frac{\partial l}{\partial L} = \frac{\partial l}{\partial G} = \frac{\partial l}{\partial H} = \frac{\partial l}{\partial g} = \frac{\partial l}{\partial h} = 0, \quad \frac{\partial l}{\partial l} = 1. \quad (369)$$

Using results of (354), (355), (358), (359), (360), (361), (363), (365), (367), (368), and (369), the transformation matrix (352) becomes as

$$\begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ dl \end{pmatrix} = \begin{pmatrix} \frac{\partial a}{\partial L} & 0 & 0 & 0 & 0 & 0 \\ \frac{\partial e}{\partial L} & \frac{\partial e}{\partial G} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & \frac{\partial I}{\partial G} & \frac{\partial H}{\partial G} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} dL \\ dG \\ dH \\ dl \\ dg \\ dh \end{pmatrix} \\ = \begin{pmatrix} \frac{2L}{eL^3} & 0 & 0 & 0 & 0 & 0 \\ \frac{\mu}{G^2} & -\frac{G}{eL^2} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & \frac{1}{G \tan I} & -\frac{1}{G \sin I} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} dL \\ dG \\ dH \\ dl \\ dg \\ dh \end{pmatrix} \quad (370)$$



$$\begin{aligned}
& \times \begin{pmatrix} \frac{2L}{eL^3} & 0 & 0 & 0 & 0 & 0 \\ \frac{\mu}{eL^3} & -\frac{G}{eL^2} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & \frac{1}{G \tan I} & -\frac{1}{G \sin I} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix} \\
& = \begin{pmatrix} 1 & & & 0 & 0 \\ & 1 & & 0 & 0 \\ & & 1 & 0 & 0 \\ & & & 1 & 0 \\ 0 & & & & 1 & 0 \\ 0 & \left(\frac{a}{r} + \frac{L^2}{G^2}\right) \sin f & 0 & 0 & 0 & \frac{a^2 \eta}{r^2} \end{pmatrix} \begin{pmatrix} \frac{2L}{eL^3} & 0 & 0 & 0 & 0 & 0 \\ \frac{\mu}{eL^3} & -\frac{G}{eL^2} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & \frac{1}{G \tan I} & -\frac{1}{G \sin I} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix} \\
& = \begin{pmatrix} \frac{2L}{eL^3} & 0 & 0 & 0 & 0 & 0 \\ \frac{\mu}{eL^3} & -\frac{G}{eL^2} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & \frac{1}{G \tan I} & -\frac{1}{G \sin I} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ \frac{G^2}{eL^3} \left(\frac{a}{r} + \frac{L^2}{G^2}\right) \sin f & -\frac{G}{eL^2} \left(\frac{a}{r} + \frac{L^2}{G^2}\right) \sin f & 0 & \frac{a^2 \eta}{r^2} & 0 & 0 \end{pmatrix} \quad (371)
\end{aligned}$$

Hence the final matrix becomes

$$\begin{aligned}
\begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ df \end{pmatrix} &= \frac{\partial(a, e, \omega, I, \Omega, f)}{\partial(L, G, H, l, g, h)} \begin{pmatrix} dL \\ dG \\ dH \\ dl \\ dg \\ dh \end{pmatrix} \\
&= \begin{pmatrix} \frac{\partial a}{\partial L} & \frac{\partial a}{\partial G} & \frac{\partial a}{\partial H} & \frac{\partial a}{\partial l} & \frac{\partial a}{\partial g} & \frac{\partial a}{\partial h} \\ \frac{\partial e}{\partial L} & \frac{\partial e}{\partial G} & \frac{\partial e}{\partial H} & \frac{\partial e}{\partial l} & \frac{\partial e}{\partial g} & \frac{\partial e}{\partial h} \\ \frac{\partial \omega}{\partial L} & \frac{\partial \omega}{\partial G} & \frac{\partial \omega}{\partial H} & \frac{\partial \omega}{\partial l} & \frac{\partial \omega}{\partial g} & \frac{\partial \omega}{\partial h} \\ \frac{\partial I}{\partial L} & \frac{\partial I}{\partial G} & \frac{\partial I}{\partial H} & \frac{\partial I}{\partial l} & \frac{\partial I}{\partial g} & \frac{\partial I}{\partial h} \\ \frac{\partial \Omega}{\partial L} & \frac{\partial \Omega}{\partial G} & \frac{\partial \Omega}{\partial H} & \frac{\partial \Omega}{\partial l} & \frac{\partial \Omega}{\partial g} & \frac{\partial \Omega}{\partial h} \\ \frac{\partial f}{\partial L} & \frac{\partial f}{\partial G} & \frac{\partial f}{\partial H} & \frac{\partial f}{\partial l} & \frac{\partial f}{\partial g} & \frac{\partial f}{\partial h} \end{pmatrix} \begin{pmatrix} dL \\ dG \\ dH \\ dl \\ dg \\ dh \end{pmatrix} \\
&= \begin{pmatrix} \frac{2L}{eL^3} & 0 & 0 & 0 & 0 & 0 \\ \frac{\mu}{eL^3} & -\frac{G}{eL^2} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & \frac{1}{G \tan I} & -\frac{1}{G \sin I} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ \frac{G^2}{eL^3} \left(\frac{a}{r} + \frac{L^2}{G^2}\right) \sin f & -\frac{G}{eL^2} \left(\frac{a}{r} + \frac{L^2}{G^2}\right) \sin f & 0 & \frac{a^2 \eta}{r^2} & 0 & 0 \end{pmatrix} \begin{pmatrix} dL \\ dG \\ dH \\ dl \\ dg \\ dh \end{pmatrix}. \quad (372)
\end{aligned}$$

Thus we have reached the conclusive partial derivatives of (317) and (318).

## Appendix C. 6 $(da, de, d\omega, dI, d\Omega, du) \rightarrow (dL, dG, dH, dl, dg, dh)$

This transformation matrix is obtained as a product of the two matrices of

$$\frac{\partial(a, e, \omega, I, \Omega, u)}{\partial(a, e, \omega, I, \Omega, l)}$$

and

$$\frac{\partial(a, e, \omega, I, \Omega, l)}{\partial(L, G, H, l, g, h)}$$

as

$$\begin{aligned} \frac{\partial(a, e, \omega, I, \Omega, u)}{\partial(L, G, H, l, g, h)} &= \frac{\partial(a, e, \omega, I, \Omega, u)}{\partial(a, e, \omega, I, \Omega, l)} \frac{\partial(a, e, \omega, I, \Omega, l)}{\partial(L, G, H, l, g, h)} \\ &= \begin{pmatrix} \frac{\partial a}{\partial a} & \frac{\partial a}{\partial e} & \frac{\partial a}{\partial \omega} & \frac{\partial a}{\partial I} & \frac{\partial a}{\partial \Omega} & \frac{\partial a}{\partial l} \\ \frac{\partial e}{\partial a} & \frac{\partial e}{\partial e} & \frac{\partial e}{\partial \omega} & \frac{\partial e}{\partial I} & \frac{\partial e}{\partial \Omega} & \frac{\partial e}{\partial l} \\ \frac{\partial \omega}{\partial a} & \frac{\partial \omega}{\partial e} & \frac{\partial \omega}{\partial \omega} & \frac{\partial \omega}{\partial I} & \frac{\partial \omega}{\partial \Omega} & \frac{\partial \omega}{\partial l} \\ \frac{\partial I}{\partial a} & \frac{\partial I}{\partial e} & \frac{\partial I}{\partial \omega} & \frac{\partial I}{\partial I} & \frac{\partial I}{\partial \Omega} & \frac{\partial I}{\partial l} \\ \frac{\partial \Omega}{\partial a} & \frac{\partial \Omega}{\partial e} & \frac{\partial \Omega}{\partial \omega} & \frac{\partial \Omega}{\partial I} & \frac{\partial \Omega}{\partial \Omega} & \frac{\partial \Omega}{\partial l} \\ \frac{\partial u}{\partial a} & \frac{\partial u}{\partial e} & \frac{\partial u}{\partial \omega} & \frac{\partial u}{\partial I} & \frac{\partial u}{\partial \Omega} & \frac{\partial u}{\partial l} \end{pmatrix} \begin{pmatrix} \frac{\partial a}{\partial L} & \frac{\partial a}{\partial G} & \frac{\partial a}{\partial H} & \frac{\partial a}{\partial l} & \frac{\partial a}{\partial g} & \frac{\partial a}{\partial h} \\ \frac{\partial e}{\partial L} & \frac{\partial e}{\partial G} & \frac{\partial e}{\partial H} & \frac{\partial e}{\partial l} & \frac{\partial e}{\partial g} & \frac{\partial e}{\partial h} \\ \frac{\partial \omega}{\partial L} & \frac{\partial \omega}{\partial G} & \frac{\partial \omega}{\partial H} & \frac{\partial \omega}{\partial l} & \frac{\partial \omega}{\partial g} & \frac{\partial \omega}{\partial h} \\ \frac{\partial I}{\partial L} & \frac{\partial I}{\partial G} & \frac{\partial I}{\partial H} & \frac{\partial I}{\partial l} & \frac{\partial I}{\partial g} & \frac{\partial I}{\partial h} \\ \frac{\partial \Omega}{\partial L} & \frac{\partial \Omega}{\partial G} & \frac{\partial \Omega}{\partial H} & \frac{\partial \Omega}{\partial l} & \frac{\partial \Omega}{\partial g} & \frac{\partial \Omega}{\partial h} \\ \frac{\partial l}{\partial L} & \frac{\partial l}{\partial G} & \frac{\partial l}{\partial H} & \frac{\partial l}{\partial l} & \frac{\partial l}{\partial g} & \frac{\partial l}{\partial h} \end{pmatrix} \\ &= \begin{pmatrix} 1 & & & 0 & 0 \\ & 1 & & & 0 \\ & & 1 & & 0 \\ & & & 1 & 0 \\ 0 & & & & 1 & 0 \\ 0 & \frac{\partial u}{\partial e} & 0 & 0 & 0 & \frac{\partial u}{\partial l} \end{pmatrix} \begin{pmatrix} \frac{\partial a}{\partial L} & 0 & 0 & 0 & 0 & 0 \\ \frac{\partial e}{\partial L} & \frac{\partial e}{\partial G} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & \frac{\partial I}{\partial G} & \frac{\partial H}{\partial G} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix} \\ &= \begin{pmatrix} 1 & & & 0 & 0 \\ & 1 & & & 0 \\ & & 1 & & 0 \\ & & & 1 & 0 \\ 0 & & & & 1 & 0 \\ 0 & \frac{\sin f}{\eta} & 0 & 0 & 0 & \frac{a}{r} \end{pmatrix} \begin{pmatrix} \frac{2L}{eL^3} & 0 & 0 & 0 & 0 & 0 \\ \frac{\mu}{eL^3} & -\frac{G}{eL^2} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & \frac{1}{G \tan I} & -\frac{1}{G \sin I} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix} \\ &= \begin{pmatrix} \frac{2L}{eL^3} & 0 & 0 & 0 & 0 & 0 \\ \frac{\mu}{eL^3} & -\frac{G}{eL^2} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & \frac{1}{G \tan I} & -\frac{1}{G \sin I} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ \frac{G^2 \sin f}{eL^3 \eta} & -\frac{G \sin f}{eL^2 \eta} & 0 & \frac{a}{r} & 0 & 0 \end{pmatrix} \\ &= \begin{pmatrix} \frac{2L}{eL^3} & 0 & 0 & 0 & 0 & 0 \\ \frac{\mu}{eL^3} & -\frac{G}{eL^2} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & \frac{1}{G \tan I} & -\frac{1}{G \sin I} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ \frac{\eta}{eL} \sin f & -\frac{\sin f}{eL} & 0 & \frac{a}{r} & 0 & 0 \end{pmatrix}. \end{aligned} \tag{373}$$

Thus finally,

$$\begin{aligned}
 \begin{pmatrix} da \\ de \\ d\omega \\ dI \\ d\Omega \\ du \end{pmatrix} &= \frac{\partial(a, e, \omega, I, \Omega, u)}{\partial(L, G, H, l, g, h)} \begin{pmatrix} dL \\ dG \\ dH \\ dl \\ dg \\ dh \end{pmatrix} \\
 &= \begin{pmatrix} \frac{\partial a}{\partial L} & \frac{\partial a}{\partial G} & \frac{\partial a}{\partial H} & \frac{\partial a}{\partial l} & \frac{\partial a}{\partial g} & \frac{\partial a}{\partial h} \\ \frac{\partial e}{\partial L} & \frac{\partial e}{\partial G} & \frac{\partial e}{\partial H} & \frac{\partial e}{\partial l} & \frac{\partial e}{\partial g} & \frac{\partial e}{\partial h} \\ \frac{\partial \omega}{\partial L} & \frac{\partial \omega}{\partial G} & \frac{\partial \omega}{\partial H} & \frac{\partial \omega}{\partial l} & \frac{\partial \omega}{\partial g} & \frac{\partial \omega}{\partial h} \\ \frac{\partial I}{\partial L} & \frac{\partial I}{\partial G} & \frac{\partial I}{\partial H} & \frac{\partial I}{\partial l} & \frac{\partial I}{\partial g} & \frac{\partial I}{\partial h} \\ \frac{\partial \Omega}{\partial L} & \frac{\partial \Omega}{\partial G} & \frac{\partial \Omega}{\partial H} & \frac{\partial \Omega}{\partial l} & \frac{\partial \Omega}{\partial g} & \frac{\partial \Omega}{\partial h} \\ \frac{\partial u}{\partial L} & \frac{\partial u}{\partial G} & \frac{\partial u}{\partial H} & \frac{\partial u}{\partial l} & \frac{\partial u}{\partial g} & \frac{\partial u}{\partial h} \end{pmatrix} \begin{pmatrix} dL \\ dG \\ dH \\ dl \\ dg \\ dh \end{pmatrix} \\
 &= \begin{pmatrix} \frac{2L}{\mu} & 0 & 0 & 0 & 0 & 0 \\ \frac{G^2}{eL^3} & -\frac{G}{eL^2} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & \frac{1}{G \tan I} & -\frac{1}{G \sin I} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ \frac{\eta}{eL} \sin f & -\frac{\sin f}{eL} & 0 & \frac{a}{r} & 0 & 0 \end{pmatrix} \begin{pmatrix} dL \\ dG \\ dH \\ dl \\ dg \\ dh \end{pmatrix}. \tag{374}
 \end{aligned}$$

Now the forward six transformations in

$$(da, de, d\omega, dI, d\Omega, df) \rightarrow (dL, dG, dH, dl, dg, dh)$$

have all been completed.

## References

- Albrow, M.D., Beaulieu, J.-P., Caldwell, A.R., DePoy, D.L., Dominik, M., Gaudi, B.S., Gould, A., Greenhill, J., Hill, K., Kane, S., Martin, R., Menzies, J., Naber, R.M., Pogge, R.W., Pollard, K.R., Sackett, P.D., Sahu, K.C., Vermaak, P., Watson, R., and Williams, A. (2000) Limits on stellar and planetary companions in microlensing event OGLE-1998-BUL-14, *Astrophys. J.*, **535**, 176–189.
- Battin, R.H. (1987) *An Introduction to The Mathematics and Methods of Astrodynamics*, American Institute of Aeronautics and Astronautics, Inc., New York.
- Bennett, D.P., Rhie, S.H., Becker, A.C., Butler, N., Dann, S., J. Kaspi, Leibowitz, E.M., Lipkin, Y., Maoz, D., Mendelson, H., Peterson, B.A., Quinn, J., Shemmer, O., Thomson, S., and Turner, S.E. (1999) Discovery of a planet orbiting a binary star system from gravitational microlensing, *nature*, **402**, 57–59.
- Boccaletti, D. and Pucacco, G. (1998) *Theory of Orbits. 2. Perturbative and geometrical methods*, Springer-Verlag, Berlin.
- Boss, A.P. (1996) Extrasolar planets, *Physics Today*, **49**, 32–38.
- Brouwer, D. and Clemence, G.M. (1961) *Methods of Celestial Mechanics*, Academic Press, New York.
- Chambers, J.E. (1999) A hybrid symplectic integrator that permits close encounters between massive bodies, *Mon. Not. R. Astron. Soc.*, **304**, 793–799.
- Chambers, J.E., Quintana, E.V., Duncan, M.J., and Lissauer, J.J. (2002) Symplectic integrator algorithms for modeling planetary accretion in binary star systems, *Astron. J.*, **123**, 2884–2894.
- Charlier, C.L. (1902) *Die Mechanik des Himmels*, Verlag Von Veit & Comp., Leipzig.
- Danby, J.M.A. (1992) *Fundamentals of Celestial Mechanics (second edition, third printing)*, Willmann-Bell Inc., Richmond, Virginia.
- Deprit, A. (1969) Canonical transformations depending on a small parameter, *Celes. Mech.*, **1**, 12–30.
- Dragt, A.J. and Finn, J.M. (1976) Lie series and invariant functions for analytic symplectic maps, *J. Math. Phys.*, **17**, 2215–2227.
- Duncan, M.J. and Lissauer, J.J. (1998) The effects of post-main-sequence solar mass loss on the stability of our planetary system, *Icarus*, **134**, 303–310.
- Duncan, M.J., Levison, H.F., and Lee, M.H. (1998) A multiple time step symplectic algorithm for integrating close encounters, *Astron. J.*, **116**, 2067–2077.
- Fukushima, T. (1998) Symmetric multistep methods revisited, in *Proc. 30th Symp. Celes. Mech.*, National Astronomical Observatory, 229–247.

- Fukushima, T. (1999) Symmetric multistep methods revisited: II. numerical experiments, in Svoren, J. and Pittich, E.M. eds., *Proc. 173rd Colloq. IAU., Evolution and source regions of asteroids and comets*, Astronomical Institute of the Slovak Academy of Sciences, 309, Tatranska Lomnica, Slovak Republic, August 24–28, 1998.
- Gladman, B., Duncan, M., and Candy, J. (1991) Symplectic integrators for long-term integrations in celestial mechanics, *Celes. Mech. Dyn. Astron.*, **52**, 221–240.
- Holman, M. and Wiegert, P.A. (1999) Long-term stability of planets in binary systems, *Astron. J.*, **117**, 621–628.
- Holman, M., Touma, J., and Tremaine, S. (1997) Chaotic variations in the eccentricity of the planet orbiting 16 Cygni B, *Nature*, **386**, 254–256.
- Hori, G. (1966) Theory of general perturbations with unspecified canonical variables, *Publ. Astron. Soc. Japan*, **18**, 287–296.
- Hori, G. (1967) Non-linear coupling of two harmonic oscillations, *Publ. Astron. Soc. Japan*, **19**, 229–241.
- Hori, G. (1970) Comparison of two perturbation theories based on the canonical transformation, *Publ. Astron. Soc. Japan*, **22**, 191–198.
- Hori, G. (1971) Theory of general perturbations for non-canonical systems, *Publ. Astron. Soc. Japan*, **23**, 567–587.
- Ito, T. and Tanikawa, K. (2002) Long-term integrations and stability of planetary orbits in our solar system, *Mon. Not. R. Astron. Soc.*, in press.
- Ito, T., Kinoshita, H., Nakai, H., and Fukushima, T. (1996) Numerical experiments to inspect the long-term stability of the planetary motion –1, in *Proc. 28th Symp. Celes. Mech.*, National Astronomical Observatory, Mitaka, Tokyo, 123–136.
- Kinoshita, H. (1998) *Celestial Mechanics and Orbital Dynamics*, University of Tokyo Press, in Japanese.
- Kinoshita, H. and Nakai, H. (1992) New method for long-term numerical integration of planetary orbits, in *Chaos, resonance and collective dynamical phenomena in the solar system*, Kluwer Academic publishers, Dordrecht, 395–406.
- Kinoshita, H., Yoshida, H., and Nakai, H. (1991) Symplectic integrators and their application to dynamical astronomy, *Celes. Mech. Dyn. Astron.*, **50**, 59–71.
- Lee, M.H., Duncan, M.J., and Levison, H.F. (1997) Variable time step integrators for long-term orbital integrations, *Astron. Soc. Pac. Conf. Ser.*, **123**, 32–37.
- Levison, H.F. and Duncan, M.J. (1994) The long-term dynamical behavior of short-period comets, *Icarus*, **108**, 18–36.



- Lichtenberg, A.J. and Lieberman, M.A. (1992) *Regular and Chaotic Dynamics*, Springer-Verlag, New York.
- Lie, S. (1888) *Theorie der Transformationgruppen I*, Teubner, Lipzig.
- Marcy, G.W. and Butler, R.P. (2000) Planets orbiting other suns, *Publ. Astron. Soc. Pac.*, **112**, 137–140.
- Marcy, G.W., Cochran, W.D., and Mayor, M. (2000) Extrasolar planets around main-sequence stars, in *Protostars & Planets IV*, The University of Arizona Press, Tucson, Arizona, 1285–1311.
- Mazeh, T., Krymolowski, Y., and Resenfeld, G. (1997) The high eccentricity of the planet orbiting 16 Cygni B, *Astrophys. J. Lett.*, **477**, L103–L106.
- Michel, P. and Valsecchi, G.B. (1997) Numerical experiments on the efficiency of second-order mixed-variable symplectic integrators for  $N$ -body problems, *Celes. Mech. Dyn. Astron.*, **65**, 355–371.
- Mikkola, S. (1997) Practical symplectic methods with time transformation for the few-body problem, *Celes. Mech. Dyn. Astron.*, **67**, 145–165.
- Mikkola, S. and Tanikawa, K. (1999) Explicit symplectic algorithms for time-transformed Hamiltonians, *Celes. Mech. Dyn. Astron.*, **74**, 287–295.
- Moriwaki, K. (2001) Stability of a planet in a binary system: MACHO-97-BLG-41, in *Proc. 33rd Symp. Celes. Mech.*, National Astronomical Observatory, 140–147.
- Moriwaki, K. and Nakagawa, Y. (2002) Stability of a planet in the binary system MACHO-97-BLG-41, *Astron. J.*, in press.
- Nagasawa, K. (1983) *Introduction to Astrodynamics*, Chijin Shokan Pub., in Japanese.
- Neri, F. (1987) Lie algebras and canonical integration, preprint.
- Plummer, H.C. (1960) *An Introductory Treatise on Dynamical Astronomy*, Dover, New York.
- Quinlan, G.D. and Tremaine, S. (1990) Symmetric multistep methods for the numerical integration of planetary orbits, *Astron. J.*, **100**, 1694–1700.
- Quintana, E.V., Lissauer, J.J., Chambers, J.E., and Duncan, M.J. (2002) Terrestrial planet formation in the  $\alpha$  Cenrauri system, *Astrophys. J.*, **576**, 982–996.
- Rauch, K.P. and Holman, M. (1999) Dynamical chaos in the Wisdom-Holman integrator: origins and solutions, *Astron. J.*, **117**, 1087–1102.
- Rhie, S.H., Bennet, D.P., Becker, A.C., Peterson, B.A., Fragile, P.C., Johnson, B.R., Quinn, J.L., Crouch, A., Gray, J., King, L., Messenger, B., Thomson, S., Bond, I.A., Abe, F., Carter, B.S., Dodd, R.J., Hearnshaw, J.B., Honda, M., Jugaku, J., Kabe, S., Kilmartin, P.M., Koribalski, B.S., Masuda, K., Matsubara, Y., Muraki, Y., Nakamura, T.,

- Nankivell, G.R., Noda, S., Rattenbury, N.J., Reid, M., Rumsey, N.J., Saito, T., Sato, H., Sato, S., Yock, P.C.M., and Yoshizawa, M. (2000) On planetary companions to the MACHO 98-BLG-35 microlens star, *Astrophys. J.*, **533**, 378–391.
- Saha, P. and Tremaine, S. (1992) Symplectic integrators for solar system dynamics, *Astron. J.*, **104**, 1633–1640.
- Saha, P. and Tremaine, S. (1994) Long-term planetary integrations with individual time steps, *Astron. J.*, **108**, 1962–1969.
- Sanz-Serna, J.M. and Calvo, M.P. (1994) *Numerical Hamiltonian Problems*, Chapman & Hall, London.
- Shniad, H. (1970) The equivalence of von Zeipel mapping and Lie transforms, *Celes. Mech.*, **2**, 114–120.
- Standish, E.M. (1990) The observational basis for JPL’s DE200, the planetary ephemerides of the astronomical almanac, *Astron. Astrophys.*, **233**, 252–271.
- Varadarajan, V.S. (1974) *Lie groups, Lie algebras and their representation*, Prentice-Hall, New Jersey.
- von Zeipel, H. (1916) Recherches sur le mouvement des petites planètes, *Ark. Astron. Mat. Fys.*, **11**, No. 1.
- Wiegert, P.A. and Holman, M.J. (1997) The stability of planets in the alpha Centauri system, *Astron. J.*, **113**, 1445–1450.
- Wisdom, J. and Holman, M. (1992) Symplectic maps for the  $n$ -body problem: stability analysis, *Astron. J.*, **104**, 2022–2029.
- Wisdom, J., Holman, M., and Touma, J. (1996) Symplectic correctors, in Marsden, J.E., Patrick, G.W., and Shadwick, W.F. eds., *Integration Algorithms and Classical Mechanics*, Vol. 10 of Fields Institute Communications, American Mathematical Society, Providence, Rhode Island, 217–244.
- Yoshida, H. (1990a) Conserved quantities of symplectic integrators for Hamiltonian systems, preprint.
- Yoshida, H. (1990b) Construction of higher order symplectic integrators, *Phys. Lett. A*, **150**, 262–268.
- Yoshida, H. (1993) Recent progress in the theory and application of symplectic integrators, *Celes. Mech. Dyn. Astron.*, **56**, 27–43.
- Yuasa, M. (1971) The comparison of Hori’s perturbation theory and von Zeipel’s theory, *Publ. Astron. Soc. Japan*, **23**, 399–403.

# Construction of Heterogeneous Computer System with GRAPE-5 and VPP5000 by Using IMPI

Mitsuru Hayashi<sup>\*</sup>, T., Ito<sup>†</sup>, E., Kokubo<sup>‡</sup>, H., Koyama<sup>†</sup>, K., Tomisaka<sup>†</sup>,  
K., Wada<sup>†</sup>, N., Uchida<sup>‡</sup>, N., Asai<sup>‡</sup>, E., Uemura<sup>‡</sup>, K., Sugimoto<sup>‡</sup>

## Abstract

We construct Heterogeneous Computer System with a special purpose computer for astrophysical many-body simulation, GRAPE-5 and a general purpose vector parallel computer VPP5000 and evaluate the performance of the system. We use IMPI to connect the computers. This study is a fundamental experiment to realize a simulation of particle-gas systems by using the characteristics of GRAPE and VPP. We carry out a calculation of self-gravity on the GRAPE-5 and a calculation of hydrodynamics on the VPP5000 to solve the physics of a self-gravitating contraction of a uniform gas sphere. The evaluation shows that the speed of the communication between GRAPE and VPP by using IMPI is slow and the calculation of self-gravity by using the straightforward approach on the GRAPE-5 occupies the most part of the calculation.

## 1 序論

重力多体専用計算機 GRAPE は数万体以上の粒子間に働く重力相互作用を解くシミュレーション研究において効率の良い計算を実現し、多くの研究成果をもたらしている。一方ベクトル型スーパーコンピュータは、格子点上の物理量の変化を時間追跡する流れの計算に代表される様な、各格子点上で同質の演算パターンが繰り返される計算に於いて効率の良い計算を実現し、天体物理学に限らず、多くの科学技術分野で様々な成果が得られている。年々、計算機能力が向上する中で、粒子計算、流体計算各々で計算を大規模化し、新たな知見を得る方法に加えて、上記の異なる計算機を連携し、従来にない計算を高効率に実現し、新たな知見を目指すことは非常に興味深い問題である。特に、銀河団ガスと銀河、AGN 近傍のガスと恒星系のダイナミクス、分子雲からの星団形成、原始惑星系における微惑星とガス等に代表される粒子-ガスの相互作用を扱う問題に上記システムの特性の活用が期待される。更に重力多体専用計算機 GRAPE を自分の重さで潰れるガスのシステム(自己重力流体システム)において、自己重力ソルバーとして採用し、特に AMR(Adaptive Mesh Refinement)、NG(Nested Grid)等の複雑な構造を持つ格子を用いた流体計算に活用することも、プログラミングが他の手法と比較して容易なことから興味深い。本研究では、上記 GRAPE と VPP の連携計算を実現するための基礎研究として、システムを構築し、具体的に自己重力が考慮された物理問題の計算を実行し、システムの性能評価を行なった。

---

<sup>\*</sup>Japan Science and Technology Corporation and National Astronomical Observatory

<sup>†</sup>National Astronomical Observatory of Japan

<sup>‡</sup>Fujitsu Limited

## 2 自己重力ソルバーの比較

上記の通り、GRAPE-VPP の連携によって最終的に実現を目指す計算は粒子-ガスシステムのシミュレーションである。一方、今回は、GRAPE-VPP で GRAPE に自己重力ソルバーの役割を持たせて実行したシミュレーションを用いて性能評価を行なった。そこで、GRAPE を自己重力ソルバーとして用いたプログラムと、他の手法による自己重力ソルバーを採用したプログラムに関して、規則的な構造を持つメッシュと、不規則な構造を持つメッシュの場合でプログラミング、パフォーマンスに関して比較する。

規則的なメッシュの場合プログラミングに関しては、GRAPE-VPP の連携は自己重力ソルバー部分のプログラミングに大きな困難はなく、容易と言って良い。一方、自己重力ソルバーとして採用される代表的な手法、FFT(Fast Fourier Transform、高速フーリエ変換)、CG(Conjugate Gradient、共役勾配法)、MG(Multi Grid、多層格子)等は高度なプログラミングを必要とする。プログラミングの容易さと言う点では GRAPE-VPP の連携にメリットはある。又、他の方法も高度なプログラミングを必要とはするが、規則的なメッシュの場合実現困難ということはない。

規則的なメッシュの場合パフォーマンスに関しては、GRAPE-汎用機の連携は直接全ての粒子に関して自己重力を計算する場合は非常に計算コストが大きくなってしまいが、遠方にある粒子はまとめて扱う Tree 法を用いることで計算コストは軽減できる。一方、FFT,CG,MG の場合は、計算コストは、自己重力部分以外の主要部分と同程度以下である。

不規則なメッシュの場合、GRAPE-VPP の連携はプログラミングは規則的なメッシュの場合同様容易であるといえる。一方、FFT,CG,MG は、規則的なメッシュより高度なプログラミングが必要となり、実現の見通しが極めて良くない場合もある。

不規則なメッシュ、特に AMR の様に位置情報もダイナミックに変動する場合の自己重力流体の問題に於いては GRAPE の特性が十分に活用され、GRAPE-VPP の連携が有効であると言える。

## 3 異機種計算機システム構築

次に構築したシステムに関して述べる。異機種並列環境を IMPI(Interoperable Message Passing Interface) を用いて構築した。IMPI は (閉じたシステムの) 並列スーパーコンピュータ等で用いられる並列化ライブラリ MPI(Message Passing Interface) を拡張したもので、アーキテクチャの異なるスーパーコンピュータ、WS クラスタ、PC クラスタ等を結合して並列計算環境を実現できる。MPI は個々のスーパーコンピュータ等で用いられる並列化ライブラリーのデファクトスタンダードとなっているが、IMPI も米国 NIST(National Institute of Standard) によって標準化を目指して開発された。

構築したハードウェアの構成は、Pentium III 1GHz の PC に Linux をインストールしたものを GRAPE-5(ピーク性能 38.4GFlops) のホストマシーンとして、ギガビットイーサで国立天文台天文学データ解析計算センターの VPP5000 に接続した。本研究に於いては、VPP5000 の 1 プロセッサを流体計算に割り当て、1 プロセッサを IMPI のサーバとして用いている。VPP5000 の 1PE(Processing Element) あたりのピーク性能は 9.6GFlops である。

## 4 GRAPE による流体の自己重力計算

流体計算の自己重力ソルバーとして GRAPE を用いる方法に関して述べる。流体計算では格子点上の流体の密度が計算される。格子点周辺には空間きざみから決定される、ある体積を持ったセルが定義される。そこで、上記 VPP 上の密度とセルの積から求められる質量が、質点として格子点上に存在すると仮定して、GRAPE の計算に必要な入力データである (粒子の) 位置情報、粒子数、質量を GRAPE のホストに転送し、GRAPE で重力を計算することで、自己重力の計算が可能とな

る(図1)。

特に、後述する自己重力収縮のシミュレーションは、Cartesian 座標を採用しており、座標は固定されている。そこで、上記シミュレーションでは、GRAPE のホストと VPP に共通の座標情報、格子点数(粒子数)を持たせたプログラミングを行なった。

実際に VPP から転送が必要なデータは密度とセルの積から求められる質量であり、GRAPE で計算された重力ポテンシャルを GRAPE から VPP に転送する(図2)。

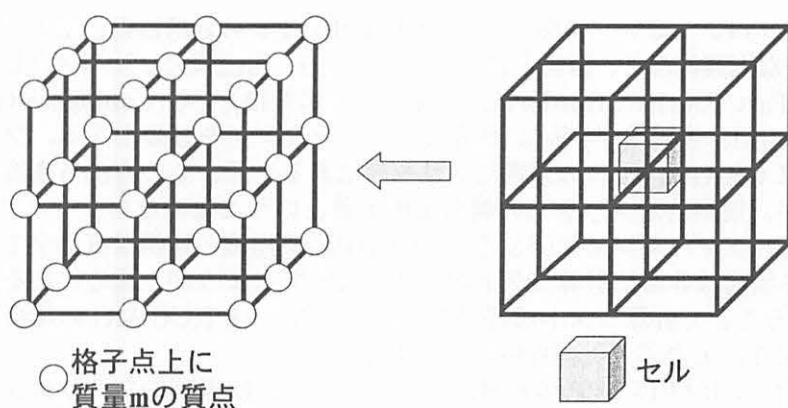


図1.粒子の質量イメージ (GRAPE計算イメージ)

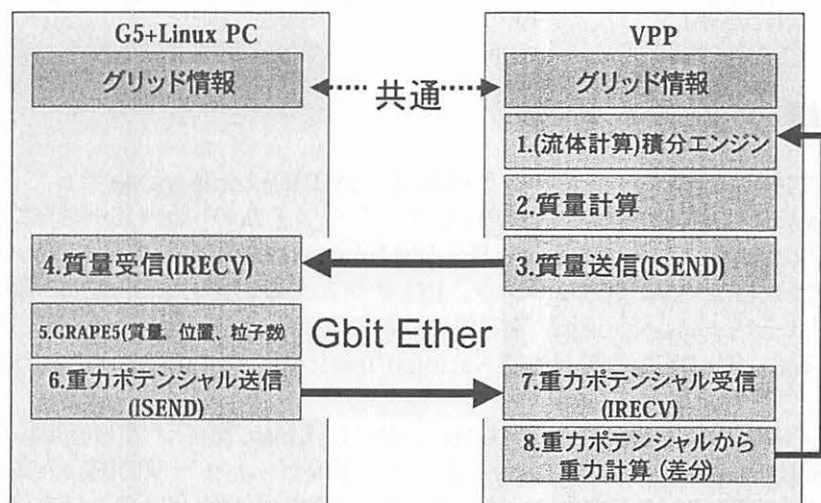


図2. データフローイメージ

## 5 一様球自己重力収縮の計算

実行した物理計算に関して述べる。初期条件は半径 1, 質量 1 の一様ガス球で、ジーンズ波長は約 0.2 である。VPP で計算される流体コードは時間、空間 2 次精度の陽的差分法 2 段階修正 Lax-Wendroff 法を採用し、上述の通り自己重力ソルバーとしては GRAPE-5 を用いた。下図は、中心密度の時間変化のとに関して、解析解と解像度の異なるシミュレーション結果を比較したものである。解像度はそれぞれ  $41^3, 49^3, 57^3$  で、解像度が向上することで、解析解に近付いて行くことが分かる。

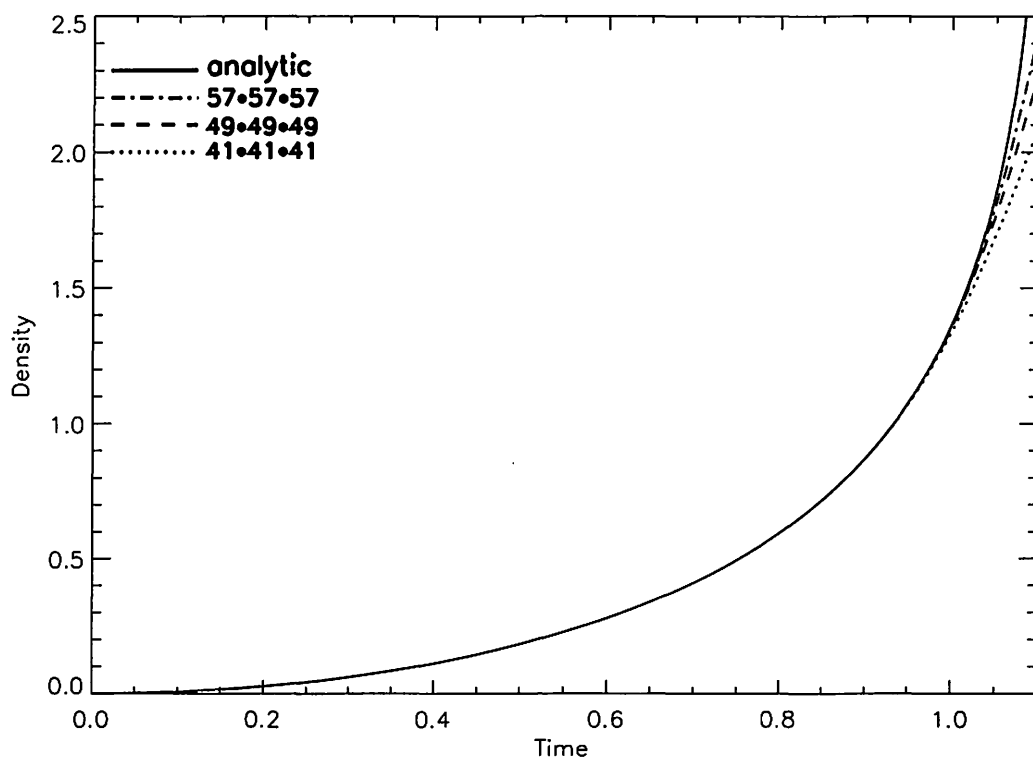


図 3. 解析解とシミュレーション結果の比較

## 6 性能評価

計算サイズと計算コストの一覧表を表 1 に示す。GRAPE-5 の計算を直接計算で実行しているため、GRAPE-5 の計算コストが非常に大きなものとなり、通信コストも、流体計算で殆どの計算コストを占める積分エンジンのコストよりも大きなものとなっている。

GRAPE-5 の計算に関しては、直接法ではなく、遠方の粒子をまとめて扱う Tree 法を採用することで、計算コストを軽減することができる。又、GRAPE-5 は GRAPE-6 に変更することで、10 倍以上のパフォーマンス向上が見込める。

通信コストに関しては、ギガビットイーサの MPI を用いた通信は、生の TCP/IP ソケット通信と比較して通信が劣化する計測結果が得られている。将来的には GRAPE のホストと VPP 間の通信は MPI ではなく、TCP/IP ベースの通信を採用することが望まれる。又、本研究の計測においては、共同利用の通常運用時で、VPP5000 プライマリ PE の負荷が大きな状況で計算実行を行なったために通信性能が低下したことも考えられる。特別にプライマリ PE の負荷がない状況で実行した計測では通常時における計測の場合よりも高い通信性能が示された。

サイズ	17 <sup>3</sup>	33 <sup>3</sup>	49 <sup>3</sup>	65 <sup>3</sup>
積分エンジン	0.007	0.02	0.06	0.13
通信	0.08	0.13	0.25	0.54
GRAPE-5	0.04	1.13	11.19	61.0
サイズ	73 <sup>3</sup>	81 <sup>3</sup>	89 <sup>3</sup>	97 <sup>3</sup>
積分エンジン	0.19	0.21	0.26	0.32
通信	0.73	1.00	1.31	1.70
GRAPE-5	121.5	226.95	398.21	665.62

表 1. 計算サイズと計算コスト一覧表

## 7 GRAPE-6 採用時の予想パフォーマンス

本研究では GRAPE-5 を用いて計算を行なったが、既に開発段階から実用段階に入った GRAPE-6 を採用した場合 (8 モジュール構成の GRAPE-6 ボードの場合)、ピーク性能の比較から予想される性能向上は GRAPE-5 の約 27 倍となる。計測結果より 129<sup>3</sup> の計算サイズの場合 GRAPE-5 の計算時間は約 3600 秒であることから、GRAPE-6 採用時の予想パフォーマンスは約 139.5 秒となる。Tree 法の採用により更に約 10 倍のパフォーマンス向上が予想され、GRAPE-6、Tree 法を採用した場合、129<sup>3</sup> の計算サイズでのパフォーマンスは約 14 秒と予想される。通信に関してはギガビットイーサを用いた場合、MPI でなく TCP/IP の生のソケット通信を行なうことで 60MB/s 程度のパフォーマンスが実現されていることが、他のプロジェクトで確認されており、129<sup>3</sup> の計算サイズの場合に適用すると双方向で 0.56 秒となる。以上より、GRAPE-6、Tree 法、生の TCP/IP ソケット通信を採用した場合、129<sup>3</sup> の計算サイズの自己重力ソルバーとしてのパフォーマンスは約 15 秒となる。一方、128<sup>3</sup> の計算サイズで数値計算ライブラリ Fujitsu SSLII VPP を用いた場合のパフォーマンスも約 15 秒であり、GRAPE-6、Tree 法、生の TCP/IP ソケット通信で自己重力ソルバーを構築することで、ベンダー提供のライブラリ程度の効率が見込める、

## 8 まとめと今後

粒子-ガス系のシミュレーション実行のための基礎研究として、重力多体専用計算機 GRAPE-5 とベクトル型並列計算機 VPP5000 の連携システムを IMPI を用いて構築し性能評価を行なうことで、上記シミュレーション実現のために有用な知見を得ることができた。特に通信に関しては、MPI 通信ではなく、TCP/IP のソケット通信を用いて通信することで、効率の良い連携計算が実現できることが予想される。GRAPE-6(ピーク性能約 1TFlops/ボード)の採用、Tree 法の採用によって、ベンダー提供の共役勾配法を用いて自己重力の計算を実行する場合と同程度の効率が期待できる。今後は、GRAPE-6 を用いた性能評価を実行、AMR を用いた自己重力問題のためのプログラム開発と性能評価の後、粒子-ガス系を解くためのアルゴリズム、プログラム開発を計画している。

# 地球の時間暦の非線型調和解析

## Non-linear harmonic analysis of the time ephemeris of the Earth

原田 渉 (東京大学)

Wataru Harada (Tokyo University)

福島登志夫 (国立天文台)

Toshio Fukushima (National Astronomical Observatory of Japan)

### abstract

A time ephemeris of the Earth is related with a relativistic time-dilation and is represented by the solar system barycentric time(TCB) and the geocentric coordinate time(TCG). At present a time ephemeris of the Earth is obtained by numerical integration (Fukushima 1995 A&A; Irwin & Fukushima 1999 A&A), but it is inconvenient for practical use. In this paper, we decomposed the results of the numerical integration into the Fourier series and the mixed secular terms. Because the frequency in the series was generally unknown we tried non-linear harmonic analysis by which the frequency was solved simultaneously. Actually we adopted the data calculated from the JPL's development ephemeris, DE405, and estimated parameters by using non-linear least square method. Our calculations reproduced the time ephemeris of the Earth more accurate than that obtained the previous analytical theory (Fairhead, Bretagnon & Lestrade 1995)

## 1 Introduction

研究の目的は非線型調和解析プログラムの確立である。つまり解析的理論があまり正確であるとは言い難く、数値積分でしか求められないようなデータがある時、このプログラムが確立されていれば、そのデータの平均値が精度良く求まるばかりでなく、数値積分では難しい長期予測も可能になる。そこで今回の物理ターゲットとして地球の時間暦を選んだ。地球の時間暦、すなわち太陽系重心座標時 TCB と地心座標時 TCG とを結びつける関係式

$$\text{TCB} - \text{TCG} = \int g(t) dt \quad , \quad g(t) \equiv \frac{1}{c^2} \left( U_E(t) + \frac{\mathbf{v}_E^2(t)}{2} \right) \quad (1)$$

は、現在、数値積分 (Fukushima 1995 A&A ; Irwin & Fukushima 1999 A&A) で求められているが、そのままでは使いづらい。ただし、 $U_E(t)$ 、 $\mathbf{v}_E(t)$  は地心における重力ポテンシャルと地球の速度で、 $t$  は TCB である。ここで平均値は  $L_c = \langle g(t) \rangle$  で与えられ、過去の研究として

$$L_c = 1.48082685594 \times 10^{-8} \pm 1. \times 10^{-17} \quad (\text{Irwin \& Fukushima 1999}) \quad (2)$$

が求められている。目標としてはこの平均値  $L_c$  の uncertainty を目安に、残差を過去の研究で行われたものよりも小さくなることを目指す。それはつまり、今回の調和解析による結果が過去の研究よりもより精度良く求められていることを意味する。



# 2 非線型最小二乗法

JPL の月惑星暦 DE405 を用いて  $g(t)$  を 3 日おきに表として求め、残差平方和  $\phi$  を

$$\phi = \sum_{n=1}^N [a_1 + a_2 t_n + \sum_{k=1}^K \{a_{2k+1} \sin(2\pi f_k t_n) + a_{2k+2} \cos(2\pi f_k t_n)\} - g(t_n)]^2 \quad (3)$$

とする。これはデータ  $g(t)$  を linear な部分 ( $a_1 + a_2 t_n$ ) と Fourier term で fitting したときの残差平方和である。この残差平方和を最小にするようなパラメーター推定を行いたい。ここで未知数はそれぞれの項の係数  $a_i$ 、Fourier term の周波数成分  $f_k$ 、そして Fourier term の項数  $K$  である。ここで  $\phi$  が最小になるには  $a_i$  の偏微分が 0 となるような  $a_i$  を求めればよい。それぞれの項の係数  $a_i$  に関して残差平方和は線型なので、 $\phi$  を  $a_i$  で微分することによって次式 (4) の正規方程式を求め、それを解くと  $a_i$  は一意的に求まる。

$$\mathbf{A} \cdot \mathbf{a} = \mathbf{b} \quad , \quad \mathbf{A}_{ij} = \frac{\partial^2 \phi}{\partial a_i \partial a_j} \quad (4)$$

しかし、周波数成分は非線型になるので一意的には求めることはできない。そこで次式 (5) のようにテイラー展開から  $\Delta \vec{f}$  を求め、それを繰り返して真値に近づける。

$$\Delta \vec{f} = -\mathbf{H}^{-1} \frac{\partial \phi}{\partial \vec{f}} \quad , \quad \mathbf{H}_{nm} = \frac{\partial^2 \phi}{\partial f_n \partial f_m} \quad (5)$$

また、未知数である Fourier term の項数  $K$  は求めることができない。そこで  $K = 1$  から 1 つずつ増やし、その都度、非線形最小二乗法を行うことを考えた。その非線形最小二乗法と全体のアルゴリズムをまとめたものが Fig.1 である。

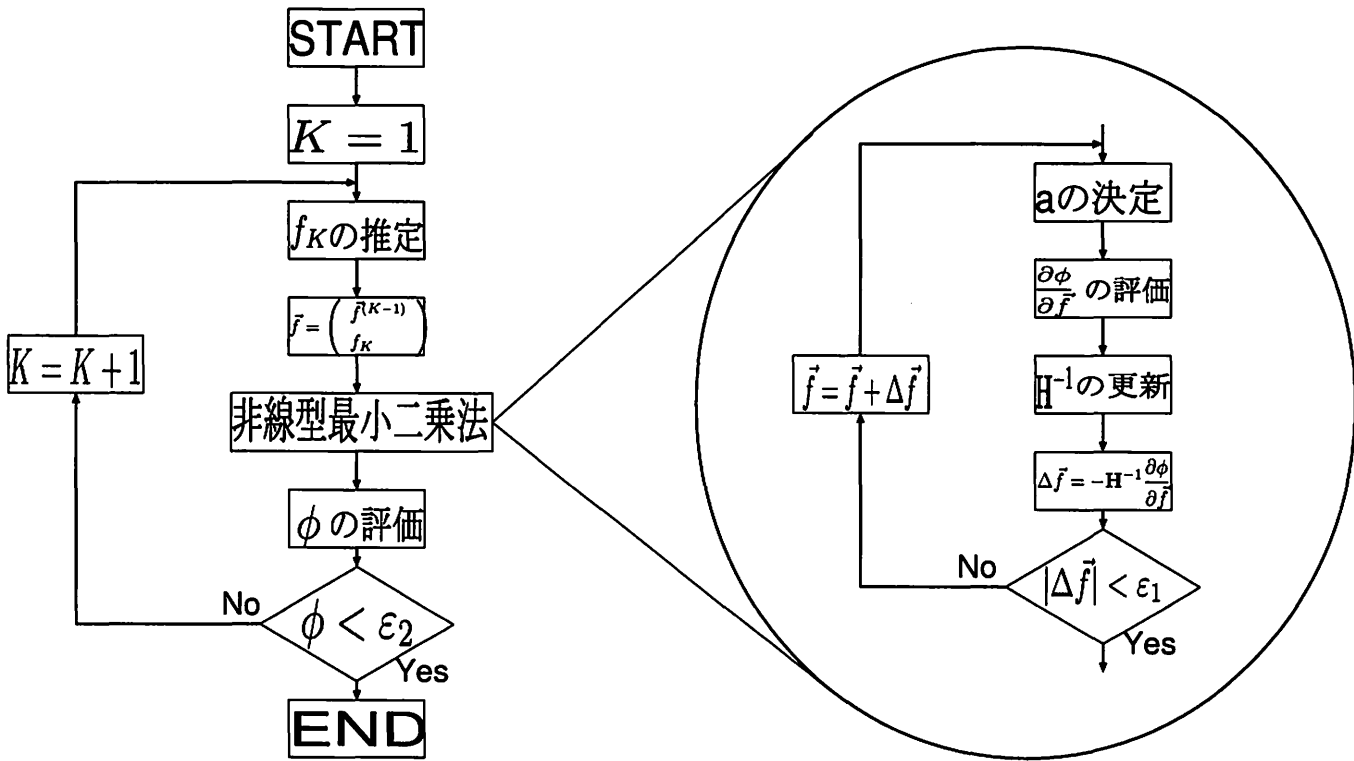


Fig.1 Algorithm of non-linear harmonic analysis

2.1  $f_K$  の推定

Fig.1 から分かるように、非線型最小二乗法を行う前には必ず  $f_K$  の推定を行わなくてはならない。精度良く推定してやるために、精度がデータ数に左右されてしまう FFT (Fast Fourier Transform) を使わずにペリオドグラムを使った。その時の最も大きい振幅を持つ周波数を新しい  $f_K$  とした。さらに精度良くするために、ピークの位置付近で放物線近似 (Brent, R.P. 1973) を使って周波数を求めている。Fig.2 は  $K = 5$  のときの残差、つまり周波数成分を 5 つ取り除いた後の残差をペリオドグラムにかけたものである。

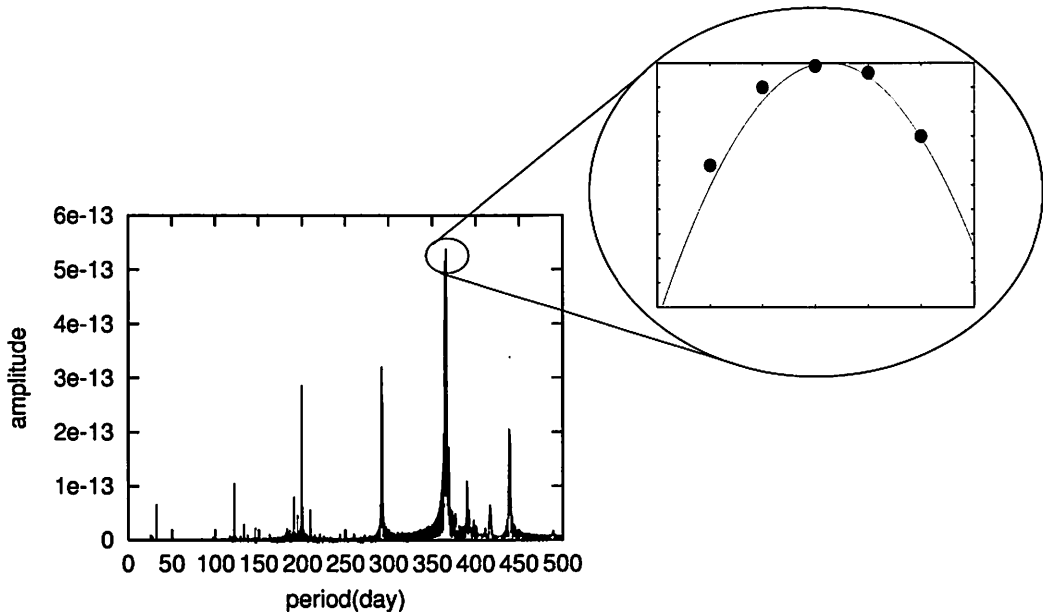


Fig.2 Estimation of  $f_K$

2.2  $H^{-1}$  の更新

上式 (5) に示した方法を逆 Hesse 法と呼ぶ。また別の方法として、周波数空間における  $\phi$  の傾きを計算して、最も傾きの大きい方向を選び、その方向で 1 次元の最小値探索をする。これを繰り返して真値に近づける方法を最急降下法という。しかし、これらの方法は周波数空間において最小値付近の形状が複雑な場合、Fig.3 のように無駄な繰り返しを行ってしまう (Fig.3 の曲線は周波数空間における  $\phi$  の contour)。今回は Hesse 行列  $H$  を容易く計算できることから、無駄な繰り返しを極力抑えることのできる準ニュートン法を使う。これは反復にしたがって Hesse 行列の逆行列の近似を高めていく方法である。代表的な準ニュートン法として DFP 法 (Davidon-Fletcher-Powell algorithm) と BFGS 法 (Broyden-Fletcher-Goldfarb-Shanno algorithm) があり、ここでは後者の BFGS 法を使った。

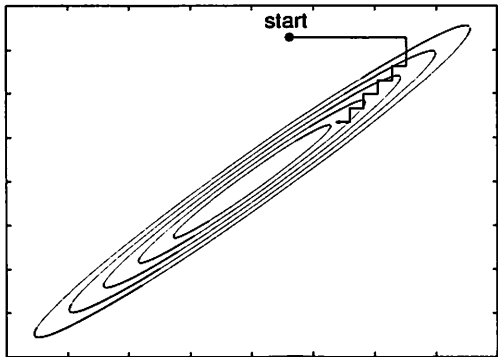


Fig.3 Iteration of frequency space

### 3 混合永年項 (Mixed secular term) 問題

#### 3.1 混合永年項の出現

周波数成分が非常に近接している場合

$$(a + \Delta a) \sin \{(\omega + \Delta\omega)t\} - a \sin(\omega t) \approx a\Delta\omega t \cos(\omega t) \quad (6)$$

となり  $t \sin$  や  $t \cos$  の項 (混合永年項) が現れてしまう。これはデータ区間を区切ってしまったために見えるもので、Fourier term のみで fitting しようとしてもなかなか収束しない。

#### 3.2 拡張ペリオドグラムの共鳴現象

$f_K$  の推定のために、今まで使ってきたペリオドグラムから混合永年項を基底関数に含む拡張ペリオドグラムを考えてみた。

$$\phi' = \sum_t \{S \sin(2\pi f t) + C \cos(2\pi f t) + S' t \sin(2\pi f t) + C' t \cos(2\pi f t) - g(t)\}^2 \quad (7)$$

つまり、上式 (7) で正規方程式を立て、そこから  $S$ 、 $C$ 、 $S'$ 、 $C'$  を決定して  $\sqrt{S^2 + C^2}$  を Fourier term の振幅、 $\sqrt{S'^2 + C'^2}$  を混合永年項の振幅とした。そして、Fourier 成分に関するペリオドグラム、混合永年項成分に関するペリオドグラムを作った。これらをまとめて拡張ペリオドグラムと呼び、それらの最大振幅を持つ周波数を見積もることによって、 $f_K$  の推定を行う。

～テスト～

$$g_{\text{test}}(t_n) = S \sin(2\pi f_1 t_n) + C' t \cos(2\pi f_2 t_n) \quad , \quad (f_1 = f_2 = 0.1) \quad (8)$$

この拡張ペリオドグラムを使って式 (8) のテストデータで確認してやると、Fig.4 に見られる共鳴現象のような偽の信号が見えた。(図①④の double peak が偽の信号)

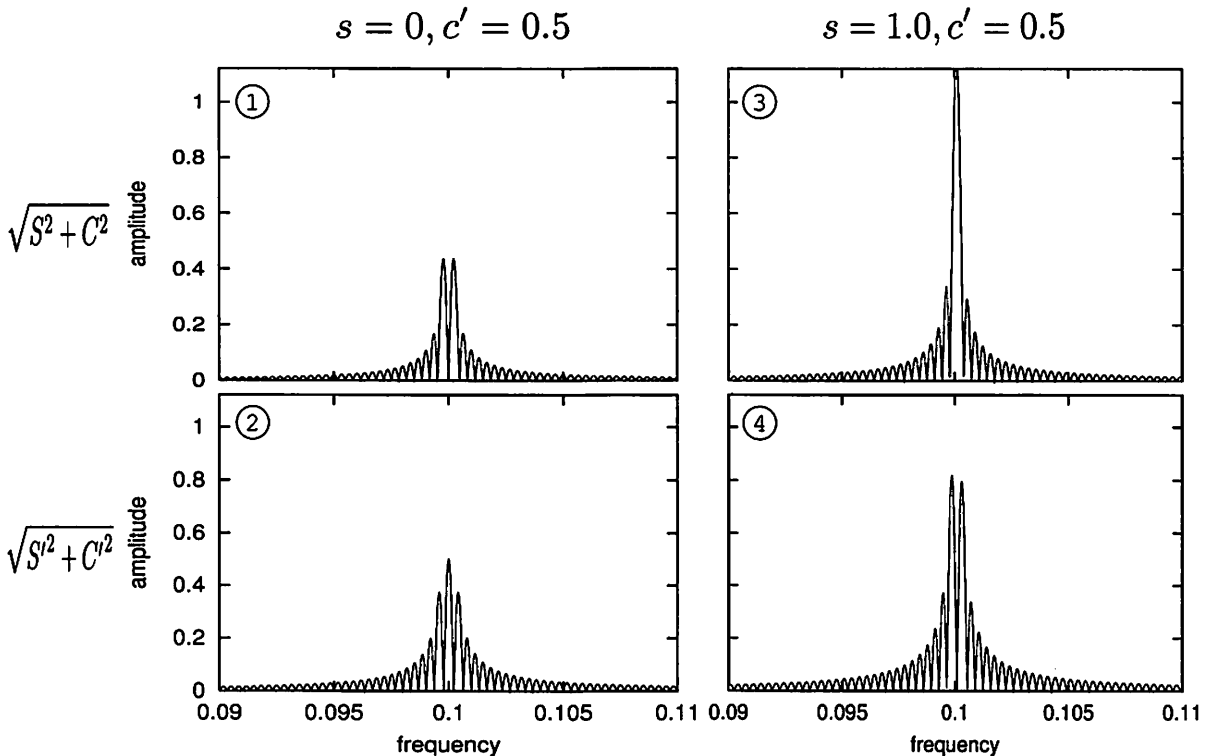


Fig.4 Test of extended periodogram

3.3 共鳴現象の解決策

この共鳴現象は単一周波数で同時にパラメーター推定することによって起きている。同一周波数で考えるなら、Fourier termの方が必ず混合永年項よりも大きいと仮定すると、先に fitting されるのは必ず Fourier termで、そのとき混合永年項のペリオドグラムに偽の信号が見えていても関係がない(図③④)。また Fourier termが fittingされた後 ( $S = 0$ ) は混合永年項に偽の信号は見られない(図①②)。よって、今回の algorithm ( $K$ を1つつ増やす)を使えば、共鳴現象による偽の信号に惑わされることはない。

4 結果

まず、時間暦の式における積分する前の関数  $g(t)$  をこの非線型調和解析プログラムに適用した。この時の得られたパラメーター  $a_i, f_j$  を手で積分してやり、解析解 (Fairhead, Bretagnon & Lestrade 1995) と数値積分との比較を行った。これら3つの中で、DE405を直接数値積分したものが最も信頼できるものである。

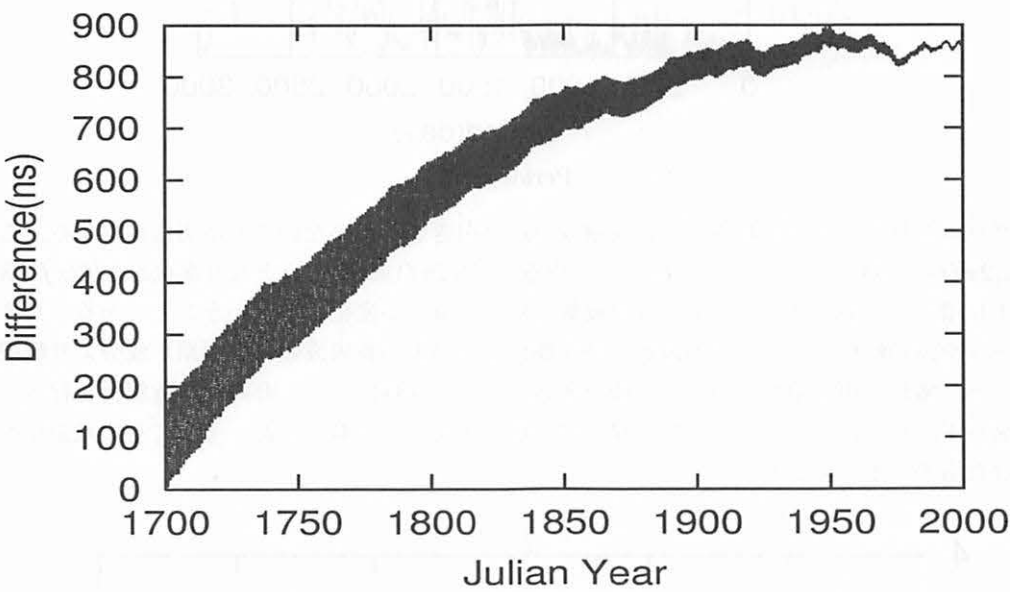


Fig.5 Comparison our results with FBL

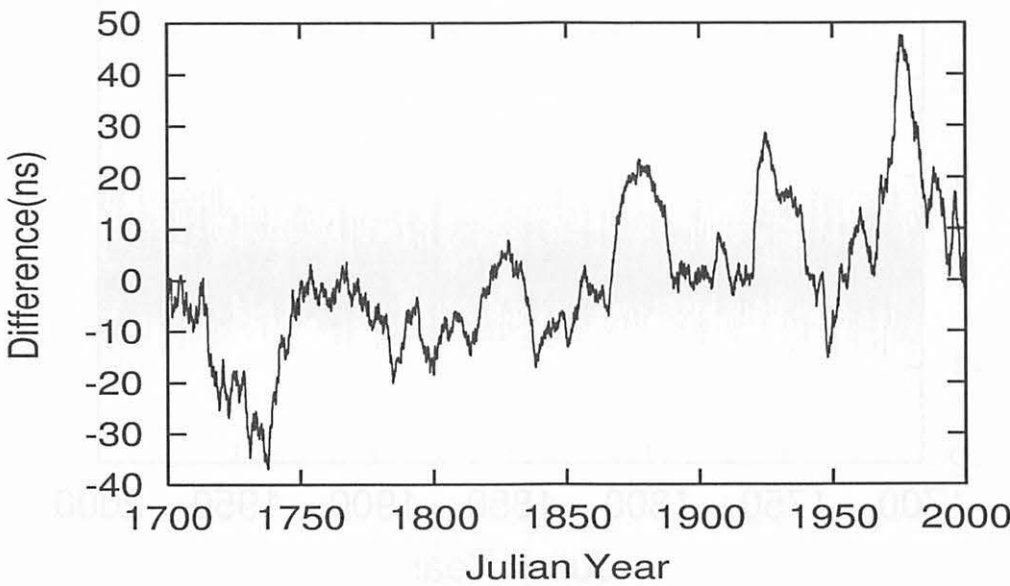


Fig.6 Comparison our results with numerical integration

これらの図 (Fig.5 Fig.6) から明らかなように解析解より精度良く求まっていることが分かる。しかし、Fig.6 を見るとまだ引ききれていない周波数成分があるようである。そこで、この比較したものの残差をペリオドグラムにかけてやると下図が見えた。

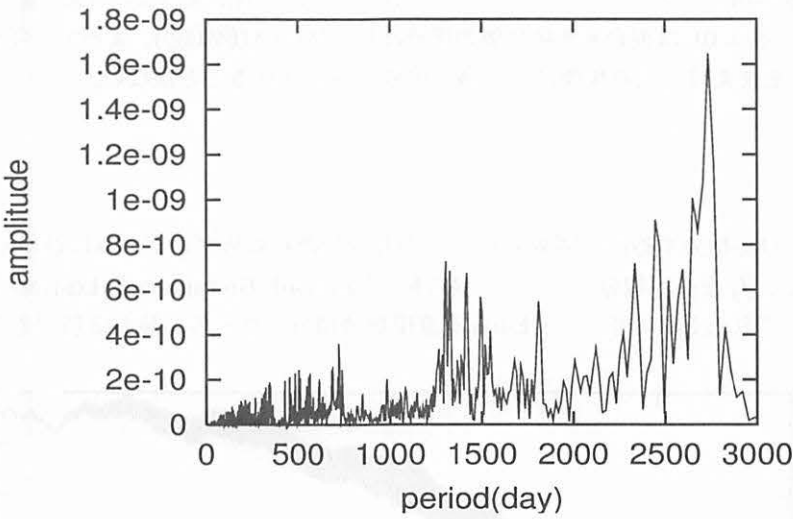


Fig.7 Peridogram

このペリオドグラムを見ると低周波数成分（長周期成分）が引ききれていないことが明らかである。これは積分前に  $S \cos(2\pi ft)$  で得られたデータは、積分してやると  $S/(2\pi f) \sin(2\pi ft)$  となり振幅が  $1/(2\pi f)$  倍されてしまう。つまり積分する前と積分した後では周波数に対して重みが変わるということである。具体的には積分前のデータを調和解析すると、短周波数（長周期）成分の方が長周波数（短周期）成分より軽視されるといえる。よって無駄な項を拾わないように積分後のデータに同様に、非調和解析解析を行った。その時の残差の図が Fig.8 である。ここでの K は積分前も積分後もほぼ同数とした。そしてその残差をペリオドグラムにかけたものが Fig.9 である。

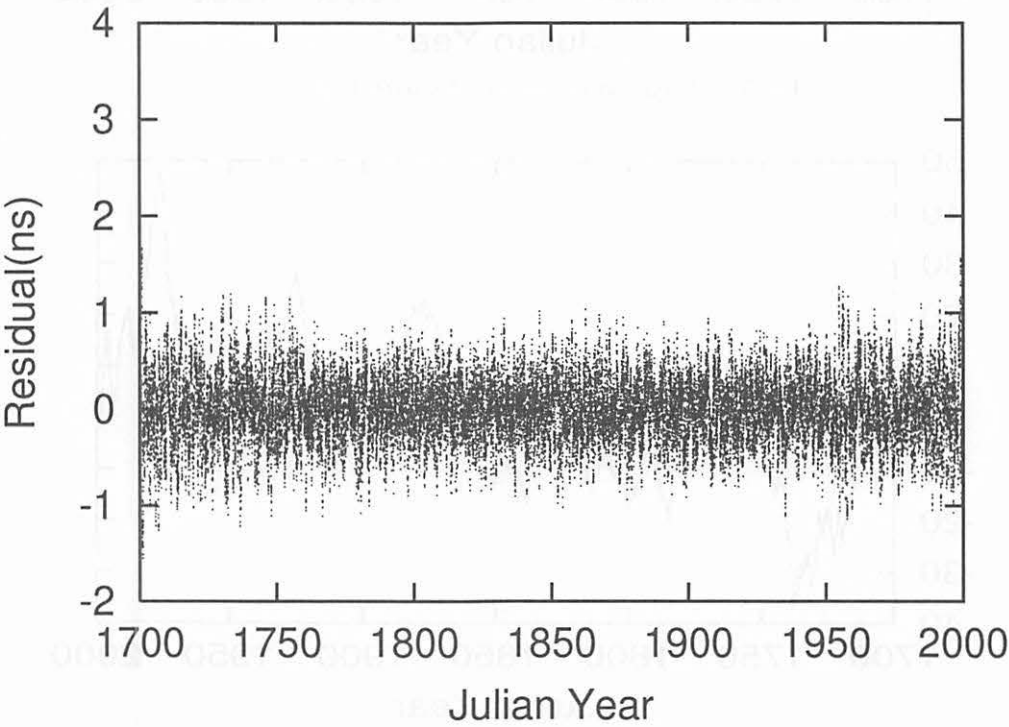


Fig.8 Residual of our analysis

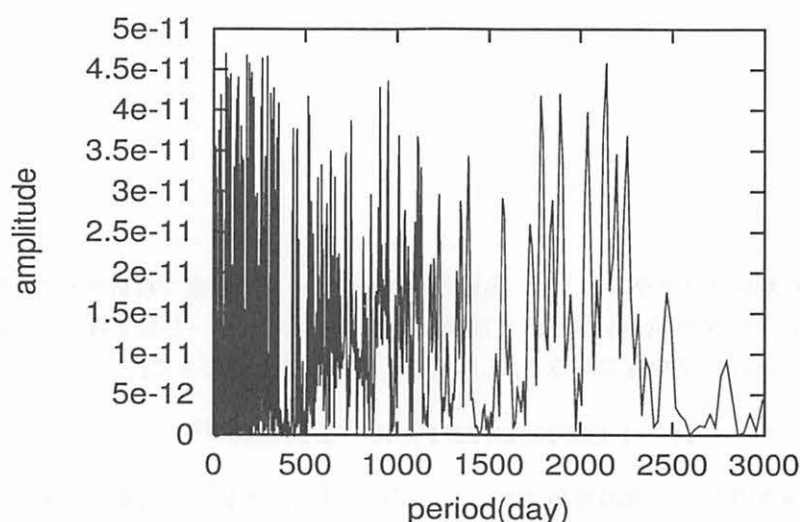


Fig.9 Periodogram

このペリオドグラムは明らかに Fig.7 より white であるといえるので、これは明らかに積分後のデータを非線型調和解析した方が少ない項数で精度良い結果を得ることができるといえる。このことは  $L_c$  の推定誤差とパラメーター数の関係からも明らかでその図は Fig.10 である。

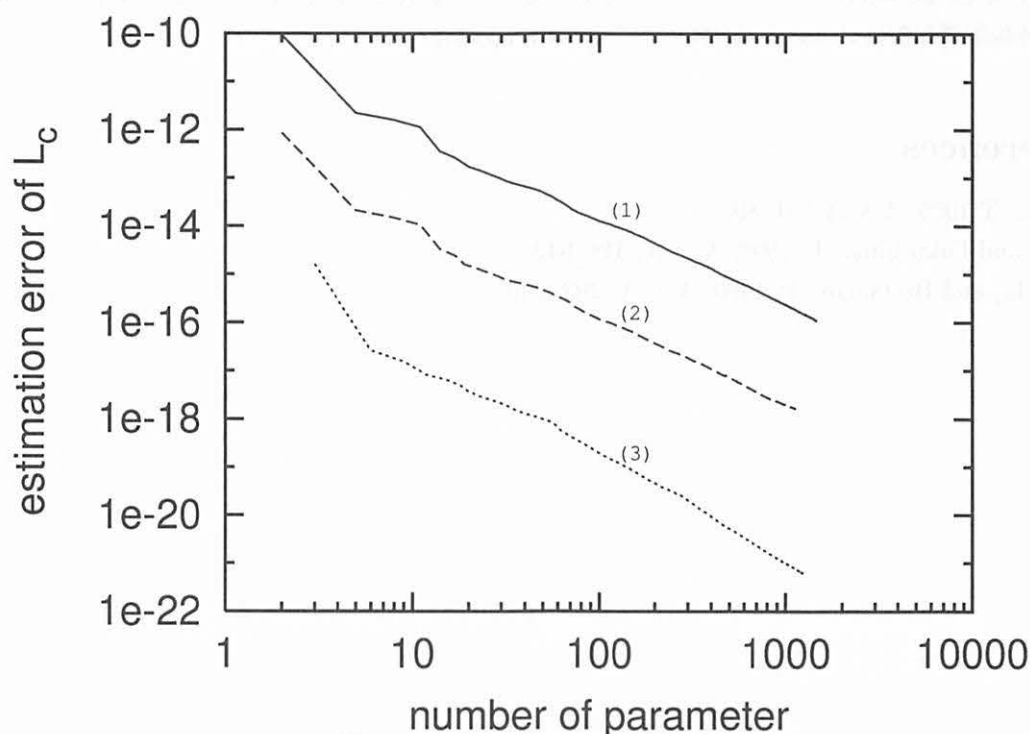


Fig.10 Estation error of  $L_c$

This figure show (1)analysis of data without mixed secular terms(full line), (2)that of data with mixed secular terms(dashed line),and (3)that of integratede data with mixed secular terms(dotted line)

この  $L_c$  の推定誤差の結果と残差の結果 (Fig.8) は、明らかに過去の研究 (Fukushima 1995 A & A ; Irwin & Fukushima 1999 A&A) よりもより精度の良い結果が得ることができた。

## 5 まとめ

### 5.1 結論

今回作った非線型調和解析プログラムにより、地球の時間暦をよりよい精度で解析することに成功した。これは、数値積分とは違いデータが Fourier term や混合永年項のパラメーターとして得られているので長期予測も可能である。また、 $g(t)$  の平均値である  $L_c$  も精度良く求めることができた。

$$L_c = 1.48082684872271 \times 10^{-8} \pm 6. \times 10^{-22}$$

付け加えるなら、混合永年項による共鳴現象に対して、プログラムに内装した拡張ペリオドグラムのアルゴリズムもまとめることができた。

### 5.2 今後の課題

まずは月・惑星暦 DE405 の全範囲（1600-2200）において、非線型調和解析の計算しなおさねばならないだろう。他には非線型調和解析プログラムの他の物理現象へ適用することが課題である。また、今回の地球の時間暦は 1 次元の調和解析であったが、物理現象によっては多次元の調和解析が必要になることもある。よってこの非線型調和解析プログラム自体のベクトル化も課題となる。

## 6 References

- Fukushima, T 1995, A & A 294, 895
- Irwin, A., and Fukushima, T. 1999, A & A, 348, 642
- Fairhead, L., and Bretagnon, P. 1990, A & A, 299, 240

# ケプラー運動に対する対称線形多段法の問題点

## Problems of Symmetric Multistep Methods for Keplerian Motion

山本 一登（総研大／国立天文台）

Tadato YAMAMOTO

tadato.yamamoto@nao.ac.jp

*Department of Astronomical Science, Graduate University for Advanced Studies,  
2-21-1, Osawa, Mitaka, Tokyo 181-8588 JAPAN*

福島 登志夫（国立天文台）

Toshio FUKUSHIMA

Toshio.Fukushima@nao.ac.jp

*National Astronomical Observatory of Japan,  
2-21-1, Osawa, Mitaka, Tokyo 181-8588 JAPAN*

### ABSTRACT

Symmetric multistep methods can avoid the linear energy error, but these methods suffer damages of resonances and instabilities at special stepsizes when integrating with non-linear equations like Keplerian motion. We consider whether there would be any method of avoiding this phenomenon.

## 1 はじめに

対称線形多段法の特徴と今日までにされてきた研究についてまとめておく。

### ◆ 主な特徴

- 線形多段法なので任意の次数の公式が作れる
- ◎ エネルギーや角運動量などの保存量の誤差がある範囲内に留まる
- × ケプラー運動のような非線形な系に使用すると刻み幅共鳴を起こすことがある
- × 離心率が大きくなるにつれて安定な刻み幅の最大値が小さくなる



## ●これまでの主な研究の流れ

- 1976 年、Lambert, J.D. and Watson, I.A.  
対称線形多段法に関する最初の論文
- 1990 年、Quinlan, G.D. and Tremaine, S.  
8、10、12、14 次の公式が導かれ、ケプラー運動に対する数値実験が行なわれた、  
しかし刻み幅共鳴の存在には気が付かなかった
- 1998 年、Fukushima, T.  
Quinlan 達よりも良い性質の公式を導くが、ケプラー運動に対してはうまくいかない
- 1999 年、Quinlan, G.D.  
刻み幅共鳴の存在を知り、その起源と対処についての研究だが、  
プレプリントのみで論文にはなっていない
- 2000 年、Arakida, H. and Fukushima, T.  
K-S 変換によりケプラー運動を調和振動子化 (正則化) することで  
刻み幅共鳴の問題は解決した、ただし N 体問題には適用が困難

K-S 変換によりケプラー運動を正則化すればこの刻み幅共鳴の問題は回避できるのだが、適用範囲が限られてしまう。今回の目的は適用範囲を広げるため、ケプラー運動のままで使用したときに刻み幅共鳴を回避および減少することができないか考えてみた。

## 2 刻み幅共鳴

刻み幅共鳴がどのように現れるのかを Quinlan and Tremaine の公式を例に見てみる。

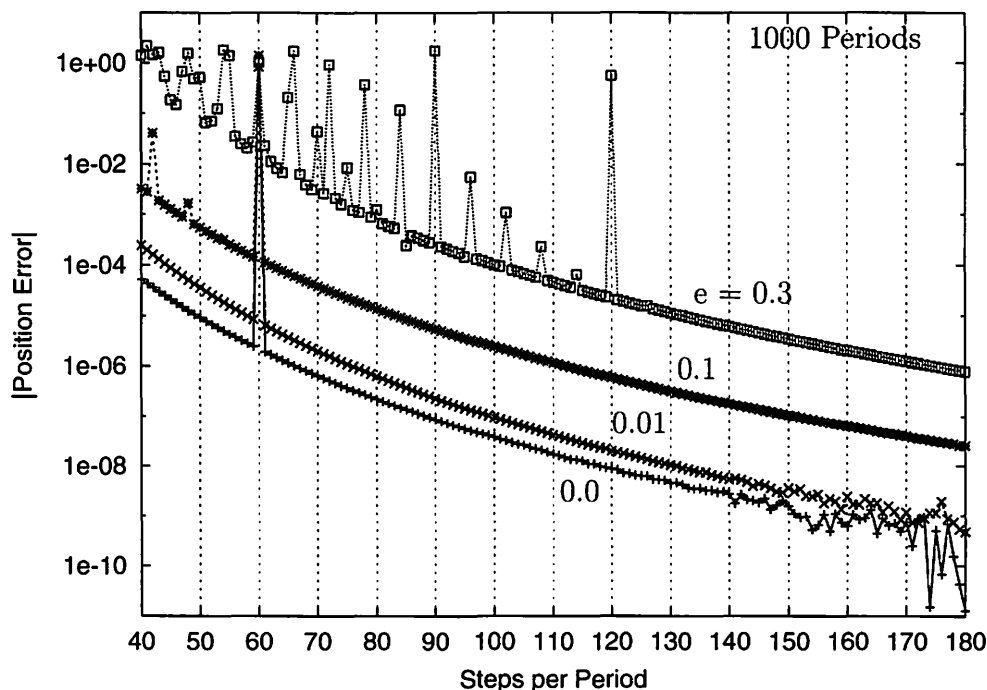


Fig.1 Quinlan and Tremaine 8th Order Keplerian Motion

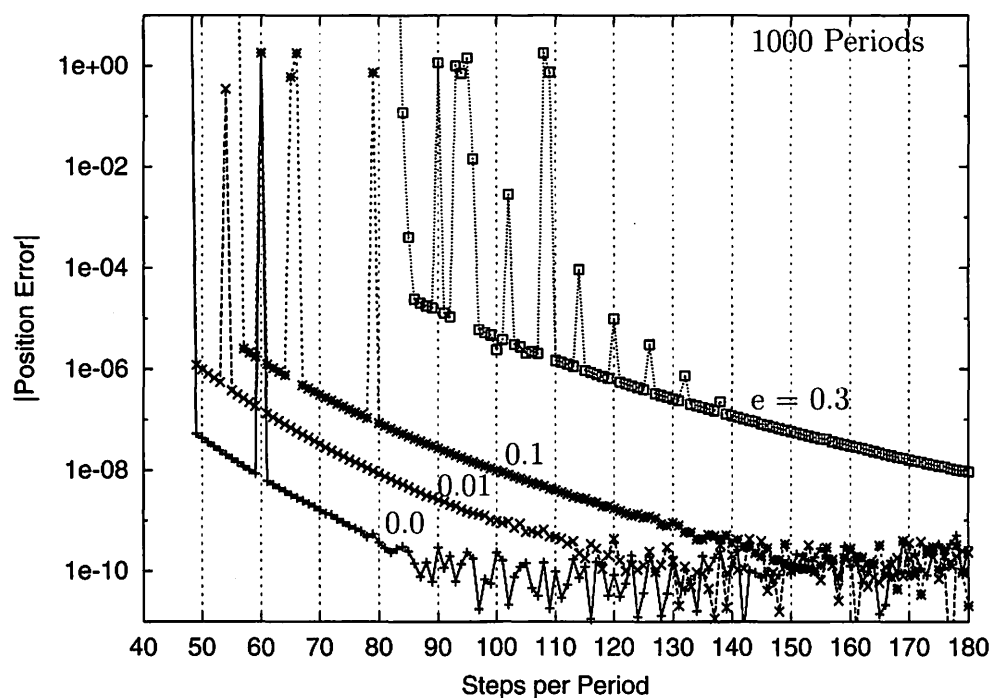


Fig.2 Quinlan and Tremaine 10th Order Keplerian Motion

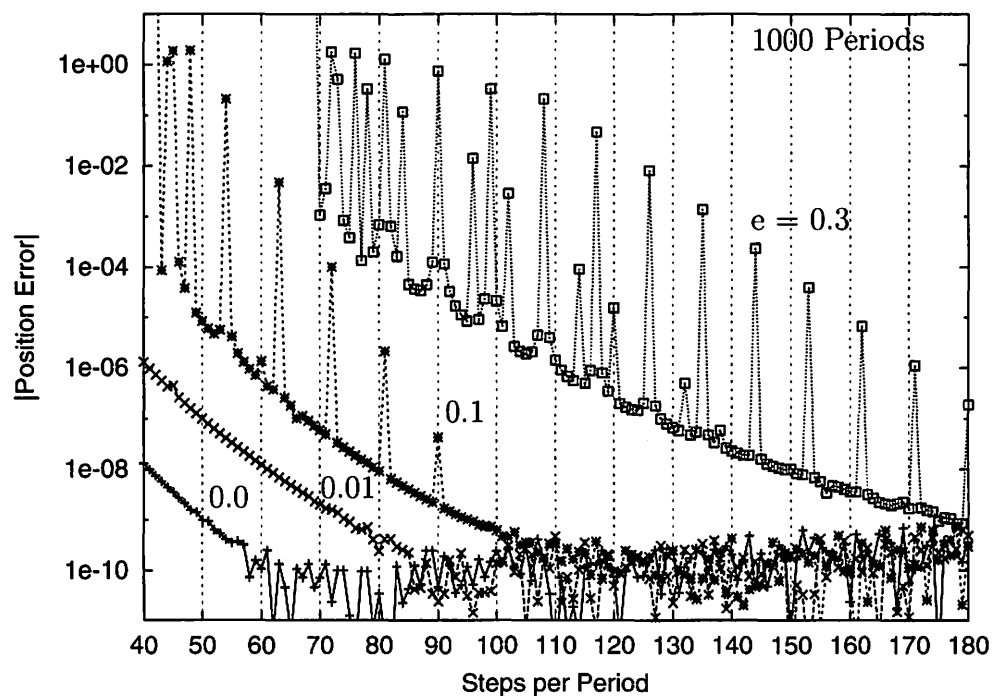


Fig.3 Quinlan and Tremaine 12th Order Keplerian Motion

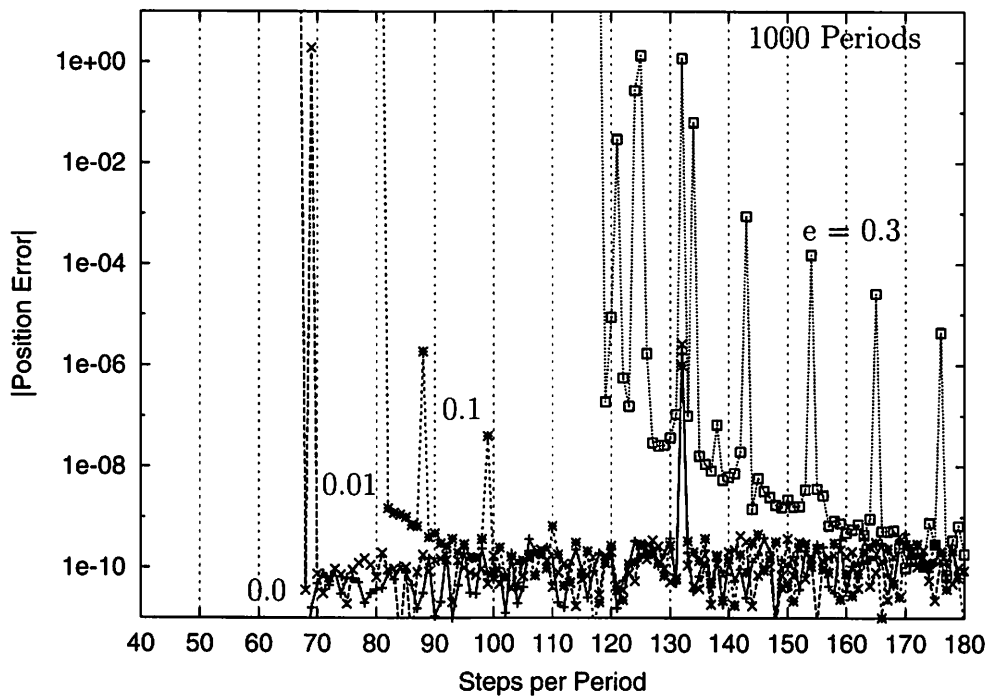


Fig.4 Quinlan and Tremaine 14th Order Keplerian Motion

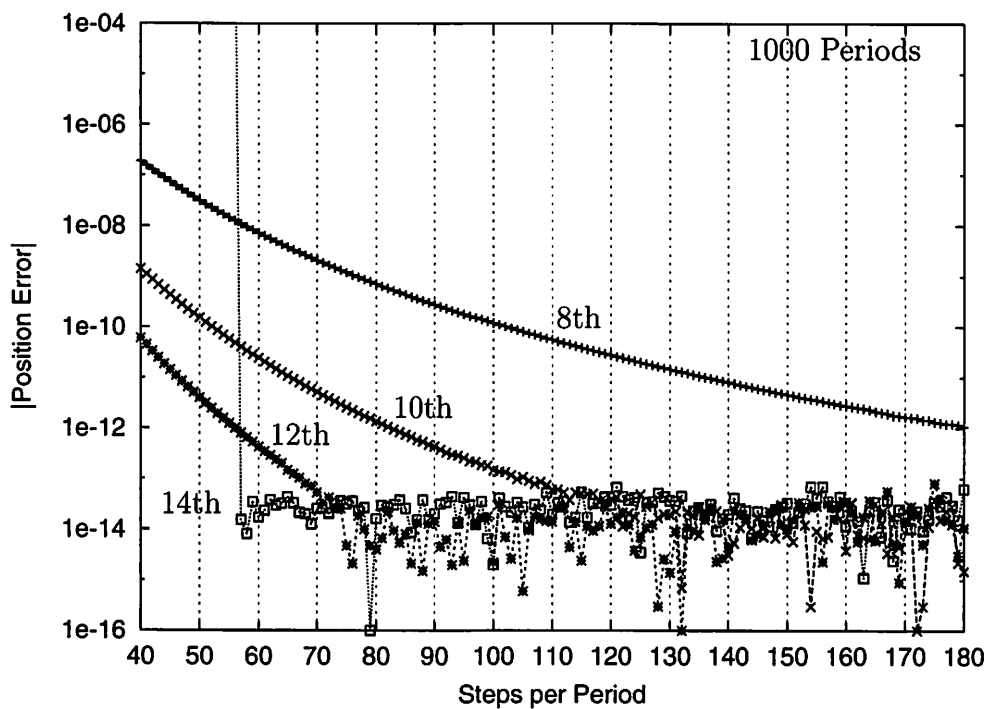


Fig.5 Quinlan and Tremaine 8th, 10th, 12th, 14th Order Harmonic Oscillator

ケプラー運動の場合、次数および離心率の大きさに依存して共鳴現象が現れるのがわかる。調和振動子の場合には共鳴現象は現れない。K-S 変換によってケプラー運動を調和振動子化した場合でも、方程式の性質が調和振動子と同じになるため共鳴現象が起こらなくなる。

### 3 対称型公式の作り方

まず、良く知られている Störmer-Cowell 型の Explicit の場合の一般系は以下のような形である。

$$x_{n+1} - 2x_n + x_{n-1} = h^2(\beta_0 f_n + \cdots + \beta_k f_{n-k}) \quad (1)$$

次数によって右辺の項の形が変わる。

一方、対称型公式の Explicit の場合は次のような形になる。

$$\begin{aligned} x_{n+1} + \alpha_0 x_n + \alpha_1 x_{n-1} + \alpha_2 x_{n-2} + \alpha_1 x_{n-3} + \alpha_0 x_{n-4} + x_{n-5} = \\ = h^2(\beta_0 f_n + \beta_1 f_{n-1} + \beta_2 f_{n-2} + \beta_1 f_{n-3} + \beta_0 f_{n-4}) \end{aligned} \quad (2)$$

右辺の項の個数 = 左辺の項の個数 - 2

係数  $\beta$  の値は母関数  $G(t)$  を使って簡単に求められる。

例) Quinlan and Tremaine の 8 次の公式

$$G(t) = \frac{1 - 2t + 2t^2 - 1t^3 + 0t^4 - 1t^5 + 2t^6 - 2t^7 + t^8}{t(\log t)^2} \quad (3)$$

これを  $t = 1$  のまわりでテーラー展開すれば係数が得られる。 $t$  の 0 次の係数が  $f_n$  の係数に対応する。

$$\begin{aligned} x_{n+1} - 2x_n + 2x_{n-1} - x_{n-2} + 0x_{n-3} - x_{n-4} + 2x_{n-5} - 2x_{n-6} + x_{n-7} = \\ = \frac{17671}{12096} f_n - \frac{3937}{2016} f_{n-1} + \frac{20483}{4032} f_{n-2} - \frac{12629}{3024} f_{n-3} \\ + \frac{20483}{4032} f_{n-4} - \frac{3937}{2016} f_{n-5} + \frac{17671}{12096} f_{n-6} \end{aligned} \quad (4)$$

次に、新しく作った公式の安定な刻み幅の最大値を解析的に調べる。

それには調和振動子  $\ddot{x} = -x = f$  を用いて

$$x_{n+1} + \alpha_0 x_n + \cdots + \alpha_0 x_{n-k+1} + x_{n-k} = \quad (5)$$

$$= -H^2(\beta_0 x_n + \beta_1 x_{n-1} + \cdots + \beta_1 x_{n-k-2} + \beta_0 x_{n-k-1}) \quad (6)$$

$x_i \Rightarrow z^i$  と置き換えて  $H^2$  の値による  $z$  の多項式の根の振舞いを考える。

初期値  $H^2 = 0$  のときの根の場所

- 対称型  
複素平面の実数軸上の 1 に重根、他は単位円上
- Störmer-Cowell  
複素平面の実数軸上の 1 に重根、0 に多重根

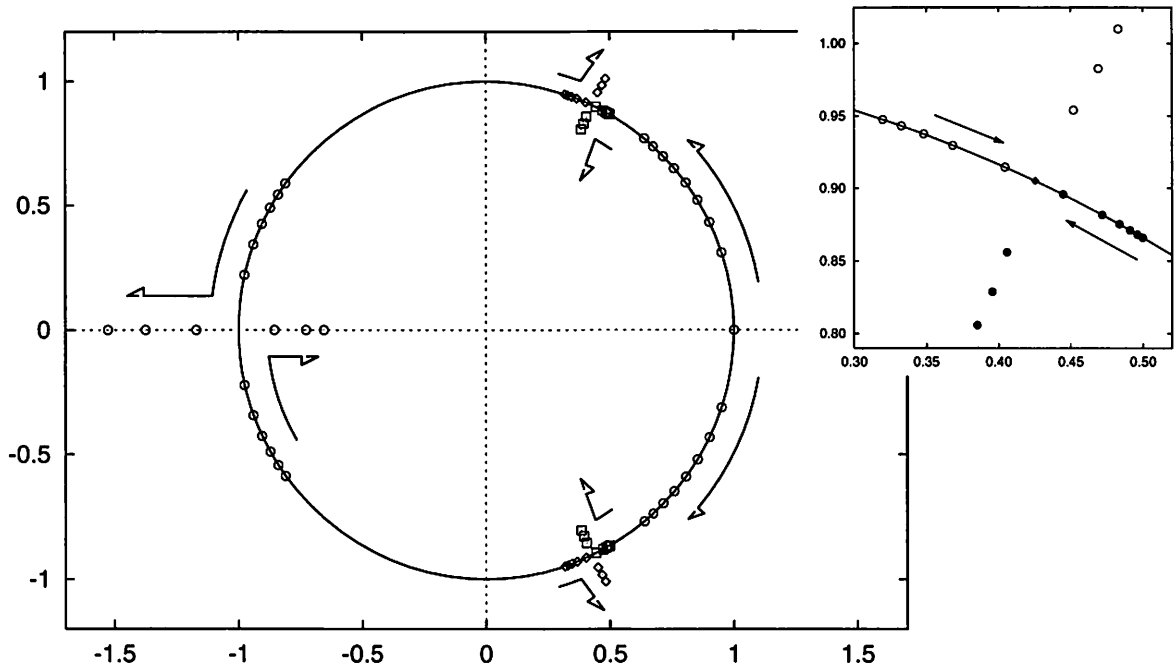


Fig.6 Quinlan and Tremaine 8th Order

Position of the roots by the value of  $H^2$  in a complex plane

$H^2$  の値を増大させていき多項式の根が一番最初に単位円の外に出るときの値が安定な刻み幅の最大値になる。Fig.6 参照

## 4 刻み幅共鳴の出現場所

Quinlan and Tremaine 8次の場合の  $H^2 = 0$  の根は

$$\begin{aligned}
 & z^8 - 2z^7 + 2z^6 - z^5 + 0z^4 - z^3 + 2z^2 - 2z + 1 = \\
 & = \left( z^2 - 2z + 1 \right) \left( z^2 - 2 \cos \left[ \frac{2\pi}{6} \right] z + 1 \right) \left( z^2 - 2 \cos \left[ \frac{2\pi}{5} \right] z + 1 \right) \left( z^2 - 2 \cos \left[ \frac{4\pi}{5} \right] z + 1 \right)
 \end{aligned} \tag{7}$$

Quinlan によると共鳴は

$$\frac{\text{Steps}}{\text{Period}} = \frac{2\pi}{\theta} \times n \quad , \quad n = 1, 2, \dots \quad , \quad \theta \text{は上の場合} \frac{2\pi}{6}, \frac{2\pi}{5}, \frac{4\pi}{5}$$

の場所に出現し  $n$  が大きくなるにつれて共鳴の大きさは小さくなっていく。

また、 $\frac{2\pi}{6} \times 6$ 、 $\frac{2\pi}{5} \times 5$  でそれぞれ1周、5、6の最小公倍数30の倍数で共鳴が起こるとなっている、これを不安定と呼び、運動方程式を離心率の冪乗で展開したものに由来するとなっている。こちらのほうが仕組みが複雑である。

## ● 数値実験

8 次の場合はフリーパラメーターは3つ、値は  $0 < \theta < 2\pi$  でそれぞれ異なるものを選ぶ。  
刻み幅の増分を  $0.1/\text{周期}$ 、離心率は  $e = 0.3$ 、積分期間は1000周で計算する。

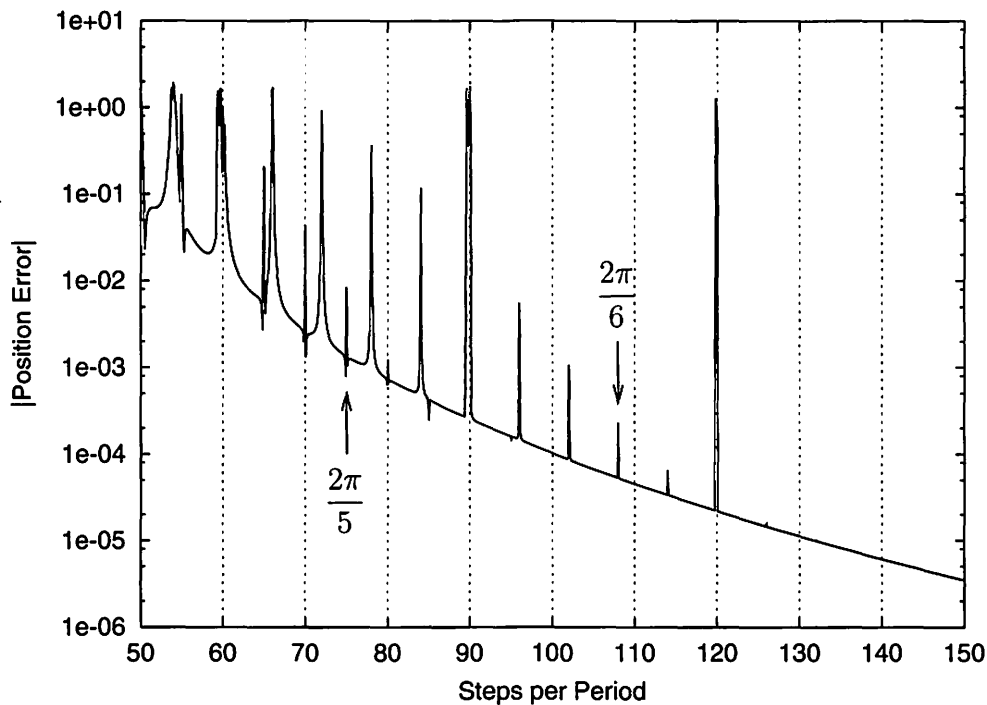


Fig.7 Example 1 Quinlan and Tremaine

60、90、120 が不安定によるもの、 $\frac{4\pi}{5}$  に由来する共鳴はこの条件では判別できない。

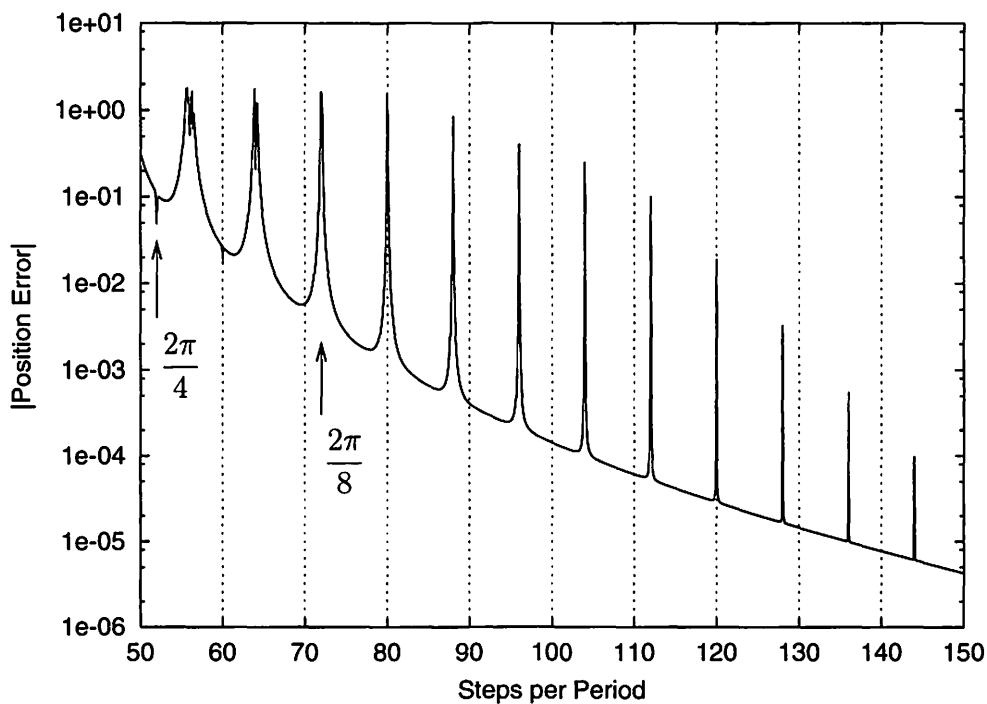


Fig.8 Example 2  $\theta = \frac{2\pi}{8}, \frac{2\pi}{4}, \frac{2\pi}{3}$

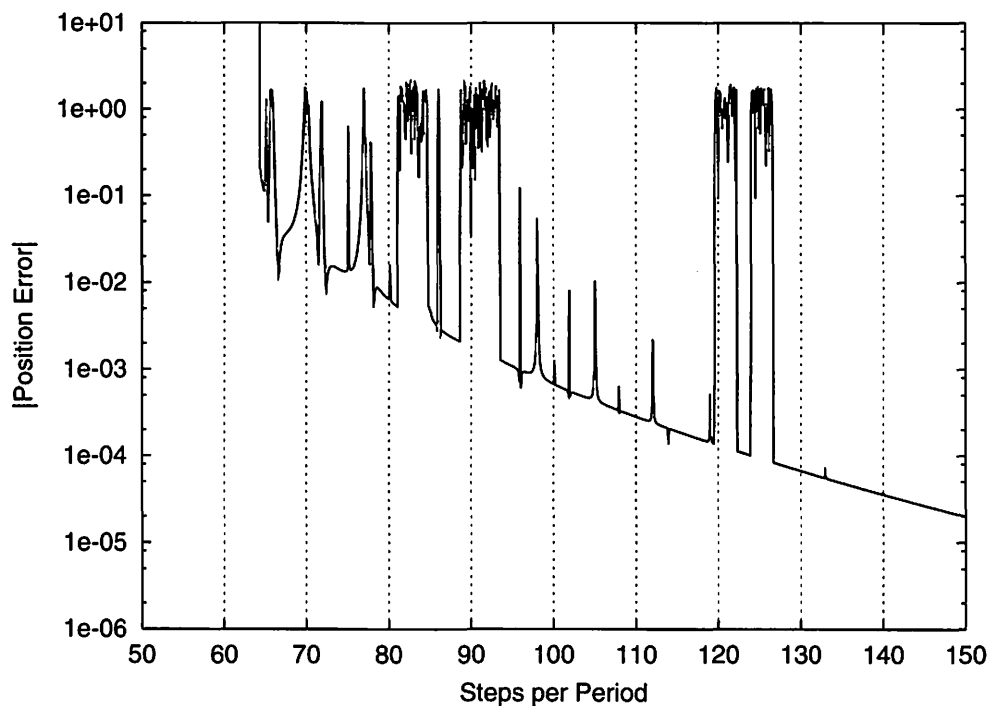


Fig.9 Example 3  $\theta = \frac{2\pi}{7}$  ,  $\frac{2\pi}{6}$  ,  $\frac{2\pi}{5}$

## 5 問題点

- 根の配置を工夫しても共鳴を除去することはできないが、共鳴と共鳴の間隔の調節は可能である。
- 無摂動のケプラー運動に対して共鳴の場所が予測できたとしても、摂動がある場合周期の変動とともに共鳴の位置が移動してしまうと使用する刻み幅の決定が困難になる。
- 今の時点で対称型の多段法を使う場合、K-S 変換が可能な問題には適用可能だが、K-S 変換せずに対称型の多段法を使うことは危険である。

## 6 今後の方針

対称型の多段法をケプラー運動に対して安心して使用できるようにするために

- 線形多段法の公式を拡張
- 要素変化法に一階の微分方程式用の多段法の適用

等を考えている。

## 参考文献

- Arakida, H. and Fukushima, T., 2000, AJ, 120, 3333
- Evans, N.W. and Tremaine, S., 1999, AJ, 118, 1888
- Fukushima, T., 1998, in Proc. 30th Symp. Celest. Mech., 229
- Fukushima, T., 1999, in IAU Colloq., 173, 309
- Hairer, H. , Nørsett, S.P. and Wanner, G., 1987, Solving Ordinary Differential Equations I, Springer, Berlin
- Lambert, J.D. and Watson, I.A., 1976, J. Inst. Maths. Applics., 18, 189
- Quinlan, G.D., 1999, astro-ph/9901136
- Quinlan, G.D. and Tremaine, S., 1990, AJ, 100, 1694



# **Response of Lifespan of Organisms to Secularly Changing Environment using a New Dynamical Model**

Toshihiro Handa<sup>1</sup>, Kiyotaka Tanikawa<sup>2</sup>, Takashi Ito<sup>2</sup>

<sup>1</sup>Institute of Astronomy, University of Tokyo

Mitaka, Tokyo 181-0015, Japan

<sup>1</sup>National Astronomical Observatory

Mitaka, Tokyo 181-8588, Japan

## **abstract**

A relation between lifespan of a species and its environment is studied using a dynamical model. We make a simple model of a species under single parameter environment. Our model has parameters which describe the dispersion of character taken over between a parent and a child, width of allowance to survive under a given environment, changing rate of the environment, and energy flux to support whole bodies in a species. We do not introduce any direct interaction between any individuals or any other species to focus our attention on lifespan by environment.

Under linearly changing environment the population of a species grows exponentially or extinct without any limitation of reproduction. With limit of energy supply the population can be stable and optimal lifespan always exists which gives the largest population.

## **1 Introduction and Background Questions**

We believe all species evolve through mutation and natural selection. Variety of species is thought to be a result of combination work by heredity, mutation and natural selection.

Using technology in molecular biology we can read DNA code itself. One of the important results in this field is that many species have a telomere in its DNA, which is reduced at a cell division and limits the number of cell divisions. The function of telomere seems irrational, because due to the function every individual of a species cannot keep itself forever even under an ideal environment; it requires reconstruction of individuals from the first stage. Why individuals of actual life form has finite lifespan? Is it possible to understand this is a result of competition between species with different lifespan under the same environment? Is there any optimal lifespan to maximize population of a species?

On the early physical and chemical Earth environment around organisms, or lifeform, has secularly changed: the amount of dangerous radiation (Boothroyd et al., 1990), total area of continents (Condie, 1989), the total amount of carbon dioxide in the atmosphere and in the sea (Tajika, 1992), the average temperature of the atmosphere and in the interior of the earth, the amount of atmospheric ozone, the spin rate of Earth (Tajika, 1992), strength and frequency of tidal force by the Moon, the gravitational flattening of the Earth, and solar insolation (Ito et al., 1993)

To survive for a long time any species must have adapted to the changing environment (Losos et al., 1997), although a species must keep its character over many generations to identify itself. A genetic system and reproduction with dispersion of genetic characters should overcome this contradictory requirement. If it is, an origin of telomere's function should be thought as a result of an evolutionary process and there should be some relations between environment and lifespan.

To access this problem we make simple evolution models of a species and study its response of the number of individuals to a change of environment. Our models consider effects of heredity and natural selection by a given environment but are not introduced any direct competition between any individual of a species. Our models are not introduced any direct effect of rival species.

In this report we construct an analytic model and its response under secular (linearly changing) environment as a first step.

## 2 Constant Reproduction Rate Model

### 2.1 Basic Assumptions

To construct a model we put the following assumptions.

**Basic assumption 1. (one-dimensional environment)** Environment around the species can be parametrized along a single dimensional axis. The parametrized environment is assigned by a parameter  $x$ , and a genetic character of an individual is parametrized along the same axis.

**Basic assumption 2. (neutral evolution)** (Kimura, 1983) A genetic character of an individual is independent of its environment. Propagation of a genetic character of individuals follows a stochastic process.

**Basic assumption 3. (Markov process)** A genetic character of an individual is set only by that of its direct parent but it deviates from parent's character as a stochastic process.

**Basic assumption 4. (statistical approach)** The character of an individual is only evaluated whether it will survive until it will reproduce its own children. The

survivability is determined only by a difference between the genetic character of an individual and an environment.

**Basic assumption 5. (characterization of a species)** The species is characterized by a property of propagation of a genetic character and survivability of individuals.

## 2.2 Assumptions and Parameters for Formulation

Under the basic assumptions mentioned above we set following parameters and functions of their relations.

**Assumption for formulation 1. (instantaneous evolution)** All parents of a species are extinct just after birth of their children. Overlapping period between two generations is negligibly short. Natural selection is done just after the birth time of individuals under the environment at that time.

**Assumption for formulation 2. (Gaussian dispersion of a genetic character)** Based on the basic assumptions of neutral evolution, Markov process, and characterization of a species, a parent with character  $x'$  bears a child with character  $x$  at a probability of  $g(x - x')$ . The probability  $g(x - x')$  is assumed to be a Gaussian with an e-folding width  $\sigma_g$ .

**Assumption for formulation 3. (Gaussian survivability)** Based on the basic assumptions of neutral evolution, Markov process, and characterization of a species, a child with character  $x$  can survive until making next generation under environment  $x_E$  with a probability  $s(x - x_E)$ . The survivability  $s(x - x')$  is assumed to be a Gaussian with an e-folding width  $\sigma_E$ .

**Assumption for formulation 4. (the same number of child for each individual)** The number of children for an individual parent, or reproduction rate, is  $c$ . Based on the basic assumption of statistical approach, the value of  $c$  is the same for all individuals of a species. Based on the basic assumption of characterization of a species  $c$  is constant for any generation.

Using them we will derive a population equation between two generations. We use the following notation.

- Population of generation  $k$  is  $P_k$ .
- Environment at the time  $t$  is  $x_E(t)$ .
- Distribution of a character for generation  $k$  is  $N(k, x)$ .
- The origin of time  $t = 0$  is set to the birth time of generation 0. The origin of environment is set to be the value of  $x_E(t)$  at  $t = 0$ , or  $x_E(t = 0) = 0$ .

- Lifespan of an individual, which is defined the time between two generations,  $\tau$ . The lifetime  $\tau$  is constant for all generations based on the basic assumption of characterization of a species. It means birth time of generation  $k$  is  $t = k\tau$ .

Using all these notations the distribution of a character for generation  $k + 1$  is given by

$$N(k + 1, x) = s(x - x_E) \int N(k, x') cg(x - x') dx'. \quad (1)$$

Population of generation  $k + 1$  is

$$\begin{aligned} P_k &= \int N(k, x) dx \\ P_{k+1} &= \int N(k + 1, x) dx. \end{aligned} \quad (2)$$

### 2.3 Recurrence relations between two generations

We derive a recurrence relations between generations  $k$  and  $k + 1$ . For the first step we derive a steady state solution under steady environment to find a basic property of solutions. We set a genetic character distribution of generation  $k$  is a Gaussian distribution at first. (Actually this is a weak constraint as shown in a later section.) If the center and e-folding width of the distribution are  $x_k$  and  $\sigma_k$ , respectively, we get

$$N(k, x - x_k) = P_k \text{Gauss}(x - x_k, \sigma_k), \quad (3)$$

where  $\text{Gauss}(x, \sigma)$  is a nomilized Gaussian defined as

$$\text{Gauss}(x, \sigma) = \frac{1}{\sqrt{\pi}\sigma} \exp\left(-\frac{x^2}{\sigma^2}\right). \quad (4)$$

Using assumptions for formulation we put

$$g(x - x') = \text{Gauss}(x - x', \sigma_g) \quad (5)$$

and

$$s(x - x_E(t)) = \sqrt{\pi}\sigma_E \text{Gauss}(x - x_E(t), \sigma_E). \quad (6)$$

Using the basic assumptions  $\sigma_E, \sigma_g$ , and  $c$  are constant.

With (5) and (6) Equation (1) gives a recurrence relation,

$$\begin{aligned} N(k + 1, x) &= s(x - x_E(t)) \int N(k, x') cg(x - x') dx' \\ &= P_k \sqrt{\pi}\sigma_E \text{Gauss}(x - x_E(t), \sigma_E) \int \text{Gauss}(x' - x_k, \sigma_k) c \text{Gauss}(x - x', \sigma_g) dx' \\ &= c P_k \frac{\sigma_E}{\sqrt{\sigma_k^2 + \sigma_g^2 + \sigma_E^2}} \exp\left(-\frac{(x_k - x_E(t))^2}{\sigma_k^2 + \sigma_g^2 + \sigma_E^2}\right) \text{Gauss}\left(x - \frac{\sigma_E^2 x_k + (\sigma_k^2 + \sigma_g^2) x_E}{\sigma_k^2 + \sigma_g^2 + \sigma_E^2}, \frac{\sqrt{\sigma_k^2 + \sigma_g^2} \sigma_E}{\sqrt{\sigma_k^2 + \sigma_g^2 + \sigma_E^2}}\right) \\ &= P_{k+1} \text{Gauss}(x - x_{k+1}, \sigma_{k+1}) \end{aligned} \quad (7)$$

This means the distribution of generation  $k+1$  is also Gaussian if that of generation  $k$  is Gaussian. And the last equality in (7) comes, therefore, from the corresponding expression of (3) for the generation  $k+1$ . Comparison of both sides of the recurrence relation of distributions gives three recurrence relations between parameters.

$$P_{k+1} = c \frac{\sigma_E}{\sqrt{\sigma_k^2 + \sigma_g^2 + \sigma_E^2}} \exp \left( -\frac{(x_k - x_E(t))^2}{\sigma_k^2 + \sigma_g^2 + \sigma_E^2} \right) P_k, \quad (8)$$

$$\sigma_{k+1} = \sigma_E \frac{\sqrt{\sigma_k^2 + \sigma_g^2}}{\sqrt{\sigma_k^2 + \sigma_g^2 + \sigma_E^2}}, \quad (9)$$

and

$$x_{k+1} = \frac{\sigma_E^2 x_k + (\sigma_k^2 + \sigma_g^2) x_E}{\sigma_k^2 + \sigma_g^2 + \sigma_E^2}. \quad (10)$$

### 3 General Property of the Constant Reproduction Rate Model

#### 3.1 Convergency of distribution

##### 3.1.1 Width of distribution

The recurrence equation of  $\sigma_k$  (9) shows that a width of the genetic character distribution, or  $\sigma_k$ , is independent of current environment  $x_E$ , and its numerical trend depends only on two generation-independent parameters  $\sigma_g$  and  $\sigma_E$ . We can easily get the general solution of Equation (9) as

$$\sigma_k^2 = \frac{1}{\left( \frac{\sigma_g^2 + \sigma_E^2 + \sigma_\infty^2}{\sigma_E^2} \right)^{2k} \left( \frac{1}{\sigma_0^2 - \sigma_\infty^2} + \frac{1}{2\sigma_\infty^2 + \sigma_g^2} \right) - \frac{1}{2\sigma_\infty^2 + \sigma_g^2}} + \sigma_\infty^2, \quad (11)$$

where

$$\sigma_\infty^2 = \frac{\sigma_g^2}{2} \left( \sqrt{1 + \left( 2 \frac{\sigma_E}{\sigma_g} \right)^2} - 1 \right), \quad \text{or} \quad \sigma_\infty = \sqrt{\frac{\sigma_g^2}{2} \left( \sqrt{1 + \left( 2 \frac{\sigma_E}{\sigma_g} \right)^2} - 1 \right)}. \quad (12)$$

This always converges to a finite value,  $\sigma_k \rightarrow \sigma_\infty$ , because the base of the power,  $\frac{\sigma_g^2 + \sigma_E^2 + \sigma_\infty^2}{\sigma_E^2} > 1$ .

It means that width of the final distribution converges to  $\sigma_\infty$ , if the initial distribution is Gaussian with any width  $\sigma_0$  or center position  $x_0$ .

For large  $k$  it is enough to only consider the case of  $\sigma_k = \sigma_\infty$ .

### 3.1.2 Center of distribution

Using the result of  $\sigma_k = \sigma_\infty$  for large  $k$  we can rewrite the recurrence relation (10) of the parameter  $x_k$ , which represents the center of a genetic character distribution for generation  $k$ , as the following:

$$x_{k+1} = \frac{\sigma_E^2 x_k + (\sigma_\infty^2 + \sigma_g^2) x_E}{\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2}. \quad (13)$$

This is reduced to

$$\begin{aligned} (\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2) x_{k+1} &= \sigma_E^2 x_k + (\sigma_\infty^2 + \sigma_g^2) x_E(t) \\ (\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2) (x_{k+1} - x_E(t)) &= \sigma_E^2 (x_k - x_E(t)). \end{aligned} \quad (14)$$

In this report we only consider the following two simplest cases on a change of the environment as a first step of studies.

#### Case of steady environment (constant $x_E(t)$ for $t$ )

A steady environment is represented by  $x_E(t) = x_E(0) = x_E$ . Using Equation (14), the parameter  $x_k - x_E$  is a geometrical series with a constant ratio  $\frac{\sigma_E^2}{\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2} < 1$ , and  $x_k$  always converges as  $x_k \rightarrow x_E(t) = x_E$ . Note that  $x_E(t)$  is a constant because of a steady environment.

#### Case of linearly changing environment ( $x_E(t)$ changes with a constant velocity)

A linearly changing environment is represented by  $x_E(t) = x_E(k\tau) = k(x_E(1) - x_E(0)) = k\dot{x}_E\tau$ . In this case Equation (14) is reduced to be

$$x_{k+1} - (k+1)\dot{x}_E\tau + \frac{\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2}{\sigma_\infty^2 + \sigma_g^2} \dot{x}_E\tau = \frac{\sigma_E^2}{\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2} \left( x_k - k\dot{x}_E\tau + \left( \frac{\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2}{\sigma_\infty^2 + \sigma_g^2} \right) \dot{x}_E\tau \right).$$

Then a general term of the series is

$$x_k - x_E(t) + \frac{\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2}{\sigma_\infty^2 + \sigma_g^2} \dot{x}_E\tau = \left( \frac{\sigma_E^2}{\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2} \right)^{k-1} \left( x_0 - x_E(0) + \frac{\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2}{\sigma_\infty^2 + \sigma_g^2} \dot{x}_E\tau \right).$$

Considering a base of the power term is  $\frac{\sigma_E^2}{\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2} < 1$ , the series of  $x_k - x_E(t)$  always converges  $x_k - x_E(t) \rightarrow -\frac{\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2}{\sigma_\infty^2 + \sigma_g^2} \dot{x}_E\tau$  at  $k \rightarrow \infty$ . That is  $x_k \rightarrow \left( k - \frac{\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2}{\sigma_\infty^2 + \sigma_g^2} \right) \dot{x}_E\tau$ .

This means that center of the distribution  $x_k$  asymptotically converges to the position with constant offset  $(1 + \frac{\sigma_E^2}{\sigma_\infty^2 + \sigma_g^2}) \dot{x}_E\tau$  from the environment  $x_E(t)$  in the case of a linearly changing environment.

Therefore we define a new notation  $x_{\text{offset}}$  by

$$x_{\text{offset}} = \frac{\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2}{\sigma_\infty^2 + \sigma_g^2} \dot{x}_E\tau = \left( 1 + \frac{\sigma_E^2}{\sigma_\infty^2 + \sigma_g^2} \right) \dot{x}_E\tau \quad (15)$$

in the case of linearly changing environment.

### 3.2 Dependence on an Initial Distribution

In the previous section we show that any Gaussian distribution of the model converges to a Gaussian distribution with the same width and the same offset from an environment at  $t \rightarrow \infty$ , or  $k \rightarrow \infty$ , if the environment changes linearly.

We will show that a convergent distribution is also the same for any practical distribution. From a view point as a transformation on the function  $N(k, t)$  the recurrence equation (7) is a linear process. If we take Dirac's delta function as a initial distribution, or  $N(0, x) = \delta(x)$  at  $k = 0$ ,  $N(1, x)$  is obviously a Gaussian. Therefore any linear combination of the delta functions, which can be called as any practical functions, converges to the same Gaussian. The converged distribution  $N(k = \infty, x)$  depends only on genetic character dispersion width  $\sigma_g$  and width of survivability  $\sigma_E$  for any initial distributions.

### 3.3 Population under Linearly Changing Environment

#### 3.3.1 Growth of population of a species

Through the discussion above we get that any actual distribution of genetic character always converges to a Gaussian distribution with  $\sigma_k = \sigma_\infty$  for large  $k$ . If the environment is steady, the center of the distribution is  $x_k = x_E = 0$  for large  $k$ . If the environment changes linearly,  $x_k - x_E(t) = x_{\text{offset}}$  for large  $k$ .

In this section we discuss population of a species  $P_k$ . Using converged values  $\sigma_k = \sigma_\infty$  and  $x_k - x_E(t) = x_{\text{offset}}$  the recurrence equation (8) is reduced to be

$$\begin{aligned} P_{k+1} &= c \frac{\sigma_E}{\sqrt{\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2}} \exp\left(-\frac{x_{\text{offset}}^2}{\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2}\right) P_k \\ &= c \frac{1}{\sqrt{\frac{\sigma_\infty^2 + \sigma_g^2}{\sigma_E^2} + 1}} \exp\left(-\frac{(\dot{x}_E \tau)^2}{\sigma_g^2}\right) P_k. \end{aligned} \quad (16)$$

This is a simple geometrical series. Growth of  $P_k$  is dependent only on a value of  $c \frac{1}{\sqrt{\frac{\sigma_\infty^2 + \sigma_g^2}{\sigma_E^2} + 1}} \exp\left(-\frac{(\dot{x}_E \tau)^2}{\sigma_g^2}\right)$ . In other words growth of  $P_k$  is categorized by the following three cases using a parameter  $c_{\text{stat}}$  as defined by

$$c_{\text{stat}} = \sqrt{1 + \frac{\sigma_\infty^2 + \sigma_g^2}{\sigma_E^2}} \exp\left(\frac{(\dot{x}_E \tau)^2}{\sigma_g^2}\right). \quad (17)$$

- In the case of  $c > c_{\text{stat}}$ : Population  $P_k$  grows exponentially, and  $P_k \rightarrow \infty$  at  $k \rightarrow \infty$ .
- In the case of  $c = c_{\text{stat}}$ : Population  $P_k$  is constant.

- In the case of  $c < c_{\text{stat}}$ : Population  $P_k$  decreases exponentially, and  $P_k \rightarrow 0$  at  $k \rightarrow \infty$ . In this case the species will be extinct.

In this model a population of a species is constant only  $c = c_{\text{stat}}$ . The system is exponentially unstable.

### 3.3.2 Relation between population and model parameters

Using expression of  $c_{\text{stat}}$  in (17) we discuss the response of population of a model species and its physical interperations.

Using Equation (12) we get

$$\frac{\sigma_\infty^2 + \sigma_g^2}{\sigma_E^2} = \frac{1}{2} \frac{\sigma_g^2}{\sigma_E^2} \left( \sqrt{1 + 4 \frac{\sigma_E^2}{\sigma_g^2}} + 1 \right), \quad (18)$$

and

$$\frac{\sigma_g^2}{\sigma_E^2} = \frac{\left( \frac{\sigma_\infty^2 + \sigma_g^2}{\sigma_E^2} \right)^2}{1 + \frac{\sigma_\infty^2 + \sigma_g^2}{\sigma_E^2}}. \quad (19)$$

Then Equation (17) is expressed by  $u = \frac{\sigma_\infty^2 + \sigma_g^2}{\sigma_E^2}$  as the following:

$$c_{\text{stat}}^2 = (u + 1) \exp \left( 2 \frac{1 + u}{u^2} \Delta u_E \right), \quad (20)$$

where  $\Delta u_E = \left( \frac{\dot{x}_E \tau}{\sigma_E} \right)^2$ .

The equation (20) clearly shows  $c_{\text{stat}}$  becomes minimum at  $\Delta u_E = 0$ , or  $\dot{x}_E \tau = 0$ , because  $u$  is independent of  $\Delta u_E = 0$ .

Using

$$\frac{\partial}{\partial \left( \frac{\sigma_g^2}{\sigma_E^2} \right)} = \frac{du}{d \left( \frac{\sigma_g^2}{\sigma_E^2} \right)} \frac{\partial}{\partial u} = \frac{(u + 1)^2}{u(u + 2)} \frac{\partial}{\partial u}$$

we get

$$\frac{\partial(c_{\text{stat}}^2)}{\partial \left( \frac{\sigma_g^2}{\sigma_E^2} \right)} = \exp \left( 2 \frac{u + 1}{u^2} \Delta u_E \right) \frac{(u + 1)^2}{u^4(u + 2)} \left( u^3 - 2\Delta u_E(u + 2)(u + 1) \right). \quad (21)$$

A right hand of Equation (21) shows that a sign of  $\frac{\partial(c_{\text{stat}}^2)}{\partial \left( \frac{\sigma_g^2}{\sigma_E^2} \right)}$  is the same as that of  $u^3 - 2\Delta u_E(u + 2)(u + 1)$ , because of  $u > 0$ . Using  $\Delta u_E \geq 0$  only one positive solution on  $u$  of

$$u^3 - 2\Delta u_E(u + 2)(u + 1) = 0 \quad (22)$$

always exists for any value of  $\Delta u_E = \frac{(\dot{x}_E \tau)^2}{\sigma_E^2} > 0$ . Setting the solution as  $u = u_{\text{opt}}$ ,  $c_{\text{stat}}^2$  is minimum at  $u = u_{\text{opt}}$ . We can include the same statement formally for the case



of  $\Delta u_E = \frac{(\dot{x}_E \tau)^2}{\sigma_E^2} = 0$ , because it gives  $c_{\text{stat}}^2 = u + 1$  from Equation (20) and  $c_{\text{stat}}$  is minimum at  $u = 0$ .

Using Equation (19), we get  $\left(\frac{\sigma_g}{\sigma_E}\right)_{\text{opt}}$  corresponding to  $u_{\text{opt}}$ , and if  $u_{\text{opt}} \geq 0$ ,  $\left(\frac{\sigma_g}{\sigma_E}\right)_{\text{opt}}$  is always positive or zero. This means that there is only one  $\frac{\sigma_g}{\sigma_E} = \left(\frac{\sigma_g}{\sigma_E}\right)_{\text{opt}} \geq 0$  which gives minimum  $c_{\text{stat}}$  for any  $\frac{\dot{x}_E \tau}{\sigma_E^2}$ .

Using Equations (18) and (19), Equation (22) is written by  $y = \sigma_g^2 / \sigma_E^2$  as

$$y^3 - 4\Delta u_E^2(y + 2) = 0,$$

or

$$\left(\frac{\sigma_g^2}{\sigma_E^2}\right)^3 - 4\left(\frac{\dot{x}_E \tau}{\sigma_E}\right)^4 \left(\frac{\sigma_g^2}{\sigma_E^2} + 2\right) = 0. \quad (23)$$

At the only one solution of Equation (23) on  $\frac{\sigma_g}{\sigma_E}$ , the reproduction rate  $c_{\text{stat}}$  is minimum.

Using Equations (12) and (17)

$$c_{\text{stat}}^2 = \left(\frac{\sigma_g^2}{2\sigma_E^2} \left(1 + \sqrt{1 + 4\frac{\sigma_E^2}{\sigma_g^2}}\right) + 1\right) \exp\left(2\frac{\dot{x}_E \tau}{\sigma_g^2}\right).$$

The partial derivative with respect to  $\sigma_E^2$  is

$$\frac{\partial(c_{\text{stat}}^2)}{\partial(\sigma_E^2)} = -\frac{1}{4\sigma_E^4 \sqrt{1 + 4\frac{\sigma_E^2}{\sigma_g^2}}} \sigma_g^2 \left(1 + \sqrt{1 + \frac{4\sigma_E^2}{\sigma_g^2}}\right)^2 \exp\left(2\frac{(\dot{x}_E \tau)^2}{\sigma_g^2}\right) < 0.$$

It means that  $c_{\text{stat}}$  always decreases if only  $\sigma_E$  increases.

Summary of parameter dependences of  $c_{\text{stat}}$  and its physical interpretation are the following:

**velocity of environment change:** Comparison between two species with the same model parameters under different velocities of environment change  $\dot{x}_E$  gives that the faster change requires the larger  $c_{\text{stat}}$ . The higher reproduction rate is required under the faster changing environment to keep the species. This is a trivial result.

**width of survivability:** Comparison between two species with different widths of survivability  $\sigma_E$  under the same conditions for other parameters gives the smaller  $\sigma_E$  requires the larger  $c_{\text{stat}}$ . Under the same environment the severer or weaker species for deviation from the environment is required the higher reproduction rate. This is also a trivial result.

**lifespan:** Comparison between two species with different lifetimes  $\tau$  under the same conditions for other parameters gives the longer  $\tau$  requires the larger  $c_{\text{stat}}$ . Under the same environment the longer-lived species requires the higher reproduction rate.

**width of genetic character dispersion:** Comparison between two species with different widths of genetic character dispersion  $\sigma_g$  under the same conditions for other parameters gives that a species with  $\sigma_g$  given by Equation (22) has the smallest  $c_{\text{stat}}$ . This means that there is a preferable value for  $\sigma_g$  to minimize  $c$ .

## 4 Adaptive Reproduction Rate Model

Our model species with the assumptions given in the section 2 are exponentially unstable. To get a stable model we modify one of the assumptions for formulation.

**Assumption for formulation 4'. (adaptive reproduction rate)** The reproduction rate  $c$  may change between different generations, although the basic assumption of characterization of a species is still valid for other parameters. It should be written as  $c_k$  to show dependence on generation  $k$ . By the basic assumption of statistical approach  $c_k$  is the same for individuals of the species in a single generation.

In this section we assume that reproduction rate  $c_k$  depends on population  $P_k$ . If population increases, reproduction rate decreases, and vice versa. It can be interpreted as an indirect effect of limited logistics.

We investigate models with the following three types of functions.

$$c_k = c_{\text{stat}} \left( 1 + r_{\text{feedback}} \frac{P_0 - P_k}{P_0} \right), \quad (24)$$

$$c_k = c_{\text{stat}} \left( \exp \left( r_{\text{feedback}} \frac{P_0 - P_k}{P_0} \right) \right), \quad (25)$$

$$c_k = c_{\text{stat}} \left( 1 + \frac{2}{\pi} \arctan \left( r_{\text{feedback}} \frac{P_0 - P_k}{P_0} \right) \right). \quad (26)$$

In these functions a parameter  $r_{\text{feedback}}$  controls strength of response.

Under linearly changing environment we estimate genetic character distributions with three types of  $c_k$ 's by numerical calculations.

The results of simulations have been reported in Tanikawa, Handa, and Ito (2000). Figures 6 and 7 in the report shows that final populations are categorized into four types; extinction, dumping to non-zero value, dumped oscillation, and chaotic oscillation. It depends on the feedback factor  $r_{\text{feedback}}$ , which is noted as  $\alpha$  in the report, for each  $\frac{\sigma_E}{\sigma_E}$  and  $\Delta u_E$ , which is noted as  $\sigma_E$  and  $\Delta x_E$  respectively in the report.

## 5 Supply Limited Model

### 5.1 Additional assumptions and modeling

The adaptive reproduction rate models can show stable population, although its population sometimes changes chaotically.

However, the relation between population and reproduction rate is excessively intimate, so we cannot get any explicit relation between lifespan and population. To get such a relation we do not take adaptive reproduction rate but add the following three assumptions for formulations on energy supply for reproduction to the original assumptions.

**Assumption for formulation 5. (constant cost for reproduction)**

Reproduction of a child by an individual parent requires a constant energy  $\epsilon$ .

The required energy comes from a stock which an individual parent has gotten during its lifespan. An actual reproduction rate of generation  $c_k$  is limited by a required energy which must be smaller than energy stock of an individual parent. Using the basic assumption of statistical approach the values of  $\epsilon$  and  $c_k$  are the same for all individuals in a species.

**Assumption for formulation 6. (constant rate to energy collection)**

An individual of a species can take energy for its own stock at the rate of  $w_k$  per unit time.

This means that an individual can make a stock of  $w_k \tau$ . The value of  $w_k$  is limited by  $w_{\text{full}}$ . Using the basic assumption of statistical approach the value of  $w_k$  is the same for all individuals in a species. Using the basic assumption of characterization of a species the value of  $w_{\text{full}}$  is constant for a species.

**Assumption for formulation 7. (limited supply)**

Environment can supply a limited energy per unit time  $W$  for total requirement by all individuals of a species.

Using the basic assumption of statistical approach this means that the value of  $w_k$  is limited by  $W/P_k$  for each individual in a species.

Using these assumptions we modify the equation of population evolution.

Using the assumption for formulation of constant rate to energy collection the total energy stock of an individual is  $w_k \tau$ . By the assumption for formulation of constant cost for reproduction the reproduction rate  $c_k$  is as

$$c_k = \frac{w_k \tau}{\epsilon}. \quad (27)$$

By the assumption for formulation of limited supply

$$w_k P_k = \min(W, w_{\text{full}} P_k),$$

$$w_k = \min\left(\frac{W}{P_k}, w_{\text{full}}\right), \quad (28)$$

$$c_k = \frac{\tau}{\epsilon} \min\left(\frac{W}{P_k}, w_{\text{full}}\right). \quad (29)$$

Using these equations we modify the recurrence equation (8) and get

$$\begin{aligned} P_{k+1} &= \frac{\sigma_E}{\sqrt{\sigma_k^2 + \sigma_g^2 + \sigma_E^2}} \exp\left(-\frac{(x_k - x_E(t))^2}{\sigma_k^2 + \sigma_g^2 + \sigma_E^2}\right) P_k c_k \\ &= \frac{\sigma_E}{\sqrt{\sigma_k^2 + \sigma_g^2 + \sigma_E^2}} \exp\left(-\frac{(x_k - x_E(t))^2}{\sigma_k^2 + \sigma_g^2 + \sigma_E^2}\right) \frac{P_k w_k \tau}{\epsilon} \\ &= \frac{\sigma_E}{\sqrt{\sigma_k^2 + \sigma_g^2 + \sigma_E^2}} \exp\left(-\frac{(x_k - x_E(t))^2}{\sigma_k^2 + \sigma_g^2 + \sigma_E^2}\right) \min(W, w_{\text{full}} P_k) \frac{\tau}{\epsilon}. \end{aligned} \quad (30)$$

## 5.2 Response under linearly changing environment

Using Equation (29),  $c_k = \frac{w_{\text{full}}\tau}{\epsilon}$  is constant if  $P_k < \frac{W}{w_{\text{full}}}$ . In this case growth of population is the same as constant reproduction rate model discussed in the section 2. To avoid the extinction a species must satisfy  $c_k \geq c_{\text{stat}}$ . Using Equation (17) this criterion is

$$\frac{w_{\text{full}}\tau}{\epsilon} \geq \sqrt{1 + \frac{\sigma_{\infty}^2 + \sigma_g^2}{\sigma_E^2}} \exp\left(\frac{(\dot{x}_E\tau)^2}{\sigma_g^2}\right). \quad (31)$$

The criterion (31) is equivalent to

$$\frac{w_{\text{full}}}{\epsilon} \frac{1}{\sqrt{\frac{\sigma_{\infty}^2 + \sigma_g^2}{\sigma_E^2} + 1}} \frac{\sigma_g}{\dot{x}_E} \geq \frac{\exp(\xi^2)}{\xi},$$

where  $\xi = \frac{\dot{x}_E}{\sigma_g}\tau$ . The right hand side of the equation  $\frac{\exp(\xi^2)}{\xi}$  has a minimum value of

$$\frac{\exp(\xi^2)}{\xi} \geq \sqrt{2e}$$

at  $\xi = \frac{1}{\sqrt{2}}$ , where  $e$  is Napier's constant. Therefore lifespan  $\tau > 0$  which satisfies the criterion (31) can exist, if

$$\frac{w_{\text{full}}}{\epsilon} \frac{1}{\sqrt{\frac{\sigma_{\infty}^2 + \sigma_g^2}{\sigma_E^2} + 1}} \frac{\sigma_g}{\dot{x}_E} \geq \sqrt{2e}. \quad (32)$$

As shown in the section 2, if the criterion (32) is satisfied,  $P_k$  always grows exponentially.

After the exponential growth population should reach  $P_k \geq \frac{W}{w_{\text{full}}}$ . In this case Equation (30) is reduced to

$$\begin{aligned} P_{k+1} &= \frac{\sigma_E}{\sqrt{\sigma_k^2 + \sigma_g^2 + \sigma_E^2}} \exp\left(-\frac{(x_k - x_E(t))^2}{\sigma_k^2 + \sigma_g^2 + \sigma_E^2}\right) W \frac{\tau}{\epsilon} \\ &= \frac{\sigma_E}{\sqrt{\sigma_{\infty}^2 + \sigma_g^2 + \sigma_E^2}} \exp\left(-\frac{(\dot{x}_E\tau)^2}{\sigma_g^2}\right) W \frac{\tau}{\epsilon}, \end{aligned} \quad (33)$$

using  $\sigma_k = \sigma_{\infty}$  and  $(x_k - x_E(t))^2 = x_{\text{offset}}^2$  for large  $k$ . The right hand side of Equation (33) is independent of  $k$ . This means population  $P_k$  is constant, and we put this value as  $P_{\text{sat}}$ . Then

$$P_{\text{sat}} = \frac{\sigma_E}{\sqrt{\sigma_{\infty}^2 + \sigma_g^2 + \sigma_E^2}} \exp\left(-\frac{(\dot{x}_E\tau)^2}{\sigma_g^2}\right) W \frac{\tau}{\epsilon}. \quad (34)$$

### 5.3 Relation between lifespan and population

We derive a relation between population  $P_k$  and lifespan  $\tau$ .

Using Equation (34) we get

$$\frac{\partial P_{\text{sat}}}{\partial \tau} = \frac{\sigma_E}{\sqrt{\sigma_\infty^2 + \sigma_g^2 + \sigma_E^2}} \frac{W}{\epsilon} \exp\left(-\frac{(\dot{x}_E \tau)^2}{\sigma_g^2}\right) \left(1 - 2\left(\tau \frac{\dot{x}_E}{\sigma_g}\right)^2\right). \quad (35)$$

This means  $P_{\text{sat}}$  has a maximum at  $\tau = \tau_{\text{opt}} = \frac{\sigma_g}{\sqrt{2}\dot{x}_E}$ .

We should check whether  $\tau = \tau_{\text{opt}}$  satisfies the criterion (31) or not. At  $\tau = \tau_{\text{opt}}$  the criterion (31) is reduced to be

$$\frac{w_{\text{full}}}{\epsilon} \frac{\sigma_g}{\sqrt{2}\dot{x}_E} \geq \sqrt{1 + \frac{\sigma_\infty^2 + \sigma_g^2}{\sigma_E^2}} \sqrt{2e},$$

and this is equivalent to the criterion (32).

The criterion (32) is deformed to

$$\frac{\epsilon}{w_{\text{full}}} \leq \frac{1}{\sqrt{2e}} \frac{1}{\sqrt{\frac{\sigma_\infty^2 + \sigma_g^2}{\sigma_E^2} + 1}} \frac{\left(\frac{\sigma_g}{\sigma_E}\right)}{\left(\frac{\dot{x}_E}{\sigma_E}\right)},$$

or

$$\frac{\epsilon}{w_{\text{full}}} \leq \frac{1}{\sqrt{2e}} \frac{1}{\sqrt{\frac{1}{2} \left(\frac{\sigma_g^2}{\sigma_E^2}\right)} \sqrt{1 + 4 \frac{\sigma_E^2}{\sigma_g^2} + 1}} \frac{\left(\frac{\sigma_g}{\sigma_E}\right)}{\left(\frac{\dot{x}_E}{\sigma_E}\right)}. \quad (36)$$

The left hand side of the criterion (36) can be evaluated by observations and the right hand side is a value related to the distribution or width along  $x$ .

Therefore we conclude that a species of which population  $P_k$  grows maximum must satisfy the criterion (36) and its lifespan is

$$\tau = \tau_{\text{opt}} = \frac{\sigma_g}{\sqrt{2}\dot{x}_E}. \quad (37)$$

The equation (37) shows relation between observable parameters and parameters along  $x$ .

We call the maximum population species as a dominated species, because it has the largest population under the same environment.

### 5.4 Properties of dominant species

For dominant species what are relations between parameters under steady population?

Because lifespan of a dominant species is given by Equation (37), the criterion (36) is reduced to

$$\frac{\epsilon}{w_{\text{full}}} \leq \frac{\tau_{\text{opt}}}{\sqrt{e}} \frac{1}{\sqrt{\frac{1}{2} \left( \frac{\sigma_g^2}{\sigma_E^2} \right) \sqrt{1 + 4 \frac{\sigma_E^2}{\sigma_g^2} + 1}}}.$$

This is reduced to

$$\frac{\sigma_g}{\sigma_E} \leq \frac{c_{\text{full}}^2 - 1}{\sqrt{e} c_{\text{full}}}, \quad (38)$$

where

$$c_{\text{full}} = \frac{s_{\text{full}} \tau_{\text{opt}}}{\epsilon}. \quad (39)$$

If all present species are dominant species, their parameter distribution on a  $c - \frac{\sigma_g}{\sigma_E}$  plane must be restricted in a region shown by the criterion (38).

## 6 Summary

Using a simple dynamical model of a species with propagation of genetic characters between generations we evaluate a behavior of its population under linearly changing environment.

Without any limitation of reproduction, population of a species grows exponentially or extinct if parameters are not well tuned; the system is exponentially unstable.

With limited energy available for reproduction of a species, its population can be stable. In this case only one optimal value of lifespan always exists which gives largest population. For such a species there are two criteria between parameters on genetic characters and observable parameters.

## References

- [1] Boothroyd, A.I., Sackmann, I.-J., and Fowler, W.A.: 1990, Our Sun. I. The standard model: successes and failures, *Astrophys. J.*, vol.360, 727-736.
- [2] Condie, K.C.: 1989, *Plate tectonics & crustal evolution*, 3rd ed. Pergamon Press, Oxford.
- [3] Ito, T., Kumazawa, M., Hamano, Y., Matsui, T., and Masuda, K.: 1993, Long term evolution of the solar insolation variation over 4Ga, *Proc. Jpn. Acad., Ser. B* **69**, 233-237(1993).
- [4] Kimura, M.: 1983, *The neutral theory of molecular evolution*, Cambridge University Press, Cambridge.

- [5] Losos, J.B., Warheit, K.I., and Schoener, T.W.: 1997, Adaptive differentiation following experimental island colonization in *Anolis* lizards, *Nature* **387**, 70-73.
- [6] Tajika, E.: 1992, *Evolution of the atmosphere and ocean of the Earth: global geochemical cycles of C, H, O, N, and S, and degassing history coupled with thermal history*, PhD Dissertation(University of Tokyo).
- [7] Tanikawa, K., Handa, T., and Ito, T., Lifespan of organisms and the secular change of the environment, in *the Proceedings of the 32nd Symposium on Celestial Mechanics*, 15-17, March, 2000, Tokyo, Japan, pp.179-193.

# Symposium Program/プログラム

## 第 1 日 11 日 (月)

14:00-15:00	受けつけ
15:00-15:10	小久保 英一郎 (国立天文台)
	はじめに
【力学系・他】	座長：荒木田 英禎
15:10-15:30	井上 猛 (京都産業大学)
	Schwarzschild 解の意味するもの
15:30-15:50	船渡 陽子 (東京大学)
	Life Expectancy of the Large Magellanic Cloud
【数値計算法】	座長: 荒木田 英禎
15:50-16:10	小久保 英一郎 (国立天文台)
	A Modified Hermite Integrator for Planetary Dynamics
【ポスター発表】	座長: 小南 淳子
16:30-17:30	ポスター 3 分発表
18:30-19:30	夕食兼懇親会 (2 階広間)
19:30-	ポスターセッション兼ウェルカムパーティー

## 第 2 日 12 日 (火)

08:00-09:00	朝食 (1 階食堂)
09:00-10:00	ポスターセッション
【招待講演】	座長: 井田 茂
10:00-10:15	井田 茂 (東京工業大学)
	メインテーマについて — いまさら力学的摩擦?
10:15-11:30	牧野 淳一郎 (東京大学)
	恒星系及び恒星系 N 体シミュレーションにおける力学的摩擦 — 現実と虚構
12:00-13:00	昼食 (1 階食堂)
13:00-14:00	ポスターセッション
14:00-15:15	田中 秀和 (東京工業大学)
	円盤重力系における力学的摩擦と動径方向移動



【惑星系形成】	座長: 岩崎 一典
15:50-16:10	小南 淳子 (東京工業大学) 原始惑星系円盤の散逸に伴う地球型惑星集積 ～円盤ガスからの力学的摩擦の影響～
16:10-16:30	小林 浩 (東京工業大学) 離心率、軌道傾斜角の大きい天体が原始惑星系星雲内でうける ガス抵抗について
16:40-17:00	武田 隆頭 (東京工業大学) 低質量な複数の衛星からなる衛星系の形成
【系外惑星系】	座長: 跡部 恵子
17:00-17:20	木下 宙 (国立天文台) 2:1 平均運動共鳴にある太陽系外型外星系の運動
18:00-20:00	夕食兼懇親会 (2 階広間)

### 第 3 日 13 日 (水)

08:00-09:00	朝食 (1 階食堂)
09:00-10:00	ポスターセッション
【銀河・恒星系力学】	座長: 出田 誠
10:00-10:20	坂本 強 (総合研究大学院大学) 銀河系の質量に対する新しい制限
10:20-10:40	樽家 篤史 (東京大学) Gravothermal Catastrophe and Tsallis' Entoropy
10:40-11:00	井口 修 (お茶の水大学) 巾的ポテンシャルによる多体系の準平衡状態
【太陽系】	座長: 吉田 二美
11:10-11:30	谷川 清隆 (国立天文台) Interactions among Planets in a Stable Solar System
11:30-11:50	Sebastian Bouquillon (国立天文台) Models for the Mercury's Rotation
11:50-12:00	伊藤 孝士 (国立天文台) おわりに

## ポスター発表 (50 音順)

跡部 恵子 (東京工業大学)	巨大惑星の摂動による地球型惑星の自転軸の傾きの時間進化
荒木田 英禎 (総合研究大学院大学)	Motion around Triangular Lagrange Points Perturbed by Other Bodies
出田 誠 (京都大学)	偏心銀河円盤とダークハロー間の力学的摩擦
伊藤 孝士 (国立天文台)	始原的なシンプレクティック数値積分法の誤差を再び簡単に考察する
岩崎 一典 (東京工業大学)	ガス円盤中での原始惑星系の安定性 -原始惑星の質量依存性について-
中井 宏 (国立天文台)	小惑星の軌道と共鳴
林 満 (総合研究大学院大学)	IMPI を用いた GRAPE-5 と VPP5000 連携計算環境の構築
原田 渉 (東京大学)	地球の時間暦の非線型調和解析
半田 利弘 (東京大学)	Response of Lifespan of Organisms to Secularly Changing Environment by a New Dynamical Model
眞崎 良光 (総合研究大学院大学)	内側天体から摂動を受けた高離心率天体の運動
松林 達史 (東京工業大学)	星団沈降プロセスにおける潮汐力と質量損失の影響
宮口 智成 (早稲田大学)	離散力学系におけるレヴィ拡散
山口 喜博 (帝京平成大学)	ねじれ写像における Non-Birkhoff 型周期軌道 円写像における Non-Birkhoff 型周期軌道
山口 義幸 (京都大学)	長距離相互作用系における緩和
山本 一登 (総合研究大学院大学)	ケプラー運動に対する線形多段法の安定領域
吉田 二美 (神戸大学)	Sub-km 小惑星のデータに力学的摩擦の徴候は見られるか?

# Author Index and Participant List/参加者リスト

AKIMOTO, Takuma (秋元琢磨), <i>Waseda University</i>	
ARAKIDA, Hideyoshi (荒木田英禎), <i>Graduate University for Advanced Studies</i> .....	265
ATOBE, Keiko (跡部恵子), <i>Tokyo Institute of Technology</i> .....	255
BOUQUILLON, Sebastien, <i>National Astronomical Observatory of Japan</i>	
FUNATO, Yoko (船渡陽子), <i>University of Tokyo</i> .....	30
HANDA, Toshihiro (半田利弘), <i>University of Tokyo</i> .....	550
HARADA, Wataru (原田渉), <i>University of Tokyo</i> .....	533
HATANAKA, Shijun (畑中至純)	
HAYASHI, Mitsuru (林満), <i>Graduate University for Advanced Studies</i> .....	528
HIGUCHI, Arika (樋口有理可), <i>Kobe University</i>	
IDA, Shigeru (井田茂), <i>Tokyo Institute of Technology</i> .....	1
IDETA, Makoto (出田誠), <i>Kyoto University</i> .....	134
IGUCHI, Osamu (井口修), <i>Ochanomizu University</i> .....	104
INOUE, Takeshi (井上猛), <i>Kyoto Sangyo University</i> .....	226
ITO, Takashi (伊藤孝士), <i>National Astronomical Observatory of Japan</i> .....	437
IWASAKI, Kazunori (岩崎一典), <i>Tokyo Institute of Technology</i> .....	193
KANAEDA, Naoko (金枝直子), <i>Ochanomizu University</i>	
KINOSHITA, Hiroshi (木下宙), <i>National Astronomical Observatory of Japan</i> .....	199
KOBAYASHI, Hiroshi (小林浩), <i>Tokyo Institute of Technology</i> .....	162
KOKUBO, Eiichiro (小久保英一郎), <i>National Astronomical Observatory of Japan</i> .....	415
KOMINAMI, Junko (小南淳子), <i>Tokyo Institute of Technology</i> .....	152
KUSAKABE, Nobuhiko (日下部展彦), <i>Tokyo Gakugei University</i>	
MAKINO, Junichiro (牧野淳一郎), <i>University of Tokyo</i> .....	2
MASAKI, Yoshimitsu (眞崎良光), <i>Graduate University for Advanced Studies</i> .....	303
MATSUBAYASHI, Tatsushi (松林達史), <i>Tokyo Institute of Technology</i> .....	112
MIYAGUCHI, Tomoshige (宮口智成), <i>Waseda University</i> .....	353
MORIWAKI, Kazumasa (森脇一匡), <i>Kobe University</i>	
NAKAI, Hiroshi (中井宏), <i>National Astronomical Observatory of Japan</i> .....	289
SAKAMOTO, Tsuyoshi (坂本強), <i>Graduate University for Advanced Studies</i> .....	38
SOTA, Yasuhide (曾田康秀), <i>Ochanomizu University</i>	
SUMITANI, Hideo (住谷秀夫), <i>Osaka College of Music</i>	
SUTO, Yasushi (須藤靖), <i>University of Tokyo</i>	
TAKEDA, Takaaki (武田隆顕), <i>Tokyo Institute of Technology</i> .....	189
TANAKA, Hidekazu (田中秀和), <i>Tokyo Institute of Technology</i> .....	25
TANIKAWA, Kiyotaka (谷川清隆), <i>National Astronomical Observatory of Japan</i> .....	236
TARUYA, Atsushi (樽家篤史), <i>University of Tokyo</i> .....	78

UMEHARA, Hiroaki (梅原広明), *Communications Research Laboratory*

YAMAGUCHI, Yoshihiro (山口喜博), *Teikyo Heisei University* ..... 359, 395

YAMAGUCHI Y., Yoshiyuki (山口義幸), *Kyoto University* ..... 407

YAMAMOTO, Tadato (山本一登), *Graduate University for Advanced Studies* ..... 541

YOSHIDA, Fumi (吉田二美), *Kobe University* ..... 330

YOSHIDA, Haruo (吉田春夫), *National Astronomical Observatory of Japan*

天体力学N体力学研究会

平成14年3月11日－13日 箱根 静雲荘