

THE ACCURATE NUMERICAL
SOLUTION OF HIGHLY OSCILLATORY
ORDINARY DIFFERENTIAL EQUATIONS

Thesis by

Robert Elmer Scheid, Jr.

In Partial Fulfillment of the Requirements
for the Degree of
Doctor of Philosophy

California Institute of Technology
Pasadena, California

1982

(Submitted March 10, 1982)

Time for you and time for me,
And time yet for a hundred indecisions,
And for a hundred visions and revisions,
Before the taking of a toast and tea.

-- T. S. Eliot

I especially must thank my advisor, Professor Heinz-Otto Kreiss, for his valuable guidance and his considerable patience. I am also grateful to the students and faculty of the Department of Applied Mathematics for many interesting discussions and many prized friendships.

Financial support for my endeavors has come from the California Institute of Technology through assistantships and from the Office of Naval Research through Contract Number N00014-80-C-0076.

And finally I must acknowledge my family, whose encouragement has been unwavering throughout. I am pleased to dedicate this thesis to the memory of my father, Robert Scheid, who died shortly after I began my studies here.

ABSTRACT

We consider systems of ordinary differential equations with rapidly oscillating solutions. Conventional numerical methods require an excessively small time step ($\Delta t = O(\epsilon h)$, where h is the step size necessary for the resolution of a smooth function of t and $\frac{1}{\epsilon}$ measures the size of the large eigenvalues of the system's Jacobian).

For the linear problem with well-separated large eigenvalues we introduce smooth transformations which lead to the separation of the time scales and computation with a large time step ($\Delta t = O(h)$). For more general problems, including systems with weak polynomial nonlinearities, we develop an asymptotic theory which leads to expansions whose terms are suitable for numerical approximation. Resonances can be detected and resolved often with a large time step ($\Delta t = O(h)$). Passage through resonance in nonautonomous systems can be resolved by a moderate time step ($\Delta t = O(\sqrt{\epsilon} h)$).

TABLE OF CONTENTS

Acknowledgments	ii
Abstract	iii
Table of Contents	iv
I Introduction	1
II The Linear Problem with Well-Separated Large Frequencies	
2.1 Oscillatory Stability	17
2.2 Normal Forms and Canonical Transformations	24
2.3 Reduction to Canonical Form	40
2.4 Numerical Methods	47
2.5 A Computational Example	54
Appendix A: The Smoothness of Eigenvalues and Eigenvectors	62
III The Nonlinear Problem without Turning Points	
3.1 Reduction to the Non-Stiff Formulation	69
3.2 Hierarchy for the Linear Problem	76
3.3 Solution by Successive Linearizations	81
3.4 Solution by Formal Expansion	90
3.5 Extensions to More General Systems	101
3.6 Extensions to Partial Differential Equations	104
3.7 Computational Examples	108
IV The Resolution of Turning Points	
4.1 Solution by Nonuniform Expansion	126
4.2 The Cost of Resolution / Turning Points of Higher Order	138
4.3 A Computational Example	142
References	152

CHAPTER I
INTRODUCTION

Mathematical modeling of a chemical, electrical, mechanical or biological process often leads to a differential system whose Jacobian has at least one eigenvalue with either a large negative real part or a large imaginary part. Even when the underlying structure is quite complicated, one generally can analyze the stiffness of such a system through the simple scalar equation:

$$(1.01) \quad \begin{cases} dy/dt = ay, & t > 0 \\ y(0) = y_0 \end{cases}$$

Case I: $\text{Re}\{-a\} \gg 1$

Case II: $|\text{Im}\{a\}| \gg 1$

Unless one is prepared to compute with an excessively small time step, most conventional numerical methods are ill-suited to the problem for reasons of stability or accuracy. For example, in table <1.10> we consider several generic schemes as applied to the system (1.01) with mesh width h .

METHOD	FORMULATION / SOLUTION	$ ah \rightarrow \infty$ ($\arg(ah) < 0$)
Forward Euler <1.10a>	$\begin{cases} v_{N+1} = v_N(1+ah) \\ v_0 = y_0 \\ v_N = (1+ah)^N v_0 \end{cases}$	$ v_N \rightarrow \infty$ ($N > 0$)
Backward Euler <1.10b>	$\begin{cases} v_{N+1}(1-ah) = v_N \\ v_0 = y_0 \\ v_N = [1/(1-ah)]^N v_0 \end{cases}$	$ v_N \rightarrow 0$ ($N > 0$)
Trapezoidal Rule <1.10c>	$\begin{cases} v_{N+1}(1-ah/2) = v_N(1+ah/2) \\ v_0 = y_0 \\ v_N = [(1+ah)/(1-ah)]^N v_0 \end{cases}$	$v_N \rightarrow (-1)^N v_0$ ($N > 0$)

table <1.10>

On considering the stiff limit ($|ah| \rightarrow \infty$ with $\arg(ah) < 0$), we find that the first method is unstable while the second and third are stable. Moreover, the solution of <1.10b> decays rapidly on the grid points, while the solution of <1.10c> can be characterized as grid oscillations. These observations do not contradict the general theory which has been developed for the non-stiff limit ($|ah| \rightarrow 0$) but rather indicate that one cannot expect convergence in the stiff limit.

Nevertheless, for case I the solution of <1.10b> is qualitatively similar to the solution of the differential equation (1.01). Much has been made of this salient feature of the backwards Euler formulation, and many schemes with similar stability properties have been proposed for stiff problems of this type (see, for example, Lambert [20] and Kreiss [18]). With the exception of a thin boundary layer, such problems have nicely behaved solutions. Our aim is a detailed numerical analysis of the highly oscillatory case (II), in which the rapid

changes are expected to persist.

Since the fundamental work of Poncaré, mathematicians studying oscillatory phenomena have developed an extensive arsenal of perturbative techniques including multi-scaling, averaging, and the near-identity transformation (see, for example, Kevorkian and Cole [16], Nayfeh [29], and Neu [28]). For the most part, these tools are difficult to implement numerically since the analytical manipulations require a competence not to be expected of a collection of FORTRAN statements; however, a number of computational schemes have also been proposed.

Many researchers have attempted to extrapolate the effects of the oscillations from grid point to grid point. For certain problems in which the high frequencies are known in advance, Gautschi [11] developed linear multistep methods which are exact for trigonometric polynomials up to a certain degree, and later Snyder and Fleming [31] proposed an aliasing technique applicable to Certaine's method for solving ordinary differential equations. Multirevolution methods [12,32] were first introduced by astronomers to calculate future satellite orbits by using some physical reference point such as a node, apogee or perigee; these ideas were further developed by Petzold [30], whose methods extrapolate the effects of the oscillations for many cycles by first calculating for one cycle near each grid point. Fatunla [8] also introduced schemes designed to follow many cycles with each time step.

Others less concerned with the details of the oscillations

have proposed filtering techniques designed to eliminate entirely the effects of the fast modes. In their study of linear problems with well-separated, slowly varying large frequencies, Amdursky and Ziv [3] used left and right eigenvectors corresponding to the high frequencies to project the solution onto the manifold of smooth components. Lindberg [21] used temporal filters to remove the grid oscillations resulting from the application of the trapezoidal rule $\langle 1.10c \rangle$. More recently Kreiss [19] has shown that for a large class of linear and nonlinear problems oscillations can be suppressed by a proper choice of initial conditions. And finally Majda [22] has demonstrated that for the linear problem time-filtered solutions have the full accuracy of the filtering method as long as the system has constant coefficients, the fast and slow scales have been separated, or the initial data have been prepared by Kreiss's technique; otherwise, the computed solutions are only first-order accurate.

Since, for many problems, the effects of the oscillations cannot be blindly suppressed or crudely approximated, a number of analytical-numerical methods have been proposed to further exploit the underlying mathematical structure. Miranker and Hoppensteadt [13,23,24] analyzed the theoretical and practical difficulties of implementing a method of averaging for such problems; however, they only executed their strategy to solve linear equations with constant coefficients. Moreover, they did not attempt to calculate the phases accurately. Amdursky and Ziv [2] also studied the linear problem with slowly varying

large frequencies by using a formulation similar to averaging. Nonlinear problems of the form

$$(1.02) \quad \begin{cases} dx/dt = (A/\varepsilon)X + g(t, X), & t > 0 \\ x(0) = X_0, & A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}, & 0 < \varepsilon \ll 1 \end{cases}$$

were studied by Miranker and van Veldhuizen [25], who introduced a Fourier expansion in the fast scale ($\tau = t/\varepsilon$). Miranker and Wabba [23,26] also analyzed such oscillators by developing a calculus of stable averaging functionals to replace the standard point functionals of analysis.

Our approach is similar to this last group in that we use analytical as well as numerical techniques to calculate the solutions accurately. While some of the latter methods require extensive preparation, we generally insist that the analytic manipulations be programmable, and also we treat linear and nonlinear systems in considerably greater generality than has previously been attempted.

In chapter II we consider the general linear problem and propose a method which separates the time scales by a smooth transformation when the large frequencies are well-separated. Using a combination of numerical and asymptotic techniques, one then can compute the solution accurately with a large time step. In chapter III we consider systems with weak polynomial nonlinearities and derive an asymptotic expansion whose terms are suitable for computation with a large time step; in addition, we demonstrate how these principles can be applied to more general systems, including partial differential equations

in which the oscillations occur temporally but not spatially. In chapter IV we discuss passage through resonance in nonautonomous systems and demonstrate how the results of chapter III can be extended to handle these circumstances.

In our consideration of linear systems with well-separated large frequencies (Chapter II), we study the partitioned real system

$$(1.03) \quad \begin{cases} \begin{pmatrix} \dot{Y}^I \\ \dot{Y}^{II} \end{pmatrix} = \begin{bmatrix} \frac{1}{\varepsilon} A_{11}(t, \varepsilon) & \frac{1}{\varepsilon} A_{12}(t, \varepsilon) \\ 0 & A_{22}(t, \varepsilon) \end{bmatrix} \begin{pmatrix} Y^I \\ Y^{II} \end{pmatrix} + \begin{pmatrix} \frac{1}{\varepsilon} F_1(t, \varepsilon) \\ F_2(t, \varepsilon) \end{pmatrix} \\ Y^I(0, \varepsilon) = Y_0^I \\ Y^{II}(0, \varepsilon) = Y_0^{II} \end{cases} \quad 0 < t < T, \quad 0 < \varepsilon \ll 1$$

Here Y^I is an n_1 -vector, and Y^{II} is an n_2 -vector. All given matrices and vectors are smooth functions of t and ε . Moreover, $\frac{1}{\varepsilon} A_{11}(t, \varepsilon)$ is block diagonal with each diagonal block a (2×2) submatrix of the form

$$(1.04) \quad \begin{bmatrix} a(t, \varepsilon) & \frac{1}{\varepsilon} b(t, \varepsilon) \\ -\frac{1}{\varepsilon} b(t, \varepsilon) & a(t, \varepsilon) \end{bmatrix}$$

The fast modes in (1.03) are completely decoupled, and Y^{II} is the solution of the smooth system

$$(1.05) \quad \begin{cases} \dot{Y}^{II} = A_{22}(t, \varepsilon) Y^{II} + F_2(t, \varepsilon) \\ Y^{II}(0, \varepsilon) = Y_0^{II} \end{cases}$$

Let $S_\varepsilon(t, s)$ be the $(n_1 \times n_1)$ solution operator of the reduced system

$$(1.06) \quad \begin{cases} Y^{\mathbf{I}} = A_{11}(t, \varepsilon) Y^{\mathbf{I}} \\ Y^{\mathbf{I}}(0, \varepsilon) = Y_0^{\mathbf{I}} \end{cases}$$

In fact $S_{\varepsilon}(t, s)$ is a block diagonal matrix with each block a (2 x 2) submatrix of the form

$$(1.07) \quad \exp(\hat{A}_{\varepsilon}(t, s)) \begin{bmatrix} \cos(\hat{B}_{\varepsilon}(t, s)/\varepsilon) & \sin(\hat{B}_{\varepsilon}(t, s)/\varepsilon) \\ -\sin(\hat{B}_{\varepsilon}(t, s)/\varepsilon) & \cos(\hat{B}_{\varepsilon}(t, s)/\varepsilon) \end{bmatrix}$$

where, in correspondence to (1.04), we have

$$(1.08) \quad \begin{cases} \hat{A}_{\varepsilon}(t, s) = \int_s^t a(r, \varepsilon) dr \\ \hat{B}_{\varepsilon}(t, s) = \int_s^t b(r, \varepsilon) dr \end{cases}$$

Now the solution for $Y^{\mathbf{I}}$ is given explicitly by

$$(1.09) \quad \begin{cases} Y^{\mathbf{I}}(t, \varepsilon) = S_{\varepsilon}(t, 0) Y_0^{\mathbf{I}} + \frac{1}{\varepsilon} \int_0^t S_{\varepsilon}(t, s) G(s, \varepsilon) ds \\ G(s, \varepsilon) = A_{12}(s, \varepsilon) Y^{\mathbf{II}}(s, \varepsilon) + F_1(s, \varepsilon) \end{cases}$$

One now can generate an asymptotic expansion for $Y^{\mathbf{I}}$ by means of integration by parts, and in Section 2.4 we show how this procedure leads to numerical techniques for approximating the solution of (1.03). Of course one does not have this form in general, and so in Sections 2.1, 2.2, and 2.3 we discuss the reduction of a general linear system to the form

$$(1.10) \quad \begin{cases} Y' = [\frac{1}{\varepsilon} A(t, \varepsilon) + \varepsilon^q B(t, \varepsilon)] Y + \frac{1}{\varepsilon} F(t, \varepsilon) \\ Y(0, \varepsilon) = Y_0 \end{cases}$$

where, after the $O(\varepsilon^q)$ terms are discarded, we have the form (1.03). For this strategy to be successful one must assume that the large frequencies are well-separated. That is, in correspondence to the block diagonal structure of $All(t, \varepsilon)$ we must have:

$$(1.11) \quad \begin{cases} \min_{\substack{0 < t < T \\ k}} |b_k(t, \varepsilon)| > K \\ \min_{\substack{0 < t < T \\ k \neq j}} | |b_k(t, \varepsilon)| - |b_j(t, \varepsilon)| | > K \end{cases}$$

where K is a positive constant independent of ε . This procedure leads to very accurate approximations with large time steps when assumption (1.11) is possible; however, the exclusion of weak nonlinearities and coalescing large frequencies considerably restricts the qualitative possibilities for highly oscillatory systems.

In Chapter III we begin our development of more general techniques. Since we intend to compute with a relatively large mesh width, it is again essential to represent the oscillatory structure analytically; otherwise, excessively small time steps will be needed to resolve these features numerically. To illustrate our approach, we first consider the scalar problem

$$(1.12) \quad \begin{cases} u' = (ia/\varepsilon)u + u^2 \\ u(0, \varepsilon) = u_0, \quad 0 < t < T \end{cases}$$

where a is a nonzero real number and ε is a small parameter ($0 < \varepsilon \ll 1$). The substitution

$$(1.13) \quad u = \exp(iat/\varepsilon) x$$

reduces the stiff system (1.12) to a formulation in which the right-hand side is bounded but rapidly oscillating:

$$(1.14) \quad \begin{cases} x' = \exp(iat/\varepsilon) x^2 \\ x(0, \varepsilon) = x_0, \quad 0 < t < T \end{cases}$$

In this introduction we refer to terms with factors such as $\exp(iat/\varepsilon)$ as oscillatory; terms without such factors are called nonoscillatory. For sufficiently small ε system (1.14) can be solved explicitly by separation of variables:

$$(1.15) \quad \begin{aligned} x &= x_0 \left\{ 1 / [1 - \varepsilon x_0 (\exp(iat/\varepsilon) - 1) / (ia)] \right\} \\ &= x_0 \sum_{k=0}^{\infty} [x_0 (\exp(iat/\varepsilon) - 1) / (ia)]^k \varepsilon^k \end{aligned}$$

For less tractable equations, of course, this method is unworkable, and solutions must be uncovered by more general techniques. On investigating the dominant balance of (1.14), we intuitively expect the rapidly oscillating terms to be less important, and accordingly in section 3.3 we demonstrate that

$$(1.16) \quad \max_{0 < t < T} |x(t, \varepsilon) - x_0| = O(\varepsilon) .$$

This analysis leads to an obvious change of variables:

$$(1.17) \quad \begin{cases} \bar{x}' = \exp(iat/\varepsilon) [(x_0^2/\varepsilon) + 2x_0\bar{x} + \varepsilon\bar{x}^2] \\ \bar{x}(0, \varepsilon) = 0, \quad x = x + \varepsilon\bar{x} \end{cases}$$

The $O(1/\varepsilon)$ oscillatory term cannot be neglected; however, after the substitution

$$(1.18) \quad \begin{cases} x = y_1(t, \varepsilon) + \tilde{\tilde{x}} \\ y_1(t, \varepsilon) = -i(x_0^2/a) \exp(iat/\varepsilon) \end{cases}$$

we have the more manageable system

$$(1.19) \quad \begin{cases} \tilde{\tilde{x}}' = (-2ix_0^2/a) \exp(i2at/\varepsilon) + 2x_0\tilde{\tilde{x}} \exp(iat/\varepsilon) \\ \quad + \varepsilon[(-ix_0^2/a) \exp(iat/\varepsilon) + \tilde{\tilde{x}}]^2 \exp(iat/\varepsilon) , \\ \tilde{\tilde{x}}(0, \varepsilon) = i(x_0^2/a) \end{cases}$$

and again by the results of Section 3.3 we can neglect oscillatory terms and $O(\varepsilon)$ terms to give

$$(1.20) \quad \max_{0 < t < T} |x(t, \varepsilon) - w_1(t)| = O(\varepsilon) ,$$

where w_1 satisfies the system:

$$(1.21) \quad \begin{cases} w_1' = 0 \\ w_1(0) = i(x_0^2/a) \end{cases}$$

Thus, the first-order approximation to the solution of (1.14) is given by

$$(1.22) \quad x = x_0 + \varepsilon(w_1(t) + y_1(t, \varepsilon)) + O(\varepsilon^2) ,$$

where $w_1(t)$ is nonoscillatory and $y_1(t, \varepsilon)$ is oscillatory.

We systematically develop this methodology for nonlinear systems in Sections 3.1, 3.2, and 3.3, where the balancing of terms is justified by a functional Newton iteration. Integration by parts yields the first oscillatory correction as in (1.18), whereupon the elimination of secondary terms determines the first nonoscillatory correction as in (1.21). When this procedure is repeated after linearization, corrections of higher order are likewise generated; the solution is then represented by an asymptotic expansion of the form

$$(1.23) \quad x(t, \varepsilon) \sim \sum_K (w_K(t) + y_K(t, \varepsilon)) \varepsilon^K ,$$

where each $w_K(t)$ is bounded and nonoscillatory and each $y_K(t, \varepsilon)$ is bounded and oscillatory. We characterize the terms of (1.23) as the solutions of equations which are easily resolved with a large time step; that is, a time step which need not be small compared with ε . Our approach is conceptually similar to the generalized method of averaging as developed by Bogoliubov and Mitropolsky [4].

Given this asymptotic form for the solution, we develop in Section 3.4 a formal procedure which generates the terms of the series so that repeated linearizations are unnecessary. For example, in the solution (1.15) of equation (1.14) the

oscillations clearly occur on the fast scale

$$(1.24) \quad \tau = t/\varepsilon ,$$

and so (1.15) can be rewritten as

$$(1.25) \quad x = x_0 \sum_{k=0}^{\infty} [x_0 (\exp(ia\tau) - 1) / (ia)]^k \varepsilon^k .$$

One might discover the first terms of this summation through a multi-time expansion of the form

$$(1.26) \quad x \sim \sum_k f_k(\tau, t) \varepsilon^k ,$$

where each $f_k(t/\varepsilon, t)$ is bounded. In general, however, the ε -dependence will not appear simply through the fast scale (1.24). For example, we consider a nonautonomous variant of (1.14)

$$(1.27) \quad \begin{cases} x' = \exp(ia(t)/\varepsilon) x^2 \\ x(0, \varepsilon) = x_0, \quad 0 < t < T \end{cases} ,$$

where $a(t)$ is a smooth real function with

$$(1.28) \quad \min_{0 < t < T} |a'(t)| > 0 .$$

As in the case of (1.14) the solution is readily obtained by separation of variables:

$$(1.29) \quad \begin{aligned} x &= x_0 \{1 / [1 - x_0 F(t, \varepsilon)]\} \\ &= x_0 \sum_{k=0}^{\infty} [x_0 F(t, \varepsilon)]^k , \end{aligned}$$

where

$$(1.30) \quad F(t, \varepsilon) = \int_0^t \exp(ia(t)/\varepsilon) dt .$$

The right-hand side of (1.30) can be integrated by parts to give an asymptotic expansion :

$$(1.31) \quad F(t, \varepsilon) \sim \int_0^t \exp(ia(t)/\varepsilon) \{ \varepsilon[-i/a'(t)] + \varepsilon^2[a''(t)/(a'(t))^3] + O(\varepsilon^3) \} .$$

And after the substitution of (1.31) into (1.29) we have an expansion of the form given in (1.23). Now the oscillatory behavior enters through the factor

$$(1.32) \quad \exp(ia(t)/\varepsilon) = \exp(ia(\varepsilon\zeta)/\varepsilon) \quad (\zeta=t/\varepsilon) ,$$

and so we do not have the form given in (1.26). Without the assumption (1.28) this procedure is unworkable because the mathematical structure of $F(t, \varepsilon)$ changes significantly over any interval where $a'(t)$ vanishes. This behavior characterizes the general theory, where the expansions first become nonuniform and eventually break down entirely due to the failure of the technique of integration by parts.

In Chapter IV we extend the theory of the preceding chapter to handle these circumstances, which can be described mathematically as a turning point or physically as a passage through resonance. Under these conditions the structure of certain terms is changing from oscillatory to nonoscillatory to oscillatory, and correspondingly the balancing of terms also

must change. Since in practice one computes simply from grid point to grid point, we need only develop a principle which can be applied locally, and as in the previous example this essentially must be a procedure for the evaluation of integrals of the form

$$(1.33) \quad \int_a^b f(t) \exp(a(t)/\mathcal{E}) dt,$$

where $a'(t)$ may now vanish at some points in the interval of interest. In this context we consider the evaluation of the integral

$$(1.34) \quad \int_0^1 \exp(it^2/2\mathcal{E}) dt,$$

which corresponds to a response to the frequency (it/\mathcal{E}) . With the introduction of the scalar parameter

$$(1.35) \quad \hat{K} > 1,$$

we can decompose (1.34) into the sum of

$$(1.36) \quad \int_0^{\hat{K}\sqrt{\mathcal{E}}} \exp(it^2/2\mathcal{E}) dt$$

and

$$(1.37) \quad \int_{\hat{K}\sqrt{\mathcal{E}}}^1 \exp(it^2/2\mathcal{E}) dt.$$

For the first integral the change of variables

$$(1.38) \quad \sqrt{\mathcal{E}} r = s$$

in the integrand gives

$$(1.39) \quad \int_0^{\hat{K}\sqrt{\mathcal{E}}} \exp(is^2/2\mathcal{E}) ds = \int_0^{\hat{K}} \exp(ir^2/2) dr.$$

Thus, using an $O(\sqrt{\varepsilon})$ stretching factor completely resolves the oscillatory effects in the neighborhood of the turning point. For the evaluation of the second integral we again use the technique of integration by parts, but here we approximate to within $O(\sqrt{\varepsilon}/\hat{K}^{2m-1})$ where m is sufficiently large so that this error can be neglected:

$$(1.40) \quad \int_{\hat{K}\sqrt{\varepsilon}}^1 \exp(is^2/2\varepsilon) ds = \frac{1}{\hat{K}\sqrt{\varepsilon}} \exp(is^2/2\varepsilon) [-\varepsilon i/s - \varepsilon^2/s^3 + \dots + \varepsilon^m a_m/s^{2m-1}] + R_m$$

$$|R_m| < \varepsilon^m A_m \left| \int_{\hat{K}\sqrt{\varepsilon}}^1 1/s^{2m} ds \right| < \tilde{A}_m (\sqrt{\varepsilon}/\hat{K}^{2m-1}) \quad [\{A_m, \tilde{A}_m, a_m\} \subset \mathbb{R}]$$

The point

$$(1.41) \quad t^* = \hat{K} \sqrt{\varepsilon}$$

then marks the transition between a fast mode which requires asymptotic analysis and a slow mode which requires full resolution, and correspondingly the value of \hat{K} determines a trade-off between the number of integrations needed in one region and the number of grid points needed in another. One might compare this procedure to a perturbation theorist's matching of inner and outer expansions; that is, after one has derived the correct balances in two disjoint regions, one postulates that the corresponding functional representations can be matched through an intermediate range in which both expansions are valid.

Our turning point arguments can resolve the transition of modes in the neighborhood of an isolated point. In general, however, if there is no well-defined separation between the frequencies of the fast and slow modes of a system, then one is not really solving a singular perturbation problem.

CHAPTER II
THE LINEAR PROBLEM WITH
WELL-SEPARATED LARGE FREQUENCIES

2.1 OSCILLATORY STABILITY

In this chapter we consider real systems of the form

$$(2.1.1) \quad \begin{cases} dx/dt = A(t, \varepsilon) X \\ X(0, \varepsilon) = X_0, \quad 0 < t < T, \quad 0 < \varepsilon \ll 1 \end{cases}$$

and

$$(2.1.2) \quad \begin{cases} dx/dt = A(t, \varepsilon) X + F(t, \varepsilon) \\ X(0, \varepsilon) = X_0, \quad 0 < t < T, \quad 0 < \varepsilon \ll 1 \end{cases}$$

where we assume:

- (i) $X(t, \varepsilon)$ and $F(t, \varepsilon)$ are n -vectors; $A(t, \varepsilon)$ is an $(n \times n)$ matrix;
- (ii) $A(t, \varepsilon) = A_1(t)/\varepsilon + A_2(t, \varepsilon)$;
- (iii) $0 < \varepsilon \ll 1$ is to be interpreted as "for sufficiently small positive ε "; for definiteness we also introduce an explicit upper bound E : $0 < \varepsilon < E$;
- (iv) $A_1(t) \in C^p(t)$; $A_2(t, \varepsilon) \in C^p(t, \varepsilon)$; $F(t, \varepsilon) \in C^p(t, \varepsilon)$; $p \geq 1$ (continuous partial derivatives of order p);
- (v) The eigenvalues of $A_1(t)$ can be divided into two sets uniformly on $0 < t < T$:

Oscillatory Set:

n_1 imaginary eigenvalues $\{\lambda_j(t)\}$ with

$$\min_j |\lambda_j(t)| > r$$

and

$$\min_{j \neq k} |\lambda_j(t) - \lambda_k(t)| > r,$$

where r is some positive real constant.

Nonoscillatory Set:

$n_2 (=n-n_1)$ identically zero eigenvalues.

Although we restrict ourselves here to real systems, these principles also can be applied to complex systems. And also our analysis applies to the system (1.03), where we allow $O(1/\epsilon)$ forcing terms in certain components. In later chapters we shall relax the restrictions of assumption (v) so as to allow the treatment of coalescing eigenvalues in the oscillatory set and the transfer of eigenvalues between the two sets. For this chapter, however, we consider only well-spaced large frequencies. $x^{(k)}$ is the k -th component of the vector x and $A^{(k)}$ is the k -th column of the matrix A .

We can write the solution of (2.1.1) in the form

$$(2.1.3) \quad X(t, \epsilon) = S_\epsilon(t, 0) X_0,$$

where $S_\epsilon(t, s)$ is the solution operator of the system (2.1.1). By Duhammel's Principle we then can write the solution of the system (2.1.2) in the form

$$(2.1.4) \quad X(t, \epsilon) = S_\epsilon(t, 0) X_0 + \int_0^t S_\epsilon(t, s) F(s, \epsilon) ds.$$

The system (2.1.1) is said to be oscillatory stable if we have an estimate for the solution operator:

$$(2.1.5) \quad |S_{\varepsilon}(t,s)| < K \exp(a(t-s)) ,$$

where K and a are real constants independent of ε ($|\cdot|$ is the maximum norm). The following example demonstrates that stability is not simply determined by the real parts of the eigenvalues.

example 2.1.1

Let $A_1(t) = \begin{bmatrix} i & 1 \\ 0 & i \end{bmatrix}$ and $A_2(t, \varepsilon) = \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}$. The eigenvalues of $A(t, \varepsilon)$ are

$$\lambda_{\pm} = i/\varepsilon \pm 1/\sqrt{\varepsilon} .$$

Thus, the solution operator cannot be bounded independently of ε , and the system is not oscillatory stable. Note that the system with $A_2(t, \varepsilon) \equiv 0$ does not have a bounded solution operator although the eigenvalues are purely imaginary.

The following lemmas follow directly from our stability requirements.

LEMMA [2.1.1]

If the system (2.1.1) is oscillatory stable, then the solution of the inhomogeneous system (2.1.2) satisfies the estimate

$$|X(t, \varepsilon)| < K [\exp(at) |X_0| + R(a, t) \max_{(s, \varepsilon)} |F(s, \varepsilon)|]$$

$$(2.1.6) \quad \text{with } R(a, t) = \begin{cases} t & , a=0 \\ (\exp(at)-1)/a & , a \neq 0 \end{cases}$$

PROOF

From (2.1.4) and (2.1.5) we have

$$|X(t, \varepsilon)| < K [\exp(at) |X_0| + \max_{(s, \varepsilon)} |F(s, \varepsilon)| \int_0^t \exp(a(t-s)) ds]$$

and the lemma follows.

LEMMA [2.1.2]

Let the matrix $B(t, \varepsilon)$ be in $C^0(t, \varepsilon)$ with

$$d = \sup_{(t, \varepsilon)} |B(t, \varepsilon)| < \infty .$$

If the system (2.1.1) is oscillatory stable, then the system

$$(2.1.7) \quad \begin{cases} dW/dt = [A(t, \varepsilon) + B(t, \varepsilon)] W \\ W(0, \varepsilon) = W_0 \quad , \quad 0 < t < T \end{cases}$$

also satisfies an estimate of the form (2.1.5):

$$(2.1.8) \quad |W(t, \varepsilon)| < K |W_0| \exp([a+Kd]t) .$$

PROOF

By Duhammel's Principle we can use the solution operator of system (2.1.1) to reformulate system (2.1.7) as

$$W(t, \varepsilon) = S_{\varepsilon}(t, 0)W_0 + \int_0^t S_{\varepsilon}(t, s)B(s, \varepsilon)W(s, \varepsilon) ds$$

By (2.1.5) we have

$$|W(t, \varepsilon)| < \exp(at)K|W_0| + Kd \int_0^t \exp(a(t-s))|W(s, \varepsilon)| ds.$$

After the substitution

$$\widetilde{W}(t, \varepsilon) = W(t, \varepsilon) \exp(-at)$$

we have:

$$|\widetilde{W}(t, \varepsilon)| < K|W_0| + Kd \int_0^t |\widetilde{W}(s, \varepsilon)| ds.$$

The result now follows directly from a fundamental inequality of stability theory (see Coddington and Levinson [6], p.37).

LEMMA [2.1.3]

Let the systems (2.1.1) and (2.1.7) be oscillatory stable, and let $S_{\varepsilon}(t, s)$ be the solution operator of (2.1.1). Then if

$$E(t, \varepsilon) = W(t, \varepsilon) - X(t, \varepsilon)$$

we have:

$$(2.1.9) \quad E(t, \varepsilon) = S_{\varepsilon}(t, 0)[W_0 - X_0] + \int_0^t S_{\varepsilon}(t, s)B(s, \varepsilon)W(s, \varepsilon) ds$$

and

$$(2.1.10) \quad |E(t, \varepsilon)| < K \exp(at) [|W_0 - X_0| + |W_0| (\exp(Kdt) - 1)].$$

PROOF

The equations for $E(t, \varepsilon)$ can be written

$$(2.1.11) \quad \begin{cases} dE/dt = A(t, \varepsilon)E + B(t, \varepsilon)W(t, \varepsilon) \\ E(0, \varepsilon) = W_0 - X_0 \end{cases}$$

and the result follows directly from Duhammel's Principle and (2.1.8).

None of these results depends on assumption (v) for its validity, but now we demonstrate that assumption (v) can be strengthened to give a sufficient condition for oscillatory stability.

THEOREM [2.1.1]

If $A_1(t)$ can be reduced to diagonal form by a smooth transformation $W(t)$, then the system (2.1.1) with is oscillatory stable.

PROOF

By Lemma [2.1.2] it is sufficient to verify the result for the system

$$\begin{cases} dx/dt = [A_1(t)/\varepsilon]X \\ X(0, \varepsilon) = X_0 \end{cases}$$

The equations for

$$\widetilde{X} = W(t)X$$

become

$$\begin{cases} d\tilde{X}/dt = (W(A_1/\varepsilon) \tilde{W}^{-1} + dW/dt \tilde{W}^{-1}) \tilde{X} \\ \tilde{X}(0, \varepsilon) = X_0 \end{cases} .$$

Again by Lemma [2.1.2] we can ignore the $O(1)$ term of the system and consider the further reduced system

$$\begin{cases} d\tilde{X}/dt = (W(A_1/\varepsilon) \tilde{W}^{-1}) \tilde{X} \\ \tilde{X}(0, \varepsilon) = \tilde{X}_0 \end{cases} ,$$

which has diagonal structure and is clearly oscillatory stable since by assumption (v) the eigenvalues of $A_1(t)$ are purely imaginary. Thus, the system (2.1.1) is oscillatory stable.

2.2 Normal Forms and Canonical Transformations

In this section we investigate the use of similarity transformations to reduce the homogeneous problem (2.1.1) to a more tractable form. The system (2.1.1) is said to be in normal form (p,q) if the matrix $A(t, \varepsilon)$ can be written in the form

$$(2.2.1) \quad A(t, \varepsilon) = \frac{1}{\varepsilon} B(t, \varepsilon) + \varepsilon^q C(t, \varepsilon),$$

where we have

$$(2.2.2) \quad \frac{1}{\varepsilon} B(t, \varepsilon) = \begin{bmatrix} \frac{1}{\varepsilon} B_{11}(t, \varepsilon) & \frac{1}{\varepsilon} B_{12}(t, \varepsilon) \\ 0 & B_{22}(t, \varepsilon) \end{bmatrix}$$

and

$$(2.2.3) \quad C(t, \varepsilon) = \begin{bmatrix} C_{11}(t, \varepsilon) & C_{12}(t, \varepsilon) \\ C_{21}(t, \varepsilon) & C_{22}(t, \varepsilon) \end{bmatrix}.$$

We assume:

- (i) $\frac{1}{\varepsilon} B_{11}(t, \varepsilon)$ is an $(n_1 \times n_1)$ block diagonal submatrix with each diagonal block a (2×2) submatrix of the form

$$(2.2.4) \quad \begin{bmatrix} a_k(t, \varepsilon) & \frac{1}{\varepsilon} b_k(t, \varepsilon) \\ -\frac{1}{\varepsilon} b_k(t, \varepsilon) & a_k(t, \varepsilon) \end{bmatrix}$$

where $a_k(t, \varepsilon)$ and $b_k(t, \varepsilon)$ are in $C^p(t, \varepsilon)$ with

$$(2.2.5) \quad \begin{cases} \min_t |b_k(t, \varepsilon)| > \tilde{b} > 0 \\ \min_{\substack{t \\ k \neq j}} | |b_k(t, \varepsilon)| - |b_j(t, \varepsilon)| | > \tilde{b} > 0 \end{cases}$$

- (ii) $B_{22}(t, \varepsilon)$ and $C_{22}(t, \varepsilon)$ are $(n_2 \times n_2)$ submatrices which are in $C^p(t, \varepsilon)$.
- (iii) $B_{12}(t, \varepsilon)$ and $C_{12}(t, \varepsilon)$ are $(n_1 \times n_2)$ submatrices which are in $C^p(t, \varepsilon)$.
- (iv) $C_{21}(t, \varepsilon)$ is an $(n_2 \times n_1)$ submatrix which is in $C^p(t, \varepsilon)$.
- (v) $C_{11}(t, \varepsilon)$ is an $(n_1 \times n_1)$ submatrix which is in $C^p(t, \varepsilon)$.
- (vi) q is a nonnegative integer and p is a positive integer.

THEOREM [2.2.1]

If the system (2.1.1) is in normal form (p, q) , then it is oscillatory stable.

PROOF

By Lemma [2.1.2] it is sufficient to consider the reduced system

$$\begin{cases} dx/dt = (B(t, \varepsilon)/\varepsilon) x \\ x(0, \varepsilon) = x_0 \end{cases} .$$

After the transformation

$$\tilde{X}(t, \varepsilon) = S(t, \varepsilon) X(t, \varepsilon),$$

where

$$S(t, \varepsilon) = \begin{bmatrix} I & (B_{11}(t, \varepsilon))^{-1} B_{12}(t, \varepsilon) \\ \emptyset & I \end{bmatrix} ,$$

we can again by Lemma [2.1.2] neglect $O(1)$ terms to give:

$$(2.2.6) \quad \begin{cases} d\tilde{X}/dt = (\tilde{B}(t, \varepsilon)/\varepsilon) \tilde{X} \\ \tilde{X}(0, \varepsilon) = \tilde{X}_0 \end{cases} ,$$

where

$$\tilde{B}(t, \varepsilon) = -\tilde{B}(t, \varepsilon)^T = \begin{bmatrix} B_{11}(t, \varepsilon) & 0 \\ 0 & 0 \end{bmatrix} .$$

The system (2.2.6) is clearly oscillatory stable since

$$d/dt(\tilde{X}^T \tilde{X}) = \tilde{X}^T (\tilde{B}(t, \varepsilon)/\varepsilon + \tilde{B}(t, \varepsilon)^T/\varepsilon) \tilde{X} = 0.$$

From the results of a later section one can prove the following result, which characterizes the eigenstructure of the oscillatory set.

THEOREM [2.2.2]

For the system (2.1.1) in normal form (p, q) there exists, correspondingly to the oscillatory set, a set $\{(u_k(t, \varepsilon), \lambda_k(t, \varepsilon))\}$ of n_1 normalized eigenpairs which are in $C^p(t, \varepsilon)$ and which satisfy

$$(2.2.7) \quad A(t, \varepsilon) u_k(t, \varepsilon) = \lambda_k(t, \varepsilon) u_k(t, \varepsilon) ,$$

where for k odd we have

$$(2.2.8) \quad \begin{cases} u_k = e_k + i e_{k+1} + \varepsilon^{2+1} \tilde{u}_k(t, \varepsilon) \\ u_{k+1} = \bar{u}_k \quad (\{\tilde{u}_k(t, \varepsilon)\} \subset C^P(t, \varepsilon)) \\ \varepsilon \lambda_k = \varepsilon a_k(t, \varepsilon) + i b_k(t, \varepsilon) + \varepsilon^{2+1} \tilde{\lambda}_k(t, \varepsilon) \\ \lambda_{k+1} = \bar{\lambda}_k \quad (\{a_k(t, \varepsilon), b_k(t, \varepsilon), \tilde{\lambda}_k(t, \varepsilon)\} \subset C^P(t, \varepsilon)) \end{cases}$$

The vectors $\{e_k\}$ are the Euclidian unit vectors.

PROOF

The theorem follows from Corollary [A1.2a] in the appendix of this chapter and the eigenstructure of the (2×2) matrix (2.2.4).

It is convenient to represent the eigenstructure in another form derived from (2.2.8). For real systems we can assume, without loss of generality that for the oscillatory set

$$(2.2.9) \quad \lambda_k(t, \varepsilon) = \bar{\lambda}_{k+1}(t, \varepsilon) = \hat{a}_k(t, \varepsilon) + i \frac{1}{\varepsilon} \hat{b}_k(t, \varepsilon) \quad (k \text{ odd}),$$

and similarly

$$(2.2.10) \quad u_k(t, \varepsilon) = \bar{u}_{k+1}(t, \varepsilon) \quad (k \text{ odd}).$$

The space spanned by $\{u_k\}$ can be equivalently represented by the $n/2$ real, independent vectors $\{v_k\}$ which are given by

$$(2.2.11) \quad \begin{cases} v_k(t, \varepsilon) = \text{Re}\{u_k(t, \varepsilon)\} & (k \text{ odd}) \\ v_{k+1}(t, \varepsilon) = \text{Im}\{u_k(t, \varepsilon)\} & (k \text{ odd}) \end{cases}$$

For convenience we define $V(t, \varepsilon)$ to be the $(n \times n/2)$ matrix whose

k-th column is $v_k(t, \varepsilon)$; furthermore, let V_0 be the $(n \times n_1)$ matrix whose k-th column is e_k . We note that

$$(2.2.12) \quad A(t, \varepsilon) V(t, \varepsilon) = V(t, \varepsilon) L(t, \varepsilon) \quad ,$$

where $L(t, \varepsilon)$ is an $(n_1 \times n_1)$ block diagonal matrix with blocks of the form given in (2.2.4). Hereafter we shall refer to the matrix $V(t, \varepsilon)$ as a canonical representation for the oscillatory set.

From the eigenstructure corresponding to the oscillatory set, we construct transformations designed for the separation of time scales.

THEOREM [2.2.3]

Assume the system (2.1.1) is in canonical form (p, q) with $(p \geq 2)$, and let $V(t, \varepsilon)$ be the canonical representation for the oscillatory set. If $S(t, \varepsilon)$ is a transformation of the form

$$(2.2.13) \quad S(t, \varepsilon) = I + \varepsilon^{q+1} T(t, \varepsilon) \quad ,$$

where $T(t, \varepsilon)$ is in $C^p(t, \varepsilon)$, and if

$$(2.2.14) \quad S(t, \varepsilon) V(t, \varepsilon) = V_0 = \begin{bmatrix} I \\ \emptyset \end{bmatrix} \quad ,$$

then the system for

$$(2.2.15) \quad \tilde{X}(t, \varepsilon) = S(t, \varepsilon) X(t, \varepsilon)$$

is in normal form $(p-1, q+1)$. $S(t, \varepsilon)$ is then called a canonical transformation.

PROOF

Since

$$(S(t, \varepsilon))^{-1} = I - \varepsilon^{q+1} T(t, \varepsilon) + O(\varepsilon^{2q+2})$$

we have

$$S(t, \varepsilon) A(t, \varepsilon) (S(t, \varepsilon))^{-1} = A(t, \varepsilon) + O(\varepsilon^q).$$

Moreover, the first columns of this matrix are given by

$$\begin{aligned} S(t, \varepsilon) A(t, \varepsilon) (S(t, \varepsilon))^{-1} V_0 &= S(t, \varepsilon) A(t, \varepsilon) V(t, \varepsilon) \\ &= S(t, \varepsilon) V(t, \varepsilon) L(t, \varepsilon) \\ &= V_0 L(t, \varepsilon) \end{aligned}$$

where $L(t, \varepsilon)$ has the form given in (2.212). Since

$$dS/dt = O(\varepsilon^{q+1}),$$

the equations for \tilde{X} are in normal form $(p-1, q+1)$.

Next we demonstrate that such transformations can be constructed from the canonical representation for the oscillatory set.

THEOREM [2.2.4]

Let the system (2.1.1) be in normal form (p, q) , and let $V(t, \varepsilon)$ be a canonical representation for the oscillatory set. Then the transformation $S(t, \varepsilon)$, whose inverse $W(t, \varepsilon)$ has its columns given by

$$(2.2.16) \quad W(t, \varepsilon)^{(k)} = \begin{cases} V(t, \varepsilon)^{(k)} & , \quad 1 < k < n_1 \\ e_k & , \quad n_1 < k < n \end{cases}$$

satisfies the requirements of Theorem [2.2.3].

PROOF

The theorem follows immediately from Theorem [2.2.2] and the definition of the canonical representation for the oscillatory set.

Let $\|\cdot\|$ be the Euclidian norm. Through the following variation of the Q-R Factorization Theorem, we can construct other forms for canonical transformations.

LEMMA [2.2.1]

If u is an n -vector such that

$$u \neq -e_k$$

and if the reflection H is defined by

$$\begin{cases} H = I - 2ww^T \\ w = (u + \sigma e_k) / \|u + \sigma e_k\| \quad , \quad \sigma = \|u\| \end{cases}$$

then

$$Hu = -\sigma e_k.$$

PROOF

We have:

$$\begin{cases} Hu = u - 2[(u + \sigma e_k)^T u / \|u + \sigma e_k\|^2] (u + \sigma e_k) \\ = -\sigma e_k \end{cases}$$

since

$$\|u + \sigma e_k\|^2 = 2(\sigma^2 + \sigma u^{(k)}) .$$

THEOREM [2.2.5]

Let the system (2.1.1) be in normal form (p, q) , and let $V(t, \varepsilon)$ be a canonical representation for the oscillatory set. Then the conditions of Theorem [2.2.3] are satisfied by a unique transformation of the form

$$(2.2.17) \quad S(t, \varepsilon) = \begin{bmatrix} R(t, \varepsilon) & \emptyset \\ \emptyset & I \end{bmatrix} Q(t, \varepsilon) ,$$

where:

(i) the $(n_1 \times n_1)$ submatrix $R(t, \varepsilon)$ has the form

$$(2.2.18) \quad R(t, \varepsilon) = I + \varepsilon^{q_H} \tilde{R}(t, \varepsilon) \quad (\tilde{R}(t, \varepsilon) \in C^p(t, \varepsilon)) ,$$

where $\tilde{R}(t, \varepsilon)$ is upper triangular;

(ii) $Q(t, \varepsilon)$ has the form

$$(2.2.19) \quad Q(t, \varepsilon) = Q_{n_1}(t, \varepsilon) Q_{n_1-1}(t, \varepsilon) \dots Q_1(t, \varepsilon) ,$$

where each $Q_k(t, \varepsilon)$ is expressible as a product of Householder transformations:

$$(2.2.20) \quad \begin{cases} Q_K(t, \varepsilon) = [I - 2e_K e_K^T] [I - 2w_K(t, \varepsilon)w_K(t, \varepsilon)^T] \\ w_K(t, \varepsilon) = e_K + \varepsilon^{q+1} z_K(t, \varepsilon) \quad (z_K(t, \varepsilon) \in C^P(t, \varepsilon)) \end{cases}$$

whereby

$$(2.2.21) \quad Q_K(t, \varepsilon) = I + \varepsilon^{2q+1} W_K(t, \varepsilon) \quad (W_K(t, \varepsilon) \in C^P(t, \varepsilon)).$$

PROOF

By Theorem [2.2.2] and (2.2.11) we have

$$(2.2.22) \quad \begin{cases} v_1(t, \varepsilon) = V(t, \varepsilon)^{(1)} \\ \quad \quad \quad = e_1 + \varepsilon^{q+1} z_1(t, \varepsilon) \\ z_1(t, \varepsilon) \in C^P(t, \varepsilon) \end{cases}$$

$Q_1(t, \varepsilon)$ now can be defined by (2.2.20) with

$$(2.2.23) \quad \begin{cases} w_1(t, \varepsilon) = (v_1(t, \varepsilon) + \sigma_1 e_1) / \|v_1(t, \varepsilon) + \sigma_1 e_1\| \\ \quad \quad \quad = e_1 + \varepsilon^{q+1} y_1(t, \varepsilon) \end{cases}$$

where

$$(2.2.24) \quad \begin{cases} \sigma_1(t, \varepsilon) = |v_1(t, \varepsilon)| \\ y_1(t, \varepsilon) \in C^P(t, \varepsilon) \end{cases}$$

to give by Lemma [2.2.1]

$$(2.2.25) \quad Q_1(t, \varepsilon) v_1(t, \varepsilon) = \sigma_1(t, \varepsilon) e_1.$$

We now proceed by induction. After the elements of $\{Q_1, Q_2, \dots, Q_{k-1}\}$ have been calculated as in (2.2.20), we define

$$\begin{aligned}
 \widetilde{v}_k(t, \varepsilon) &= Q_{k-1}(t, \varepsilon) Q_{k-2}(t, \varepsilon) \dots Q_1(t, \varepsilon) v(t, \varepsilon)^{(k)} \\
 (2.2.26) \qquad &= e_k + \varepsilon^{q+1} z_k(t, \varepsilon)
 \end{aligned}$$

and by (2.2.11), Theorem [2.2.2], and our inductive hypothesis we have $z_k(t, \varepsilon) \in C^p(t, \varepsilon)$. Next we construct the projection of \widetilde{v}_k onto the span of $\{e_k, e_{k+1}, \dots, e_n\}$:

$$\begin{aligned}
 \widetilde{\widetilde{v}}_k(t, \varepsilon) &= \widetilde{v}_k(t, \varepsilon) - \sum_{j=1}^{k-1} (e_j^T \widetilde{v}_k) e_j \\
 (2.2.27) \qquad &= e_k + \varepsilon^{q+1} \widetilde{z}_k(t, \varepsilon)
 \end{aligned}$$

where $\widetilde{z}_k(t, \varepsilon) \in C^p(t, \varepsilon)$. And now we define $Q_k(t, \varepsilon)$ as in (2.2.20) with

$$(2.2.28) \quad \begin{cases} w_k(t, \varepsilon) = (\widetilde{\widetilde{v}}_k(t, \varepsilon) + \sigma_k e_k) / \|\widetilde{\widetilde{v}}_k(t, \varepsilon) + \sigma_k e_k\| \\ = e_k + \varepsilon^{q+1} y_k(t, \varepsilon) \end{cases}$$

where

$$(2.2.29) \quad \begin{cases} \sigma_k(t, \varepsilon) = \|\widetilde{\widetilde{v}}_k(t, \varepsilon)\| \\ y_k(t, \varepsilon) \in C^p(t, \varepsilon) \end{cases}$$

to give by Lemma [2.2.1]

$$(2.2.30) \quad \begin{cases} Q_k(t, \varepsilon) e_j = e_j & (j < k) \\ Q_k(t, \varepsilon) \widetilde{\widetilde{v}}_k(t, \varepsilon) = \sigma_k(t, \varepsilon) e_k \end{cases}$$

Thus, with $Q(t, \varepsilon)$ defined by (2.2.19) we have

$$(2.2.31) \quad Q(t, \varepsilon) V(t, \varepsilon) = \begin{bmatrix} U(t, \varepsilon) \\ \emptyset \end{bmatrix},$$

where the $(n_1 \times n_1)$ submatrix $U(t, \varepsilon)$ is upper triangular.

Furthermore, we can write

$$(2.2.32) \quad U(t, \varepsilon) = I + \varepsilon^{q+l} \tilde{U}(t, \varepsilon) \quad [\tilde{U}(t, \varepsilon) \in C^p(t, \varepsilon)]$$

since $Q(t, \varepsilon)$ and $V(t, \varepsilon)$ also have this property. With

$$(2.2.33) \quad R(t, \varepsilon) = U(t, \varepsilon)^{-1}$$

in (2.2.17) we have the desired result.

The smoothness of the transformations of Theorem [2.2.4] and Theorem [2.2.5] leads to the following useful result which characterizes the effect of perturbations of $V(t, \varepsilon)$.

THEOREM [2.2.6]

Let $V(t, \varepsilon)$ in the previous theorems be replaced by the perturbation

$$(2.2.34) \quad \tilde{V}(t, \varepsilon, \delta) = V(t, \varepsilon) + \delta F(t, \varepsilon),$$

where $F(t, \varepsilon) \in C^r(t, \varepsilon)$ ($r \leq p$) and δ is a small scalar parameter. Then for sufficiently small δ the derived transformations of Theorems [2.2.4] and [2.2.5] are likewise perturbed:

$$(2.2.35) \quad \left\{ \begin{array}{l} \widetilde{S}(t, \varepsilon, \delta) = S(t, \varepsilon) + f(t, \varepsilon, \delta) \quad \text{in (2.2.16)} \\ \widetilde{R}(t, \varepsilon, \delta) = R(t, \varepsilon) + g(t, \varepsilon, \delta) \quad \text{in (2.2.18)} \\ \widetilde{w}_k(t, \varepsilon, \delta) = w_k(t, \varepsilon) + h_k(t, \varepsilon, \delta) \quad \text{in (2.2.21)} \\ \{f, g, h_k\} \subset C^r(t, \varepsilon, \delta) \end{array} \right. .$$

PROOF

The theorem follows from the smoothness of the transformations, which are given explicitly by ((2.2.16) - (2.2.32)).

By the straightforward application of Theorem [2.2.3] we can transform the basic equations to an improved normal form; by the repeated application of this procedure we can reduce the original system to a more manageable formulation.

THEOREM [2.2.7]

Let the system (2.1.1) be in normal form (p,q), and let the positive integer r be less than p. Then there exists a transformation $S(t, \varepsilon)$, expressible as a product of canonical transformations, such that the equations for

$$(2.2.36) \quad \widetilde{X}(t, \varepsilon) = S(t, \varepsilon) X(t, \varepsilon)$$

can be written as

$$(2.2.37) \quad \left\{ \begin{array}{l} d\widetilde{X}/dt = A(t, \varepsilon) \widetilde{X} \\ X(0, \varepsilon) = X_0 \end{array} \right. ,$$

where the system (2.2.37) is in canonical form (p-r, q+r) with

$$(2.2.38) \quad \widetilde{A}(t, \varepsilon) = \widetilde{B}(t, \varepsilon)/\varepsilon + \varepsilon^{q+r} \widetilde{C}(t, \varepsilon).$$

Furthermore, if $Y(t, \varepsilon)$ is the solution of the system

$$(2.2.39) \quad \begin{cases} dY/dt = [B(t, \varepsilon)/\varepsilon] Y \\ Y(0, \varepsilon) = Y_0 \end{cases},$$

then

$$(2.2.40) \quad \max_t |\widetilde{X}(t, \varepsilon) - Y(t, \varepsilon)| = O(X_0 - Y_0) + O(\varepsilon^{2+r}).$$

PROOF

The theorem follows directly from Theorem [2.2.3], Theorem [2.2.4] (or Theorem [2.2.5]), and Lemma [2.1.3].

The reduction of system (2.1.1) to system (2.2.39) represents the major result of this chapter, for the presence of the small parameter in the latter formulation no longer places serious practical limitations on our ability to represent or to approximate the solution. Since the time scales have been effectively separated, the now apparent singular nature of the problem can be treated by asymptotic techniques. For example, consider the system (2.2.39) with $B(t, \varepsilon)$ as in (2.2.2) and with the decomposition

$$(2.2.41) \quad Y = \begin{pmatrix} Y^I \\ Y^{II} \end{pmatrix}$$

where Y^I is n_1 -dimensional and Y^{II} is n_2 -dimensional. Clearly Y^I can be represented as the solution of the system

$$(2.2.42) \quad \begin{cases} dY^{\mathbb{I}}/dt = B_{22}(t, \varepsilon) Y^{\mathbb{I}} \\ Y^{\mathbb{I}}(0, \varepsilon) = Y_0^{\mathbb{I}} \end{cases},$$

where the coefficients depend smoothly on the small parameter.

The components of $Y^{\mathbb{I}}$ are handled in pairs which correspond to the block diagonal structure of $B_{11}(t, \varepsilon)$; therefore, without loss of generality we can assume that $Y^{\mathbb{I}}$ is two-dimensional. The appropriate equations have the form:

$$(2.2.43) \quad d/dt \begin{pmatrix} u \\ v \end{pmatrix} = \begin{bmatrix} a(t, \varepsilon) & b(t, \varepsilon)/\varepsilon \\ -b(t, \varepsilon)/\varepsilon & a(t, \varepsilon) \end{bmatrix} \begin{pmatrix} u \\ v \end{pmatrix} + \begin{pmatrix} f_1(t, \varepsilon)/\varepsilon \\ f_2(t, \varepsilon)/\varepsilon \end{pmatrix}$$

with the appropriate initial conditions

$$(2.2.44) \quad \begin{pmatrix} u(0, \varepsilon) \\ v(0, \varepsilon) \end{pmatrix} = \begin{pmatrix} u_0 \\ v_0 \end{pmatrix}.$$

The change of variables

$$(2.2.45) \quad \begin{cases} \begin{pmatrix} \tilde{u} \\ \tilde{v} \end{pmatrix} = \begin{pmatrix} u \\ v \end{pmatrix} \exp(-\hat{A}_\varepsilon(t, 0)) \\ \tilde{t} = \hat{B}_\varepsilon(t, 0) \end{cases},$$

where

$$(2.2.46) \quad \begin{cases} \hat{A}_\varepsilon(t, \tau) = \int_\tau^t a(r, \varepsilon) dr \\ \hat{B}_\varepsilon(t, \tau) = \int_\tau^t b(r, \varepsilon) dr \end{cases},$$

gives for the homogeneous problem

$$(2.2.47) \quad \begin{cases} d/d\tilde{t} \begin{pmatrix} \tilde{u} \\ \tilde{v} \end{pmatrix} = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \begin{pmatrix} \tilde{u} \\ \tilde{v} \end{pmatrix} , \\ \tilde{u}(0, \varepsilon) = \tilde{u}_0 , \quad \tilde{v}(0, \varepsilon) = \tilde{v}_0 \end{cases} ,$$

which has the solution operator

$$(2.2.48) \quad \tilde{S}_\varepsilon(\tilde{t}, \tilde{\tau}) = \begin{bmatrix} \cos(\tilde{t}-\tilde{\tau}) & \sin(\tilde{t}-\tilde{\tau}) \\ -\sin(\tilde{t}-\tilde{\tau}) & \cos(\tilde{t}-\tilde{\tau}) \end{bmatrix} .$$

Therefore, the solution operator for the system ((2.2.43)-(2.2.44)) is

$$(2.2.49) \quad S_\varepsilon(t, \tau) = \exp(\hat{A}_\varepsilon(t, \tau)) \begin{bmatrix} \cos(\hat{B}_\varepsilon(t, \tau)/\varepsilon) & \sin(\hat{B}_\varepsilon(t, \tau)/\varepsilon) \\ -\sin(\hat{B}_\varepsilon(t, \tau)/\varepsilon) & \cos(\hat{B}_\varepsilon(t, \tau)/\varepsilon) \end{bmatrix} ,$$

and the solution of the original system is explicitly given by Duhammel's Principle (2.1.4). For the first component we have

$$(2.2.50) \quad \begin{aligned} u(t, \varepsilon) = & \exp(\hat{A}_\varepsilon(t, 0)) \cos(\hat{B}_\varepsilon(t, 0)/\varepsilon) u_0 \\ & + \exp(\hat{A}_\varepsilon(t, 0)) \sin(\hat{B}_\varepsilon(t, 0)/\varepsilon) v_0 \\ & + \int_0^t \exp(\hat{A}_\varepsilon(t, s)) \cos(\hat{B}_\varepsilon(t, s)/\varepsilon) f_1(s, \varepsilon)/\varepsilon ds \\ & + \int_0^t \exp(\hat{A}_\varepsilon(t, s)) \sin(\hat{B}_\varepsilon(t, s)/\varepsilon) f_2(s, \varepsilon)/\varepsilon ds . \end{aligned}$$

An asymptotic expansion for the integrals easily can be generated by integration by parts. For example, we have:

$$(2.2.51) \quad \begin{aligned} \int_0^t \cos(\hat{B}_\varepsilon(t, s)/\varepsilon) f(s, \varepsilon)/\varepsilon ds = & \\ & \Big|_0^t \{ \sin(\hat{B}_\varepsilon(t, s)) [-f(s, \varepsilon)/b(s, \varepsilon)] \} \\ & + \Big|_0^t \{ \varepsilon \cos(\hat{B}_\varepsilon(t, s)) [(d/ds) (f(s, \varepsilon)/b(s, \varepsilon))] [1/b(s, \varepsilon)] \} \\ & + \dots + O(\varepsilon^m) \end{aligned}$$

and

$$\begin{aligned}
 & \int_0^t \sin(\widehat{B}_\varepsilon(t,s)/\varepsilon) f(s,\varepsilon)/\varepsilon \, ds = \\
 & \quad \int_0^t \{ \cos(\widehat{B}_\varepsilon(t,s)) [f(s,\varepsilon)/b(s,\varepsilon)] \} \\
 (2.2.52) \quad & + \int_0^t \varepsilon \{ \sin(\widehat{B}_\varepsilon(t,s)) [(d/ds) (f(s,\varepsilon)/b(t,\varepsilon))] [1/b(s,\varepsilon)] \} \\
 & + \dots + O(\varepsilon^m)
 \end{aligned}$$

In subsequent sections we shall use these theoretical results to justify algorithms suitable for computation with large time steps.

2.3 REDUCTION TO CANONICAL FORM

The purpose of this section is to reduce the general problem (2.1.1) to the normal form described in Section 2.2. By Theorem [A1.2] of Appendix A, there exist n_1 smooth eigenpairs which correspond to the oscillatory set:

$$(2.3.1) \quad \{[\lambda_k(t, \varepsilon), u_k(t, \varepsilon)]\},$$

where

$$(2.3.2) \quad \begin{cases} \{\lambda_k(t, \varepsilon)\} \subset C^p(t, \varepsilon) \\ \{u_k(t, \varepsilon)\} \subset C^p(t, \varepsilon) \end{cases} .$$

As in the previous section we assume, without loss of generality,

$$(2.3.3) \quad \begin{cases} \lambda_k(t, \varepsilon) = \overline{\lambda_{k+1}(t, \varepsilon)} & , k \text{ odd} \\ u_k(t, \varepsilon) = \overline{u_{k+1}(t, \varepsilon)} & , k \text{ odd} \end{cases} ,$$

and so the eigenspace can be equivalently represented by $\{v_k\}$, where

$$(2.3.4) \quad \begin{cases} v_k(t, \varepsilon) = \operatorname{Re}\{u_k(t, \varepsilon)\} & , k \text{ odd} \\ v_{k+1}(t, \varepsilon) = \operatorname{Im}\{u_k(t, \varepsilon)\} & , k \text{ odd} \end{cases} .$$

As in the previous section we define the $(n \times n_1)$ matrix $V(t, \varepsilon)$ by

$$(2.3.5) \quad V(t, \varepsilon)^{(k)} = v_k(t, \varepsilon)$$

and refer to $V(t, \varepsilon)$ as a canonical representation for the oscillatory set. Furthermore, let V_0 be the $(n \times n_1)$ matrix whose k -th column is e_k . Theorem [2.2.3] now must be replaced by the following assumption.

Assumption (2.3a)

If there exists a transformation $S(t, \varepsilon)$ such that

$$(2.3.6) \quad \begin{cases} \{S(t, \varepsilon), S(t, \varepsilon)^{-1}\} \subset C^p(t, \varepsilon) \\ S(t, \varepsilon)V(t, \varepsilon) = V_0 \end{cases}$$

then the equations for

$$(2.3.7) \quad \tilde{X}(t, \varepsilon) = S(t, \varepsilon)X(t, \varepsilon)$$

are in normal form $(p-1, 0)$.

By (2.3.6) the structure of the first n_1 columns is clearly that of (2.2.1), and also the first n_1 rows of the last $n_2 (=n-n_1)$ columns have the appropriate form; thus, the assumption simply insures that no $O(1/\varepsilon)$ terms appear in the lower $(n_2 \times n_2)$ diagonal submatrix, which must correspond to the nonoscillatory set. If unbounded terms do exist in this submatrix, then smooth perturbations of the system can change the number of large eigenvalues.

example (2.3a)

The system

$$(2.3.8) \quad dx/dt = \begin{bmatrix} 0 & 1/\varepsilon \\ 0 & 0 \end{bmatrix} x, \quad x(0, \varepsilon) = x_0$$

has eigenvalues

$$(2.3.9) \quad \lambda_{\pm} = 0$$

while the perturbed system

$$(2.3.10) \quad dx/dt = \begin{bmatrix} 0 & 1/\varepsilon \\ \delta & 0 \end{bmatrix} x, \quad x(0, \varepsilon) = x_0$$

has eigenvalues

$$(2.3.11) \quad \tilde{\lambda}_{\pm} = \pm \sqrt{\delta/\varepsilon}.$$

Moreover, neither system is oscillatory stable since neither has a bounded solution operator.

In Section 2.1 two methods for the construction of such transformations were demonstrated. The method of Theorem [2.2.4] seems not so useful here since one cannot conveniently augment the matrix $V(t, \varepsilon)$ so as to span the underlying n -dimensional space; however, we have the following generalization of Theorem [2.2.5].

THEOREM [2.3.1]

Let $V(t, \varepsilon)$ be a canonical representation for the oscillatory set; then there exists a locally valid canonical transformation $S(t, \varepsilon)$, where

$$(2.3.12) \quad S(t, \varepsilon) = \begin{bmatrix} R(t, \varepsilon) & \emptyset \\ \emptyset & I \end{bmatrix} Q(t, \varepsilon);$$

the $(n_1 \times n_1)$ submatrix $R(t, \varepsilon)$ is upper triangular with

$$(2.3.13) \quad R(t, \varepsilon) \in C^P(t, \varepsilon)$$

and

$$(2.3.14) \quad Q(t, \varepsilon) = P Q_{n_1}(t, \varepsilon) Q_{n_1-1}(t, \varepsilon) \dots Q_1(t, \varepsilon),$$

where P is a constant permutation matrix and

$$(2.3.15) \quad Q_k(t, \varepsilon) = [I - 2e_{p_k} e_{p_k}^T] [I - 2w_k(t, \varepsilon) w_k(t, \varepsilon)^T]$$

with

$$(2.3.16) \quad \begin{cases} w_k(t, \varepsilon) \in C^P(t, \varepsilon) \\ 1 \leq p_k \leq n_1 \end{cases} .$$

PROOF

Essentially we repeat the constructive argument of Theorem [2.2.5]. Let p_1 be chosen so that the p_1 -component of $v_1(t, \varepsilon)$ is locally bounded away from zero. Without loss of generality, we assume that for some \tilde{T}_1 ,

$$(v_1(t, \varepsilon))^{(p_1)} > \emptyset \quad (0 \leq t \leq \tilde{T}_1 \leq T)$$

since otherwise a renormalization is possible. By Lemma [2.2.1] $Q_1(t, \varepsilon)$ is given by (2.3.15) with

$$\begin{cases} w_1(t, \varepsilon) = (v_1(t, \varepsilon) + \sigma_1 e_{p_1}) / \|v_1(t, \varepsilon) + \sigma_1 e_{p_1}\| \\ \sigma_1(t, \varepsilon) = \|v_1(t, \varepsilon)\| \end{cases}$$

We proceed by induction. After the elements of $\{Q_1, Q_2, \dots, Q_{k-1}\}$ have been determined, we define

$$\widetilde{v}_k(t, \varepsilon) = Q_{k-1}(t, \varepsilon) Q_{k-2}(t, \varepsilon) \dots Q_1(t, \varepsilon) v_k(t, \varepsilon),$$

where by our inductive hypothesis

$$\widetilde{v}_k(t, \varepsilon) \in C^p(t, \varepsilon).$$

$\widehat{\widetilde{v}}_k(t, \varepsilon)$ is defined as the projection of this vector onto the complement of the span of $\{e_{p_1}, e_{p_2}, \dots, e_{p_{k-1}}\}$:

$$\widehat{\widetilde{v}}_k(t, \varepsilon) = \widetilde{v}_k(t, \varepsilon) - \sum_{j=1}^{k-1} (e_{p_j}^T \widetilde{v}_k(t, \varepsilon)) e_{p_j}.$$

We choose p_k and \widetilde{T}_k such that

$$\widehat{\widetilde{v}}_k(t, \varepsilon)^{(p_k)} > 0 \quad (0 \leq t \leq \widetilde{T}_k \leq \widetilde{T}_{k-1}),$$

and now with $Q_k(t, \varepsilon)$ as in (2.3.15) with

$$\begin{cases} w_k(t, \varepsilon) = (\widehat{\widetilde{v}}_k(t, \varepsilon) + \sigma_k e_{p_k}) / \|\widehat{\widetilde{v}}_k(t, \varepsilon) + \sigma_k e_{p_k}\| \\ \sigma_k(t, \varepsilon) = \|\widehat{\widetilde{v}}_k(t, \varepsilon)\| \end{cases}$$

we have

$$\begin{cases} Q_k(t, \varepsilon) e_{p_j} = e_{p_j} \quad (j < k) \\ Q_k(t, \varepsilon) \widehat{\widetilde{v}}_k(t, \varepsilon) = \sigma_k(t, \varepsilon) e_{p_k} \end{cases}$$

Here $\tilde{U}_k(t, \varepsilon)$ cannot vanish because of the linear independence of $\{v_k(t, \varepsilon)\}$.

This construction gives

$$\tilde{Q}(t, \varepsilon)V(t, \varepsilon) = P^T \begin{bmatrix} U(t, \varepsilon) \\ \\ \emptyset \end{bmatrix},$$

where P is a constant permutation matrix, the $(n_1 \times n_1)$ submatrix $U(t, \varepsilon)$ is upper triangular, and

$$\tilde{Q}(t, \varepsilon) = Q_{n_1}(t, \varepsilon)Q_{n_1-1}(t, \varepsilon)\dots Q_1(t, \varepsilon);$$

moreover, since the left-hand side of this equation is in $C^p(t, \varepsilon)$ the right-hand side is also in $C^p(t, \varepsilon)$. Thus we have for $0 \leq t \leq \tilde{T}_{n_1}$ the form given by (2.3.12) with

$$R(t, \varepsilon) = U(t, \varepsilon)^{-1}.$$

Singularities in the unitary transformations arise when $\|v_k + e_{p_k}\|$ becomes small; then a small change in $\tilde{v}_k(t, \varepsilon)$ may mean a large change in $Q_k(t, \varepsilon)$. In fact the transformation is not even continuous in a neighborhood of

$$(2.3.17) \quad \|\tilde{v}_k(t, \varepsilon) + e_{p_k}\| = 0.$$

Note that in Theorem [2.3.5] $V(t, \varepsilon)$ is normalized so that

$$(2.3.18) \quad v_k(t, \varepsilon) \approx e_k,$$

and so singularities are avoided.

The following smoothness result, which is analagous to Theorem [2.2.6], will also be useful.

THEOREM [2.3.2]

Let $V(t, \varepsilon)$ in the previous theorem be replaced by the perturbation

$$(2.3.19) \quad \left\{ \begin{array}{l} \widetilde{V}(t, \varepsilon, \delta) = V(t, \varepsilon) + \delta F(t, \varepsilon) \\ F(t, \varepsilon) \in C^r(t, \varepsilon) \quad (r < p) \end{array} \right.$$

where δ is a small scalar parameter. Then for sufficiently small δ the derived transformations of Theorem [2.3.1] are likewise perturbed:

$$(2.3.20) \quad \left\{ \begin{array}{l} \widetilde{R}(t, \varepsilon, \delta) = R(t, \varepsilon) + \delta g(t, \varepsilon, \delta) \\ \widetilde{w}_k(t, \varepsilon, \delta) = w_k(t, \varepsilon) + h_k(t, \varepsilon, \delta) \\ \{g, h_k\} \subset C^r(t, \varepsilon, \delta) \end{array} \right. .$$

PROOF

The theorem follows from the smoothness of the transformations, which are given explicitly in the proof of Theorem [2.3.1].

2.4 NUMERICAL METHODS

The theory developed in the previous sections leads very naturally to the formulation of algorithms which are well-suited to highly oscillatory problems. In practice, the solutions of differential equations must be approximated through the techniques of discretization. In this context it is sufficient to consider the homogeneous problem; thus, we study the following reformulation of (2.1.1):

$$(2.4.1) \quad \begin{cases} dX/dt = A(t, \varepsilon) X \\ X(0, \varepsilon) = X_0, \quad 0 < t < T \end{cases}$$

where the system satisfies assumption (2.2a) and where the solution must be approximated on the grid $\{t_k\}$ with

$$(2.4.2) \quad \begin{cases} h_k = (t_{k+1} - t_k) \\ h = \max_k |h_k| \\ t_0 = 0, \quad t_N = T \end{cases}$$

The system (2.4.1) can be analytically transformed to the normal form (p, q) by Theorem [2.2.7] and Theorem [2.3.1] by which we have:

$$(2.4.3) \quad \begin{cases} A_0 = A(t, \varepsilon) \\ A_{N+1} = S_N A_N S_N^{-1} + S_N' S_N^{-1} \end{cases}$$

Here S_N , the canonical transformation associated with A_N , can be explicitly constructed from the canonical representation vectors

that correspond to the oscillatory set of A_N . Thus, from (2.2.12) one must calculate the solutions of the equation

$$(2.4.4) \quad G(A_N(t, \varepsilon), L, V) = A_N(t, \varepsilon)V - VL = \emptyset,$$

where V is normalized in some convenient way and where

$$(2.4.5) \quad \begin{cases} |A_N| = O(1/\varepsilon) \\ |L| = O(1/\varepsilon) \end{cases} .$$

Since, by Theorem [A1.2], these quantities depend smoothly on the coefficients of A_N , one equivalently can obtain V by solving the system

$$(2.4.6) \quad G(\varepsilon A_N(t, \varepsilon), \varepsilon L, V) = \emptyset,$$

where V is normalized as before. If the eigenvalue problem is well-conditioned, as it must be for the normal form, this system can be solved by a Newton iteration to an accuracy within roundoff error.

Suppose the grid values for A_N are given to $O(\rho_N/\varepsilon)$, where ρ_N is a measure of the relative error. By Corollary [A1.2a] the eigenvectors which correspond to the oscillatory set of the perturbed system are changed only by $O(\rho_N)$, and likewise by Theorem [2.2.6] and Theorem [2.3.2] the constructed transformations are changed only by $O(\rho_N)$. Thus, we can assume that, for a given vector v , $S_N v$ and $S'_N v$ can be calculated to within relative error $O(\rho_N)$. And also we can approximate S'_N by a linear operator:

$$(2.4.7) \quad L\{S_{N+O}(\rho_N), t, h\} = S'_N(t) + O(h^m) + O(\rho_N/h).$$

Then by (2.4.3) it is reasonable to assume that the error in A_{N+1} is

$$(2.4.8) \quad \rho_{N+1}/\varepsilon = O(\rho_N/\varepsilon + \rho_N/h + h^m).$$

If we take

$$(2.4.9) \quad \rho_0 = O(\delta),$$

where δ is a measure of the relative roundoff error, then for $N > 1$ we have by induction :

$$(2.4.10) \quad \rho_N/\varepsilon = O((1 + (\varepsilon/h)^N)(\delta/\varepsilon) + (1 + (\varepsilon/h)^{N-1})h^m)$$

as

$$(2.4.11) \quad \left\{ \begin{array}{l} \delta \rightarrow 0 \\ \varepsilon \rightarrow 0 \\ h \rightarrow 0 \end{array} \right. .$$

Thus, after these transformations have been carried out, we can assume that the system (2.4.1) is in normal form (p, q) with a grid error given by $O(\rho/\varepsilon)$. Setting $(\rho=0)$ then corresponds to the analysis without the consideration of roundoff error. We propose to approximate the solution on the grid by the following Linear Solver of the form $(p, q, m_1, m_2, m_3, m_4)$:

I

The $O(\varepsilon^q)$ terms of the system in normal form (p, q) are discarded

to give on the grid points a system of the form

$$(2.4.12) \quad \begin{cases} d/dt \begin{pmatrix} Y^I \\ Y^{II} \end{pmatrix} = \begin{bmatrix} \frac{1}{\varepsilon} B_{11}(t, \varepsilon) & \frac{1}{\varepsilon} B_{12}(t, \varepsilon) \\ 0 & B_{22}(t, \varepsilon) \end{bmatrix} \begin{pmatrix} Y^I \\ Y^{II} \end{pmatrix} \\ Y^I(0, \varepsilon) = Y^I_0 \\ Y^{II}(0, \varepsilon) = Y^{II}_0 \end{cases} \quad 0 < t < T$$

Y^I is an n_1 -dimensional vector; Y^{II} is an n_2 -dimensional vector; the coefficients are given by

$$(2.4.13) \quad \begin{cases} \frac{1}{\varepsilon} B_{11}(t, \varepsilon) = \frac{1}{\varepsilon} B_{11}(t, \varepsilon) + O(\rho/\varepsilon) \\ B_{22}(t, \varepsilon) = B_{22}(t, \varepsilon) + O(\rho/\varepsilon) \\ \frac{1}{\varepsilon} B_{12}(t, \varepsilon) = \frac{1}{\varepsilon} B_{12}(t, \varepsilon) + O(\rho/\varepsilon) \end{cases},$$

where the first right-hand-side terms are in $C^p(t, \varepsilon)$. By Lemma [2.1.3] this truncation gives rise to an error

$$(2.4.14) \quad \sigma_I = O(\varepsilon^2).$$

II

The system for Y^{II} is solved approximatively by some self-starting method which is accurate to within

$$(2.4.15) \quad \sigma_{II} = O(h^m) + O(\rho/(h\varepsilon))$$

with

$$(2.4.16) \quad m_1 < (p-1).$$

(see Fröberg [10], p. 260)

III

The system for Y^I is solved by approximating the terms of the

asymptotic expansion derived in Section 2.2. $O(\varepsilon^2)$ terms are ignored, and the solution operator (2.2.49) is estimated by

$$(2.4.17) \begin{cases} L1\{a+O(\rho/\varepsilon), t, \tau, \varepsilon, h\} = \hat{A}_\varepsilon(t, \tau) + O(h^{m_2}) + O(\rho/\varepsilon) \\ L2\{b/\varepsilon+O(\rho/\varepsilon), t, \tau, \varepsilon, h\} = \hat{B}_\varepsilon(t, \tau)/\varepsilon + O(h^{m_3}/\varepsilon) + O(\rho/\varepsilon) \end{cases}$$

where the discrete operators $L1$ and $L2$ approximate the integrals given in (2.2.49) with

$$(2.4.18) \quad \max \{m_2, m_3\} < p .$$

The t -derivatives in (2.2.51) and (2.2.52) are also estimated by discrete operators such as, for example,

$$(2.4.19) \begin{cases} L3\{g+O(\rho/\varepsilon h), t, h\} = (d/dt)[g(t, \varepsilon)] + O(h^{m_4}) + O(\rho/\varepsilon h^2) \\ L4\{g+O(\rho/\varepsilon h), t, h\} = (d^2/dt^2)[g(t, \varepsilon)] + O(h^{m_4-1}) + O(\rho/\varepsilon h^3) \\ L5\{g+O(\rho/\varepsilon h), t, h\} = (d^3/dt^3)[g(t, \varepsilon)] + O(h^{m_4-2}) + O(\rho/\varepsilon h^4) \end{cases}$$

with

$$(2.4.20) \quad m_4 < p - 1 .$$

(See Fröberg [10]: p. 190, p. 195) By the form of the expansions (2.2.51) and (2.2.52), these discretizations give rise to a cumulative error

$$(2.4.21) \begin{cases} \sigma_{III} = O(\varepsilon^2 + h^{m_2} + h^{m_3}/\varepsilon + \rho/\varepsilon h) + \\ O(h^{m_4}(1 + (\varepsilon/h)^{q-1}) + (\rho/\varepsilon h)(1 + (\varepsilon/h)^{q-1})) \end{cases}$$

This procedure gives an approximation of the form

$$(2.4.22) \quad Y(t, \varepsilon) = G(t, \varepsilon) + \sum_K C_K(t, \varepsilon) f_K(t, \varepsilon),$$

where $G(t, \varepsilon)$ and $\{C_k(t, \varepsilon)\}$ are approximations to smooth vectors and the $\{f_k(t, \varepsilon)\}$ are approximations to oscillatory functions of the forms

$$(2.4.23) \quad \sin(\widehat{B}_\varepsilon(t, \theta)/\varepsilon)$$

and

$$(2.4.24) \quad \cos(\widehat{B}_\varepsilon(t, \theta)/\varepsilon).$$

By our analysis of the algorithm we have:

THEOREM [2.4.1]

Let $X(t, \varepsilon)$ be the solution of the system (2.4.1), which is in normal form (p, q) , and let $Y(t, \varepsilon)$ be the approximation by a linear solver of the form $(p, q, m_1, m_2, m_3, m_4)$ on the grid (2.4.2), where the system coefficients are given to within $O(\rho/\varepsilon)$. We then have:

$$(2.4.25) \quad \max_t |X - Y| = O(\sigma_I + \sigma_{II} + \sigma_{III}),$$

where σ_I , σ_{II} , and σ_{III} are given respectively by (2.4.14), (2.4.15), and (2.4.21).

The composite error (2.4.25) gives much information as to the relative importance of the approximations which are made; in particular, the critical terms arise from roundoff error in the transformed system and the frequency approximations for the fast modes (2.4.17). A concern for the balance between the various contributions to the total error must influence any realistic

implementation of these methods.

2.5 A Computational Example

In this section we describe a simple implementation of the techniques of Section 2.4. We consider the following system, which describes the motion of two weakly coupled simple harmonic oscillators:

$$(2.5.1) \quad \left\{ \begin{array}{l} X' = \begin{bmatrix} 0 & k_1(t)/\varepsilon & a(t) & 0 \\ -k_1(t)/\varepsilon & 0 & 0 & a(t) \\ b(t) & 0 & 0 & k_2(t)/\varepsilon \\ 0 & b(t) & -k_2(t)/\varepsilon & 0 \end{bmatrix} X \\ X(0, \varepsilon) = [1, 0, 1, 0]^T \quad T_0 < t < T_1 \end{array} \right. ,$$

where

$$(2.5.2) \quad \left\{ \begin{array}{ll} k_1(t) = 1.5 + .9t - .2t^2 & k_2(t) = 2.7 + .5t \\ a(t) = \sin(t) & b(t) = 1 - .5t \\ T_0 = 0 & T_1 = 5 \\ \varepsilon = .01 & \end{array} \right. .$$

We propose to solve these equations to within an accuracy of three or four significant figures by using the step size

$$(2.5.3) \quad h = .05$$

If the calculations are done with standard single-precision accuracy, then in the notation of Section 2.4 we have

$$(2.5.4) \quad \begin{cases} \rho \sim 10^{-7} \\ \rho/\varepsilon \sim 10^{-5} \end{cases}$$

and so the effects of the roundoff error can be ignored. The system is in normal form (p,q) with

$$(2.5.5) \quad \begin{cases} p = +\infty \\ q = 0 \end{cases};$$

after transforming the equations to normal form $(p-2,q+2)$ we can continue as in the previous section by dropping the $O(\varepsilon^2)$ terms. This procedure is modeled by tableau in table <2.5.1>.

NORMAL FORM		GRID POINTS		
		T_0	...	T_1
(a)	(p,q)	x x x x x x x x	...	x x x x x x
(b)	$(p-1,q+1)$	x x x x x x	...	x x x x
(c)	$(p-2,q+2)$	x x x x x x	...	x x x x

table <2.5.1>

In <2.5.1a> we use grid values of the system in its original form; two extra points have been added at the end points of the interval so that centered differences can be used throughout to estimate the derivatives. At the grid points the canonical representation vectors are calculated by a Newton iteration applied to equation (2.4.2); we note that the computed values at

$t=t_N$ can be used as the initial guess for $t=t_{N+1}$ since, by Theorem [2.2.2], this estimate gives the correct value to within $O(h)$. At the grid point $t=t_N$ we calculate the canonical transformation

$$(2.5.6) \quad T_N = R_N Q_N$$

where, as in Theorem [2.2.5], R_N is upper triangular and Q_N is orthonormal. For this implementation, however, we express Q_N as a product of plane rotations rather than reflections. We note that for

$$(2.5.7) \quad Q = \begin{bmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{bmatrix}$$

we have

$$(2.5.8) \quad \left\{ \begin{aligned} Q'Q^* &= \begin{bmatrix} -\sin(\theta) & \cos(\theta) \\ -\cos(\theta) & -\sin(\theta) \end{bmatrix} \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix} \theta' \\ &= \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \theta' \end{aligned} \right.$$

To estimate the derivatives of the transformation we use a fourth order centered approximation [see Fröberg [10], p. 192]:

$$(2.5.9) \quad h L\{f, N\} = (2/3)(f_{N+1} - f_{N-1}) + (1/12)(f_{N-2} - f_{N+2}) \\ = h f'(t_N) + O(h^4)$$

Thus, the transformation can be carried out at the grid points to bring the system to normal form $(p-1, q+1)$:

$$(2.5.10) \quad \begin{cases} A_N \rightarrow \widetilde{A}_N = T_N A_N T_N^{-1} + T'_N T_N^{-1} \\ T'_N T_N^{-1} = O(\varepsilon) \end{cases} .$$

This procedure can be repeated at the grid points to transform the system to normal form (p-2,q+2):

$$(2.5.11) \quad \begin{cases} \widetilde{A}_N \rightarrow \widetilde{\widetilde{A}}_N = \widetilde{T}_N \widetilde{A}_N \widetilde{T}_N^{-1} + \widetilde{T}'_N \widetilde{T}_N^{-1} \\ \widetilde{T}'_N \widetilde{T}_N^{-1} = O(\varepsilon^2) \end{cases} .$$

Here we need not carry out the differentiation since the $O(\varepsilon^2)$ terms are to be discarded; thus, no extra points are needed at the ends of the interval. The initial guess for the Newton iteration now comes from the values of the system in normal form (p-1,q+1) at the same grid point since, by Theorem [2.2.2], the canonical representation vectors and the rescaled eigenvalues change only by $O(\varepsilon)$.

We note that there are no smooth components to be approximated as in (2.4.12). The necessary approximations for the integrals in (2.4.17) are made by a sixth-order Newton-Cotes formula [see Fröberg [10], p. 198]:

$$(2.5.12) \quad \begin{aligned} M\{f, N, h\} &= [7(f_N + f_{N+4}) + 32(f_{N+1} + f_{N+3}) + 12(f_{N+2})](h/90) \\ &= \int_{t_N}^{t_{N+4}} f \, dt + O(h^6) \end{aligned} .$$

Thus, the error in the integration of the frequencies is $O(h^6/\varepsilon)$. In the nomenclature of Section 4.4, we have outlined a linear solver of the form (p-2,q+2,2,..,6,6,..).

All calculations were done with single-precision accuracy on a VAX11/780 computer. In figure <2.5.1> we plot the leading order energy functions of the two oscillators:

$$(2.5.13) \quad \begin{cases} E_1(t) = (k_1(t)^2)/2 & (x_1^2 + x_2^2) \\ E_2(t) = (k_2(t)^2)/2 & (x_3^2 + x_4^2) \end{cases}$$

The amplitudes and phases have been linearly interpolated between grid points. And in table <2.5.1> we compare the computed grid values with the accepted function values, which were computed with double-precision accuracy by means of a fourth order Runge-Kutta scheme with a time step

$$(2.5.14) \quad h = 10^{-4} .$$

In the table we list

$$(2.5.15) \quad \begin{cases} \text{ERR1}(t) = \max_i \{|x_i(t) - \widetilde{x}_i(t)|\} \\ \text{ERR2}(t) = \max_i \{|(E_i(t) - \widetilde{E}_i(t))/\widetilde{E}_i(t)|\} \end{cases}$$

where $\{x_i\}$ and $\{E_i\}$ are the computed function values and where $\{\widetilde{x}_i\}$ and $\{\widetilde{E}_i\}$ are the accepted function values. We have

$$(2.5.16) \quad \max_i \max_t |x_i(t)| < 1.1 .$$

The analysis of this chapter has demonstrated that when the large frequencies of a linear system are well-separated the fast modes can be essentially decoupled; thus, in figure <2.5.1> we see that the energies behave rather independently. In Chapter

III we find that the introduction of nonlinearities considerably enriches the qualitative possibilities. For example, in Section 3.7 we study a coupled system where, although the leading order problem is oscillatory stable, the nonlinearities induce an internal resonance which results in a sequence of energy exchanges between the two oscillators. Our assumptions also preclude the possibility of passage through resonance since our approximation techniques break down when the large frequencies coalesce; under these circumstances the turning point arguments of Chapter IV must be utilized.

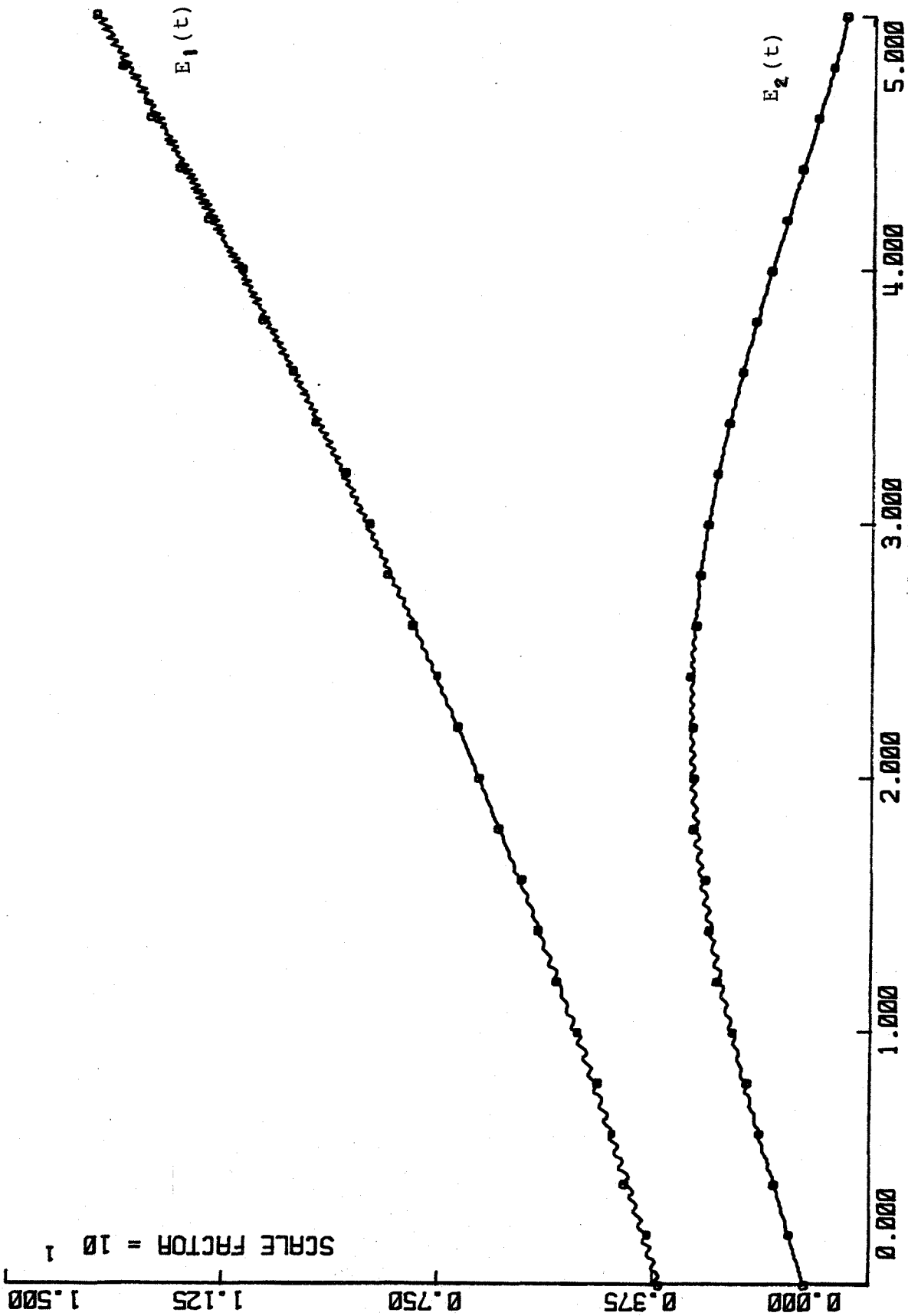


figure <2.5.1>

N	T_N	$E_1(T_N)$	$E_2(T_N)$	ERR1(T_N)	ERR2(T_N)
0	0.00	0.11250E+01	0.36450E+01	0.0E+00	0.0E+00
4	0.20	0.13929E+01	0.38615E+01	0.3E-04	0.1E-04
8	0.40	0.16821E+01	0.42653E+01	0.6E-04	0.2E-04
12	0.60	0.19306E+01	0.44863E+01	0.8E-04	0.5E-04
16	0.80	0.21575E+01	0.47505E+01	0.1E-03	0.6E-04
20	1.00	0.24134E+01	0.51131E+01	0.1E-03	0.7E-04
24	1.20	0.26682E+01	0.54840E+01	0.1E-03	0.7E-04
28	1.40	0.28324E+01	0.58010E+01	0.1E-03	0.8E-04
32	1.60	0.28919E+01	0.61043E+01	0.1E-03	0.9E-04
36	1.80	0.31073E+01	0.64937E+01	0.1E-03	0.1E-03
40	2.00	0.31090E+01	0.68469E+01	0.1E-03	0.2E-03
44	2.20	0.31315E+01	0.72284E+01	0.1E-03	0.2E-03
48	2.40	0.31700E+01	0.75886E+01	0.1E-03	0.2E-03
52	2.60	0.30731E+01	0.80339E+01	0.1E-03	0.2E-03
56	2.80	0.29934E+01	0.84461E+01	0.2E-03	0.2E-03
60	3.00	0.28831E+01	0.87763E+01	0.2E-03	0.2E-03
64	3.20	0.27164E+01	0.91938E+01	0.2E-03	0.2E-03
68	3.40	0.25307E+01	0.97375E+01	0.2E-03	0.2E-03
72	3.60	0.23057E+01	0.10128E+02	0.2E-03	0.2E-03
76	3.80	0.20723E+01	0.10650E+02	0.2E-03	0.2E-03
80	4.00	0.17970E+01	0.10992E+02	0.2E-03	0.2E-03
84	4.20	0.15418E+01	0.11598E+02	0.2E-03	0.2E-03
88	4.40	0.12674E+01	0.12095E+02	0.2E-03	0.3E-03
92	4.60	0.99624E+00	0.12595E+02	0.2E-03	0.2E-03
96	4.80	0.73748E+00	0.13094E+02	0.2E-03	0.2E-03
100	5.00	0.50035E+00	0.13550E+02	0.2E-03	0.2E-03

figure <2.5.1>

CHAPTER IIAPPENDIX ATHE SMOOTHNESS OF EIGENVALUES AND EIGENVECTORS

Although the study of eigenvalues and eigenvectors is usually developed within an algebraic framework, the tools of analysis are often more effective in the treatment of perturbations of the spectrum. Indeed, the smoothness results will follow from the validity of local expansions which are derived by the standard tools of analysis. The first basic result is the continuity of the eigenvalues.

THEOREM [A1.1]

Let $L(t, \mathcal{E})$ be an $(n \times n)$ matrix which is continuous on some domain. The eigenvalues of $L(t, \mathcal{E})$ can be represented by an unordered n -tuple

$$(A1.1) \quad G(t, \mathcal{E}) = \{g_1(t, \mathcal{E}), \dots, g_n(t, \mathcal{E})\},$$

where we define the distance between two such n -tuples as

$$(A1.2) \quad \text{dist}\{G^1, G^2\} = \min \max_i |g_i^1 - g_i^2|$$

where the min is taken over all possible orderings of the elements of the n -tuples. The eigenvalues are continuous at each point (t_0, \mathcal{E}_0) ; that is,

$$(A1.3) \quad \left\{ \begin{array}{l} (t, \mathcal{E}) \rightarrow (t_0, \mathcal{E}_0) \\ \implies \\ \text{dist}\{G(t, \mathcal{E}), G(t_0, \mathcal{E}_0)\} \rightarrow 0 \end{array} \right.$$

PROOF

This fundamental result is most easily interpreted algebraically as it is really a theorem about the roots of the polynomial

$$(A1.4) \quad |L(t, \varepsilon) - \lambda I| = 0 .$$

See, for example, Franklin [9].

The existence of smooth functions which represent the eigenvalues is quite another matter; however, in certain circumstances such parameterizations can be demonstrated.

THEOREM [A1.2]

Consider the $(n \times n)$ matrix $L(t, \varepsilon)$ where:

- (i) $L(t, \varepsilon) = A_1(t) + A_2(t, \varepsilon);$
- (ii) $A_1(t)$ and $A_2(t, \varepsilon)$ are $(n \times n)$ matrices which are in $C^p(t, \varepsilon)$ for $0 < t < T$ and $0 < \varepsilon < 1;$
- (iii) The eigenvalues of $A_1(t)$ can be divided into two groups uniformly on $0 < t < T :$

Oscillatory Set:

n_1 imaginary eigenvalues $\{\lambda_k(t)\}$ with

$$\min_k |\lambda_k(t)| > r$$

and

$$\min_{k \neq j} |\lambda_k(t) - \lambda_j(t)| > r ,$$

where r is some positive real constant.

Nonoscillatory Set:

$n_2 (=n-n_1)$ identically zero eigenvalues.

Then there exist, correspondingly to the oscillatory set,

- (i) n_1 scalar functions $\{\lambda_k(t, \varepsilon)\} (\subset C^p(t, \varepsilon))$ which represent eigenvalues of $L(t, \varepsilon)$;
- (ii) n_1 matrix functions $\{P_k(t, \varepsilon)\} (\subset C^p(t, \varepsilon))$ which represent one-dimensional eigenprojections of $L(t, \varepsilon)$; and
- (iii) n_1 vector functions $\{v_k(t, \varepsilon)\} (\subset C^p(t, \varepsilon))$ which correspond to normalized eigenvectors of $L(t, \varepsilon)$.

Moreover, these functions can be represented by local expansions which depend smoothly on the coefficients of $L(t, \varepsilon)$.

PROOF

First we argue that the eigenvalues of the oscillatory set can be represented by n_1 continuous functions

$$(A1.5) \quad \{\lambda_k(t, \varepsilon)\}.$$

To achieve a unique parameterization one can number the eigenvalues of the oscillatory set in ascending order with respect to their imaginary parts. By assumption (iii) the ordering of these eigenvalues is preserved while the eigenvalues of the nonoscillatory set have imaginary parts which are $o(1)$. The continuity of these functions follows from Theorem [A1.1].

The resolvent matrix of L , defined by

$$(A1.6) \quad R(L, \lambda) = [L - \lambda I]^{-1},$$

is well-defined for $\gamma \in P(L)$, the resolvent set of L . A resolvent matrix theory for linear operators on finite-dimensional spaces has been developed, especially for the case where L is an analytic function of a single variable (see, for example, Kato [15]). Here we outline a simplified theory for the case where L is a smooth function of its arguments.

For the case where L is constant one has the local representation formulas for the oscillatory set:

$$(A1.7) \quad P_k = (-1/2\pi i) \int_{\Gamma_k} R(L, \gamma) d\gamma$$

and

$$(A1.8) \quad \lambda_k = \text{trace}\{LP_k\} ,$$

where Γ_k is a simple closed curve enclosing λ_k but not enclosing or passing through λ_j for $j \neq k$.

We now summarize some of the basic results of the resolvent theory (for proofs see Kato [15]). P_k is said to be the projection for the eigenvalue λ_k ; that is, P_k is an $(n \times n)$ matrix which projects elements of the underlying linear space onto the one-dimensional eigenspace associated with the eigenvalue which is enclosed by Γ_k . (A1.8) gives an analytic representation for the eigenvalue enclosed by Γ_k . The most fundamental properties of the projections are

$$(A1.9) \quad P_k P_j = \delta_{kj} P_k$$

and

$$(A1.10) \quad LP_k = P_k L = P_k LP_k .$$

Finally with suitable choices for w and m we can represent the k -th normalized eigenvector by

$$(A1.11) \quad \begin{cases} v_k = P_k w / \| (P_k w) \| \\ (P_k w)^{(m)} > 0 \end{cases}$$

Thus, for each value of its arguments $L(t, \varepsilon)$ has, correspondingly to the oscillatory set, a set of n_1 triples, each comprised of eigenvalue, eigenprojection, and eigenvector; moreover, by (A1.5) these triples can be parameterized so that the eigenvalues are represented by n_1 continuous functions of the independent variables. The smoothness of these triples can be analyzed through (A1.7), where L is now considered as a smooth function.

Let Γ_k be a simple closed curve which encloses $\lambda_k(t_0, 0)$ but no other eigenvalue of L_0 . Near $(t, \varepsilon) = (t_0, 0)$ we have:

$$(A1.12) \quad L(t, \varepsilon) = L_0 + \tilde{L}(t, \varepsilon) = L_0 + O(|t - t_0| + \varepsilon).$$

Now we construct the so-called second Neumann series for the resolvent:

$$(A1.13) \quad \begin{aligned} R(L(t, \varepsilon), \gamma) &= [L(t, \varepsilon) - \gamma I]^{-1} \\ &= R(\gamma) [I + \tilde{L}(t, \varepsilon) R(\gamma)]^{-1} \quad \text{where } R(\gamma) = [L_0 - \gamma I]^{-1} \\ &= R(\gamma) \sum_{j=0}^{\infty} [-\tilde{L}(t, \varepsilon) R(\gamma)]^j \quad \text{for } |\tilde{L}(t, \varepsilon)| < \max_{\gamma \in \Gamma_k} |R(\gamma)|^{-1}. \end{aligned}$$

For sufficiently small $[|t - t_0| + \varepsilon]$ we can expand the resolvent in some neighborhood of $(t, \varepsilon) = (t_0, 0)$, and then by (A1.7) and

(A1.8) we have:

$$(A1.14) \quad P_K(t, \varepsilon) = \sum_i \sum_j P_{Kij} (t-t_0)^i (\varepsilon)^j + o((t-t_0)^P) + o(\varepsilon^P)$$

and

$$(A1.15) \quad \lambda_K(t, \varepsilon) = \sum_i \sum_j \lambda_{Kij} (t-t_0)^i (\varepsilon)^j + o((t-t_0)^P) + o(\varepsilon^P).$$

And likewise we have a similar local expansion for the normalized eigenvectors. Thus, the functional representations for the eigenvalues, eigenprojections, and eigenvectors are all in $C^P(t, \varepsilon)$.

COROLLARY [A1.2a]

Let $K(t, \varepsilon, \delta)$ be the $(n \times n)$ matrix defined by

$$(A1.16) \quad K(t, \varepsilon, \delta) = L(t, \varepsilon) + \delta F(t, \varepsilon),$$

where δ is a small scalar parameter, $F(t, \varepsilon)$ is an $(n \times n)$ matrix which is in $C^r(t, \varepsilon)$ ($r \leq p$), and $L(t, \varepsilon)$ is as given in the previous theorem. For sufficiently small δ there exist, correspondingly to the oscillatory set of $L(t, \varepsilon)$, the following n_l eigenvalues, eigenprojections, and eigenvectors of $K(t, \varepsilon, \delta)$:

$$(A1.17) \quad \begin{cases} \tilde{\lambda}_K(t, \varepsilon, \delta) = \lambda_K(t, \varepsilon) + \delta f_1^K(t, \varepsilon, \delta) \\ \tilde{P}_K(t, \varepsilon, \delta) = P_K(t, \varepsilon) + \delta f_2^K(t, \varepsilon, \delta) \\ \tilde{v}_K(t, \varepsilon, \delta) = v_K(t, \varepsilon) + \delta f_3^K(t, \varepsilon, \delta), \end{cases}$$

where $\{P_K\}$, $\{\lambda_K\}$, and $\{v_K\}$ are as in the previous theorem and where each of $\{f_i^K\}$ is in $C^r(t, \varepsilon, \delta)$.

PROOF

The results follow from the substitution of $K(t, \epsilon, \delta)$ for $L(t, \epsilon)$ in (A1.7), (A1.8), and (A1.11).

CHAPTER IIITHE NONLINEAR PROBLEM WITHOUT TURNING POINTS3.1 REDUCTION TO THE NON-STIFF FORMULATION

In this chapter we consider the system

$$(3.1.1) \quad \begin{cases} Z' = A(t)/\varepsilon Z + H(Z, t) \\ Z(0, \varepsilon) = \hat{Z}_0; \quad 0 < \varepsilon \ll 1; \quad 0 < t < T \end{cases}$$

where:

- (i) \hat{Z}_0 is independent of ε ;
- (ii) $H(Z, t)$ has components which are polynomial in the components of Z with t -dependent coefficients in $C^p(t)$ ($p \geq 0$, p continuous derivatives);
- (iii) $A(t)$ is in diagonal form:

$$(3.1.2) \quad \begin{cases} A(t) = \text{diag}(\lambda_k(t)) \quad (\text{Re}\{\lambda_k(t)\} \leq 0) \\ \Delta(t)^{(k)} = \lambda_k(t) \in C^p(t) \end{cases}$$

Here $\Delta(t)^{(k)}$ is the k -th component of the vector $\Delta(t)$.

If $A(t)$ is not in diagonal form, then, provided the reduced system

$$(3.1.3) \quad \begin{cases} Z' = A(t)/\varepsilon Z \\ Z(0) = \hat{Z}_0 \end{cases}$$

is oscillatory stable (Section 2.1), one can construct an associated canonical transformation $T(t)$ which diagonalizes $A(t)$

(Section 2.3). The system for

$$(3.1.4) \quad \widetilde{Z} = T(t)Z$$

is then in the form (3.1.1) with p replaced by $(p-1)$. By the constructive argument of Theorem [2.2.2] this procedure will work even in the case of repeated eigenvalues if there exists a corresponding set of smooth eigenvectors. The system now can be transformed to a formulation in which the coefficients are bounded but some are rapidly oscillating.

(example 3.1a)

The following system describes the motion of an unforced oscillator with cubic damping:

$$(3.1.5) \quad \begin{cases} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}' = \begin{bmatrix} 0 & 1/\varepsilon \\ -1/\varepsilon & 0 \end{bmatrix} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} + \begin{pmatrix} 0 \\ -(z_1)^3 \end{pmatrix} \\ \begin{pmatrix} z_1(0, \varepsilon) \\ z_2(0, \varepsilon) \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad 0 < t < T, \quad 0 < \varepsilon \ll 1 \end{cases}$$

After the change of variables

$$(3.1.6) \quad \begin{cases} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = S \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = 1/2 \begin{bmatrix} 1 & 1 \\ -i & i \end{bmatrix} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = \begin{pmatrix} (u_1 + u_2)/2 \\ i(u_1 - u_2)/2 \end{pmatrix} \\ \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = S^{-1} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = \begin{bmatrix} 1 & i \\ 1 & -i \end{bmatrix} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = \begin{pmatrix} z_1 + iz_2 \\ z_1 - iz_2 \end{pmatrix} \end{cases}$$

the equations become

$$(3.1.7) \quad \begin{cases} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}' = \begin{bmatrix} -i/\varepsilon & 0 \\ 0 & i/\varepsilon \end{bmatrix} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} + \begin{pmatrix} F(u_1, u_2) \\ -F(u_1, u_2) \end{pmatrix} \\ \begin{pmatrix} u_1(0, \varepsilon) \\ u_2(0, \varepsilon) \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad 0 < t < T \end{cases}$$

where

$$(3.1.8) \quad \begin{cases} F(u_1, u_2) = 1/8 u_1^3 - 3/8 u_1^2 u_2 + 3/8 u_1 u_2^2 - 1/8 u_2^3 \\ u_2 = \bar{u}_1 \end{cases}$$

The system (3.1.7) has the form (3.1.1). Next we make the change of variables

$$(3.1.9) \quad \begin{cases} u_1 = \exp(-it/\varepsilon) x_1 \\ u_2 = \exp(it/\varepsilon) x_2 \end{cases}$$

and obtain

$$(3.1.10) \quad \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}' = \begin{pmatrix} \exp(it/\varepsilon) F(\exp(-it/\varepsilon)x_1, \exp(it/\varepsilon)x_2) \\ -\exp(-it/\varepsilon) F(\exp(-it/\varepsilon)x_1, \exp(it/\varepsilon)x_2) \end{pmatrix}$$

by which we have:

$$(3.1.11) \quad \begin{cases} x_1' = -3/8 x_1^2 x_2 + \\ \quad 1/8 \exp(-2it/\varepsilon) x_1^3 + 3/8 \exp(2it/\varepsilon) x_1 x_2^2 + \\ \quad -1/8 \exp(4it/\varepsilon) x_2^3 \\ x_1(0, \varepsilon) = 1, \quad x_2 = \bar{x}_1, \quad 0 < t < T \end{cases}$$

For the general case, let $S_\varepsilon(t,s)$ be the solution operator of the reduced system (3.1.3) where $A(t)$ satisfies (3.1.2); in fact, we can write:

$$(3.1.12) \quad S_\varepsilon(t,s) = \text{diag}\left\{\exp\left(\int_s^t \lambda_k(\tau) d\tau\right)/\varepsilon\right\}.$$

Thus, with the change of variables

$$(3.1.13) \quad X = S_\varepsilon(t,0)Z$$

we reach a system in which the coefficients are bounded but some are rapidly oscillating:

$$(3.1.14) \quad \begin{cases} X' = G(X,t,\varepsilon) \\ \quad = g_{\text{I}}(X,t,\varepsilon) + g_{\text{II}}(X,t) + f_{\text{I}}(t,\varepsilon) + f_{\text{II}}(t) \\ X(0,\varepsilon) = \hat{X}_0; \quad 0 < \varepsilon \ll 1; \quad 0 < t < T \end{cases}$$

Here the forms of the coefficients are given by:

- (i) $\hat{X}_0 = \hat{Z}_0$ is independent of ε ;
- (ii) $g_{\text{I}}(X,t,\varepsilon)^{(i)} = \sum_j a_{ij}(t) \exp[B_{ij}(t)/\varepsilon] p_{ij}(X)$;
- (iii) $g_{\text{II}}(X,t)^{(i)} = \sum_j d_{ij}(t) q_{ij}(X)$;
- (iv) $f_{\text{I}}(t,\varepsilon)^{(i)} = \sum_j c_{ij}(t) \exp[G_{ij}(t)/\varepsilon]$;
- (v) $f_{\text{II}}(t)^{(i)} = h_i(t)$;
- (vi) $\{a_{ij}(t), d_{ij}(t), c_{ij}(t), h_i(t)\} \subset C^p(t)$;
 $p_{ij}(X)$ and $q_{ij}(X)$ are monomials of positive degree in the components of X ; $B_{ij}(t)/\varepsilon$ and $G_{ij}(t)/\varepsilon$ can be represented by n -vector scalar products of the form

$$(3.1.15) \quad N^T P(t)/\varepsilon$$

where N is a constant n -vector with integral components and $P(t)$ is given by

$$(3.1.16) \quad P(t)^{(i)} = \int_0^t \lambda_i(s) ds.$$

By (3.1.2) and (3.1.16) the t -derivative of the expression (3.1.15) is

$$(3.1.17) \quad (N^T P(t)/\epsilon)' = N^T \Lambda(t)/\epsilon.$$

The entries of $\Lambda(t)/\epsilon$ are called fundamental frequencies while the relevant terms of the form (3.1.17) are called secondary frequencies. We also impose the following restriction on all relevant secondary frequencies:

$$(3.1.18) \quad |N^T \Lambda(t)| \geq K > 0,$$

where K is some positive constant. (K/ϵ) is then a measure of the stiffness of the system.

This assumption arises out of the necessity for some concrete specification as to the meaning of "fast oscillations"; moreover, this restriction must be maintained in subsequent levels of analysis. Functions which have the form given by (iv) and which satisfy (3.1.18) are called strictly oscillatory (of class p); functions of the form given by (v) are called strictly nonoscillatory (of class p). In this chapter the subscript I designates a strictly oscillatory function, and the subscript II designates a strictly nonoscillatory function. The system (3.1.14) is said to be in non-stiff oscillatory form.

If for some relevant

$$(3.1.19) \quad \sigma(t) = N^T P(t)$$

we have

$$(3.1.20) \quad \sigma'(t) = N^T \Lambda(t) \equiv 0,$$

then the corresponding term can be reclassified as nonoscillatory; however, if $\sigma'(t)$ vanishes at some isolated point, then the approximation must be made as a turning point calculation. Essentially the underlying structure of some terms is changing from oscillatory to nonoscillatory to oscillatory, and correspondingly the balancing of the terms also must change. We develop this procedure in chapter IV. If, as in example (3.3a), the entries of $\Lambda(t)$ are integral constants, this difficulty cannot occur since all possible secondary frequencies must satisfy (3.1.18) or (3.1.20) with

$$(3.1.21) \quad K = 1 \quad .$$

The strength of (3.1.18) also allows us to define the leading order antiderivative of any oscillatory function through the linear operator

$$(3.1.22) \quad \mathcal{L}\{c(t) \exp[B(t)/\epsilon]\} = [c(t)/B'(t)] \exp[B(t)/\epsilon]$$

since

$$(3.1.23) \quad \epsilon \mathcal{L}\{c(t) \exp[B(t)/\epsilon]\}' = c(t) \exp[B(t)/\epsilon] + O(\epsilon).$$

Our aim is to characterize the solution of the system (3.1.14) in terms of these concepts. Thus, the function $f(t, \epsilon)$ is said to be decomposable if

$$(3.1.24) \quad f(t, \varepsilon) = \sum_K \varepsilon^K (W_K(t) + Y_K(t, \varepsilon)),$$

where each $Y_K(t, \varepsilon)$ is strictly oscillatory and each $W_K(t)$ is strictly nonoscillatory. And likewise $f(t, \varepsilon)$ is said to be decomposable to $O(\varepsilon^m)$ if

$$(3.1.25) \quad f(t, \varepsilon) = \sum_K \varepsilon^K (W_K(t) + Y_K(t, \varepsilon)) + O(\varepsilon^{m+1}),$$

where each $Y_K(t, \varepsilon)$ is strictly oscillatory and each $W_K(t)$ is strictly nonoscillatory. The characterization of the solution of (3.1.14) in terms of such an asymptotic expansion stands as the major goal of this chapter; the breakdown of this decomposability principle will correspond to a violation of (3.1.18), whereupon turning point techniques must be used.

3.2 Hierarchy for the Linear Problem

Since our treatment is based on a functional Newton iteration, We first discuss the linear problem

$$(3.2.1) \quad \begin{cases} X' = A_{\text{I}}(t, \varepsilon)X + A_{\text{II}}(t)X + f_{\text{I}}(t, \varepsilon) + f_{\text{II}}(t) \\ X(0, \varepsilon) = \hat{X}_0; \quad 0 < t < T, \quad 0 < \varepsilon \ll 1 \end{cases}$$

where the subscript I denotes a strictly oscillatory function and the subscript II denotes a strictly nonoscillatory function. We begin with a rather standard result on the stability of ordinary differential equations.

LEMMA [3.2.1]

The system

$$(3.2.2) \quad \begin{cases} Y' = B(t)Y + h(t) \\ Y(0) = Y_0, \quad t > 0 \end{cases}$$

where

$$(3.2.3) \quad \begin{cases} B(t) \in C^0(t) \\ h(t) \in C^0(t) \\ \max_t |B(t)| = K \end{cases}$$

has a bounded solution operator $S(t, s)$, by which we have:

$$(3.2.4) \quad \begin{cases} |S(t, s)| < \exp[K(t-s)] \quad (t > s) \\ Y(t) = S(t, 0)Y_0 + \int_0^t S(t, s)h(s)ds \end{cases}$$

PROOF

We derive a bound for the solution operator from the homogeneous version of (3.2.2). Let

$$\begin{cases} w(t) = |Y(t)| \\ w = |Y_0| \end{cases} ;$$

then

$$Y(t) = Y_0 + \int_0^t B(s)Y(s) ds$$

implies

$$w(t) < w_0 + K \int_0^t w(s) ds,$$

and by a fundamental lemma of stability theory (Coddington and Levinson [6], page 37) we have

$$w(t) \leq w_0 \exp[Kt] \quad (t > 0).$$

This guarantees the bound for the solution operator. The solution of the inhomogeneous equation is given by Duhammel's Principle.

LEMMA [3.2.2]

Consider the two systems

$$(3.2.5) \quad \begin{cases} Y' = B(t)Y + h(t, \varepsilon) \\ Y(0, \varepsilon) = Y_0, \quad 0 < t < T \end{cases}$$

and

$$(3.2.6) \quad \begin{cases} Z' = B(t)Z + h(t, \varepsilon) + \varepsilon \tilde{h}(t, \varepsilon) \\ Z(0, \varepsilon) = Y_0, \quad 0 < t < T \end{cases}$$

where each vector or matrix is a bounded continuous function of t , and ε is a sufficiently small positive real number. Then we have:

$$(3.2.7) \quad \begin{cases} \max_t |Y(t, \varepsilon) - Z(t, \varepsilon)| = O(\varepsilon) \\ \max_t |Y'(t, \varepsilon) - Z'(t, \varepsilon)| = O(\varepsilon) \end{cases}$$

PROOF

By the previous lemma both systems have bounded solution operators. The system for

$$R = Y - Z$$

then has the form given by (3.2.2) with

$$\begin{cases} R' = B(t)R - \varepsilon \tilde{h}(t, \varepsilon) \\ R(0, \varepsilon) = 0 \end{cases}$$

and therefore we have (3.2.7).

Thus, to achieve leading-order accuracy one simply ignores certain terms of the system. This principle leads to the following useful result concerning the system (3.2.1).

THEOREM [3.2.1]

Let $X(t, \varepsilon)$ be the solution of (3.2.1) and let $V(t)$ be the solution of the system

$$(3.2.8) \quad \begin{cases} V' = A_{\text{II}}(t)V + f_{\text{II}}(t) \\ V(0) = \hat{X}_0 \end{cases}$$

then

$$(3.2.9) \quad \begin{cases} \max_t |X(t, \varepsilon) - V(t)| = O(\varepsilon) \\ \max_t |X'(t, \varepsilon) - V'(t) - \varepsilon F'(t, \varepsilon)| = O(\varepsilon) \\ F(t, \varepsilon) = \mathcal{L}\{f_{\text{I}} + A_{\text{I}}V\} \end{cases}$$

PROOF

Let

$$Z(t, \varepsilon) = X(t, \varepsilon) - V(t).$$

Then the equations for $Z(t, \varepsilon)$ are

$$\begin{cases} Z' = (A_{\text{I}}(t, \varepsilon) + A_{\text{II}}(t))Z + (f_{\text{I}}(t, \varepsilon) + A_{\text{I}}(t, \varepsilon)V(t)) \\ Z(0, \varepsilon) = 0 \end{cases}$$

Since $V(t)$ is strictly nonoscillatory both forcing terms are strictly oscillatory. Then by (3.1.23) the equations for

$$\tilde{Z} = Z - \varepsilon(F(t, \varepsilon) - F(0, \varepsilon))$$

have the form

$$\begin{cases} \tilde{z}' = (A_{\underline{I}}(t, \varepsilon) + A_{\underline{II}}(t))\tilde{z} + \varepsilon h(t, \varepsilon) \\ \tilde{z}(0, \varepsilon) = 0 \end{cases}$$

and so by Lemma [3.2.2]:

$$\begin{cases} \max_t |z(t, \varepsilon)| = o(\varepsilon) \\ \max_t |z'(t, \varepsilon) - \varepsilon F'(t, \varepsilon)| = o(\varepsilon) \end{cases}$$

3.3 SOLUTION BY SUCCESSIVE LINEARIZATIONS

By using the results of the previous section, we now generate an asymptotic expansion for the solution of the system (3.1.14), which is in non-stiff oscillatory form. If a linearization technique is to be successful one must have a suitable value for the initial approximation.

ASSUMPTION [3.3a]

In correspondence to the system (3.1.14), the reduced system

$$(3.3.1) \quad \begin{cases} v' = g_{\mathbf{II}}(v, t) + f_{\mathbf{II}}(t) \\ v(0) = \hat{x}_0, \quad 0 < t < T \end{cases}$$

is well-posed and has a bounded solution in $C^{p+1}(t)$.

Given this assumption, we can define the $(n \times n)$ matrices

$$(3.3.2) \quad \begin{cases} A_{\mathbf{I}}(t, \varepsilon) = g_{\mathbf{I}}(X, t, \varepsilon)_{\mathbf{X}} \big|_{X=V} = g_{\mathbf{I}}(V, t, \varepsilon)_{\mathbf{X}} \\ A_{\mathbf{II}}(t) = g_{\mathbf{II}}(X, t)_{\mathbf{X}} \big|_{X=V} = g_{\mathbf{II}}(V, t)_{\mathbf{X}} \end{cases}$$

where $A_{\mathbf{I}}(t, \varepsilon)$ is strictly oscillatory of class p and $A_{\mathbf{II}}(t)$ is strictly nonoscillatory of class p . Here the notation $g(X)_{\mathbf{X}}$ indicates the Jacobian of the vector function, and $g(X)^{(i)}$ indicates the i -th component. We also define the operators:

$$(3.3.3) \begin{cases} G(X, t, \varepsilon) \equiv g_{\text{I}}(X, t, \varepsilon) + g_{\text{II}}(X, t) + f_{\text{I}}(t, \varepsilon) + f_{\text{II}}(t) \\ M(X) \equiv G(X, t, \varepsilon) - X' \end{cases}$$

For a positive integer m , \tilde{X}^m is said to be an ε^m -approximate solution of (3.1.14) if

$$(3.3.4) \begin{cases} \tilde{X}^m(t, \varepsilon) = \sum_{k=0}^m \varepsilon^k (W_k(t) + Y_k(t, \varepsilon)) + \varepsilon^{m+1} Y_{m+1}(t, \varepsilon) \\ \tilde{X}^m(0, \varepsilon) = \tilde{X}_0 + O(\varepsilon^{m+1}) \\ M(\tilde{X}^m) = O(\varepsilon^{m+1}) \end{cases}$$

where each Y_k is strictly oscillatory and each W_k is strictly nonoscillatory. By means of Assumption (3.3a) we can immediately demonstrate the existence of such an approximate solution.

THEOREM [3.3.1] The function \tilde{X}^0 , which is given by

$$(3.3.5) \begin{cases} \tilde{X}^0(t, \varepsilon) = W_0(t) + \varepsilon Y_1(t, \varepsilon) \\ W_0(t) = V(t) \\ Y_1(t, \varepsilon) = \mathcal{L}\{g_{\text{I}}(V, t, \varepsilon) + f_{\text{I}}(t, \varepsilon)\} \end{cases}$$

is an ε^0 -approximate solution of the system (3.1.14). $W_0(t)$ is strictly nonoscillatory of class $(p+1)$, and $Y_1(t, \varepsilon)$ is strictly oscillatory of class p .

PROOF

To verify (3.3.4) we consider

$$M(\tilde{X}^0) = G(\tilde{X}^0, t, \varepsilon) - (\tilde{X}^0)'$$

By (3.3.1) and (3.1.23) we have:

$$\begin{aligned}(\tilde{X}^0)' &= V' + g_{\mathbf{I}}(V, t, \varepsilon) + f_{\mathbf{I}}(t, \varepsilon) + O(\varepsilon) \\ &= G(V, t, \varepsilon) + O(\varepsilon)\end{aligned}$$

and also a simple Taylor expansion gives

$$G(\tilde{X}^0, t, \varepsilon) = G(V, t, \varepsilon) + O(\varepsilon).$$

Therefore, \tilde{X}^0 is an ε^0 -approximate solution of the system (3.1.14).

We now demonstrate that an ε^m -approximate solution actually approximates the exact solution of the system.

THEOREM [3.3.2]

If \tilde{X}^m is an ε^m -approximate solution of the system (3.1.14), then

$$(3.3.6) \quad \max_t |\tilde{X}^m(t, \varepsilon) - X(t, \varepsilon)| = O(\varepsilon^{m+1}),$$

where $X(t, \varepsilon)$ is the solution of (3.1.14).

PROOF

Consider the $(n+1)$ -dimensional space

$$\mathcal{D} = \{(x, t) \ni |x - \tilde{X}^m| < \delta, 0 < t < T\},$$

where δ is some arbitrary positive constant. By our assumptions on \tilde{X}^m and the system (3.1.14), we conclude that for some positive constants K_1 and K_2 we have:

- (i) $G(x, t, \varepsilon)$ is continuous in \mathcal{D} ;
- (ii) $|G(x, t, \varepsilon)| < K_1$ in \mathcal{D} ;
- (iii) $|G(x_1, t, \varepsilon) - G(x_2, t, \varepsilon)| < K_2|x_1 - x_2|$ in \mathcal{D} .

Note that the Lipschitz condition (iii) is guaranteed even though

the t -derivatives of $G(x, t, \varepsilon)$ are unbounded as $\varepsilon \rightarrow 0$. For example, we have:

$$|\exp(it/\varepsilon)x_1^2 - \exp(it/\varepsilon)x_2^2| < 2\delta |x_1 - x_2|.$$

We now introduce a sequence of Picard iterates:

$$\begin{cases} x_0 = \tilde{X}^m \\ x_{N+1} = \hat{X}_0 + \int_0^t G(x_N, t, \varepsilon) dt. \end{cases}$$

Since \tilde{X}^m satisfies the equation to within $O(\varepsilon^{m+1})$ we have

$$\begin{aligned} |x_1 - x_0| &= |\tilde{X}^m - \hat{X}_0 - \int_0^t G(\tilde{X}^m, t, \varepsilon) dt| \\ &< R \varepsilon^{m+1} \end{aligned}$$

where R is a positive constant. Provided the successive iterates are all in \mathcal{D} , we have by induction

$$|x_{N+1} - x_N| < \varepsilon^{m+1} R (K_2 t)^N / N! ,$$

and thus for all positive N

$$\max_t |x_N - x_0| < \varepsilon^{m+1} R \exp(K_2 T).$$

Therefore, for sufficiently small ε all iterates remain in \mathcal{D} . By the uniform convergence of the iteration, we have the existence of a unique continuously differentiable function X which satisfies

$$\begin{cases} M(X) = 0 \\ X(0, \varepsilon) = \hat{X}_0 \\ \max_t |\tilde{X}^m - X| = O(\varepsilon^{m+1}). \end{cases}$$

COROLLARY [3.3.2]

The system (3.1.14) with assumption (3.3a) is well-posed with

$$(3.3.7) \quad \max_t |V(t) - X(t, \varepsilon)| = O(\varepsilon),$$

where $X(t, \varepsilon)$ is the solution of (3.1.14).

PROOF

Since an ε^0 -approximate solution is given by Theorem [3.3.1], the error estimate follows from Theorem [3.3.2].

By using Theorem (3.3.2) we now extend the result of Corollary [3.3.2] to obtain higher order approximations.

THEOREM [3.3.3]

Consider the system (3.1.14). Let \tilde{X}^m be an ε^m -approximate solution where, for ($k > 1$), W_k is strictly nonoscillatory of class $-(p+2-k)$ and Y_k is strictly oscillatory of class $(p+1-k)$. Let $M(\tilde{X}^m)$ be decomposable to $O(\varepsilon^{m+1})$ with the form

$$(3.3.8) \quad M(\tilde{X}^m) = \varepsilon^{m+1} f_{\text{I}} + \varepsilon^{m+1} f_{\text{II}} + O(\varepsilon^{m+2}),$$

where f_{I} is strictly oscillatory of class $(p-m-1)$ and f_{II} is strictly nonoscillatory of class $(p-m)$. Then an ε^{m+1} -approximate solution of the system is given by

$$(3.3.9) \quad \tilde{X}^{m+1} = \tilde{X}^m + \varepsilon^{m+1} W_{m+1} + \varepsilon^{m+2} Y_{m+1},$$

where

$$(3.3.10) \quad \begin{cases} W_{m+1}' = A_{II}(t)W_{m+1} + f_{II}(t) \\ W_{m+1}(\theta) = -Y_{m+1}(\theta, \varepsilon) \\ Y_{m+2}(t, \varepsilon) = \mathcal{L}\{f_I + A_I W_{m+1}\} \end{cases}$$

Here W_{m+1} is strictly nonoscillatory of class $(p-m-1)$, and Y_{m+2} is strictly oscillatory of class $(p-m-1)$. Thus, we have:

$$(3.3.11) \quad \begin{cases} |X - \tilde{X}^{m+1}| = O(\varepsilon^{m+2}) \\ \tilde{X}^{m+1} = \sum_{k=0}^{m+1} \varepsilon^k (W_k + Y_k) + \varepsilon^{m+2} Y_{m+2} \end{cases}$$

Moreover, $M(\tilde{X}^{m+1})$ has the form (3.3.8) with m replaced by $(m+1)$ if

(1) $M(\tilde{X}^m)$ is decomposable to $O(\varepsilon^{m+2})$;

(2) $A_I(t, \varepsilon)Y_{m+2}$ is decomposable;

and

(3) $[(G(X, t, \varepsilon)_X)(W_{m+1})]_X (W_I + Y_I) \Big|_{X=V}$ is decomposable.

PROOF

By Theorem [3.3.2] we have

$$\begin{cases} X = \tilde{X}^m + \varepsilon^{m+1} Z \\ Z(\theta, \varepsilon) = -Y_{m+1}(\theta, \varepsilon) \end{cases},$$

where $Z(t, \varepsilon)$ is a continuously differentiable function of t , and X is the solution of (3.1.14). The differential equation then can be written as

$$\begin{aligned} (\tilde{X}^m + \varepsilon^{m+1} Z)' &= G(\tilde{X}^m + \varepsilon^{m+1} Z) \\ &= G(\tilde{X}^m + \varepsilon^{m+1} Z) - G(\tilde{X}^m) + G(\tilde{X}^m). \end{aligned}$$

We now carry out a linearization which is equivalent to a functional Newton iteration:

$$\begin{cases} Z' = [G(\tilde{X}^m + \varepsilon^{m+1} Z) - G(\tilde{X}^m)]/\varepsilon^{m+1} + M(\tilde{X}^m)/\varepsilon^{m+1} \\ \quad = [A_{\mathbf{I}}(t, \varepsilon) + A_{\mathbf{II}}(t)]Z + f_{\mathbf{I}}(t, \varepsilon) + f_{\mathbf{II}}(t) + O(\varepsilon) \end{cases}$$

where $A_{\mathbf{I}}$ and $A_{\mathbf{II}}$ are given by (3.3.2). By Lemma [3.2.2] and Theorem [3.2.1] we have

$$\begin{cases} \max_t |Z - (W_{m+1} + \varepsilon Y_{m+2})| = O(\varepsilon) \\ \max_t |Z' - (W_{m+1} + \varepsilon Y_{m+2})'| = O(\varepsilon) \end{cases}$$

where

$$\begin{cases} W_{m+1}' = A_{\mathbf{II}}(t)W_{m+1} + f_{\mathbf{II}}(t) \\ W_{m+1}(\emptyset) = -Y_{m+1}(\emptyset, \varepsilon) \\ Y_{m+2}(t, \varepsilon) = \mathcal{L}\{f_{\mathbf{I}} + A_{\mathbf{I}}W_{m+1}\} \end{cases}$$

Since

$$S_{\varepsilon}(\emptyset, \emptyset) = I$$

by (3.1.12), the initial condition for W_{m+1} is actually independent of ε . By our construction \tilde{X}^{m+1} and $(\tilde{X}^{m+1})'$ approximate X and X' respectively to within $O(\varepsilon^{m+2})$. Thus, we have

$$\begin{aligned} M(\tilde{X}^{m+1}) &= M(\tilde{X}^{m+1}) - M(X) \\ &= [G(\tilde{X}^{m+1}, t, \varepsilon) - G(X, t, \varepsilon)] + [(\tilde{X}^{m+1})' - X'] \\ &= O(\varepsilon^{m+2}) \end{aligned}$$

Since

$$\widetilde{X}^{m+1}(\theta, \varepsilon) = \widehat{X}_0 + \varepsilon^{m+2} Y_{m+2}(\theta, \varepsilon)$$

we conclude that \widetilde{X}^{m+1} is an ε^{m+1} -approximate solution of the system (3.1.14); (3.3.11) then follows from Theorem [3.3.2]. We have:

$$\begin{aligned} M(\widetilde{X}^{m+1}) &= M(\widetilde{X}^m) + M(\widetilde{X}^{m+1}) - M(\widetilde{X}^m) \\ &= M(\widetilde{X}^m) - \varepsilon^{m+1} W_{m+1}'(t) - \varepsilon^{m+2} Y_{m+2}'(t, \varepsilon) + \\ &\quad \varepsilon^{m+1} [G(\widetilde{X}^m, t, \varepsilon)]_X W_{m+1}(t) + \\ &\quad \varepsilon^{m+2} [G(\widetilde{X}^m, t, \varepsilon)]_X Y_{m+2}(t, \varepsilon) + \\ &\quad O(\varepsilon^{2m+2}) \\ &= M(\widetilde{X}^m) - \\ &\quad \varepsilon^{m+1} W_{m+1}'(t) - \varepsilon^{m+2} Y_{m+2}'(t, \varepsilon) + \\ &\quad \varepsilon^{m+1} [A_I + A_{II}] W_{m+1}(t) + \\ &\quad \varepsilon^{m+2} [A_I + A_{II}] Y_{m+2}(t, \varepsilon) + \\ &\quad \varepsilon^{m+2} [[G(X, t, \varepsilon)]_X W_{m+1}]_X [W_{m+1} + Y_{m+2}] \Big|_{X=Y} + \\ &\quad O(\varepsilon^{m+3}) \end{aligned}$$

By our specification for Y_{m+2} and W_{m+1} , the second and third lines of the last expression for $M(\widetilde{X}^{m+1})$ combine with the strictly $O(\varepsilon^{m+1})$ terms of $M(\widetilde{X}^m)$ to give a term of the form

$$\varepsilon^{m+2} \widehat{f}_I(t, \varepsilon)$$

where \widehat{f}_I is strictly oscillatory of class $(p-m-2)$.

The three additional assumptions of the theorem guarantee that the other terms are likewise decomposable. The smoothness conditions for these terms are satisfied since the terms must be

algebraic combinations of terms which meet those requirements. Thus, provided the conditions of the theorem are met, we have decomposability to the next order.

3.4 Solution by Formal Expansion

Under the assumptions of the preceding section, one can approximate the solution of the system (3.1.1) by an expansion whose terms are solutions of equations which can be treated by standard numerical techniques. In principle Theorem [3.3.3] can be applied repeatedly until the decomposability argument breaks down. However, provided the asymptotic form of the solution is guaranteed, one can generate the corresponding smooth equations by a formal procedure that is well-suited to computational implementation since the analytic manipulations, although cumbersome, are simply the Taylor expansions of polynomials. First we introduce a fast time scale

$$(3.4.1) \quad \begin{cases} \mathcal{T} = t/\varepsilon \\ X' = X_t + (1/\varepsilon)X_{\mathcal{T}} \end{cases}$$

and we accordingly reformulate the system (3.1.6):

$$(3.4.2) \quad \begin{cases} X = g_{\text{I}}(X, t, \mathcal{T}, \varepsilon) + g_{\text{II}}(X, t) + f_{\text{I}}(t, \mathcal{T}, \varepsilon) + f_{\text{II}}(t) \\ X(0, \varepsilon) = \hat{X}_0, \quad 0 < \varepsilon \ll 1, \quad 0 < t < T \end{cases}$$

with the appropriate modifications of the conditions on the coefficients:

- (i) \hat{X}_0 is independent of ε ;
- (ii) $g_{\text{I}}(X, t, \mathcal{T}, \varepsilon)^{(i)} = \sum_j a_{ij}(t) \exp[B_{ij}(\varepsilon \mathcal{T})/\varepsilon] p_{ij}(X)$;
- (iii) $g_{\text{II}}(X, t)^{(i)} = \sum_j d_{ij}(t) q_{ij}(X)$;
- (iv) $f_{\text{I}}(t, \mathcal{T}, \varepsilon)^{(i)} = \sum_j c_{ij}(t) \exp[G_{ij}(\varepsilon \mathcal{T})/\varepsilon]$;

$$(v) \quad f_{II}(t)^{(i)} = h_i(t);$$

$$(vi) \quad \{a_{ij}(t), d_{ij}(t), c_{ij}(t), h_i(t)\} \subset C^p(t);$$

$p_{ij}(X)$ and $q_{ij}(X)$ are monomials of positive degree in the components of X ; $B_{ij}(\epsilon\tau)/\epsilon$ and $G_{ij}(\epsilon\tau)/\epsilon$ can be represented by n -vector dot products of the form

$$(3.4.3) \quad N^T P(\epsilon\tau)/\epsilon$$

where N is a constant n -vector with integral components and

$$(3.4.4) \quad \begin{cases} P(\epsilon\tau)^{(i)} = \int_0^{\epsilon\tau} \lambda_i(s) ds \\ P'(\epsilon\tau) = \Lambda(\epsilon\tau) \\ |N^T \Lambda(\epsilon\tau)| > K \end{cases}$$

for some positive K .

To extend the nomenclature of the previous section we call terms of the form (iii) strictly oscillatory (of class p), and we call terms of the form (iv) strictly nonoscillatory (of class p). As before subscript I denotes a strictly oscillatory function and subscript II denotes a strictly nonoscillatory function. The leading order antiderivative of an oscillatory function is given by the linear operator

$$(3.4.5) \quad \mathcal{L}\{c(t) \exp[B(\epsilon\tau)/\epsilon]\} = (c(t)/B'(t)) \exp[B(\epsilon\tau)/\epsilon].$$

To maintain our formalism we must insist that \mathcal{J} -dependence occur only as in (iv). Thus, we must interpret the \mathcal{J} -derivative of an oscillatory function by the rule,

$$(3.4.6) \quad \frac{\partial}{\partial \mathcal{J}} (c(t) \exp[B(\epsilon\tau)/\epsilon]) = c(t) B'(t) \exp[B(\epsilon\tau)/\epsilon],$$

whereby we have:

$$(3.4.7) \quad \begin{cases} \frac{\partial}{\partial \tau} \tilde{\mathcal{L}}\{c(t) \exp[B(\varepsilon \tau)/\varepsilon]\} = c(t) \exp[B(\varepsilon \tau)/\varepsilon] \\ \tilde{\mathcal{L}}\{\frac{\partial}{\partial \tau} c(t) \exp[B(\varepsilon \tau)/\varepsilon]\} = c(t) \exp[B(\varepsilon \tau)/\varepsilon] \end{cases}$$

Clearly this procedure can not correspond to a traditional multi-scaling argument; however, we are only attempting to derive a set of formal rules which mimic the balancing arguments of Section 3.3. The following assumptions give justification to our methodology.

ASSUMPTION [3.4A]

The reduced system

$$(3.4.8) \quad \begin{cases} V' = g_{\text{II}}(V, t) + f_{\text{II}}(t) \\ V(0) = \hat{X}_0 \end{cases}$$

is well posed and has a bounded solution in $C^{p+1}(t)$.

ASSUMPTION [3.4B]

Theorem [3.3.3] can be successively applied to system (3.1.14) to give an asymptotic expansion for $X(t, \varepsilon)$:

$$(3.4.9) \quad \begin{cases} X = \sum_{k=0}^{m-1} X_k \varepsilon^k + \varepsilon^m y_{m+1}(t, \tau, \varepsilon) + o(\varepsilon^{m+1}) \quad (m < p) \\ X_k = Y_k(t, \tau, \varepsilon) + W_k(t) \end{cases}$$

These assumptions assure that our formalism will not break down since our procedure is really a reworking of the decomposability argument of Theorem [3.3.3]. The substitution of the expansion into the differential equation gives:

$$(3.4.10) \quad \left\{ \begin{array}{l} (X_{0t} + X_{1\tau}) + \varepsilon(X_{1t} + X_{2\tau}) + \varepsilon^2(X_{2t} + X_{3\tau}) + \\ = \\ g_{\mathbf{I}}(X_0 + \varepsilon X_1 + \varepsilon^2 X_2 + \dots, t, \tau, \varepsilon) + f_{\mathbf{I}}(t, \tau, \varepsilon) + \\ g_{\mathbf{II}}(X_0 + \varepsilon X_1 + \varepsilon^2 X_2 + \dots, t) + f_{\mathbf{II}}(t) \end{array} \right.$$

We now adopt a formal procedure to solve the system. First we expand each monomial as a power series in ε ; then we balance successive powers of ε . Given that on the k -th level we have previously determined X_0, X_1, \dots, X_{k-1} and Y_k , we balance the $O(\varepsilon^k)$ terms by the following rules:

- I. Determine $W_k(t)$ to eliminate all terms which are strictly nonoscillatory. This is essentially a secularity condition designed to eliminate powers of τ in the expansion. The appropriate initial conditions are:

$$(3.4.11) \quad \left\{ \begin{array}{l} W_0(0) = \hat{X}_0 \\ \vdots \\ W_k(0) = -Y_k(0, 0, \varepsilon) \quad (k > 0) \end{array} \right.$$

- II. Determine $Y_{k+1}(t, \tau, \varepsilon)$ by (3.4.7) to balance the remaining terms.

For the expansion of the monomial $p_{ij}(X)$ we have

$$(3.4.12) \quad \begin{cases} p_{ij}(X_0 + \varepsilon X_1 + \dots) = \\ p_{ij}(X_0) + \varepsilon p_{ij1}(X_0, X_1) + \varepsilon^2 p_{ij2}(X_0, X_1, X_2) + \dots \end{cases}$$

where

$$(3.4.13) \quad \begin{cases} p_{ijk}(X_0, X_1, \dots, X_k) = \\ [p_{ij}(X_0)]_X X_k + \tilde{p}_{ijk}(X_0, X_1, \dots, X_{k-1}) \end{cases}$$

with \tilde{p}_{ijk} a polynomial in the components of its arguments. After a similar expansion for $q_{ij}(X)$ we have the following expression for the i -th component of the right-hand side of (3.4.10):

$$(3.4.14) \quad \begin{cases} g_{\text{I}}(X_0, t, \mathcal{J}, \varepsilon)^{(i)} + g_{\text{II}}(X_0, t)^{(i)} + f_{\text{I}}(t, \mathcal{J}, \varepsilon)^{(i)} + f_{\text{II}}(t)^{(i)} + \\ \sum_{\mathcal{F}} \sum_j \{ \varepsilon^k a_{ij}(t) \exp(B_{ij}(\varepsilon \mathcal{J})/\varepsilon) [p_{ij}(X_0)]_X X_k \} + \\ \sum_{\mathcal{F}} \sum_j \{ \varepsilon^k a_{ij}(t) \exp(B_{ij}(\varepsilon \mathcal{J})/\varepsilon) \tilde{p}_{ijk}(X_0, X_1, \dots, X_{k-1}) \} + \\ \sum_{\mathcal{F}} \sum_j \{ \varepsilon^k d_{ij}(t) [q_{ij}(X_0)]_X X_k \} + \\ \sum_{\mathcal{F}} \sum_j \{ \varepsilon^k d_{ij}(t) \tilde{q}_{ijk}(X_0, X_1, \dots, X_{k-1}) \} \end{cases}$$

Now by assumption $[p_{ij}(X_0)]_X$ and $[q_{ij}(X_0)]_X$ are smooth vector functions and therefore can be combined with the other smooth components. In correspondence to (3.3.2) let $A_{\text{I}}(t, \mathcal{J}, \varepsilon)$ and $A_{\text{II}}(t)$ be $(n \times n)$ matrices whose (m_1, m_2) components are given by

$$(3.4.15) \quad \begin{cases} [A_{\text{I}}(t, \mathcal{J}, \varepsilon)]^{(m_1, m_2)} = \sum_j a_{m_1 j}(t) \exp[B_{ij}(\varepsilon \mathcal{J})/\varepsilon] [p_{ij}(X_0)]_X^{(m_2)} \\ [A_{\text{II}}(t)]^{(m_1, m_2)} = \sum_j d_{m_1 j}(t) [q_{m_1 j}(X_0)]_X^{(m_2)} \end{cases}$$

By our inductive hypothesis we have already determined X_0, X_1, \dots, X_{k-1} , and Y_k when we are ready to balance the $O(\varepsilon^k)$ terms. Thus, \tilde{p}_{ijk} and \tilde{q}_{ijk} are in principle known functions at

this stage, and, with reference to assumption [3.4B], we decompose these into oscillatory and nonoscillatory terms. Considering only the $O(\varepsilon^k)$ terms of (3.4.14), we have

$$(3.4.16) \quad \left\{ \begin{array}{l} \sum_j a_{ij}(t) \exp[B_{ij}(\varepsilon\tau)/\varepsilon] \tilde{p}_{ijK}(X_0, X_1, \dots, X_{K-1}) + \\ \sum_j d_{ij}(t) \tilde{q}_{ijK}(X_0, X_1, \dots, X_{K-1}) \\ = \\ f_{I_K}(t, \tau, \varepsilon)^{(i)} + f_{II_K}(t)^{(i)} \end{array} \right.$$

where f_{I_K} and f_{II_K} depend implicitly on X_0, X_1, \dots, X_{K-1} . And finally by combining (3.4.14), (3.4.15), and (3.4.16) with the observation

$$(3.4.17) \quad (X_K)_\tau = (Y_K(t, \tau, \varepsilon) + W_K(t))_\tau = Y_{K\tau} ,$$

we have:

$$(3.4.18) \quad \left\{ \begin{array}{l} (X_{0\tau} + Y_{1\tau}) + \varepsilon(X_{1\tau} + Y_{2\tau}) + \varepsilon^2(X_{2\tau} + Y_{3\tau}) + \dots \\ = \\ g_I(X_0, t, \tau, \varepsilon) + g_{II}(X_0, t) + f_I(t, \tau, \varepsilon) + f_{II}(t) + \\ \sum_K \varepsilon^K \{ A_I(t, \tau, \varepsilon) X_K + A_{II}(t) X_K + f_{IK}(t, \tau, \varepsilon) + f_{IIK}(t) \} \end{array} \right.$$

We now balance the successive powers of ε by using the formal rules and the given assumptions. For the strictly $O(1)$ terms we have:

$$(3.4.19) \quad \left\{ \begin{array}{l} 0 = [-X_{0\tau} + g_{II}(X_0, t) + f_{II}(t)] + \\ \\ [-Y_{1\tau} + g_I(X_0, t, \tau, \varepsilon) + f_I(t, \tau, \varepsilon)] \end{array} \right.$$

and, therefore, by rules I and II and by assumption [3.4A]:

$$(3.4.20) \quad \begin{cases} X_0 = V(t) \\ Y_1 = \tilde{\mathcal{L}}\{g_{\mathbf{I}}(V, t, \mathcal{V}, \varepsilon) + f_{\mathbf{I}}(t, \mathcal{V}, \varepsilon)\} \end{cases}$$

And likewise for the strictly $O(\varepsilon^k)$ terms we have:

$$(3.4.21) \quad \begin{cases} \emptyset = [-W_{k,t} + A_{\mathbf{I}}W_k + f_{\mathbf{I}k}] + \\ [-Y_{k+1,\mathcal{V}} + A_{\mathbf{I}}W_k + A_{\mathbf{I}}Y_k + -Y_{k,t} + f_{\mathbf{I}k}] + \\ \{A_{\mathbf{I}}Y_k\} \end{cases}$$

where $f_{\mathbf{I}k}$, $f_{\mathbf{I}k}$, and Y_k are determined while W_k is not yet determined. From the functional forms it is clear that $A_{\mathbf{I}}W_k$, $A_{\mathbf{I}}Y_k$, and $Y_{k,t}$ are oscillatory and so can be absorbed into $f_{\mathbf{I}}$. $A_{\mathbf{I}}Y_k$ may contain oscillatory as well as nonoscillatory components but by Assumption [3.4B] must be decomposable. Thus, we have:

$$(3.4.22) \quad \begin{cases} \emptyset = [-W_{k,t} + A_{\mathbf{I}}W_k + \hat{f}_{\mathbf{I}k}] + \\ [-Y_{k+1,\mathcal{V}} + \hat{f}_{\mathbf{I}k}] \end{cases}$$

where $\hat{f}_{\mathbf{I}k}$ is completely determined while $\hat{f}_{\mathbf{I}k}$ depends implicitly on W_k . By rule I we have

$$(3.4.23) \quad \begin{cases} W_k' = A_{\mathbf{I}}W_k + \hat{f}_{\mathbf{I}k} \\ W_k(\emptyset) = -Y_k(\emptyset, \varepsilon) \end{cases}$$

and since $\hat{f}_{\mathbf{I}k}$ is now determined we can use rule II to determine Y_{k+1} :

$$(3.4.24) \quad Y_{k+1} = \tilde{\mathcal{L}}\{\hat{f}_{\mathbf{I}k}\}.$$

In this manner we can successively generate the equations which determine the terms of the expansion.

Example (3.4a)

We return to the nonlinear oscillator of example (3.1a) and calculate the leading-order approximation:

$$(3.4.25) \quad \left\{ \begin{array}{l} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \tilde{V} + \varepsilon(\tilde{W}_1 + \tilde{Y}_1) + \varepsilon^2 \tilde{Y}_2 + O(\varepsilon^2) \\ \quad \quad \quad = \tilde{V} + \varepsilon(\tilde{W}_1 + \tilde{Y}_1) + O(\varepsilon^2) \\ \tilde{V} = \begin{pmatrix} V \\ \bar{V} \end{pmatrix}, \quad \tilde{W}_1 = \begin{pmatrix} W_1 \\ \bar{W}_1 \end{pmatrix}, \quad \tilde{Y}_1 = \begin{pmatrix} Y_1 \\ \bar{Y}_1 \end{pmatrix} \end{array} \right. .$$

In the notation of this section we have:

$$(3.4.26) \quad \left\{ \begin{array}{l} g_{\text{I}}(X, t, \zeta, \varepsilon)^{(1)} = 1/8 \exp(-2i\zeta) x_1^3 + 3/8 \exp(2i\zeta) x_1 x_2^2 \\ \quad \quad \quad - 1/8 \exp(4i\zeta) x_2^3 \\ g_{\text{I}}(X, t, \zeta, \varepsilon)^{(2)} = \bar{g}_{\text{I}}(X, t, \zeta, \varepsilon)^{(1)} \\ g_{\text{II}}(X, t)^{(1)} = -3/8 x_1^2 x_2 \\ g_{\text{II}}(X, t)^{(2)} = \bar{g}_{\text{II}}(X, t)^{(1)} \end{array} \right. ,$$

and the appropriate version of (3.4.18) is

$$(3.4.27) \quad \left\{ \begin{array}{l} \tilde{V}_t + \tilde{Y}_{1\zeta} + \varepsilon(\tilde{Y}_{1t} + \tilde{W}_{1\zeta} + \tilde{Y}_{2\zeta}) + O(\varepsilon^2) \\ \quad \quad \quad = \\ g_{\text{I}}(\tilde{V}, t, \varepsilon) + g_{\text{II}}(\tilde{V}, t) + \\ \varepsilon(A_{\text{I}}\tilde{W}_1 + A_{\text{II}}\tilde{W}_1 + A_{\text{I}}\tilde{Y}_1 + A_{\text{II}}\tilde{Y}_1) + O(\varepsilon^2) \end{array} \right.$$

The first row of A_{I} is the row vector $A_{\text{I}}(t, \zeta, \varepsilon)$, where

$$(3.3.28) \quad \mathcal{A}_{\underline{I}}(t, \zeta, \varepsilon)^T = \begin{bmatrix} [3/8 \exp(-2i\zeta)v^2 + 3/8 \exp(2i\zeta)\bar{v}^2] \\ [3/4 \exp(2i\zeta)v\bar{v} - 3/8 \exp(4i\zeta)\bar{v}^2] \end{bmatrix},$$

and likewise the first row of $A_{\underline{II}}(t)$ is the row vector $\mathcal{A}_{\underline{II}}(t)$, where

$$(3.4.29) \quad \mathcal{A}_{\underline{II}}(t)^T = \begin{bmatrix} -3/4 \bar{v} \\ -3/8 v^2 \end{bmatrix}.$$

By Rule I we determine $V(t)$ as the solution of

$$(3.4.30) \quad \begin{cases} V' = -3/8 v^2 \bar{v} \\ V(0) = 1 \end{cases}$$

and thus $V(t)$ must be real ($V = \bar{V}$). By Rule II we have:

$$(3.4.31) \quad \begin{cases} Y_1(t, \zeta, \varepsilon) = \tilde{\mathcal{L}}\{1/8 \exp(-2i\zeta)v^3 + 3/8 \exp(2i\zeta)v^3 \\ \quad - 1/8 \exp(4i\zeta)v^3\} \\ \quad = i/16 \exp(-2i\zeta)v^3 - 3i/16 \exp(2i\zeta)v^3 \\ \quad \quad + i/32 \exp(4i\zeta)v^3 \end{cases}$$

With the exception of $A_{\underline{I}} \tilde{V}_{\underline{I}}$ all $O(\varepsilon)$ terms of (3.4.27) are either strictly oscillatory or strictly nonoscillatory. We have:

$$(3.4.32) \quad \begin{cases} \mathcal{A}_{\underline{I}} \begin{pmatrix} Y_1 \\ \bar{Y}_1 \end{pmatrix} = F(V) + \langle \text{strictly oscillatory terms} \rangle \\ F(V) = 27i/256 v^5 \end{cases}$$

and thus by Rule I the system for W_1 is

$$(3.4.33) \quad \begin{cases} W_1' = A_{II} \begin{pmatrix} W_1 \\ \bar{W}_1 \end{pmatrix} + F(V) \\ W_1(0) = -Y_1(0,0,\varepsilon) = 3i/32 \end{cases}$$

In terms of the original variables we have

$$(3.4.34) \quad \begin{cases} z_1 = \operatorname{Re}\{\exp(-i\mathcal{Y}) x_1\} \Big|_{\mathcal{Y} = \frac{t}{\varepsilon}} \\ z_2 = \operatorname{Im}\{\exp(-i\mathcal{Y}) x_1\} \Big|_{\mathcal{Y} = \frac{t}{\varepsilon}} \end{cases}$$

We now compare the results of our procedure with the multi-scale approximation as given by Kevorkian and Cole [16, p.123]:

$$(3.4.35) \quad \begin{cases} z_1 = (1 + 3t/4)^{-\frac{1}{2}} \cos(t/\varepsilon) + \\ \quad \varepsilon (3t + 4)^{-\frac{1}{2}} [3/8 (3t + 4)^{-1} + 15/32] \sin(t/\varepsilon) + \\ \quad \varepsilon [1/4 (3t + 4)^{-\frac{3}{2}}] \sin(3t/\varepsilon) + O(\varepsilon^2) \end{cases}$$

The solution of (3.4.30) is

$$(3.4.36) \quad V = (1 + 3t/4)^{-\frac{1}{2}}.$$

After the substitution

$$(3.4.37) \quad W_1 = i w_1,$$

where, by (3.4.33), $w_1(t)$ is real and must satisfy the system

$$(3.4.38) \quad \begin{cases} w_1' = A_{II} \begin{pmatrix} w_1 \\ -w_1 \end{pmatrix} + (27/256) V^5 \\ w_1(0) = 3/32 \end{cases}$$

we can rewrite the first equality of (3.4.34) as

$$(3.4.39) \quad \left\{ \begin{aligned} z_1 &= \operatorname{Re}\{\exp(-it/\varepsilon)V\} + \varepsilon \operatorname{Re}\{\exp(-it/\varepsilon)Y_1\} \\ &\quad + \varepsilon \operatorname{Re}\{i \exp(-it/\varepsilon)w_1\} \\ &= V \cos(t/\varepsilon) + [3/16 V^3 + w_1] \sin(t/\varepsilon) + \\ &\quad [1/32 V^3] \sin(3t/\varepsilon) \end{aligned} \right. .$$

The representations (3.4.35) and (3.4.39) are equivalent if

$$(3.4.40) \quad w_1 = 15/32 (3t + 4)^{-1/2} - 9/8 (3t + 4)^{-3/2} .$$

By inspection one easily can verify that (3.4.40) gives the solution of (3.4.39), and so the representations are indeed equivalent.

Our procedure also can be used to generate higher order corrections. Since the fundamental frequencies of the system (3.1.5) are

$$(3.4.41) \quad \{i/\varepsilon, -i/\varepsilon\}$$

our decomposability principle guarantees the existence of an asymptotic expansion in powers of ε to arbitrarily high order (cf. (3.1.21)).

3.5 EXTENSIONS TO MORE GENERAL SYSTEMS

The most obvious extension of the system (3.1.1) is

$$(3.5.1) \quad \begin{cases} Z' = A(t)/\varepsilon Z + H(Z, t, \varepsilon) \\ Z(0, \varepsilon) = \hat{Z}_0(\varepsilon), \quad 0 < t < T \end{cases}$$

where the t -dependent functions of $H(X, t)$ in (3.1.1) have been replaced by power series in ε with t -dependent coefficients and the initial condition \hat{Z}_0 has been replaced by a power series in ε . Both the Newton iteration theory and the formal expansion theory are essentially unchanged since our principle for determining the leading-order solution of a linear problem is still valid. In general, however, one does not have such an explicit ε -dependence and is content with form (3.1.1), which has the artificial small parameter only in the large component of the system. From our point of view a formalism is useful in so far as it mimics the balancing arguments of the linearization theory which justifies our procedure. Thus, one can absorb an $O(\varepsilon)$ term into an $O(1)$ term without compromising the integrity of the method.

Of somewhat more importance is the system

$$(3.5.2) \quad \begin{cases} Z' = A(t, \varepsilon)/\varepsilon Z + H(Z, t, \varepsilon) \\ Z(0, \varepsilon) = \hat{Z}_0(\varepsilon), \quad 0 < t < T \end{cases}$$

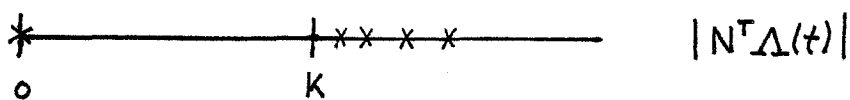
where the t -dependent functions of $A(t)$ also have been replaced by power series in ε with t -dependent coefficients. For

theoretical purposes this system is, of course, equivalent to (3.5.1); however, in practice the ε -dependence is not explicitly given, and so this reduction is not possible.

The theory for system (3.5.2) must be developed with some modifications in the requirements for secondary frequencies. For system (3.3.1) we insist that for some positive K all relevant terms of the form

$$(3.5.3) \quad N^T \underline{\Lambda}(t) \quad (A(t) = \text{diag}\{\underline{\Lambda}(t)^{(c)}\})$$

fall uniformly into one of two distinct subsets of the real line:

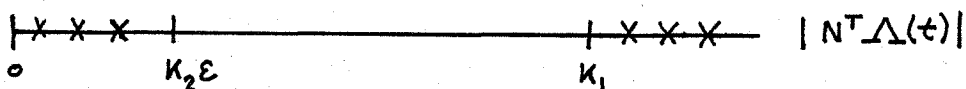


Thus, terms which correspond to frequencies in the subset satisfying (3.1.18) are called strictly oscillatory while, if any $|N^T \underline{\Lambda}(t)|$ vanishes identically, the corresponding term can be reclassified as strictly nonoscillatory. (K/ε) then measures the stiffness of the problem and the n -th term of the derived asymptotic expansion is $O((K/\varepsilon)^n)$.

For the system (3.5.2) all relevant terms of the form

$$(3.5.4) \quad N^T \underline{\Lambda}(t, \varepsilon) \quad (A(t, \varepsilon) = \text{diag}\{\underline{\Lambda}(t, \varepsilon)^{(c)}\})$$

must likewise fall uniformly into one of two well-separated subsets of the real line:



Terms whose frequencies fall into the subset given by

$$(3.5.5) \quad |N^T \underline{\Lambda}(t, \epsilon)| \geq K_1$$

are called strictly oscillatory, and terms whose frequencies fall into the subset given by

$$(3.5.6) \quad |N^T \underline{\Lambda}(t, \epsilon)| \leq K_2 \epsilon$$

are redefined as strictly nonoscillatory. Here K_1 and K_2 are positive real numbers which measure the stiffness of the problem. That is, the frequencies satisfying (3.5.6) are resolved as smooth functions, and the frequencies satisfying (3.5.5) are treated asymptotically by an expansion whose n -th term is $O((K_1/\epsilon)^n)$.

Essentially a frequency must be classifiable as fast or slow if a technique based on the separation of modes is to be successful. The violation of this principle in a neighborhood of an isolated point can be treated by the turning point arguments of the next chapter. In general, however, if there is no well-defined separation between the fast and slow modes of the system (3.5.2), then one is not really solving a singular perturbation problem.

3.6 EXTENSIONS TO PARTIAL DIFFERENTIAL EQUATIONS

The study of oscillatory phenomena frequently leads to systems of partial differential equations in which, as one might expect from physical considerations, the oscillations occur temporally but not spatially. We consider the n -dimensional vector system:

$$(3.6.1) \quad \left\{ \begin{array}{l} U_t = \frac{1}{\epsilon} P_0 \left(\frac{\partial}{\partial \underline{x}} \right) U \\ \quad + (P_1(\underline{x}, t, U, \epsilon, \frac{\partial}{\partial \underline{x}}) + B(\underline{x}, t, U, \epsilon)) U + F(\underline{x}, t, \epsilon) \\ \\ U(\underline{x}, 0) \text{ is given; } t > 0; \quad 0 < \epsilon \ll 1 \\ \\ U(\underline{x}, t) \text{ is } 2\pi\text{-periodic in each component of } \underline{x} \end{array} \right.$$

where

(i) P_0 and P_1 are first order differential operators:

$$(3.6.2) \quad \left\{ \begin{array}{l} P_0 \left(\frac{\partial}{\partial \underline{x}} \right) = \sum_{j=1}^n A_j \frac{\partial}{\partial x_j}, \quad A_j = A_j^* \\ \\ P_1(\underline{x}, t, U, \epsilon, \frac{\partial}{\partial \underline{x}}) = \sum_{j=1}^n B_j(\underline{x}, t, U, \epsilon) \frac{\partial}{\partial x_j}, \quad B_j = B_j^* \end{array} \right. ;$$

the elements of $\{B_j, B\}$ are smooth matrix functions of their arguments; the U -dependence is polynomial in the components of U , and the x -dependence is 2π -periodic in each component;

(ii) $F(\underline{x}, t, \epsilon)$ is a smooth vector function of its arguments and 2π -periodic in each spatial component.

For convenience we shall suppress the ϵ -dependence in the arguments of the dependent variable. Browning and Kreiss [5]

have demonstrated that there exists a time interval

$$(3.6.2) \quad 0 < t < T,$$

where T is independent of ε , but dependent on $\|U(\cdot, 0)\|_p$ and $\max_{0 < s < T} \|F(\cdot, s, \varepsilon)\|_p$, such that

$$(3.6.3) \quad \|U(\cdot, t)\|_p < \text{const} \left(\|U(\cdot, 0)\|_p + \max_{0 < s < T} \|F(\cdot, s, \varepsilon)\|_p \right).$$

Here we define

$$(3.6.4) \quad \begin{cases} \|U(\cdot, t)\|_p = \sum_{|i| \leq p} \|\partial_1^{i_1} \cdots \partial_n^{i_n} U(\cdot, t)\| \\ |i| = \sum_{j=1}^n |i_j| \\ \partial_j = \frac{\partial}{\partial x_j} \end{cases}$$

and we require

$$(3.6.5) \quad p \geq \lfloor n/2 \rfloor + 2,$$

where $\lfloor \cdot \rfloor$ is the integer-part function. Thus the spatial dependency is smooth, and so we can decompose the dependent variable in space and then solve the stiff problem in time.

To illustrate this approach, we consider a scalar version of (3.6.1):

$$(3.6.6) \quad \begin{cases} u_t = \frac{1}{\varepsilon} u_x + r(u, x, t) \\ u(x, t) = u(x + 2\pi, t); \quad 0 < t < T \\ u(x, 0) = f(x) \end{cases}$$

The Fourier decomposition of u is

$$(3.6.7) \quad u(x, t) = \frac{1}{2} \sum_{n=-\infty}^{\infty} c_n(t) \exp(inx),$$

where

$$(3.6.8) \quad \left\{ \begin{array}{l} c_n(t) = \frac{1}{\pi} \int_0^{2\pi} u(x,t) \exp(-inx) dx \\ c_n(0) = \frac{1}{\pi} \int_0^{2\pi} f(x) \exp(-inx) dx = c_n^0 \end{array} \right.$$

From (3.6.6), (3.6.7), and (3.6.8) one can derive equations for the t -dependent coefficients:

$$(3.6.9) \quad \left\{ \begin{array}{l} c_n'(t) = (n/\varepsilon) i c_n(t) + \sum_{j=1}^l \alpha_{jn}(t) p_j(c_0, c_1, c_{-1}, \dots) \\ c_n(0) = c_n^0; \quad -\infty < n < \infty \\ 0 < t < T \end{array} \right.$$

Here each $\alpha_{jn}(t)$ is a smooth coefficient and each $p_j(c_0, \dots)$ is a monomial in its arguments.

Given the a priori bounds on the solution's spatial derivatives, one can use (3.6.8) to achieve bounds on the coefficients of the high modes and thereby justify the truncation of (3.6.9) to a finite system. Let

$$(3.6.10) \quad C(t) = [c_0, c_1, c_{-1}, \dots, c_m, c_{-m}]^T.$$

Then by ignoring all terms dependent on the coefficients of the higher modes we reach a system of the form

$$(3.6.11) \quad \left\{ \begin{array}{l} C' = \frac{i}{\varepsilon} \Lambda C + H(C, t, \varepsilon) \\ C(0) = C_0; \quad 0 < t < T \end{array} \right.$$

System (3.6.11) is in the canonical form of Chapter III (3.1.1); moreover, since the entries of the diagonal

matrix Λ are integral constants, Theorem [3.3] guarantees the existence of an asymptotic expansion in powers of ϵ to higher orders (cf. (3.1.21)). This procedure is likewise applicable to the more general system (3.6.1), where the introduction of a multiple Fourier expansion

$$(3.6.12) \quad U(x, t) = \frac{1}{2^n} \sum_{m_1, \dots, m_n} \xi_{m_1, \dots, m_n}(t) \exp(i(m_1 x_1 + \dots + m_n x_n))$$

yields a coupled system of highly oscillatory ordinary differential equations for the time dependent vectors

$$(3.6.13) \quad \xi_{m_1, \dots, m_n}(t) \quad .$$

3.7 Computational Examples

In this section we describe possible implementations for problems in nonlinear oscillations. For illustrative purposes we include the resulting algebra although, as we have previously noted, the analytic manipulations are conceptually simple enough to be computationally feasible.

A simple mass spring system with small damping can be modeled by the equation for a Rayleigh oscillator:

$$(3.7.1) \quad \begin{cases} d^2y/d\tilde{t}^2 + y = \epsilon [dy/d\tilde{t} - (1/3) (dy/d\tilde{t})^3] \\ y(0) = 1, \quad y'(0) = 0 \\ 0 < \tilde{t} < T/\epsilon, \quad 0 < \epsilon \ll 1 \end{cases} .$$

T is independent of ϵ , and ϵ is the ratio of the characteristic time for one oscillation to the characteristic damping time [see Kevorkian and Cole [17]]. After the change of variables

$$(3.7.2) \quad \begin{cases} z_1 = y, \quad z_2 = dy/dt \\ t = \epsilon \tilde{t} \end{cases} ,$$

we can rewrite the system as

$$(3.7.3) \quad \begin{cases} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}' = \begin{bmatrix} 0 & 1/\epsilon \\ -1/\epsilon & 0 \end{bmatrix} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} + \begin{pmatrix} 0 \\ z_2 - (z_2)^3/3 \end{pmatrix} \\ \begin{pmatrix} z_1(0, \epsilon) \\ z_2(0, \epsilon) \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad 0 < t < T, \quad 0 < \epsilon \ll 1 \end{cases} ,$$

and, as in example (3.4a), the system can be reduced to diagonal form; thus, after the change of variables

$$(3.7.4) \quad \begin{cases} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = S \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ -i & i \end{bmatrix} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = \begin{pmatrix} (u_1 + u_2)/2 \\ i(u_1 - u_2)/2 \end{pmatrix} \\ \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = S^{-1} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = \begin{bmatrix} 1 & i \\ 1 & -i \end{bmatrix} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = \begin{pmatrix} z_1 + iz_2 \\ z_1 - iz_2 \end{pmatrix} \end{cases}$$

the equations become

$$(3.7.5) \quad \begin{cases} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}' = \begin{bmatrix} -i/\varepsilon & 0 \\ 0 & i/\varepsilon \end{bmatrix} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} + \begin{pmatrix} -f(u_1, u_2) \\ f(u_1, u_2) \end{pmatrix} \\ \begin{pmatrix} u_1(0, \varepsilon) \\ u_2(0, \varepsilon) \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad 0 < t < T, \quad 0 < \varepsilon \ll 1 \\ f(u_1, u_2) = -\frac{1}{2} u_1 + \frac{1}{2} u_2 + \frac{1}{24} u_2^3 - \frac{1}{8} u_2^2 u_1 + \frac{1}{8} u_2 u_1^2 - \frac{1}{24} u_1^3 \end{cases}$$

The system (3.7.5) has the form (3.1.1). Next we make the change of variables

$$(3.7.6) \quad \begin{cases} u_1 = \exp(-it/\varepsilon) x_1 \\ u_2 = \exp(it/\varepsilon) x_2 \end{cases}$$

and obtain

$$(3.7.7) \quad \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}' = \begin{pmatrix} -\exp(it/\varepsilon) f(\exp(-it/\varepsilon)x_1, \exp(it/\varepsilon)x_2) \\ \exp(-it/\varepsilon) f(\exp(-it/\varepsilon)x_1, \exp(it/\varepsilon)x_2) \end{pmatrix},$$

and so we have:

$$(3.7.8) \quad \left\{ \begin{array}{l} x_1' = 1/2 x_1 - 1/8 x_2 x_1^2 + \\ \quad - 1/2 \exp(2it/\varepsilon) x_2 - 1/24 \exp(4it/\varepsilon) x_2^3 + \\ \quad - 1/8 \exp(2it/\varepsilon) x_2^2 x_1 + 1/24 \exp(-2it/\varepsilon) x_1^3 \\ x(0, \varepsilon) = 1, \quad x_2 = \bar{x}_1, \quad 0 < t < T \end{array} \right.$$

Thus the original two-dimensional real system has been replaced by a one-dimensional complex system. Following the formalism of Section 3.4, we now introduce an asymptotic expansion in powers of ε :

$$(3.7.9) \quad \left\{ \begin{array}{l} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \tilde{V} + \varepsilon(\tilde{W}_1 + \tilde{Y}_1) + O(\varepsilon^2) \\ \tilde{V} = \begin{pmatrix} V \\ \bar{V} \end{pmatrix}, \quad \tilde{W}_1 = \begin{pmatrix} W_1 \\ \bar{W}_1 \end{pmatrix}, \quad \tilde{Y}_1 = \begin{pmatrix} Y_1 \\ \bar{Y}_1 \end{pmatrix} \end{array} \right.$$

Here the fast time scale is

$$(3.7.10) \quad \zeta = t/\varepsilon,$$

and in the notation of the section we have:

$$(3.7.11) \quad \left\{ \begin{array}{l} g_{\mathbf{I}}(X, t, \zeta, \varepsilon)^{(1)} = -1/2 \exp(2i\zeta) x_2 - 1/24 \exp(4i\zeta) x_2^3 + \\ \quad 1/8 \exp(2i\zeta) x_2^2 x_1 + 1/24 \exp(-2i\zeta) x_1^3 \\ g_{\mathbf{I}}(X, t, \zeta, \varepsilon)^{(2)} = \overline{g_{\mathbf{I}}}(X, t, \zeta, \varepsilon)^{(1)} \\ g_{\mathbf{II}}(X, t)^{(1)} = 1/2 x_1 - 1/8 x_2 x_1^2 \\ g_{\mathbf{II}}(X, t)^{(2)} = \overline{g_{\mathbf{I}}}(X, t)^{(1)} \end{array} \right.$$

The appropriate version of (3.4.18) is

$$(3.7.12) \quad \left\{ \begin{array}{l} \widetilde{V}_t + \widetilde{Y}_{1\zeta} + \varepsilon(\widetilde{Y}_{1t} + \widetilde{W}_{1t} + \widetilde{Y}_{2\zeta}) + O(\varepsilon^2) \\ = \\ g_{\text{I}}(\widetilde{V}, t, \zeta, \varepsilon) + g_{\text{II}}(\widetilde{V}, t) + \\ \varepsilon(A_{\text{I}}\widetilde{W}_1 + A_{\text{II}}\widetilde{W}_1 + A_{\text{I}}\widetilde{Y}_1 + A_{\text{II}}\widetilde{Y}_1) + O(\varepsilon^2) \end{array} \right.$$

The first row of A_{I} is the row vector $A_{\text{I}}(t, \zeta, \varepsilon)$ where

$$(3.7.13) \quad A_{\text{I}}(t, \zeta, \varepsilon)^T = \begin{bmatrix} [1/8 \exp(2i\zeta)\bar{V}^2 + 1/8 \exp(-2i\zeta)V^2] \\ [-1/2 \exp(2i\zeta) - 1/8 \exp(4i\zeta)\bar{V}^2 \\ + 1/4 \exp(2i\zeta)\bar{V}V] \end{bmatrix},$$

and likewise the first row of $A_{\text{II}}(t)$ is the row vector $A_{\text{II}}(t)$ where

$$(3.7.14) \quad A_{\text{II}}(t)^T = \begin{bmatrix} 1/2 - 1/4 \bar{V}V \\ -1/8 V^2 \end{bmatrix}.$$

By Rule I we determine $V(t)$ as the solution of

$$(3.7.15) \quad \left\{ \begin{array}{l} V' = 1/2 V - 1/8 \bar{V}V^2 \\ V(0) = 1 \end{array} \right.$$

and by Rule II Y_1 is given by

$$(3.7.16) \quad \left\{ \begin{array}{l} Y_1(t, \zeta, \varepsilon) = \mathcal{L}\{-1/2 \exp(2i\zeta)\bar{V} - 1/24 \exp(4i\zeta)\bar{V}^3 + \\ 1/8 \exp(2i\zeta)\bar{V}^2V + 1/24 \exp(-2i\zeta)V^3\} \\ = i/4 \exp(2i\zeta)\bar{V} + i/96 \exp(4i\zeta)\bar{V}^3 + \\ -i/16 \exp(2i\zeta)\bar{V}^2V + i/48 \exp(-2i\zeta)V^3 \end{array} \right.$$

With the exception of $A_{\text{I}} Y$ all $O(\epsilon)$ terms of (3.7.12) are either strictly oscillatory or strictly nonoscillatory. We have:

$$(3.7.17) \quad \begin{cases} A_{\text{I}} \left(\frac{Y_1}{\bar{Y}_1} \right) = R(V) + \langle \text{strictly oscillatory terms} \rangle \\ R(V) = i/8 V + 3i/256 \bar{V}^2 V^3 - i/16 \bar{V} V^2 \end{cases}$$

and thus by Rule I the system for W_1 is

$$(3.7.18) \quad \begin{cases} W_1' = A_{\text{II}} \left(\frac{W_1}{\bar{W}_1} \right) + R(V) \\ W_1(0) = -Y_1(0, 0, \epsilon) = -7i/32 \end{cases}$$

By using (3.7.10), we can restore the full t -dependence of the coefficients; in terms of the original variables we then have

$$(3.7.19) \quad \begin{cases} z = \text{Re}\{\exp(-i\tau) x_1\} \Big|_{\tau = \frac{t}{\epsilon}} \\ z = \text{Im}\{\exp(-i\tau) x_1\} \Big|_{\tau = \frac{t}{\epsilon}} \end{cases}$$

Using this analysis, we outline an approximation scheme based on standard discretization techniques. First we ignore $O(\epsilon^2)$ terms of the expansion. Using a step size h , we approximate the solution V of the equation (3.7.15) by means of a fourth-order Runge-Kutta scheme (see Lambert [20], p.126); then Y_1 is given explicitly by (3.7.16), and with the same step size h we approximate W_1 from (3.7.18) by means of a forward Euler predictor followed by two trapezoidal rule correctors (see Lambert [20], p. 85). The total error for this approximation is then:

$$(3.7.20) \quad O(h^4) + O(\epsilon^2) + O(\epsilon h^2) \quad .$$

Computations were done with single-precision accuracy on a VAX11/780 for the case,

$$(3.7.21) \quad \begin{cases} \epsilon = .01 \\ h = .1 \end{cases} ,$$

and therefore, by (3.7.20), one might expect a grid error of approximately 10^{-3} .

In phase space the solutions of (3.7.3) approach a stable limit cycle of approximate radius two. Thus, in figure <3.7.1> we plot the amplitude

$$(3.7.22) \quad A(t) = [z_1^2 + z_2^2]^{1/2}$$

as a function of the rescaled time variable, and also we plot the curve of

$$(3.7.23) \quad z_1(t) = y(t/\epsilon)$$

in figure <3.7.2>. In both cases we have linearly interpolated the amplitudes and the phases between grid points. In table <3.7.1> we compare the computed grid values with the accepted function values, which were computed with double-precision accuracy by means of a fourth order Runge-Kutta scheme with a time step

$$(3.7.24) \quad h = 10^{-4} .$$

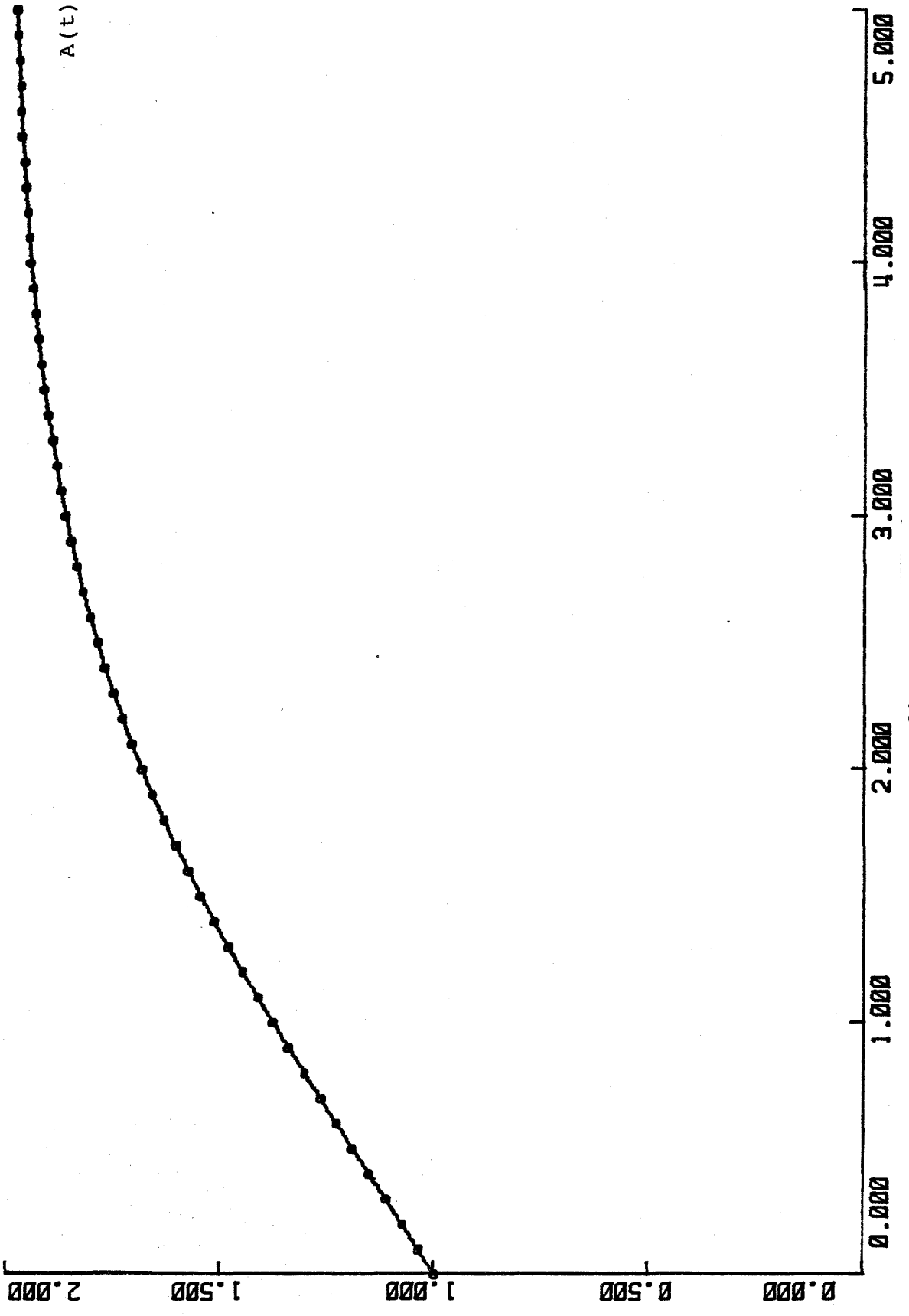


figure <3.7.1>

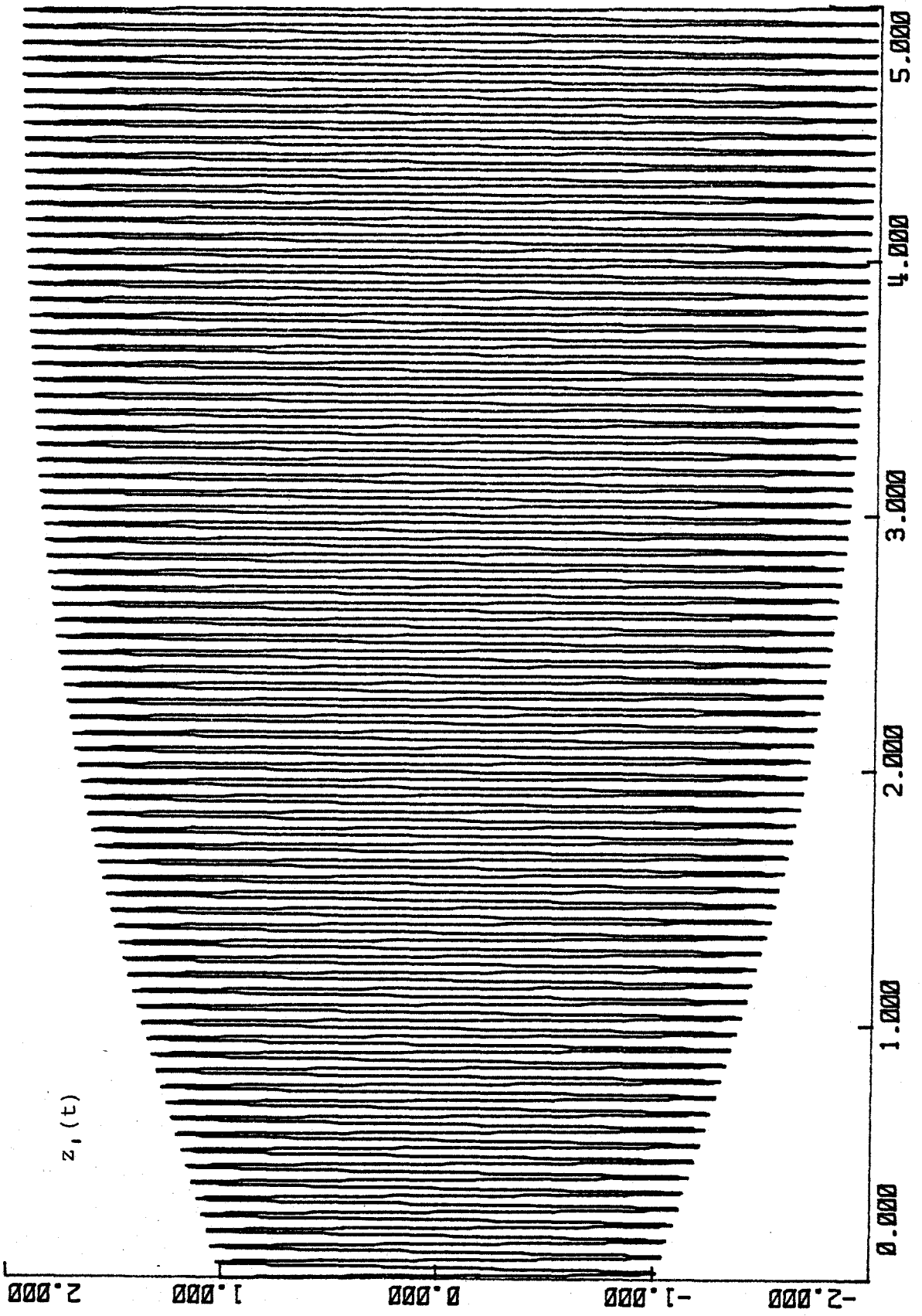


figure <3.7.2>

N	T_N	$A(T_N)$	$Z_1(T_N)$	ERR(A) (absolute errors)	ERR(Z_1) (absolute errors)
0	0.00	1.00000000	1.00000000	0.0E+00	0.0E+00
1	0.10	1.0361134	-0.8686084	0.8E-06	0.1E-06
2	0.20	1.0747031	0.4352194	0.7E-06	0.1E-05
3	0.30	1.1145539	0.1760627	0.5E-06	0.3E-05
4	0.40	1.1541607	-0.7716992	0.1E-05	0.4E-05
5	0.50	1.1920174	1.1502984	0.2E-05	0.2E-05
6	0.60	1.2283492	-1.1698117	0.3E-05	0.4E-05
7	0.70	1.2661535	0.7996659	0.1E-05	0.1E-05
8	0.80	1.3050122	-0.1398034	0.4E-05	0.1E-04
9	0.90	1.3432373	-0.6051558	0.1E-05	0.1E-04
10	1.00	1.3802726	1.1908320	0.4E-05	0.2E-05
11	1.10	1.4145483	-1.4131837	0.6E-05	0.6E-05
12	1.20	1.4483656	1.1782711	0.4E-05	0.1E-04
13	1.30	1.4831481	-0.5411782	0.3E-05	0.2E-05
14	1.40	1.5163764	-0.3039049	0.7E-05	0.3E-05
15	1.50	1.5486927	1.0847387	0.4E-05	0.5E-06
16	1.60	1.5789385	-1.5404052	0.8E-05	0.2E-05
17	1.70	1.6066599	1.5069457	0.8E-05	0.2E-04
18	1.80	1.6350559	-0.9760583	0.3E-05	0.3E-04
19	1.90	1.6617370	0.1064526	0.1E-04	0.4E-05
20	2.00	1.6864364	0.8245143	0.4E-06	0.3E-05
21	2.10	1.7109314	-1.5126638	0.9E-05	0.6E-05
22	2.20	1.7324642	1.7257520	0.1E-04	0.1E-04
23	2.30	1.7534370	-1.3801000	0.9E-05	0.6E-05
24	2.40	1.7739080	0.5748035	0.6E-05	0.3E-04
25	2.50	1.7909012	0.4346920	0.9E-05	0.3E-04
26	2.60	1.8083004	-1.3218288	0.1E-04	0.2E-05
27	2.70	1.8247517	1.7961023	0.1E-04	0.1E-04
28	2.80	1.8390280	-1.6950778	0.1E-04	0.1E-04
29	2.90	1.8539859	1.0405447	0.5E-05	0.6E-04
30	3.00	1.8656223	-0.0386439	0.1E-04	0.8E-05
31	3.10	1.8758882	-0.9876618	0.4E-05	0.5E-04
32	3.20	1.8878965	1.7061471	0.1E-04	0.2E-04
33	3.30	1.8976505	-1.8811355	0.2E-04	0.2E-04
34	3.40	1.9076006	1.4482372	0.1E-04	0.5E-04
35	3.50	1.9163814	-0.5415334	0.7E-05	0.9E-05
36	3.60	1.9214646	-0.5470619	0.6E-05	0.7E-04
37	3.70	1.9287945	1.4660136	0.1E-04	0.3E-05
38	3.80	1.9361542	-1.9188806	0.2E-04	0.1E-04
39	3.90	1.9421252	1.7552749	0.1E-04	0.4E-04
40	4.00	1.9489839	-1.0224489	0.6E-05	0.1E-04
41	4.10	1.9519360	-0.0444038	0.1E-04	0.5E-04
42	4.20	1.9546887	1.1006080	0.8E-05	0.5E-04
43	4.30	1.9602964	-1.8064740	0.1E-04	0.1E-04
44	4.40	1.9641300	1.9338502	0.2E-04	0.2E-04
45	4.50	1.9688038	-1.4372377	0.1E-04	0.2E-04
46	4.60	1.9718472	0.4744309	0.6E-05	0.5E-04
47	4.70	1.9715779	0.6431262	0.1E-05	0.7E-04
48	4.80	1.9746908	-1.5553792	0.1E-04	0.3E-04
49	4.90	1.9780447	1.9700294	0.2E-04	0.1E-04
50	5.00	1.9805924	-1.7512197	0.1E-04	0.5E-04

TABLE <3.7.1>

One likewise can apply these techniques to systems of coupled nonlinear oscillators. The following extensively analyzed system is taken from the theory of stellar orbits in a galaxy (see, for example, Contopoulos [7], Hori [14], and Kevorkian and Cole [16]):

$$(3.7.25) \quad \left\{ \begin{array}{l} d^2 r_1 / d\tilde{t}^2 + a^2 r_1 = \varepsilon r_2^2 \\ d^2 r_2 / d\tilde{t}^2 + b^2 r_2 = 2\varepsilon r_1 r_2 \\ r_1(0, \varepsilon) = 1, r_2(0, \varepsilon) = 1 \\ r_1'(0, \varepsilon) = 0, r_2'(0, \varepsilon) = 0 \\ 0 < \varepsilon \ll 1, 0 < \tilde{t} < T/\varepsilon \end{array} \right.$$

Here r_1 stands for the radial displacement of the orbit of a star from a reference circular orbit, and r_2 stands for the deviation of the orbit from the galactic plane. With the change of variables

$$(3.7.26) \quad \left\{ \begin{array}{l} Z = [z_1, z_2, z_3, z_4]^T = [r_1, r_1'/a, r_2, r_2'/b]^T \\ t = \varepsilon \tilde{t} \end{array} \right.$$

we have:

$$(3.7.27) \quad \left\{ \begin{array}{l} Z' = (1/\varepsilon) \begin{bmatrix} 0 & a & 0 & 0 \\ -a & 0 & 0 & 0 \\ 0 & 0 & 0 & b \\ 0 & 0 & -b & 0 \end{bmatrix} Z + \begin{pmatrix} 0 \\ (1/a) z_3^2 \\ 0 \\ (2/b) z_1 z_3 \end{pmatrix} \\ Z(0, \varepsilon) = [1, 0, 1, 0]^T, \quad 0 < t < T, \quad 0 < \varepsilon \ll 1 \end{array} \right.$$

By a transformation similar to (3.7.4), we can reduce the system to diagonal form. Thus after the change of variables

$$(3.7.28) \quad \begin{cases} z = \tilde{S} U \\ U^T = [u_1, u_2, u_3, u_4] \\ \tilde{S} = \begin{bmatrix} S & \emptyset \\ \emptyset & S \end{bmatrix}, S = (1/2) \begin{bmatrix} 1 & 1 \\ -i & i \end{bmatrix} \end{cases}$$

the equations become:

$$(3.7.29) \quad \begin{cases} U' = (i/\varepsilon) \begin{bmatrix} -a & \emptyset & \emptyset & \emptyset \\ \emptyset & a & \emptyset & \emptyset \\ \emptyset & \emptyset & -b & \emptyset \\ \emptyset & \emptyset & \emptyset & b \end{bmatrix} U + \begin{pmatrix} f_1(U) \\ -f_1(U) \\ f_2(U) \\ -f_2(U) \end{pmatrix} \\ U(0, \varepsilon) = [1, 1, 1, 1]^T, \quad 0 < t < T, \quad 0 < \varepsilon \ll 1 \\ f_1(U) = i (u_3 + u_4)^2 / 4a \\ \quad = (i/4a) [u_3^2 + 2u_3u_4 + u_4^2] \\ f_2(U) = i (u_1 + u_2)(u_3 + u_4) / 2b \\ \quad = (i/2b) [u_1u_3 + u_1u_4 + u_2u_3 + u_2u_4] \end{cases}$$

As in (3.7.6), we now factor out the leading order oscillatory behavior by means of the transformation

$$(3.7.30) \quad \begin{cases} U = T(t, \varepsilon) X \\ X^T = [x_1, x_2, x_3, x_4] \\ T(t, \varepsilon) = \text{diag}[\exp(-iat/\varepsilon), \exp(iat/\varepsilon), \exp(-ibt/\varepsilon), \exp(ibt/\varepsilon)] \end{cases}$$

and obtain the system:

$$(3.7.31) \quad \left\{ \begin{array}{l} X' = \begin{bmatrix} \exp(iat/\varepsilon) f_1(T(t,\varepsilon) X) \\ -\exp(-iat/\varepsilon) f_1(T(t,\varepsilon) X) \\ \exp(ibt/\varepsilon) f_2(T(t,\varepsilon) X) \\ -\exp(-ibt/\varepsilon) f_2(T(t,\varepsilon) X) \end{bmatrix} \\ X(0,\varepsilon) = [1,1,1,1]^T, \quad 0 < t < T, \quad 0 < \varepsilon \ll 1 \end{array} \right.$$

From the structure of the transformations we have

$$(3.7.32) \quad \left\{ \begin{array}{l} x_2 = \bar{x}_1 \\ x_4 = \bar{x}_3 \end{array} \right. ,$$

and therefore, similarly to the first example, we have replaced the original four-dimensional real system with a two-dimensional complex system. Once again, in the spirit of Section 3.4, we introduce an asymptotic expansion in powers of ε :

$$(3.7.33) \quad \left\{ \begin{array}{l} x = \tilde{V} + \varepsilon(\tilde{W}_1 + \tilde{Y}_1) + O(\varepsilon^2) \\ \tilde{V}^T = [V_A, \bar{V}_A, V_B, \bar{V}_B] \\ \tilde{W}_1^T = [W_A, \bar{W}_A, W_B, \bar{W}_B] \\ \tilde{Y}_1^T = [Y_A, \bar{Y}_A, Y_B, \bar{Y}_B] \end{array} \right. ;$$

here the fast scale is again (3.7.10). The most interesting resonances occur for the case

$$(3.7.34) \quad a = 2b ,$$

and so we consider the parameter valuations

$$(3.7.35) \quad \begin{cases} a = 1 \\ b = .5 \end{cases}$$

System (3.7.31) then has the form (3.4.1) with

$$(3.7.36) \quad \begin{cases} g_{\text{H}}^{(1)} = (i/2) x_3 x_4 \exp(i\zeta) + (i/4) x_4^2 \exp(2i\zeta) \\ g_{\text{H}}^{(2)} = \overline{g_{\text{H}}^{(1)}} \\ g_{\text{H}}^{(3)} = i x_1 x_3 \exp(-i\zeta) + i x_2 x_3 \exp(i\zeta) \\ \quad + i x_2 x_4 \exp(2i\zeta) \\ g_{\text{H}}^{(4)} = \overline{g_{\text{H}}^{(3)}} \\ g_{\text{H}}^{(5)} = (i/4) x_3^2 \\ g_{\text{H}}^{(6)} = \overline{g_{\text{H}}^{(5)}} \\ g_{\text{H}}^{(7)} = i x_1 x_4 \\ g_{\text{H}}^{(8)} = \overline{g_{\text{H}}^{(7)}} \end{cases}$$

We proceed as in the first example by applying the balancing arguments of Section 3.4 to equation (3.7.12). The first and third rows of A_{II} are given by the (2 x 4) matrix A_{II}^{T} , where

$$(3.7.37) \quad A_{\text{II}}^{\text{T}} = \begin{bmatrix} \emptyset & i \overline{V_{\text{B}}} \\ \emptyset & \emptyset \\ (i/2) V_{\text{B}} & \emptyset \\ \emptyset & i V_{\text{A}} \end{bmatrix}$$

The first and third rows of A_{I} are given by the (2 x 4) matrix A_{I} , where

$$(3.7.38) \quad \mathcal{A}_{\mathcal{I}}^T = \begin{bmatrix} 0 & i V_B \exp(-i\mathcal{I}) \\ 0 & i V_B \exp(i\mathcal{I}) + i \bar{V}_B \exp(2i\mathcal{I}) \\ (i/2) \bar{V}_B \exp(i\mathcal{I}) & i V_A \exp(-i\mathcal{I}) + i \bar{V}_A \exp(i\mathcal{I}) \\ (i/2) V_B \exp(i\mathcal{I}) + (i/2) \bar{V}_B \exp(2i\mathcal{I}) & i \bar{V}_A \exp(2i\mathcal{I}) \end{bmatrix}$$

\tilde{V} is then determined as the solution of

$$(3.7.39) \quad \begin{cases} \begin{pmatrix} V_A \\ V_B \end{pmatrix}' = \begin{pmatrix} (i/4) V_B^2 \\ (i/2) \bar{V}_B V_A \end{pmatrix} \\ \begin{pmatrix} V_A(\theta, \mathcal{E}) \\ V_B(\theta, \mathcal{E}) \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \theta < t < T \end{cases}$$

and then, as in (3.7.16), \tilde{Y}_1 is given explicitly by

$$(3.7.40) \quad \begin{pmatrix} Y_A \\ Y_B \end{pmatrix} = \begin{pmatrix} (1/2) V_B \bar{V}_B \exp(i\mathcal{I}) + (1/8) \bar{V}_B^2 \exp(2i\mathcal{I}) \\ - V_B V_A \exp(-i\mathcal{I}) + V_B \bar{V}_A \exp(i\mathcal{I}) + (1/2) \bar{V}_B \bar{V}_A \exp(2i\mathcal{I}) \end{pmatrix}$$

With the exception of $\mathcal{A}_{\mathcal{I}} \tilde{Y}_1$ all $O(\mathcal{E})$ terms of (3.7.12) are either strictly oscillatory or strictly nonoscillatory. We have:

$$(3.7.41) \quad \mathcal{A}_{\mathcal{I}} \begin{pmatrix} Y_A \\ \bar{Y}_A \\ Y_B \\ \bar{Y}_B \end{pmatrix} = \begin{pmatrix} (i/4) V_A V_B \bar{V}_B \\ (i/2) V_A \bar{V}_A V_B + (9i/8) V_B^2 \bar{V}_B \\ + \langle \text{oscillatory terms} \rangle \end{pmatrix}$$

and thus the system for \tilde{W}_1 is:

$$(3.7.42) \quad \left\{ \begin{array}{l} \begin{pmatrix} W_A \\ W_B \end{pmatrix}' = \begin{pmatrix} (i/2) V_B W_B \\ i \bar{V}_B W_A + i V_A \bar{W}_B \end{pmatrix} \\ \quad + \begin{pmatrix} (i/4) V_A V_B \bar{V}_B \\ (i/2) V_A \bar{V}_A V_B + (9i/8) V_B^2 \bar{V}_B \end{pmatrix} \\ \begin{pmatrix} W_A(0) \\ W_B(0) \end{pmatrix} = \begin{pmatrix} -Y_A(0,0,\varepsilon) \\ -Y_B(0,0,\varepsilon) \end{pmatrix} = \begin{pmatrix} (-5/8) \\ (-1/2) \end{pmatrix} \\ \quad 0 < t < T \end{array} \right. .$$

Once again the full t -dependence of the system can be restored by (3.7.10). In the original system variables we have:

$$(3.7.43) \quad \left\{ \begin{array}{l} z_1 = \operatorname{Re}\{\exp(-i\zeta) x_1\} \Big|_{\zeta=t/\varepsilon} \\ z_2 = \operatorname{Im}\{\exp(-i\zeta) x_1\} \Big|_{\zeta=t/\varepsilon} \\ z_3 = \operatorname{Re}\{\exp(-i\zeta/2) x_3\} \Big|_{\zeta=t/\varepsilon} \\ z_4 = \operatorname{Im}\{\exp(-i\zeta/2) x_3\} \Big|_{\zeta=t/\varepsilon} \end{array} \right. .$$

Now we can apply the same approximation techniques to this system, which is characterized by the energy integral

$$(3.7.44) \quad \left\{ \begin{array}{l} E(t) = (a^2/2) (z_1^2 + z_2^2) + (b^2/2) (z_3^2 + z_4^2) - \varepsilon z_1 z_3^2 \\ \quad = E_0 \end{array} \right.$$

The sharing of this energy between the two oscillators is illustrated in figure <3.7.3>, where we have plotted the leading order energy functions

$$(3.7.45) \quad \left\{ \begin{array}{l} \widetilde{E}_1(t) = (a^2/2) [z_1^2 + z_2^2] \\ \widetilde{E}_2(t) = (b^2/2) [z_3^2 + z_4^2] \end{array} \right. .$$

Once again we have linearly interpolated the amplitudes and phases between the grid points. In table <3.7.2> we compare the computed grid values with the accepted function values, which were computed with double-precision accuracy by means of a fourth-order Runge-Kutta scheme with a time step

$$(3.7.46) \quad h = 10^{-4}.$$

The grid error, somewhat larger than in the first example, is mainly due to the truncation of the asymptotic expansion. Indeed, decreasing ϵ by a factor of (.1) caused the grid error to fall by a factor of (.01).

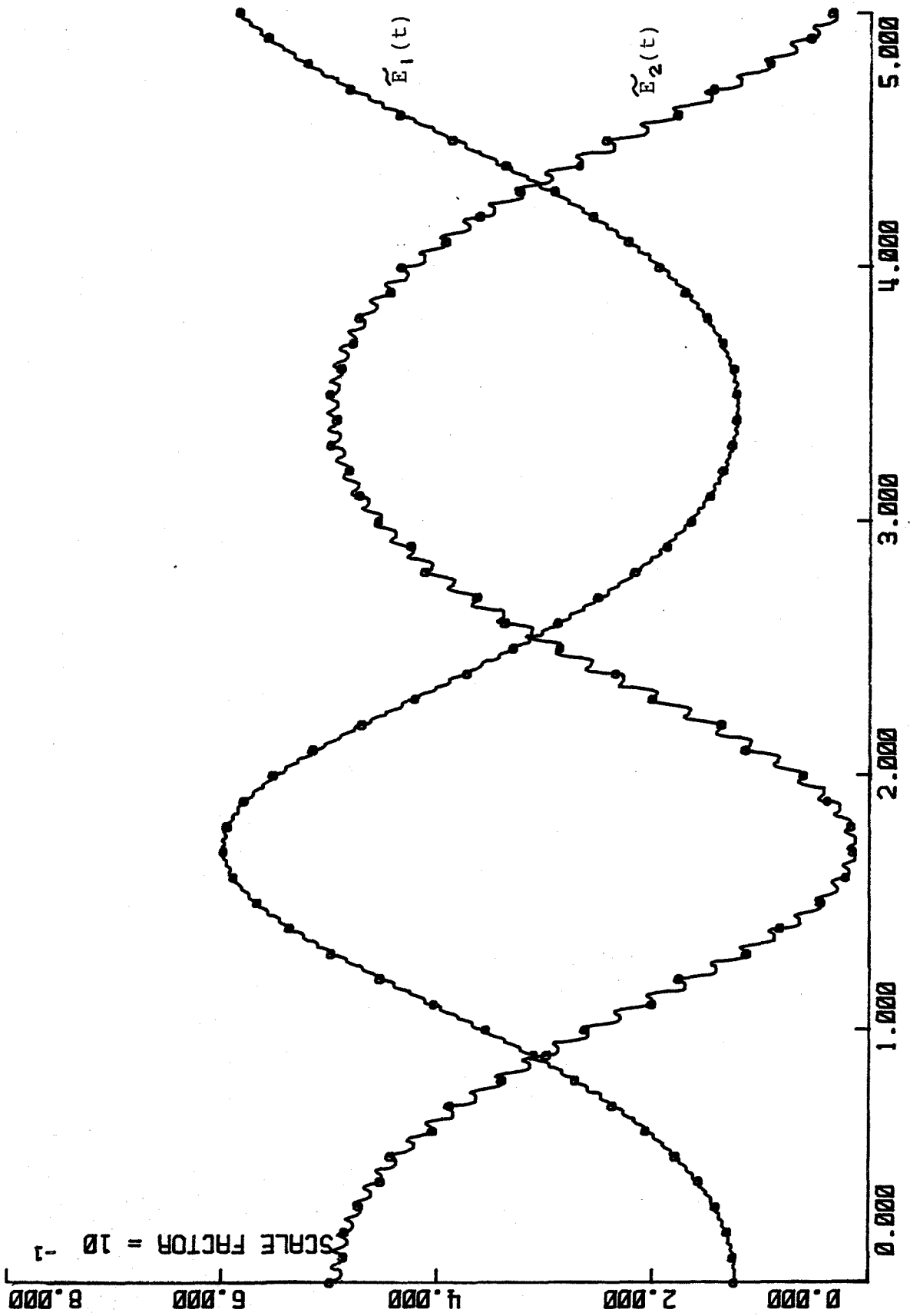


figure <3.7.3>

N	T_N	$\tilde{E}_1(T_N)$	$\tilde{E}_2(T_N)$	$E(T_N)$	ERR(\tilde{E}_1) (absolute errors)	ERR(\tilde{E}_2) (absolute errors)
0	0.00	0.500000	0.125000	0.615000	0.0E+00	0.0E+00
1	0.10	0.488006	0.126700	0.615007	0.1E-03	0.6E-04
2	0.20	0.487005	0.132210	0.615005	0.7E-04	0.2E-04
3	0.30	0.473004	0.142310	0.615006	0.5E-04	0.6E-05
4	0.40	0.453260	0.158210	0.615120	0.2E-03	0.6E-04
5	0.50	0.443740	0.180190	0.614990	0.3E-05	0.1E-04
6	0.60	0.405770	0.207550	0.615006	0.1E-03	0.6E-04
7	0.70	0.388960	0.238890	0.615005	0.1E-03	0.8E-04
8	0.80	0.341830	0.273160	0.615008	0.4E-04	0.3E-04
9	0.90	0.300006	0.311870	0.615008	0.2E-03	0.1E-03
10	1.00	0.264120	0.355660	0.615006	0.5E-04	0.4E-05
11	1.10	0.202230	0.404630	0.615270	0.2E-03	0.4E-04
12	1.20	0.177790	0.454410	0.615100	0.2E-03	0.1E-03
13	1.30	0.114940	0.499960	0.615420	0.8E-04	0.3E-03
14	1.40	0.084008	0.538900	0.614920	0.7E-04	0.2E-03
15	1.50	0.046270	0.569350	0.615560	0.7E-05	0.6E-03
16	1.60	0.023290	0.591250	0.615450	0.2E-04	0.4E-03
17	1.70	0.016780	0.600760	0.615580	0.4E-04	0.6E-03
18	1.80	0.019002	0.597540	0.616760	0.7E-04	0.2E-02
19	1.90	0.040630	0.581320	0.616120	0.2E-03	0.9E-03
20	2.00	0.062960	0.554850	0.617860	0.3E-03	0.3E-02
21	2.10	0.115640	0.517500	0.616740	0.5E-03	0.1E-02
22	2.20	0.139004	0.472006	0.617360	0.7E-03	0.2E-02
23	2.30	0.202920	0.423210	0.618130	0.1E-02	0.2E-02
24	2.40	0.236680	0.374770	0.618580	0.2E-02	0.1E-02
25	2.50	0.288560	0.331940	0.620560	0.3E-02	0.3E-02
26	2.60	0.339810	0.291006	0.619740	0.3E-02	0.2E-02
27	2.70	0.365410	0.253740	0.620270	0.4E-02	0.2E-02
28	2.80	0.413100	0.219004	0.619960	0.4E-02	0.1E-02
29	2.90	0.426240	0.189006	0.620330	0.5E-02	0.6E-03
30	3.00	0.456590	0.166420	0.621390	0.5E-02	0.1E-02
31	3.10	0.474300	0.148990	0.621350	0.6E-02	0.7E-03
32	3.20	0.483990	0.137650	0.621990	0.6E-02	0.1E-02
33	3.30	0.501800	0.129670	0.621490	0.6E-02	0.5E-03
34	3.40	0.495740	0.125300	0.621320	0.6E-02	0.1E-03
35	3.50	0.502590	0.125003	0.621470	0.6E-02	0.3E-03
36	3.60	0.492240	0.128390	0.621200	0.6E-02	0.2E-03
37	3.70	0.481850	0.138009	0.621990	0.7E-02	0.4E-03
38	3.80	0.475290	0.152560	0.621640	0.7E-02	0.2E-03
39	3.90	0.446560	0.172950	0.621640	0.6E-02	0.2E-03
40	4.00	0.436530	0.197680	0.621290	0.6E-02	0.2E-03
41	4.10	0.394950	0.225004	0.620000	0.6E-02	0.6E-03
42	4.20	0.363640	0.258250	0.620580	0.5E-02	0.2E-03
43	4.30	0.326880	0.295130	0.619840	0.5E-02	0.2E-03
44	4.40	0.271220	0.340047	0.620850	0.5E-02	0.9E-03
45	4.50	0.246580	0.389210	0.620300	0.4E-02	0.1E-02
46	4.60	0.179590	0.437750	0.619005	0.3E-02	0.8E-03
47	4.70	0.146460	0.484280	0.618630	0.2E-02	0.1E-02
48	4.80	0.094110	0.523660	0.617820	0.2E-02	0.1E-02
49	4.90	0.055550	0.560370	0.618920	0.1E-02	0.3E-02
50	5.00	0.035740	0.586940	0.618950	0.8E-03	0.3E-02

TABLE <3.7.2>

CHAPTER IVTHE RESOLUTION OF TURNING POINTS4.1 Solution by Nonuniform Expansion

As previously illustrated, the fundamental difficulty of the turning point approximation is the transition of frequencies from the fast scale ($\tau=t/\varepsilon$), which is treated asymptotically, to the $\tau=t/\sqrt{\varepsilon}$ scale, which is resolved by an $O(\sqrt{\varepsilon})$ time step. Here we assume that all turning points are of the simplest type.

ASSUMPTION [4.1a]

All relevant secondary frequencies are bounded away from zero except in the neighborhood of a turning point such as s , where for a vanishing frequency we have:

$$(4.1.1) \quad N^T \Lambda(t)/\varepsilon = a (t-s)/\varepsilon + O((t-s)^2/\varepsilon) \quad (a \neq 0).$$

In the neighborhood of a turning point we must extend the notion of oscillatory and nonoscillatory to accommodate this transition; thus, given the parameters $\{\varepsilon, K_1, K_2\}$, where

$$(4.1.2) \quad \begin{cases} 0 < \varepsilon \ll 1 \\ K_1 > 0 \\ K_2 > 1 \end{cases} ,$$

we classify a secondary frequency $(N^T \Lambda(t)/\varepsilon)$ in the neighborhood of a turning point by introducing the following subdivisions:

oscillatory range:

$$K_1/\varepsilon < |N^T \Lambda(t)/\varepsilon|$$

nonoscillatory range:

$$|N^T \Lambda(t)/\varepsilon| < K_2/\sqrt{\varepsilon}$$

(4.1.3)

quasi-oscillatory range:

$$K_2/\sqrt{\varepsilon} < |N^T \Lambda(t)/\varepsilon| < K_1/\varepsilon$$



In this chapter terms with frequencies in the oscillatory, nonoscillatory, and quasi-oscillatory ranges are denoted respectively by the subscripts I, II, and III. The intention here is to identify the strictly oscillatory terms, discussed in Chapter III, with the oscillatory range of frequencies since the same asymptotic principles are applied to each, and also we identify the strictly nonoscillatory terms of Chapter III with the nonoscillatory range of frequencies since these terms are completely resolved as smooth functions. The quasi-oscillatory range then gives the transition region between these two ranges. Since we need only apply our principle locally -- that is, from grid point to grid point -- we can assume, in correspondence to the decomposability principle of Chapter III, that the classification of relevant frequencies is uniform on the interval of interest.

ASSUMPTION [4.1b]

On the interval $\mathcal{J} = [T_1, T_2]$ all relevant secondary frequencies fall uniformly into the oscillatory, nonoscillatory, and quasi-oscillatory ranges, where the parameters (4.1.2) have been specified.

We accordingly reformulate the system (3.1.14):

$$(4.1.4) \quad \left\{ \begin{array}{l} X' = G(X, t, \varepsilon) \\ \quad = g_{\text{I}}(X, t, \varepsilon) + g_{\text{II}}(X, t, \varepsilon) + g_{\text{III}}(X, t, \varepsilon) + \\ \quad \quad f_{\text{I}}(t, \varepsilon) + f_{\text{II}}(t, \varepsilon) + f_{\text{III}}(t, \varepsilon) \\ X(T_1, \varepsilon) = \widehat{X}_0, \quad T_1 \leq t \leq T_2 \end{array} \right.$$

where we assume:

- (i) \widehat{X}_0 is independent of ε ;
- (ii) $g_{\text{I}}(X, t, \varepsilon)^{(i)} = \sum_j a_{ij}(t) \exp[B_{ij}(t)/\varepsilon] p_{ij}(X)$;
- (iii) $g_{\text{II}}(X, t, \varepsilon)^{(i)} = \sum_j \widetilde{a}_{ij}(t) \exp[\widetilde{B}_{ij}(t)/\varepsilon] \widetilde{p}_{ij}(X)$;
- (iv) $g_{\text{III}}(X, t, \varepsilon)^{(i)} = \sum_j \widetilde{\widetilde{a}}_{ij}(t) \exp[\widetilde{\widetilde{B}}_{ij}(t)/\varepsilon] \widetilde{\widetilde{p}}_{ij}(X)$;
- (v) $f_{\text{I}}(t, \varepsilon)^{(i)} = \sum_j c_{ij}(t) \exp[G_{ij}(t)/\varepsilon]$;
- (vi) $f_{\text{II}}(t, \varepsilon)^{(i)} = \sum_j \widetilde{c}_{ij}(t) \exp[\widetilde{G}_{ij}(t)/\varepsilon]$;
- (vii) $f_{\text{III}}(t, \varepsilon)^{(i)} = \sum_j \widetilde{\widetilde{c}}_{ij}(t) \exp[\widetilde{\widetilde{G}}_{ij}(t)/\varepsilon]$;
- (viii) All t -dependent functions are in $C^\infty(t)$; all X -dependent functions are polynomials in the components of X ;
- (ix) Assumptions [4.1a] and [4.1b] hold; terms with frequencies in the oscillatory, nonoscillatory, and quasi-oscillatory ranges are respectively given by the subscripts I, II, and III.

The extra smoothness conditions are for convenience since more

integrations are needed to resolve the effects of the oscillations in the quasi-oscillatory range. We also insist that the system (4.1.4) satisfy the appropriate modification of Assumption [3.3a]:

ASSUMPTION [4.1c]

In correspondence to the system (4.1.4), the reduced system

$$(4.1.5) \quad \begin{cases} V' = g_{\text{II}}(V, t, \varepsilon) + f_{\text{II}}(t, \varepsilon) \\ V(T_1, \varepsilon) = \hat{X}_0, \quad T_1 < t < T_2 \end{cases}$$

is well-posed and has a bounded solution in $C^\infty(t)$.

In Chapter III we introduced the linear operator (3.1.21), which essentially returned the result of one integration by parts of the indefinite integral of a rapidly oscillating function. In the neighborhood of a turning point the technique of integration by parts is still effective, but it is convenient to carry out the integration completely because the simple order relation (3.1.22) degenerates as the turning point is approached, and weak singularities appear in the expansion coefficients. For example we consider

$$(4.1.6) \quad \begin{cases} \int_{\alpha}^{\beta} a(t) \exp[it^2/2\varepsilon] dt \\ \sqrt{\varepsilon}K_2 < \alpha < \beta < K_1 \end{cases}$$

Integration by parts yields:

$$(4.1.7) \left\{ \begin{array}{l} \int_{\alpha}^{\beta} \exp[it^2/2\varepsilon] a(t) dt \\ = \int_{\alpha}^{\beta} \exp[it^2/2\varepsilon] \left(\sum_{k=1}^m \varepsilon^k b_k(t)/t^{2k-1} \right) \\ + R_m(\alpha, \beta) \end{array} \right.$$

where

$$(4.1.8) \left\{ \begin{array}{l} R_m(\alpha, \beta) = \varepsilon^m \int_{\alpha}^{\beta} \exp(it^2/2\varepsilon) \tilde{b}_m(t)/t^{2m} dt \\ \{b_k(t), \tilde{b}_m(t)\} \subset C^{\infty}(t) \end{array} \right.$$

From (4.1.8) we conclude:

$$(4.1.9) \left\{ \begin{array}{l} \varepsilon^k \int_{\alpha}^{\beta} \exp(it^2/2\varepsilon) b_k(t)/t^{2k-1} = O(\sqrt{\varepsilon}/K_2^{2k-1}) \\ |R_m(\alpha, \beta)| < \varepsilon^m M \int_{\alpha}^{\beta} 1/t^{2m} dt = O(\sqrt{\varepsilon}/K_2^{2m-1}) \end{array} \right.$$

where, by our smoothness assumption and (4.1.2), we can choose m so large that the error term is insignificant. Accordingly we define a linear operator which carries out this integration procedure on functions with quasi-oscillatory frequencies:

$$(4.1.10) \left\{ \begin{array}{l} \hat{\mathcal{B}}\{f_{\text{III}}\} = h_{\text{III}} \\ \sqrt{\varepsilon} h_{\text{III}}' = f_{\text{III}} \end{array} \right.$$

although in practice the approximation is made only to within $O(1/K_2^{2m-1})$. In like manner the antiderivative of a term with an oscillatory frequency is given by

$$(4.1.11) \quad \begin{cases} \hat{\mathcal{L}}\{f_{\mathbf{I}}\} = h_{\mathbf{I}} \\ \varepsilon h_{\mathbf{I}}' = f_{\mathbf{I}} \end{cases} .$$

And again in practice one only carries out this procedure to within a certain accuracy.

Analogously to (3.3.4) we say that \tilde{X} is an $\varepsilon^{m/2}$ -approximate solution of (4.1.4) if

$$(4.1.12) \quad \begin{cases} \tilde{X} = \sum_K (Q_{K/2}(t, \varepsilon) + W_{K/2}(t, \varepsilon) + Y_{K/2}(t, \varepsilon)) \varepsilon^K \\ M(\tilde{X}) = G(\tilde{X}, t, \varepsilon) - \tilde{X}' = O(\varepsilon^{\frac{m+1}{2}}) \\ \tilde{X}(T_1, \varepsilon) = \hat{X}_0 + O(\varepsilon^{\frac{m+1}{2}}) \end{cases}$$

where the frequencies of $W_{K/2}$, $Y_{K/2}$, and $Q_{K/2}$ are respectively nonoscillatory, oscillatory, and quasi-oscillatory. Each of these functions is bounded, but, because of the nonuniformity of the expansion, some may be $O(\sqrt{\varepsilon})$. The basic theory developed in Chapter III now can be extended. We present the appropriate analogues of Theorems [3.2.1], [3.3.1], [3.3.2], and [3.3.3].

THEOREM [4.1.1]

If Assumptions [4.1a] and [4.1b] are valid for the linear system

$$(4.1.13) \quad \begin{cases} Y' = [A_{\mathbf{I}} + A_{\mathbf{II}} + A_{\mathbf{III}}]Y + f_{\mathbf{I}} + f_{\mathbf{II}} + f_{\mathbf{III}} \\ Y(T_1, \varepsilon) = \hat{Y}_0, \quad T_1 < t < T_2 \end{cases}$$

then we have:

$$(4.1.14) \quad \begin{cases} \max_t |Y - \tilde{Y}| = O(\varepsilon^{\frac{1}{2}}) \\ \max_t |Y' - \tilde{Y}'| = O(\varepsilon^{\frac{1}{2}}) \end{cases}$$

where

$$(4.1.15) \quad \begin{cases} \tilde{Y} = V + \varepsilon \hat{\mathcal{L}}\{A_I V + f_I\} + \sqrt{\varepsilon} \hat{B}\{A_{III} V + f_{III}\} \\ \begin{cases} V' = A_{II} V + f_{II} \\ V(T_1) = \hat{Y}_0 \end{cases} \end{cases}$$

PROOF

After the change of variables

$$Z(t, \varepsilon) = X(t, \varepsilon) - V(t),$$

the equations for $Z(t, \varepsilon)$ are

$$\begin{cases} Z' = [A_I + A_{II} + A_{III}]Z + (f_I + A_I V) + (f_{III} + A_{III} V) \\ Z(T_1, \varepsilon) = 0, \quad T_1 < t < T_2 \end{cases}$$

The equations for

$$\tilde{Z} = Z - \varepsilon \hat{\mathcal{L}}\{f_I + A_I V\} - \sqrt{\varepsilon} \hat{B}\{f_{III} + A_{III} V\}$$

have the form

$$\begin{cases} \tilde{Z}' = [A_I + A_{II} + A_{III}]Z + \sqrt{\varepsilon} h(t, \varepsilon) \\ \tilde{Z}(T_1, \varepsilon) = \sqrt{\varepsilon} \hat{Z}_0, \quad T_1 < t < T_2 \end{cases}$$

and so by Lemma [3.2.2] we have:

$$\begin{cases} \max_{\mathcal{I}} |\tilde{Z}| = o(\varepsilon^{\frac{1}{2}}) \\ \max_{\mathcal{I}} |\tilde{Z}'| = o(\varepsilon^{\frac{1}{2}}) \end{cases}$$

THEOREM [4.1.1]

Let m be a nonnegative integer. If \tilde{X} is an $\varepsilon^{m/2}$ -approximate solution of the system (4.4.1), then

$$(4.1.16) \quad \max_{\mathcal{I}} |X(t, \varepsilon) - \tilde{X}(t, \varepsilon)| = o(\varepsilon^{\frac{m+1}{2}}).$$

Moreover, the existence of an \mathcal{C}^0 -approximate solution is guaranteed.

PROOF

The proof of this theorem is very similar to the proof of Theorem [3.3.2]; that is, by a Picard iteration one can show that if

$$\begin{cases} M(Z) = o(\varepsilon^{\frac{m+1}{2}}) \\ Z(T_1, \varepsilon) = \hat{X}_0 + o(\varepsilon^{\frac{m+1}{2}}) \end{cases},$$

then

$$\max_{\mathcal{I}} |Z - X| = o(\varepsilon^{\frac{m+1}{2}}),$$

where X is the solution of (4.1.4). Using Assumption [4.1c], one can demonstrate the existence of an \mathcal{C}^0 -approximate solution in the manner of Theorem [3.3.1].

THEOREM [4.1.3]

Let \tilde{X} be an $\varepsilon^{\frac{m}{2}}$ -approximate solution of the system (4.1.4), and assume that Assumptions [4.1a] and [4.1b] can be applied to the $O(\varepsilon^{\frac{m+1}{2}})$ terms of $M(\tilde{X})$ to give:

$$(4.1.17) \quad M(\tilde{X}) = \varepsilon^{\frac{m+1}{2}} [f_{\text{I}} + f_{\text{II}} + f_{\text{III}}] + O(\varepsilon^{\frac{m+2}{2}}).$$

Then we have an $\varepsilon^{\frac{m+1}{2}}$ -approximate solution of (4.1.4):

$$(4.1.18) \quad \tilde{\tilde{X}} = \tilde{X} + \varepsilon^{\frac{m+1}{2}} R.$$

Here $R(t, \varepsilon)$ is determined by

$$(4.1.19) \quad \begin{cases} R = \tilde{V} + \varepsilon \hat{\mathcal{L}}\{f_{\text{I}} + A_{\text{I}}V\} + \sqrt{\varepsilon} \hat{\mathcal{B}}\{f_{\text{III}} + A_{\text{III}}V\} \\ \tilde{V}' = A_{\text{II}}\tilde{V} + f_{\text{II}} \\ \tilde{V}(T_1) = (\hat{X}_0 - \tilde{X}(T_1, \varepsilon)) \end{cases}$$

where

$$(4.1.20) \quad \begin{cases} A_{\text{I}} = g_{\text{I}}[X, t, \varepsilon]_X \big|_{X=V} \\ A_{\text{II}} = g_{\text{II}}[X, t, \varepsilon]_X \big|_{X=V} \\ A_{\text{III}} = g_{\text{III}}[X, t, \varepsilon]_X \big|_{X=V} \end{cases}.$$

Here $V(t)$ is the approximation guaranteed by Assumption [4.1c].

PROOF

The proof of this theorem is very similar to the proof of Theorem [3.3.3], where the decomposability principle replaces Assumption [4.1b].

As in Chapter III one also can develop a formal expansion theory based of the Taylor expansions of polynomials. Here,

however, there are three time scales to consider.

In practice one must solve the equations for the nonoscillatory behavior approximatively by standard numerical techniques, and so the step size control is crucial for any implementation. As illustrated by (4.1.7), the integration procedure leads to the introduction of mild singularities which must be resolved as smooth functions when they arise in the differential equations for the nonoscillatory terms. Again it is sufficient to consider the case of a turning point at $t=0$; by the analysis of (4.1.7), the induced singularities must be of the form

$$(4.1.21) \quad \begin{cases} \varepsilon^{K-\frac{1}{2}}/t^{2K-1} \\ K_2\sqrt{\varepsilon} < |t| < K, \end{cases}$$

Using the stretching factor

$$(4.1.22) \quad S(t) = |it/\varepsilon|/(K_1/\varepsilon) = |t|/K_1,$$

we locally can change the independent variable by

$$(4.1.23) \quad \begin{cases} t = t^*(1 + Ms) \\ M = t^*/(|t^*|K_1) \\ t = t^* \text{ at } s=0 \end{cases}$$

to give

$$(4.1.24) \quad \varepsilon^{K-\frac{1}{2}}/t^{2K-1} = [\varepsilon^{K-\frac{1}{2}}/(t^*)^{2K-1}] \frac{1}{(1 + Ms)^{2K-1}},$$

and for

$$(4.1.25) \quad K_1 = 1$$

the expression (4.1.24) is a smooth function of s in a neighborhood of

$$(4.1.26) \quad s = 0.$$

Thus for a given secondary frequency

$$(4.1.27) \quad N_k \Lambda(t)/\varepsilon$$

in the quasi-oscillatory range, we define the stretching factor $S_k(t)$ by

$$(4.1.28) \quad \begin{cases} S_k(t) = |(N_k^T \Lambda(t)/\varepsilon)/(K_1/\varepsilon)| \\ \\ = |N_k^T \Lambda(t)/K_1| \end{cases}$$

As the frequency passes through the quasi-oscillatory range, we have

$$(4.1.29) \quad \begin{cases} S_k(t) = 1 & \text{for } |N_k^T \Lambda(t)/\varepsilon| = K_1/\varepsilon \\ \\ S_k(t) \rightarrow (K_2/K_1)\sqrt{\varepsilon} & \text{as } |N_k^T \Lambda(t)/\varepsilon| \rightarrow K_2/\sqrt{\varepsilon} \end{cases}$$

If more than one quasi-oscillatory frequencies are present, we define the uniform stretching factor by

$$(4.1.30) \quad \tilde{S}(t) = \min_k |S_k(t)|,$$

where the minimization is taken over all relevant quasi-oscillatory frequencies.

If the step size

$$(4.1.31) \quad \bar{h} = h \bar{S}(t)$$

is used, the induced singularities are resolved adequately by the analysis of (4.1.24); moreover, as secondary frequencies pass from the oscillatory range through the quasi-oscillatory range the step size is reduced to $O(\sqrt{\epsilon}h)$. Since frequencies in the nonoscillatory range must be resolved by an $O(\sqrt{\epsilon}h)$ time step, the resolution can be made without an unduly abrupt adjustment in the mesh width. The appropriate step size for the problem then depends primarily on the smallest quasi-oscillatory frequency and the largest nonoscillatory frequency.

4.2 The Cost of Resolution / Turning Points of Higher Order

The cost of resolution can be measured by the number of points used. Consider a term with a secondary frequency satisfying Assumption (4.1a) at the turning point $t=0$:

$$(4.2.1) \quad \begin{cases} N^T \Delta(t)/\varepsilon = iat/\varepsilon + O(t^2/\varepsilon) \\ a \in \mathbb{R} \setminus \{0\} \end{cases} .$$

If we ignore the $O(t^2/\varepsilon)$ correction, the regions for the frequencies are:

$$(4.2.2) \quad \begin{cases} \text{oscillatory:} & K_1/|a| < |t| \\ \text{quasi-oscillatory:} & K_1 K_2 \sqrt{\varepsilon} / |a| < |t| < K_1/|a| . \\ \text{nonoscillatory:} & |t| < K_1 K_2 \sqrt{\varepsilon} / |a| \end{cases}$$

If we assume that all other components of the system are smooth functions of t and that h is the step size required to resolve these functions, then $O(1/h)$ points will be needed for the oscillatory range. For the quasi-oscillatory range the local density of points is

$$(4.2.3) \quad 1/(hS(t)),$$

where the stretching factor $S(t)$ is given by

$$(4.2.4) \quad S(t) = |at|/K_1 .$$

Thus, the number of points required is

$$(4.2.5) \quad \left\{ \begin{array}{l} 2 \int_{\alpha}^{\beta} K_1 / (h|a|t) dt = O(|\log(\varepsilon)|/h) \\ (\alpha = K_1/|a|, \beta = K_1 K_2 \sqrt{\varepsilon'}/|a|) \end{array} \right.$$

For the nonoscillatory region the local density of points is $O(1/(\sqrt{\varepsilon}h))$, and so the number of points needed is

$$(4.2.6) \quad O(1/(\sqrt{\varepsilon}h)) O(\sqrt{\varepsilon'}) = O(1/h).$$

Therefore, the cost of resolution is:

$$(4.2.7) \quad O(|\log(\varepsilon)|/h)$$

The restriction on the order of the turning point is superfluous. For example, let a secondary frequency with turning point $t=0$ be given by

$$(4.2.8) \quad \left\{ \begin{array}{l} N^T \Delta(t)/\varepsilon = i a t^{p+1} / ((p+1)\varepsilon) \\ a \in \mathbb{R} \setminus \{0\}, \quad p \in \mathbb{N} \setminus \{0,1\} \end{array} \right.$$

One could develop a theory based on an expansion in powers of $(\varepsilon^{\frac{1}{p+1}})$ as we did for the restricted case $p=1$. In practice though one would like to compute with a consistent method for all turning points, and so we apply the technique of Section 4.1 to (4.2.8). The stretching factor is then

$$(4.2.9) \quad S(t) = |a t^p| / K_1$$

and the frequency ranges are:

$$(4.2.14) \quad O(1/(h \varepsilon^{\frac{p-1}{2p}})).$$

Thus, the penalty for using our procedure to resolve turning points of higher order is not a loss of accuracy but rather a modest increase in the number of points required.

4.3 A Computational Example

In this section we describe a possible implementation of the turning point procedure. We consider the system

$$(4.3.1) \quad \begin{cases} y' = \exp(it^2/2\varepsilon) y^2 \\ y(-2)=1, \quad -2 < t < 2, \quad \varepsilon = .01 \end{cases}$$

which one might have derived from the system

$$(4.3.2) \quad \begin{cases} z' = (it/\varepsilon)z + z^2 \\ z(-2)=\exp(2i/\varepsilon), \quad -2 < t < 2, \quad \varepsilon = .01 \end{cases}$$

after the usual change of variables. We propose to solve equation (4.3.1) to within three or four digits of accuracy.

The fundamental frequency of the system is given by

$$(4.3.3) \quad \{it/\varepsilon\} .$$

For the stiffness parameters (4.1.2) we choose:

$$(4.3.4) \quad \begin{cases} K_1 = 1.0 \\ K_2 = 5.0 \end{cases}$$

and thus, by (4.1.3), the ranges for the fundamental frequency are:

$$(4.3.5) \quad \begin{cases} \text{(a) I oscillatory range: } 1.0 < |t| \\ \text{(b) II nonoscillatory range: } |t| < 0.5 \\ \text{(c) III quasi-oscillatory range: } 0.5 < |t| < 1.0 \end{cases} .$$

When (it/ϵ) is in the nonoscillatory range, all components of the system are resolved as smooth functions, and no asymptotic analysis is needed. When the fundamental frequency is in the oscillatory or quasi-oscillatory range, then our asymptotic techniques introduce other relevant secondary frequencies of the form:

$$(4.3.6) \quad 2it/\epsilon, 3it/\epsilon, \dots$$

but by our parameter valuations these must be in the oscillatory range.

In a given subinterval where the classification is uniform, the problem is of the form

$$(4.3.7) \quad \begin{cases} y' = \exp(it^2/2\epsilon) y^2 \\ y(t_0) = y_0, \quad t_0 < t < t_1 \end{cases}$$

By the analysis of Chapter III we have for the case (4.3.5a):

$$(4.3.8) \quad \begin{cases} y \sim y_0 + \epsilon (-iy_0^2/t \exp(it^2/2\epsilon) + w_1) + O(\epsilon^2) \\ w = iy_0/t \exp(it^2/2\epsilon) \Big|_{t=t_0} \end{cases}$$

and in the case (4.3.5b) all coefficients are resolved as smooth functions. For (4.3.5c), however, the analysis is more complicated. As in (4.3.8) we determine the expansion to within an error of $O(\epsilon^2)$, and also we need only calculate the integration operators ((4.1.10)-(4.1.11)) to within the same accuracy; thus, we ignore terms of the size $O(\epsilon^2)$, $O(\epsilon^{1/2}/K_2^5)$, or

$O(\varepsilon / K_2^3)$.

The simplest way to generate the asymptotic expansion is to introduce the time scales

$$(4.3.9) \quad \left\{ \begin{array}{l} \tau = t/\varepsilon^{1/2} \\ \zeta = t/\varepsilon \\ d/dt = \frac{\partial}{\partial t} + \frac{1}{\varepsilon} \frac{\partial}{\partial \zeta} + \frac{1}{\varepsilon^{1/2}} \frac{\partial}{\partial \tau} \end{array} \right.$$

and the formal expansion

$$(4.3.10) \quad \left\{ \begin{array}{l} Y \sim \sum_k X_{k/2} \varepsilon^{k/2} \\ X_{k/2} = W_{k/2}(t, \varepsilon) + \tilde{Y}_{k/2}(\zeta, \varepsilon) + \tilde{Q}_{k/2}(\tau, \varepsilon) \\ \tilde{Y}_{k/2}(\zeta, \varepsilon) = \tilde{Y}_{k/2}(t/\varepsilon, \varepsilon) = Y_{k/2}(t, \varepsilon) \\ \tilde{Q}_{k/2}(\tau, \varepsilon) = \tilde{Q}_{k/2}(t/\varepsilon, \varepsilon) = Q_{k/2}(t, \varepsilon) \\ Y_0(t, \varepsilon) \equiv \emptyset, Q_0(t, \varepsilon) \equiv \emptyset, Y_{1/2}(t, \varepsilon) \equiv \emptyset \end{array} \right. .$$

Then when the $O(\varepsilon^{k/2})$ terms are to be balanced, we have on the left-hand side

$$(4.3.11) \quad d/dt W_{k/2}(t, \varepsilon) + d/d\zeta \tilde{Y}_{\frac{k+2}{2}}(\zeta, \varepsilon) + d/d\tau Q_{\frac{k+1}{2}}(\tau, \varepsilon),$$

and the right-hand side will depend on

$$(4.3.12) \quad \{W_{j/2}, Y_{j/2}, Q_{j/2} : j < k\} ;$$

moreover, the dependence on the $(j=k)$ -values will be linear for $k > 1$. After the right-hand side is decomposed according to Assumption [4.1b], the t -dependence is restored in the left-hand side. The terms corresponding to each frequency range are then matched to guarantee the balance on this level. We have:

$$(4.3.13) \left\{ \begin{array}{l} (a) \ \varepsilon^0: W_{0,t} + Q_{1/2,T} + Y_{1,T} = (X_0)^2 \exp(it^2/2\varepsilon) \\ (b) \ \varepsilon^{1/2}: W_{1/2,t} + Q_{1,T} + Y_{3/2,T} = 2X_0 X_{1/2} \exp(it^2/2\varepsilon) \\ (c) \ \varepsilon: W_{1,t} + Q_{3/2,T} + Y_{2,T} = (X_{1/2}^2 + 2X_0 X_1) \exp(it^2/2\varepsilon) \\ (d) \ \varepsilon^{3/2}: W_{3/2,t} + Q_{2,T} + Y_{5/2,T} = (2X_{3/2} X_0 + 2X_{1/2} X_1) \exp(it^2/2\varepsilon) \end{array} \right.$$

By following the procedure outlined above we have from (4.3.13a):

$$(4.3.14) \left\{ \begin{array}{l} \left\{ \begin{array}{l} W_{0,t} = 0 \\ W_0(0, \varepsilon) = Y_0 \end{array} \right. \Rightarrow W_0 \equiv w_0 = Y_0 \\ \varepsilon Y_{1,t} = 0 \Rightarrow Y_1 = \hat{\mathcal{L}}\{0\} = 0 \\ \varepsilon^{1/2} Q_{1/2,t} = w_0^2 \exp(it^2/2\varepsilon) \\ \Rightarrow \left\{ \begin{array}{l} Q_{1/2} = \hat{\mathcal{B}}\{w_0^2 \exp(it^2/2\varepsilon)\} \\ = q_{1/2} \exp(it^2/2\varepsilon) \\ q_{1/2}(t, \varepsilon) = w_0^2 [\varepsilon^{1/2}/it - \varepsilon^{3/2}/t^3] + O(1/K_2^5) \end{array} \right. \end{array} \right.$$

From (4.3.13b) we have:

$$(4.3.15) \left\{ \begin{array}{l} \left\{ \begin{array}{l} W_{1/2,t} = 0 \\ W_{1/2}(0, \varepsilon) = -Q_{1/2}(0, \varepsilon) \end{array} \right. \Rightarrow W_{1/2} \equiv w_{1/2} = -Q_{1/2}(0, \varepsilon) \\ \varepsilon Y_{3/2,t} = 2w_0 q_{1/2} \exp(it^2/\varepsilon) \\ \Rightarrow \left\{ \begin{array}{l} Y_{3/2} = \hat{\mathcal{L}}\{2w_0 q_{1/2} \exp(it^2/\varepsilon)\} \\ = (-w_0^3) [\varepsilon^{1/2}/t^2 + \varepsilon^{3/2}/it^4] \exp(it^2/\varepsilon) \\ + O(1/K_2^3) \end{array} \right. \\ \varepsilon^{1/2} Q_{1,t} = 2w_0 w_{1/2} \exp(it^2/2\varepsilon) \\ \Rightarrow \left\{ \begin{array}{l} Q_1 = \hat{\mathcal{B}}\{2w_0 w_{1/2} \exp(it^2/2\varepsilon)\} \\ = q_1 \exp(it^2/2\varepsilon) \\ q_1 = 2w_0 w_{1/2} (\varepsilon^{1/2}/it) + O(1/K_2^3) \end{array} \right. \end{array} \right.$$

From (4.3.13c) we have:

$$(4.3.16) \left\{ \begin{array}{l} W_{1,t} = 0 \\ W_1(\theta, \varepsilon) = -Q_1(\theta, \varepsilon) \\ \varepsilon^{1/2} Q_{3/2,t} = \exp(it^2/2\varepsilon) [w_{1/2}^2 + 2w_0 w_1] \\ \Rightarrow \begin{cases} Q_{3/2} = \hat{B} \{ [w_{1/2}^2 + 2w_0 w_1] \exp(it^2/2\varepsilon) \} \\ = q_{3/2} \exp(it^2/2\varepsilon) \\ q_{3/2} = [w_{1/2}^2 + 2w_0 w_1] (\varepsilon^{1/2}/it) + O(1/K_2^3) \end{cases} \end{array} \right. \Rightarrow W_1 \equiv w_1 = -Q_1(\theta, \varepsilon)$$

And from (4.3.13d) we have:

$$(4.3.17) \left\{ \begin{array}{l} W_{3/2,t} = 0 \\ W_{3/2}(\theta, \varepsilon) = -Q_{3/2}(\theta, \varepsilon) - Y_{3/2}(\theta, \varepsilon) \\ \Rightarrow W_{3/2} \equiv w_{3/2} = -Q_{3/2}(\theta, \varepsilon) - Y_{3/2}(\theta, \varepsilon) \end{array} \right.$$

For (4.3.5a) and (4.3.5b) all smooth components are constants, and so we do not really require any points; however, we mark grid points as might be required for the resolution of a smooth function. For the cases (4.3.5a) and (4.3.5c) we respectively use the step sizes

$$(4.3.18) \quad h = .1$$

and

$$(4.3.19) \quad h = .05$$

Thus, if a fourth-order Runge-Kutta scheme were used to resolve smooth components in (4.3.5a), the error would be comparable to the truncation error from the asymptotic expansion. For (4.3.5c) we have from (4.1.28) the stretching factor

$$(4.3.20) \quad \begin{cases} S(t) = |t| \\ 0.5 < S(t) < 1.0 \end{cases} .$$

In practice one might adopt a strategy which leads to adjusting the mesh only by a factor of (1/2) or (2); thus, the step size (4.3.19) is reasonable.

For (4.3.5b) we fully resolve the coefficients of the equation by means of a fourth-order Runge-Kutta solver with a time step sufficient so that the the estimated local truncation error is bounded above by

$$(4.3.21) \quad 10^{-3} h ,$$

and so the global error is bounded by

$$(4.3.22) \quad (10^{-3} h) O(\sqrt{N}/h) = 10^{-3} O(\sqrt{N}) .$$

Thus, depending on how well or how poorly the bound (4.3.21) is achieved, we accordingly adjust the step size by a factor of (1/2) or (2). One can estimate the local truncation error by comparing the results of two increments with the results of one double-increment over the same interval (see Lambert [20], p. 130).

All calculations were done with single-precision accuracy on a VAX11/780 computer. Plots of

$$(4.3.23) \quad u = \text{Re}\{y\} - .8$$

and

$$(4.3.24) \quad v = \text{Im}\{y\}$$

illustrate the passage through resonance in figure <4.3.1>, where we have linearly interpolated the amplitudes and phases between grid points. In table <4.3.1> we compare the computed grid values with the accepted function values, which were computed with double-precision accuracy by means of a fourth-order Runge-Kutta scheme with a time step

$$(4.3.25) \quad h = 10^{-4} .$$

In the table we list

$$(4.3.26) \quad \text{ERR}(t) = \max [\text{Re}\{y-\tilde{y}\}, \text{Im}\{y-\tilde{y}\}] ,$$

where $y(t)$ and $\tilde{y}(t)$ are respectively the computed and the accepted function values.

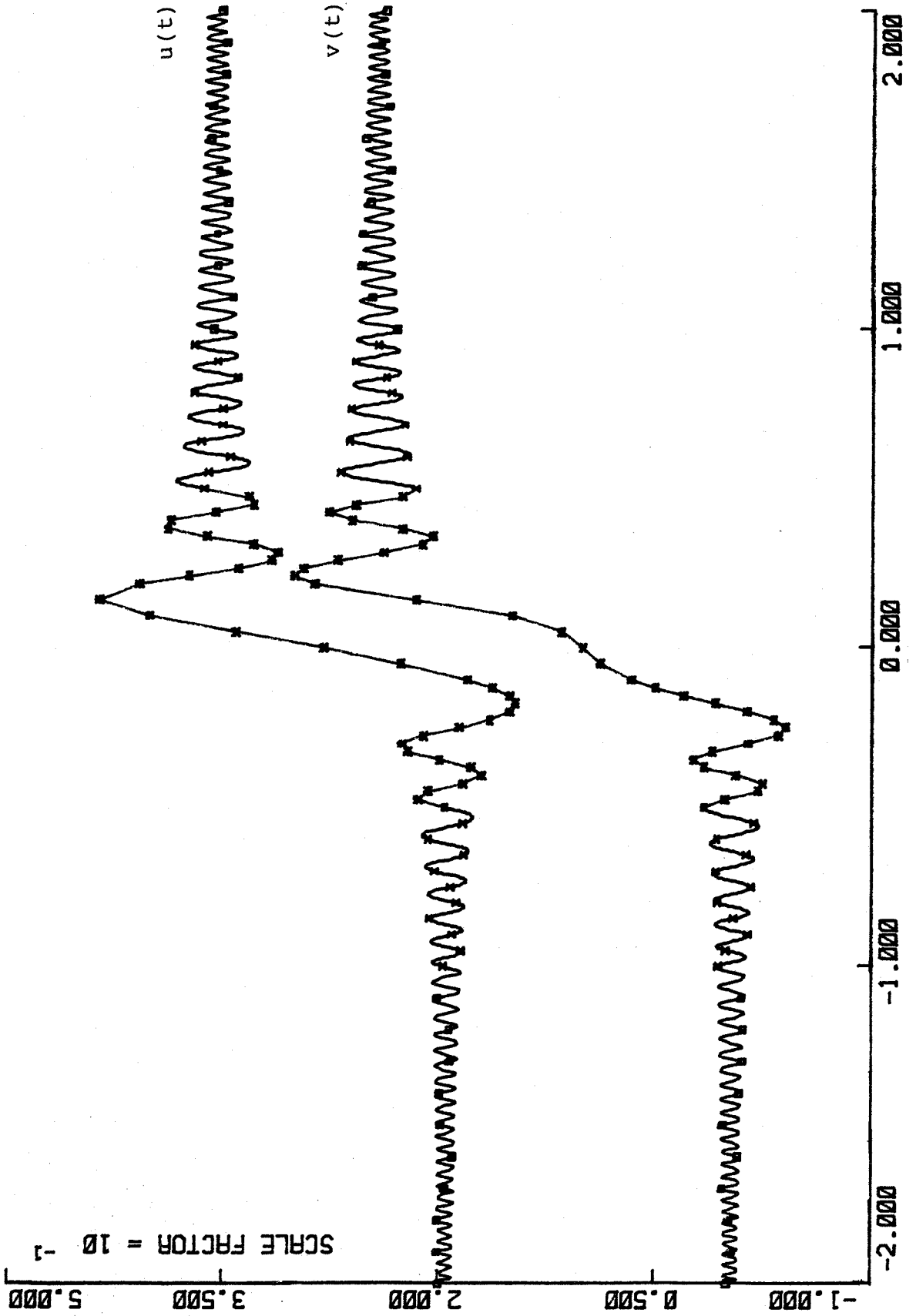


figure <4.3.1>

RANGE	h	N	T_N	$\text{Re}\{y(T_N)\}$	$\text{Im}\{y(T_N)\}$	$\text{ERR}(T_N)$
I	0.100	0	-2.0000	1.00000	0.00000	0.0E+00
I	0.100	1	-1.9000	1.00084	-0.00318	0.2E-04
I	0.100	2	-1.8000	1.00107	-0.00129	0.9E-05
I	0.100	3	-1.7000	0.99571	0.00345	0.2E-04
I	0.100	4	-1.6000	0.99113	-0.00677	0.1E-03
I	0.100	5	-1.5000	0.99938	0.00308	0.1E-04
I	0.100	6	-1.4000	0.99973	-0.00829	0.1E-03
I	0.100	7	-1.3000	0.99319	-0.00973	0.2E-03
I	0.100	8	-1.2000	0.99352	-0.01050	0.2E-03
I	0.100	9	-1.1000	1.00222	-0.00871	0.1E-03
I	0.100	10	-1.0000	0.99826	0.00721	0.4E-04
III	0.050	11	-0.9500	0.98625	0.00204	0.4E-04
III	0.050	12	-0.9000	0.99165	-0.01260	0.4E-04
III	0.050	13	-0.8500	1.00740	-0.00262	0.5E-04
III	0.050	14	-0.8000	0.98892	0.00797	0.5E-04
III	0.050	15	-0.7500	0.99317	-0.01534	0.5E-04
III	0.050	16	-0.7000	1.00420	0.00904	0.6E-04
III	0.050	17	-0.6500	0.98374	-0.01168	0.7E-04
III	0.050	18	-0.6000	1.00840	0.00842	0.7E-04
III	0.050	19	-0.5500	0.98498	-0.01673	0.1E-03
III	0.050	20	-0.5000	0.99741	0.01743	0.1E-03
II	0.025	21	-0.4750	1.01621	0.00280	0.1E-03
II	0.025	22	-0.4500	1.00870	-0.02025	0.1E-03
II	0.025	23	-0.4250	0.98518	-0.02272	0.1E-03
II	0.025	24	-0.4000	0.97171	-0.00420	0.1E-03
II	0.025	25	-0.3750	0.97913	0.01744	0.1E-03
II	0.025	26	-0.3500	1.00153	0.02518	0.1E-03
II	0.025	27	-0.3250	1.02276	0.01222	0.1E-03
II	0.025	28	-0.3000	1.02715	-0.01301	0.1E-03
II	0.025	29	-0.2750	1.01229	-0.03369	0.1E-03
II	0.025	30	-0.2500	0.98836	-0.03937	0.1E-03
II	0.025	31	-0.2250	0.96656	-0.03041	0.1E-03
II	0.025	32	-0.2000	0.95287	-0.01222	0.1E-03
II	0.025	33	-0.1750	0.94870	0.00979	0.1E-03
II	0.025	34	-0.1500	0.95326	0.03179	0.1E-03
II	0.025	35	-0.1250	0.96497	0.05150	0.1E-03
II	0.025	36	-0.1000	0.98220	0.06779	0.1E-03
II	0.050	37	-0.0500	1.02771	0.08963	0.1E-03
II	0.050	38	0.0000	1.08236	0.10200	0.1E-03

table <4.3.1a>

RANGE	h	N	T_N	$\text{Re}\{y(T_N)\}$	$\text{Im}\{y(T_N)\}$	$\text{ERR}(T_N)$
II	0.050	39	0.0500	1.14302	0.11667	0.1E-03
II	0.050	40	0.1000	1.20317	0.15097	0.1E-03
II	0.050	41	0.1500	1.23784	0.21743	0.1E-03
II	0.050	42	0.2000	1.20994	0.28779	0.1E-03
II	0.025	43	0.2250	1.17553	0.30218	0.1E-03
II	0.025	44	0.2500	1.14090	0.29565	0.1E-03
II	0.025	45	0.2750	1.11767	0.27181	0.1E-03
II	0.025	46	0.3000	1.11363	0.24000	0.1E-03
II	0.025	47	0.3250	1.13092	0.21313	0.1E-03
II	0.025	48	0.3500	1.16308	0.20580	0.1E-03
II	0.025	49	0.3750	1.19055	0.22681	0.1E-03
II	0.025	50	0.4000	1.18836	0.26220	0.1E-03
II	0.025	51	0.4250	1.15735	0.27764	0.1E-03
II	0.025	52	0.4500	1.13042	0.25888	0.1E-03
II	0.025	53	0.4750	1.13437	0.22740	0.1E-03
II	0.025	54	0.5000	1.16518	0.21792	0.1E-03
III	0.050	55	0.5500	1.16186	0.26997	0.2E-03
III	0.050	56	0.6000	1.14675	0.22434	0.2E-03
III	0.050	57	0.6500	1.16692	0.26425	0.1E-03
III	0.050	58	0.7000	1.15234	0.22559	0.2E-03
III	0.050	59	0.7500	1.15238	0.26273	0.2E-03
III	0.050	60	0.8000	1.17186	0.23485	0.2E-03
III	0.050	61	0.8500	1.14233	0.23848	0.1E-03
III	0.050	62	0.9000	1.15608	0.26028	0.1E-03
III	0.050	63	0.9500	1.17215	0.24437	0.1E-03
III	0.050	64	1.0000	1.15907	0.23093	0.2E-03
I	0.100	65	1.1000	1.14544	0.24934	0.5E-03
I	0.100	66	1.2000	1.15570	0.25631	0.3E-03
I	0.100	67	1.3000	1.15654	0.25550	0.3E-03
I	0.100	68	1.4000	1.14887	0.25014	0.4E-03
I	0.100	69	1.5000	1.15540	0.23566	0.2E-03
I	0.100	70	1.6000	1.16079	0.25278	0.2E-03
I	0.100	71	1.7000	1.16034	0.23715	0.2E-03
I	0.100	72	1.8000	1.15089	0.24039	0.3E-03
I	0.100	73	1.9000	1.15017	0.24295	0.3E-03
I	0.100	74	2.0000	1.15296	0.23930	0.2E-03

table <4.3.1b>

REFERENCES

- [1] V. Amdursky and A. Ziv, On the numerical treatment of stiff highly-oscillatory systems, IBM Israel Scientific Center Technical Report No. 15, Haifa, 1974.
- [2] _____, The numerical treatment of linear highly oscillatory O.D.E. systems by reduction to non-oscillatory type, IBM Israel Scientific Center Report No. 39, Haifa, 1976.
- [3] _____, On the numerical solution of stiff linear systems of the oscillatory type, SIAM J. Appl. Math., 33 (1977), pp. 593-606.
- [4] N. N. Bogoliubov and Y. A. Mitropolsky, Asymptotic Methods in the Theory of Nonlinear Oscillations, Gordon and Breach Science Publishers Inc., New York, 1961.
- [5] G. Browning and H.-O. Kreiss, Problems with different time scales for partial differential equations, SIAM J. Appl. Math., to appear.
- [6] E. A. Coddington and N. Levinson, Theory of Ordinary Differential Equations, McGraw-Hill Book Co., New York, 1955.
- [7] G. Contopoulos, A third integral of motion in a galaxy, Zeitschrift für Astrophysik, 49 (1960), pp. 273-291.
- [8] S. O. Fatunla, Numerical integrators for stiff and highly oscillatory differential equations, Math. Comp., 34 (1980), pp. 373-390.
- [9] J. N. Franklin, Matrix Theory, Prentice-Hall, Inc., Englewood Cliffs, New Jersey, 1968.
- [10] C.-E. Fröberg, Introduction to Numerical Analysis, 2nd edition, Addison-Wesley Publishing Co., Reading, Massachusetts, 1969.
- [11] W. Gautschi, Numerical integration of ordinary differential equations based on trigonometric polynomials, Numer. Math., 3 (1961), pp. 381-397.
- [12] O. F. Graff and D. G. Bettis, Modified multirevolution integration methods for satellite orbit computation, Celestial Mechanics, 11 (1975), pp. 433-448.
- [13] F. C. Hoppensteadt and W. L. Miranker, Differential equations having rapidly changing solutions: analytic methods for weakly nonlinear systems, J. Diff. Equ., 22 (1976), pp. 237-249.

- [14] G. Hori, Non-linear coupling of two harmonic oscillators, Publications of the Astronomical Society of Japan, 19 (1967), pp. 229-241.
- [15] T. Kato, Perturbation Theory for Linear Operators, Springer-Verlag New York Inc., 1966.
- [16] J. Kevorkian and J. D. Cole, Perturbation Methods in Applied Mathematics, Springer-Verlag New York Inc., 1981.
- [17] _____, Uniformly valid asymptotic approximations for certain non-linear differential equations, in Nonlinear Differential Equations and Nonlinear Mechanics (J. P. LaSalle, Editor), Academic Press Inc., New York, 1963.
- [18] H.-O. Kreiss, Difference methods for stiff ordinary differential equations, SIAM J. Numer. Anal., 15 (1978), pp. 21-58.
- [19] _____, Problems with different time scales for ordinary differential equations, SIAM J. Numer. Anal., 16 (1979), pp. 980-998.
- [20] J. D. Lambert, Computational Methods in Ordinary Differential Equations, John Wiley and Sons Ltd., Chichester, England, 1973.
- [21] B. Lindberg, On smoothing and extrapolation for the trapezoidal rule, BIT, 11 (1971), pp. 29-52.
- [22] G. Majda, Filtering techniques for oscillatory stiff O.D.E.'s, SIAM J. Numer. Anal., to appear.
- [23] W. L. Miranker, Numerical Methods for Stiff Equations and Singular Perturbation Problems, D. Reidel Publishing Co., Dordrecht, Holland, 1981.
- [24] W. L. Miranker and F. Hoppensteadt, Numerical methods for stiff systems of differential equations related with transistors, tunnel diodes, etc., Lecture Notes in Computer Science 10, Springer-Verlag New York Inc., 1974.
- [25] W. L. Miranker and M. van Veldhuizen, The method of envelopes, Math. Comp., 32 (1978), pp. 453-498.
- [26] W. L. Miranker and G. Wabba, An averaging method for the stiff highly oscillatory problem, Math. Comp., 30 (1976), pp. 383-399.
- [27] A. Nadeau, J. Guyard, and M. R. Feix, Algebraic-numerical method for the slightly perturbed harmonic oscillator, Math. Comp., 28 (1974), pp. 1057-1066.

- [28] J. C. Neu, The method of near-identity transformations and its applications, SIAM J. of Appl. Math., 38 (1980), pp. 189-200.
- [29] A. H. Nayfeh, Perturbation Methods, John Wiley and Sons Ltd., New York, 1973.
- [30] L. R. Petzold, An efficient numerical method for highly oscillatory ordinary differential equations, SIAM J. Numer. Anal., 18 (1981), pp. 455-479.
- [31] A. D. Snyder and G. C. Fleming, Approximation by aliasing with applications to "Certain" stiff differential equations, Math. Comp., 28 (1974), pp. 465-473.
- [32] C. E. Velez, Numerical integration of orbits in multirevolution steps, NASA Technical Note D-5915, Goddard Space Flight Center, Greenbelt, Maryland, 1970.