# COLLOCATION FOR SINGULAR PERTURBATION PROBLEMS II: LINEAR FIRST ORDER SYSTEMS WITHOUT TURNING POINTS

by U. Ascher\* and R. Weiss\*\* Technical Report 82-4 May 1982

### Abstract

We consider singularly perturbed linear boundary value problems for ODES, with variable coefficients, but without turning points. Convergence results are obtained for collocation schemes based on Gauss and Lobatto points, showing that highly accurate numerical solutions for these problems can be obtained at a very reasonable cost using such schemes, provided that appropriate meshes are used. The implementation of the numerical schemes and the practical construction of corresponding meshes are discussed.

These results extend those of a previous paper which deals with systems with constant coefficients.

Subject classification: 65L10

<sup>\*</sup>Department of Computer Science, The University of British Columbia, Vancouver, Canada V6T 1W5. The research of this author was supported in part under NSERC grant A4306.

<sup>\*\*</sup>Institut fur Angewandte und Numerische Mathematik, Technische Universitat Wien, 1040 Wien, Gusshausstrasse 27-29, Austria.

#### 1. Introduction

In part I of this work [2], (hereinafter referred to as "Part I") we have considered the numerical solution of singularly perturbed boundary value ordinary differential equations with constant coefficients. Our attention was focused on symmetric collocation schemes, which include the midpoint (or box) and the trapezoidal difference schemes as special cases. We have shown that such schemes can be used to compute highly accurate numerical solutions at a very reasonable cost, provided that appropriate meshes are used. Such a mesh consists, in general, of three parts: Two fine grids near the boundaries, to cover the possible two layer regions, and a coarser grid in between.

Similar results for the variable coefficient case are obtained in Weiss [9] for the trapezoidal and midpoint schemes. The eigenvalues of the "fast component" part of the differential equations are assumed to stay away from the imaginary axis for all values of the independent variable. Thus, turning points are excluded from the discussion. In the passage from constant to variable coefficients, the analysis had to be extended significantly.

In this paper we extend the results of the two papers mentioned above to include convergence results for the collocation schemes based on Gauss and Lobatto points for linear two-point boundary value problems which have a uniformly bounded inverse and which are restricted as in [9]. In addition, we describe an implementation of these schemes, discuss practical mesh construction and demonstrate our results numerically.

The general problem considered in this paper is of order n+m, with n equations singularly perturbed,

(1.1) 
$$\varepsilon y' = A_{11}(t,\varepsilon)y + A_{12}(t,\varepsilon)z + f_1(t)$$
  $0 \le t$ 

(1.2)  $z' = A_{21}(t,\varepsilon)y + A_{22}(t,\varepsilon)z + f_{2}(t)$ 

plus the boundary conditions for  $x(t) = \begin{pmatrix} y(t) \\ \tilde{z}(t) \end{pmatrix}$ (1.3)  $B_0 x(0) + B_1 x(1) = \beta$ .

The assumption (2.3) below on the eigenvalues of  $A_{11}$  plus the other regularity assumptions lead to the conclusion that the solution of (1.1) - (1.3) consists of a smooth curve away from the boundaries, possibly connected at each end to the boundary by a thin transition layer. As was pointed out in Part I, with Gauss or Lobatto schemes these boundary layer solutions must be approximated accurately, because otherwise layer errors would propagate throughout the entire interval of integration. The meshes used for collocation thus consist of three parts: Two fine grids near each boundary, with maximum mesh spacing  $h_L \leq Ke$  for a suitable constant K, connected to a much sparser mesh away from the boundaries with minimum mesh spacing  $\underline{h} \gg \varepsilon$ . The determination of the sparse mesh is based on the accuracy needed in the approximation of the reduced solution. The total number of mesh points required to meet a given error tolerance can be made to be independent of  $\varepsilon$ .

Of course, the mesh described above becomes highly nonuniform for very small  $\varepsilon$ . However, higher order collocation methods can handle such nonuniformity, see Part I and Ascher, Pruess and Russell [1]. Thus they are preferable to convergence acceleration methods in this context.

Following a short section where some results on the analytic solution of (1.1) - (1.3) are gathered for later use is Section 3, where the numerical schemes, their implementation and properties and the convergence results are presented. In §3.1 we describe a careful implementation of the collocation schemes which uses local-unknowns elimination, resulting with a well-conditioned system of linear equations (3.59) of a familiar sparse

≤ 1

structure, independent of the order of the scheme. This implementation is used both for the analysis and for the numerical calculations in following sections.

In §3.2 we consider a transformation of the dependent variables, needed for the analysis. Whereas in Part I this transformation commutes with the collocation operator, here it does not, and the resulting residue is shown to be sufficiently small in norm so that it can be considered as a small perturbation in regions where the mesh is dense, i.e. in boundary layer regions.

In §3.3 the mesh is described, together with the general collocation solution decomposition on each of its three parts. Then, in §3.4, our convergence results are stated. Theorem 3.1 summarizes the results for the layer regions near the boundaries while theorem 3.2 describes our results in the region away from the boundaries. Theorem 3.3 then states the combined results of the previous two theorems on the entire interval. This theorem is essentially the same as theorem 5.3 of Weiss [9]. An outline of its proof is followed by some remarks on the practical calculation of solutions by these schemes. Finally, we discuss the construction of the mesh near a boundary when there is no boundary layer because of a particularly lucky choice of boundary values. Such "luck" occurs frequently in practical problems. We find the argument in Kreiss and Kreiss [4] in this regard incomplete.

Sections 4 and 5 are devoted almost entirely to the proofs of theorems 3.1 and 3.2, respectively. In §4 we also discuss the layer mesh construction and show that the number of mesh points needed to achieve overall accuracy  $\delta$  for any  $\varepsilon$ ,  $0 < \varepsilon \leq \varepsilon_1$ , is  $0(\delta^{-P})$ , where p is the order of superconvergence of the method, defined in (3.42). This, provided that the mesh

defined in (3.46), (3.47) is used. If a uniform layer mesh is used instead then the number of mesh points needed is  $O(-\delta^{-1/P} \ln \delta)$ . But the actual advantage of (3.46), (3.47) over a uniform layer mesh is more significant than these bounds would indicate; see table 4.2 of Part I.

It is interesting to note that, perhaps contrary to one's first intuition, the analysis for the "long" interval away from the boundaries, where the solution varies slowly, is much more gruelling than the analysis for the layer intervals, where the solution varies very rapidly. In fact, the solution in the layer is dominated by a rapidly decreasing exponential and so its form is very smooth and simple to approximate, provided that we have a layer mesh with step sizes proportional to  $\varepsilon$ , affecting a stretching transformation. Indeed, it is the simple, exponential form of the layer solution which enables us to come up with the a-priori error equidistributing mesh (3.46), (3.47), whereas in general such meshes can be constructed only adaptively. Markowich and Ringhofer [6] had a similar success with problems on infinite intervals.

In §6, we seal this paper with a numerical example demonstrating our theoretical results.

### 2. Analytic preliminaries

In this section we mention some analytic results needed in the sequel and develop some notation. Since this section covers the same ground as §2 of Weiss [9], we allow ourselves to omit some details here.

Consider the linear problem (1.1), (1.2) where  $A_{ij} = A_{ij}(t,\epsilon)$  and  $f_i = f_i(t,\epsilon)$  are assumed, for simplicity, to be in  $C^{\infty}([0,1] \times [0,\epsilon_0])$  for some  $\epsilon_0 > 0$ ,  $1 \le i$ ,  $j \le 2$ . Further, assume that

(2.1) 
$$A_{11}(t,0) = E(t) \Lambda(t) E^{-1}(t)$$

with E  $\varepsilon$  C<sup>∞</sup>[0,1],

(2.2) 
$$\Lambda(t) = \text{diag} \{\lambda_1(t), \dots, \lambda_n(t)\}$$

and

(2.3) 
$$re(\lambda_{i}(t)) \begin{cases} < 0, & i = 1, ..., n \\ > 0, & i = n_{+1}, ..., n \end{cases}$$
  $te[0,1]$ 

Let  $n_+:= n-n_-$ , and denote  $\Lambda_-(t) = \text{diag} \{\lambda_1(t), \dots, \lambda_{n_-}(t)\}, \Lambda_+(t) = \text{diag} \{\lambda_{n_-+1}(t), \dots, \lambda_{n_-}(t)\}.$ 

We wish to decouple the slow components z from the fast ones and to diagonalize the remaining system for y. This is possible for a system with constant coefficients, see Part I; here, we can only almost get it. With L(t) a smooth solution to

(2.4) 
$$\varepsilon L' = -LA_{11} + \varepsilon (A_{22}L - LA_{12}L) + A_{21}$$

define the transformation

$$(2.5) \qquad \begin{pmatrix} u \\ v \\ v \\ z \end{pmatrix} = \begin{pmatrix} E^{-1} & 0 \\ -\varepsilon L & I \end{pmatrix} \begin{pmatrix} y \\ z \\ z \end{pmatrix}$$

(See Weiss [9] for justification). The system (2.1) - (2.2) is then transformed into

(2.6) 
$$\varepsilon u' = (\Lambda + \varepsilon B_{11})u + B_{22}v + g_1$$

(2.7)  $v' = B_{22}v + g_{22}v$ 

where  $B_{11}$ ,  $B_{12}$ ,  $g_1$ ,  $g_2$  are smooth functions of t and  $\epsilon$ .

For the transformed system (2.6) - (2.7), a desirable representation of the solution is obtained [5]: Writing it compactly as

(2.8) Hw = g  
with w(t) = 
$$\begin{pmatrix} u(t) \\ \tilde{v}(t) \end{pmatrix}$$
,  $g(t) = \begin{pmatrix} g_1(t) \\ \tilde{g}_2(t) \end{pmatrix}$ , and introducing  
the maps P\_ $\varepsilon$   $\mathbb{R}^{n-xn}$  an P<sub>+</sub>  $\varepsilon$   $\mathbb{R}^{n+xn}$  defined by  
(2.9) P\_ $_{-\sim}^{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$ ,  $P_{+\sim}^{x} = \begin{pmatrix} x_n \\ \vdots \\ x_n \end{pmatrix}$ ,  $x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$ 

we have

Theorem 2.1. The system (2.8) subject to boundary conditions

(2.10)  $P_{-u}(0) = n_{e} \in \mathbb{R}^{n}$ ,  $P_{+u}(1) = n_{+} \in \mathbb{R}^{n}$ ,  $v(0) = n_{0} \in \mathbb{R}^{m}$ has a unique solution which satisfies

 $(2.11) ||w|| \le \text{const} (||g|| + ||n_{-}|| + ||n_{+}|| + ||n_{0}||),$ 

provided that  $\varepsilon$  is sufficiently small,  $0 < \varepsilon \le \varepsilon_1$ . Also, for any  $q \ge 0$ there is a particular solution  $\underset{\sim p}{w}(t) = \begin{pmatrix} u_p(t) \\ v_p(t) \end{pmatrix}$  of (2.8) which satisfies (2.12)  $\sum_{\substack{j=0\\j=0}}^{q} \left| \left| \frac{d^j w_p}{dt^j} \right| \right| \le \text{const}, \quad 0 < \varepsilon \le \varepsilon_1.$ 

Now, define matrix solutions  $W_{-}$ ,  $W_{+}$  and  $W_{0}$  to the homogeneous problem (2.8) with g = 0 as follows:

(i) 
$$W_{-} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$
,  $v_{-} \in \mathbb{R}^{n \times n}$ , where  $v_{-}$  satisfies  
 $\varepsilon v_{-}' = (\Lambda + \varepsilon B_{11})v_{-}$ ,  $P_{-}v_{-}(0) = I$ ,  $P_{+}v_{-}(1) = 0$   
(ii)  $W_{+} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ ,  $v_{+} \in \mathbb{R}^{n \times n}$ , where  $v_{+}$  satisfies  
 $\varepsilon v_{+}' = (\Lambda + \varepsilon B_{11})v_{+}$ ;  $P_{-}v_{+}(0) = 0$ ,  $P_{+}v_{+}(1) = I$ 

(iii) 
$$W_0 = \begin{pmatrix} U_0 \\ V_0 \end{pmatrix}$$
,  $U_0 \in \mathbb{R}^{n \times m}$ ,  $V_0 \in \mathbb{R}^{m \times m}$ , where

 $HW_{0} = 0; V_{0}(0) = I, P_{-}\cup_{0}(0) = S_{-}(\varepsilon), P_{+}\cup_{0}(1) = S_{+}(\varepsilon)$ and  $S_{-} \in \mathbb{R}^{n_{-}\times m}, S_{+} \in \mathbb{R}^{n_{+}\times m}$  can be chosen by theorem 2.1 such that  $(2.13) \qquad \sum_{i=0}^{q} |\frac{d^{i}W_{0}}{dt^{i}}|| \leq \text{const.}$ 

Then we obtain the desired representation of the general solution to (2.6) - (2.7):

Theorem 2.3 Any solution w of (2.8) can be written as

(2.14)  $w = W_{-\zeta_{-}} + W_{+\zeta_{+}} + W_{0\zeta_{0}} + w_{p}$ with  $\zeta_{-} \in \mathbb{R}^{n_{-}}, \zeta_{+} \in \mathbb{R}^{n_{+}}$  and  $\zeta_{0} \in \mathbb{R}^{m_{-}}$ . The (smooth) particular solution satisfies (2.12), and the matrices  $W_{-}, W_{+}$  and  $W_{0}$ , defined above, have the asymptotic expansions

From the expansions (2.15) it is clear that

(2.16) 
$$\begin{array}{l} \upsilon_{-0}(t/\varepsilon) = \exp\left(\Lambda_{-}(0) t/\varepsilon\right) \\ \upsilon_{+0}(\frac{t-1}{\varepsilon}) = \exp\left(\Lambda_{+}(1) (t-1)/\varepsilon\right). \end{array}$$

Consider now the linear boundary value problem (1.1) - (1.3). The boundary conditions are transformed by (2.5) into a similar form for w(0) and w(1) and substituting the representation (2.14) into these boundary conditions we get

(2.17) 
$$M(\varepsilon)\zeta = \hat{\beta}(\varepsilon)$$
,  
where  $\zeta = (\zeta_{-}, \zeta_{+}, \zeta_{0})^{T}$  and the matrix M has an expansion  
(2.18)  $M(\varepsilon) = \sum_{j=0}^{q} M_{j}\varepsilon^{j} + O(\varepsilon^{q+1})$ .

We assume that  $M_0$  is nonsingular. This is equivalent to assuming that the boundary value problem (1.1) - (1.3) is well posed, i.e. for  $\varepsilon$  small enough

(2.19)  $||_{\tilde{z}}^{\tilde{y}}|| \leq \text{const} (||_{f_2}^{f_1}|| + ||_{\beta}||)$ with the constant independent of  $\varepsilon$ .

It is clear that the preceding representation of the general solution of (2.8) can be made analogously on any interval  $[t,\bar{t}] \in [0,1]$  with the solution matrices appropriately defined. In particular, in (2.15), (2.16),  $\underline{t}$  would replace 0 and  $\overline{t}$  would replace 1. Denoting by  $(\upsilon_{-})_{\ell}$  and  $(\upsilon_{+})_{\ell}$  the  $\ell$ -th columns of  $\upsilon_{-}$  and  $\upsilon_{+}$  we get

$$(2.20) \qquad \left| \left| \frac{d^{j}(\upsilon_{-})_{\ell}}{dt^{j}} \right| \right| \leq \operatorname{const} \varepsilon^{-j} [\exp\{\operatorname{re}(\lambda_{\ell}(\underline{t}))(\underline{t}-\underline{t})/\varepsilon\} + 0(\varepsilon)], \\ \underline{t} \leq \underline{t} \leq \overline{t}, \ \ell = 1, \dots, n_{-} \qquad j = 0, 1, \dots, q \\ (2.21) \qquad \left| \left| \frac{d^{j}(\upsilon_{+})_{\ell}}{dt^{j}} \right| \right| \leq \operatorname{const} \varepsilon^{-j} [\exp\{\operatorname{re}(\lambda_{\ell}(\overline{t}))(\underline{t}-\overline{t})/\varepsilon\} + 0(\varepsilon)], \\ \underline{t} \leq \underline{t} \leq \overline{t}, \ \ell = n + 1, \dots, n, \ j = 0, 1, \dots, q. \end{cases}$$

#### 3. Numerical solutions and their convergence

#### 3.1. The numerical schemes and their implementation

In section 3 of Part I we have presented some classes of collocation methods and discussed their equivalent Runge-Kutta formulation and some of their properties. Here we mention only some of these details again and rely on familiarity with Part I for the rest.

A collocation procedure under consideration is completely determined in terms of k points ( $k \ge 1$ ),

(3.1)  $0 \le \rho_1 < \ldots < \rho_k \le 1$ ,

which we take to be either Gauss or Lobatto points, and a mesh

(3.2) 
$$\Delta: 0 = t_1 < t_2 < \dots < t_N < t_{N+1} = 1$$
$$h_i := t_{i+1} - t_i, 1 \le i \le N, h := \max\{h_i, 1 \le i \le N\}$$

On a given mesh  $\Delta$ , the collocation solution  $\underset{\sim}{x_{\Delta}}(t) = \begin{pmatrix} y_{\Delta}(t) \\ z_{\Delta}(t) \end{pmatrix}$  to (1.1) - (1.3) is a continuous piecewise polynomial vector function of degree at most k satisfying the boundary conditions (1.3) and the differential equations (1.1), (1.2) at the collocation points

(3.3)  $t_{ij}:=t_i + h_i \rho_j$  i = 1, ..., N, j = 1, ..., k.Inside each subinterval  $[t_i, t_{i+1}]$ , the polynomials  $y_{\Delta}(t)$  and  $z_{\Delta}(t)$  can be represented in terms of the values

(3.4)  $y_{i} := y_{\Delta}(t_{i}), \quad z_{i} := z_{\Delta}(t_{i}), \quad 1 \le i \le N+1$  $y_{ij} := y_{\Delta}(t_{ij}), \quad z_{ij} := z_{\Delta}(t_{ij}) \quad 1 \le i \le N, \quad 1 \le j \le k$ 

(strictly speaking, for Lobatto points some additional derivative values are required as well), which satisfy the difference equations (3.5)  $\frac{\varepsilon}{h_{i}} (y_{ij} - y_{i}) = \sum_{\substack{k=1 \ j \neq i}}^{k} \hat{a}_{jk} (A_{11}(t_{ik},\varepsilon)y_{ik} + A_{12}(t_{ik},\varepsilon)z_{ik} + f_{1}(t_{ik}))$  $1 \le i \le k$ 

(3.6) 
$$\frac{1}{h_i} (z_{ij} - z_i) = \sum_{l=1}^{K} \hat{a}_{jl} (A_{21}(t_{il}, \epsilon) y_{il} + A_{22}(t_{il}, \epsilon) z_{il} + f_2(t_{il}))$$

For Lobatto points,  $\rho_k = 1$  and  $\rho_1 = 0$ . Thus  $y_{i+1} = y_{ik}$ ,  $z_{i+1} = z_{ik}$  and equations (3.5), (3.6) are trivial for j=1. For Gauss points,  $\rho_k < 1$ ,  $\rho_1 > 0$ and we extend the range of j in (3.5), (3.6) to include j = k+1 as well, with  $y_{i+1} = y_{i,k+1}$ ,  $z_{i+1} = z_{i,k+1}$  and  $\hat{a}_{k+1,k} = \hat{b}_k$ ,  $k=1,\ldots,k$ ; see §3 of Part I for the definitions of the constants  $\hat{a}_{jk}$ ,  $\hat{b}_k$ , as well as the matrices  $\hat{A}$  and  $\bar{A}$ used later.

In the sequel, we shall adhere to the following notational convention, used already above. The collocation approximation to a function  $\psi(t)$  is denoted by  $\psi_{\Delta}(t)$ . Its values at mesh points are  $\psi_i$ ,  $1 \le i \le N+1$ , and those at collocation points are  $\psi_{ij}$ ,  $1 \le i \le N$ ,  $1 \le j \le k$ . Also,  $\psi^c$  will denote the vector formed by the restriction of  $\psi(t)$  to  $\Delta \cup \{t_{ij}; 1 \le i \le N, 1 \le j \le k\}$ . As well, c, K and  $c_j$ , j = 0,1,2... will denote constants independent of  $\varepsilon$ and  $\Delta$ .

Next, we describe a particular, careful implementation of the collocation schemes which is used both for the numerical calculations reported in §6 and for the analysis in §5. The differential equations (1.1), (1.2) are written as one system

(3.7) x' = A(t)x + f(t)

for which the numerical method is written in Runge-Kutta form

(3.8)  $x_{i+1} = x_i + h_i \sum_{j=1}^{k} \hat{b}_{j-ij}$   $1 \le i \le N$  $1 \le j \le k$ 

(3.9) 
$$F_{ij} = x_{\Delta}'(t_{ij}) = A(t_{ij})x_{ij} + f(t_{ij}) = A(t_{ij})(x_i + h_i \sum_{\ell=1}^{\kappa} \hat{a}_{j\ell}F_{i\ell}) + f(t_{ij}).$$

The unknowns  $F_{ij}$  (or  $x_{ij}$ ) for each interval  $[t_i, t_{i+1}]$  are local and can be eliminated locally. (This is sometimes referred to as "condensation of parameters" - see Ascher, Pruess and Russell [1]). We choose to locally

1.

eliminate the  $F_{ij}$  in case that  $\rho_k < 1$ , and the  $x_{ij}$  in case that  $\rho_k = 1$ . These choices avoid unnecessary loss of digits due to cancellation error, as can be readily verified for the example y' = y/ $\varepsilon$  + 1/ $\varepsilon$  with 0 <  $\varepsilon$  << 1.

Consider <u>Gauss points</u> first. Equations (3.9) can be written as (n+m)k linear equations

(3.10)  $J_{i}F_{i} = R_{i}$ 

where

$$F_{i} = (F_{i1}, \dots, F_{ik})^{T}, \quad R_{i} = C_{A_{i}}^{X} + f_{i}$$
(3.11)  

$$C_{A} = \begin{bmatrix} A(t_{i1}) \\ \vdots \\ A(t_{ik}) \end{bmatrix}, \quad f_{i} = \begin{pmatrix} f(t_{i1}) \\ \vdots \\ f(t_{ik}) \end{pmatrix}$$
(3.12)  

$$J_{i} = I - h_{i} \begin{bmatrix} \hat{a}_{11}A(t_{i1}) & \hat{a}_{12}A(t_{i1}) & \dots & \hat{a}_{1k}A(t_{i1}) \\ \hat{a}_{21}A(t_{i2}) & \hat{a}_{22}A(t_{i2}) & \dots & \hat{a}_{2k}A(t_{i2}) \\ \vdots \\ \hat{a}_{k1}A(t_{ik}) & \dots & \dots & \hat{a}_{kk}A(t_{ik}) \end{bmatrix} = I - h_{i}D_{A}(\hat{A}\otimes I)$$

in which I stands for an identity matrix of the appropriate dimension (n+m or (n+m)k) and  $D_A = \text{diag } \{A(t_{i1}), \ldots, A(t_{ik})\}$ . (The dependence on i is suppressed in  $C_A$  and  $D_A$ ). Introducing for notational purposes the (n+m)x(n+m)k matrix

(3.13)  $\hat{B} = [\hat{b}_1 I, \dots, \hat{b}_k I]$ 

we can write (3.8) as

(3.14) 
$$\chi_{i+1} = \Gamma_i \chi_i + g_i$$
  $1 \le i \le N$ 

where

(3.15) 
$$\Gamma_i = I + h_i \hat{B} J_i^{-1} C_A$$
,  $g_i = h_i \hat{B} J_i^{-1} f_i$ 

The difference equations (3.14) together with the boundary equations corresponding to (1.3)

$$(3.16) \qquad B_0 X_1 + B_1 X_{N+1} = \beta$$

form a set of (N+1)(n+m) linear equations whose size and structure are independent of k for the solution values at the mesh points.

For <u>Lobatto points</u> we perform a similar elimination of local parameters, but now our parameters are  $x_i = x_{i1}, x_{i2}, \dots, x_{ik-1}, x_{ik} = x_{i+1}$ . Instead of (3.8), (3.9) we write, as in (3.5), (3.6)

$$(3.17) \quad 1/h_{i}(x_{ij} - x_{i}) = \sum_{\ell=1}^{K} \hat{a}_{j\ell}(A(t_{i\ell})x_{i\ell} + f(t_{i\ell})) \qquad 2 \le j \le k$$

and this can be written as (n+m)(k-1) linear equations

where

$$\bar{x}_{i} = (x_{i2}, \dots, x_{ik})^{T}, \quad \bar{R}_{i} = (\bar{R}_{i2}, \dots, \bar{R}_{ik})^{T}$$
(3.19)

$$(3.20) \quad \overline{J}_{i} = [I + h_{i}\hat{a}_{j1}A(t_{i})]x_{i} + h_{i}\sum_{\ell=1}^{K} \hat{a}_{j\ell}f(t_{i\ell})$$

$$(3.20) \quad \overline{J}_{i} = I - h_{i}\begin{bmatrix}\hat{a}_{22}A(t_{i2}) & \dots & \hat{a}_{2k}A(t_{ik})\\ \vdots \\ \hat{a}_{k2}A(t_{i2}) & \dots & \hat{a}_{kk}A(t_{ik})\end{bmatrix} = I - h_{i}(\overline{A}\otimes I)\overline{D}_{A}$$

with  $\overline{D}_A$  = diag {A(t<sub>12</sub>),...,A(t<sub>ik</sub>)} and  $\overline{A}$  is a nonsingular matrix, as in (3.14) of Part I.

Since  $\rho_k = 1$ ,  $x_{i+1}$  is obtained as the last n+m rows of  $\bar{J}_i^{-1}R_i$ . Partitioning  $\bar{J}_i^{-1}$  into blocks of size (n+m)x(n+m),  $\bar{J}_i^{-1} = ((\bar{J}_i^{-1})_{jl})_{j,l=2}^k$ , we get difference equations of the form (3.14) where now, instead of (3.15),

(3.21) 
$$\Gamma_{i} = \sum_{\ell=2}^{k} (\bar{J}_{i}^{-1})_{k\ell} [I + h_{i}\hat{a}_{\ell}] A(t_{i})], g_{i} = h_{i} \sum_{\ell=2}^{k} (\bar{J}_{i}^{-1})_{k\ell} \sum_{s=1}^{k} \hat{a}_{\ell} f(t_{is}).$$

An advantage of the difference equations (3.14), (3.16), obtained both for Gauss and for Lobatto points, is that even when some rows of A(t) of (3.7) depend on  $1/\epsilon$  and  $\epsilon \ll h_i$ , the components of  $\Gamma_i$  and  $g_i$ remain bounded and are constructed accurately.

# 3.2 Transformation of variables

Consider the linear problem (1.1) - (1.3) and the transformed system (2.6), (2.7). Since the latter has a structure more amenable to analysis, we will rely on it in parts of our treatment. However, we stress that the actual numerical procedure is applied to (1.1), (1.2) and not to (2.6), (2.7).

In the constant coefficient case, the operators of collocation and the transformation (2.5) commute. Here they do not, in general. Thus, if we define vector functions  $u_{\Delta}(t)$ ,  $v_{\Delta}(t)$  by

$$(3.22) \qquad \begin{pmatrix} u \\ \tilde{v} \\ \omega \\ \omega \end{pmatrix} = \begin{pmatrix} E^{-1} & 0 \\ -\varepsilon L & I \end{pmatrix} \begin{pmatrix} y \\ \tilde{z} \\ \omega \\ \omega \end{pmatrix}$$

then  $\underset{\sim\Delta}{u}$ ,  $\underset{\sim\Delta}{v}$  collocate the transformed equations, but are not necessarily piecewise polynomials of degree at most k. Correspondingly, applying the transformation (3.22) to the difference equations (3.5), (3.6) we obtain

$$(3.23) \quad \frac{\varepsilon}{h_{i}} (u_{ij} - u_{i}) = \sum_{\ell=1}^{k} \hat{a}_{j\ell} \{ [\Lambda(t_{i\ell}) + \varepsilon B_{11}(t_{i\ell})] u_{i\ell} + B_{12}(t_{i\ell}) v_{i\ell} + g_{1}(t_{i\ell}) \} + \frac{\varepsilon}{h_{i}} R_{ij} \qquad 1 \le i \le N$$

(3.24) 
$$\frac{1}{h_{i}} (v_{ij} - v_{i}) = \sum_{\ell=1}^{K} \hat{a}_{j\ell} \{B_{22}(t_{i\ell})v_{i\ell} + g_{2}(t_{i\ell})\} + \frac{1}{h_{i}} \sum_{ij}^{S} ij$$

where  $\epsilon/h_i \underset{ij}{R_{ij}}$  and  $1/h_i \underset{ij}{S_{ij}}$  consist of linear operators acting on  $\underset{il}{u_{il}}, \underset{il}{v_{il}},$  $l=1,\ldots,k$  for Lobatto points and  $l=1,\ldots,k+1$  for Gauss points, and inhomogeneities. We now show that their norms are  $O(h_i)$  and so they can be dealt with as small perturbations when  $h_i$  is small.

$$\frac{\text{Lemma 3.1}}{\binom{2}{1}}: \text{ For each } i, 1 \le i \le N,$$

$$\binom{3.25}{\binom{\epsilon/h_i R_{ij}}{1/h_i \tilde{s}_{ij}}} = {}^{h_i \{ \phi_{ij} [u_{i1}, \dots, u_{iq}; v_{i1}, \dots, v_{iq}] + }}{\underset{\psi_{ij} [g_1(t_{i1}), \dots, g_1(t_{iq}), g_2(t_{i1}), \dots, g_2(t_{iq})] \}}{}}$$

where  $\phi_{ij}$ ,  $\psi_{ij}$  are bounded linear operators,

 $||\phi_{ij}||, ||\psi_{ij}|| \le c,$   $1 \le i \le N, 1 \le j \le q$ with q = k for Lobatto points, q = k+1 for Gauss points.

Those readers who wish to skip the proof of this lemma can do so without loss of continuity.

<u>Proof</u>: Writing  $u'_{\Delta}(t)$  and  $v'_{\Delta}(t)$  in terms of their polynomial interplants of order k on  $[t_i, t_{i+1}]$  we get

$$u_{\Delta}'(t) = \sum_{j=1}^{k} u_{\Delta}'(t_{ij}) L_{j}(\frac{t-t_{i}}{h_{i}}) + \frac{1}{k!} u_{\Delta}^{(k+1)}(\zeta_{t}) \prod_{j=1}^{k} (t-t_{ij})$$
$$t_{i} \leq \zeta_{t} \leq t_{i+1}$$

where  $L_{i}$  are the Lagrange polynomials. Integrating,

$$u_{ij} - u_{i} = h_{i} \sum_{\ell=1}^{k} u_{\Delta}'(t_{i\ell}) \hat{a}_{j\ell} + \frac{1}{k!} \int_{t_{i}}^{t_{ij}} u_{\Delta}^{(k+1)}(z_{t}) \prod_{\ell=1}^{k} (t-t_{i\ell}) dt$$

and so, by (3.23),

(3.26) 
$$R_{ij} = \frac{1}{k!} \int_{t_i}^{t_{ij}} u_{\Delta}^{(k+1)}(z_t) \int_{\ell=1}^{k} (t-t_{i\ell}) dt$$

with a similar expression for  $S_{ij}$ ,  $v_{\Delta}$  replacing  $u_{\Delta}$ .

Next, since  $y_{\Delta}$  and  $z_{\Delta}$  are polynomials of degree at most k on  $[t_i, t_{i+1}]$ , by the transformation (3.22),

(3.27) 
$$u_{\Delta}^{(k+1)}(\tau) = \sum_{\nu=1}^{k} {\binom{k+1}{\nu}} (E^{-1}(\tau))^{(\nu)} y_{\Delta}^{(k+1-\nu)}(\tau)$$

 $t_i \leq \tau \leq t_{i+1}$ 

(3.28) 
$$v_{\Delta}^{(k+1)}(\tau) = -\varepsilon \sum_{\nu=1}^{k} {\binom{k+1}{\nu}} L^{(\nu)}(\tau) y_{\Delta}^{(k+1-\nu)}(\tau)$$

and

(3.29) 
$$\varepsilon \underbrace{y}_{\Delta}^{(\ell)}(\tau) = \varepsilon \frac{d^{\ell-1}}{d\tau^{\ell-1}} \underbrace{(y_{\Delta}^{\prime}(\tau))}_{j=1} = \underbrace{\sum_{j=1}^{k} \varepsilon y_{\Delta}^{\prime}(t_{ij})}_{d\tau^{\ell-1}} \frac{d^{\ell-1}}{d\tau^{\ell-1}} L_{j}(\frac{\tau - t_{i}}{h_{i}}).$$

Replacing the vectors  $\epsilon y'_{\Delta}(t_{ij})$  through the collocation equations (3.9) and

(3.30) 
$$||u^{(k+1)}|| \leq h_i \phi/\epsilon, ||v^{(k+1)}_{\Delta}|| \leq h_i \phi$$

(3.31)  $\phi = ch_i^{-k} \max\{||u_{ij}||, ||v_{ij}||, ||g_1(t_{ij})||, ||g_2(t_{ij})||, 1 \le j \le q\}$ Finally, substituting (3.30), (3.31) into (3.26) and the corresponding expression for  $S_{ij}$ , the desired result (3.25) is obtained.

QED.

# 3.3 The mesh and the decomposition of numerical solutions

The meshes considered in this paper have the following structure. Near the boundaries, the step sizes  $h_i$  are comparable to  $\epsilon$ . Specifically, there are given numbers 0 <  $N_0$ ,  $N_1$  < N and constants  $K_0$ ,  $K_1$ , such that

 $h_i \leq K_0 \varepsilon$ ,  $i = 1, \dots, N_0$ 

(3.32)

 $h_{i} \leq K_{1}\epsilon$ ,  $i = N-N_{1} + 1, ..., N$ .

In between, much larger step sizes may be used,

i.e.  $h_i >> \epsilon$ ,  $i = N_0 + 1, ..., N-N_1$ . We will assume for convenience of notation that h, the largest step size, occurs away from the boundaries. Such a mesh is depicted in figure 1 below.



Figure 1: The mesh

For brevity of notation we set  $\underline{i} = N_0 + 1$ ,  $\overline{i} = N - N_1 + 1$  and write (3.33)  $t_{\overline{i}} = T_0 \varepsilon$ ,  $t_{\overline{i}} = 1 - T_1 \varepsilon$ 

Our next step is to write down a decomposition representation to the discrete solution of (3.23), (3.24), similar to the representation (2.14) for the analytic solution. Moreover, we write down such a representation for each of the three parts of the mesh.

We write the system (3.23), (3.24), in analogy to (2.8) as (3.34a)  $(H_{W_{\Delta}})(t_{ij}) = g(t_{ij})$  j = 1, ..., q, i = 1, ..., Nor in shorthand as

(3.34b)  $H_{\Delta \overset{w}{\sim} \Delta} = \overset{g}{}_{\Delta}$ , where  $\underset{\sim}{w}(t) = \begin{pmatrix} u_{\Delta}(t) \\ \vdots \\ y_{V}(t) \end{pmatrix}$  is in the class  $S_{\Delta}$  of functions defined by (3.22) with  $y_{\Delta}(t)$ ,  $z_{\Delta}(t)^{\Delta}$  continuous piecewise polynomial vector functions of degree at most k. Let

(3.35) 
$$W^{M}_{\Delta S}(t) = \begin{pmatrix} U^{M}(t) \\ \Delta S \\ V^{M}_{\Delta S}(t) \end{pmatrix} \in \mathbb{R}^{(n+m)\times n} s, t \in [t_{M_{0}}, t_{M_{1}}]$$

be matrix valued functions with columns in class  $S_{\Delta}$ . Here s stands for -, + or 0,  $n_0$ := m, M stands for I, II or III, to denote the three mesh regions considered, and so  $I_0 = 1$ ,  $I_1 = \underline{i} = II_0$ ,  $II_1 = \overline{i} = III_0$ ,  $III_1 = N+1$ . <u>On the interval [0,  $T_0 \in$ ]</u>: Define  $W_{\Delta s}^{I}$ ,  $w_{p\Delta}^{I} \in S_{\Delta}$  as follows (omitting the superscript I):

$$(3.36a) \qquad H_{\Delta}W_{\Delta-} = 0, P_{-}U_{\Delta}(0) = I, P_{+}U_{\Delta}(T_{0}\varepsilon) = 0, V_{\Delta}(0) = 0;$$

(3.36b) 
$$H_{\Delta}W_{\Delta+} = 0, P_{\Delta+}(0) = 0, P_{+}U_{\Delta+}(T_{0}\varepsilon) = I, V_{\Delta+}(0) = 0;$$

$$(3.36c) \qquad H_{\Delta}W_{\Delta 0} = 0, P_{\Delta}U_{\Delta 0}(0) = S_{\epsilon}(\epsilon), P_{+}U_{\Delta 0}(T_{0}\epsilon) = P_{+}U_{0}(T_{0}\epsilon), V_{\Delta 0}(0) = I;$$

(3.36d) 
$$H_{\Delta \sim P\Delta} = g_{\Delta}, P_{\Delta \sim P\Delta}(0) = P_{\Delta \sim P}(0), P_{\Delta \sim P\Delta}(T_{0}\varepsilon) = P_{\Delta \sim P}(T_{0}\varepsilon),$$
$$v_{P\Delta}(0) = v_{P}(0)$$

The general solution of (3.34) on the left layer mesh is written formally as:

$$(3.37) \qquad \begin{split} & \underset{\Delta}{\text{w}}_{\Delta}^{I} = \underset{\Delta}{\text{w}}_{\Delta}^{I} \underbrace{\zeta}_{a}^{I} + \underset{\Delta}{\text{w}}_{\Delta}^{I} \underbrace{\zeta}_{a}^{I} + \underset{0}{\text{w}}_{0}^{I} \underbrace{\zeta}_{0}^{I} + \underset{P}{\text{w}}_{P\Delta}^{I} \\ & \text{with } \underbrace{\zeta}_{a}^{I} \in \mathbb{R}^{n}, \underbrace{\zeta}_{a}^{I} \in \mathbb{R}^{n}, \underbrace{\zeta}_{0} \in \mathbb{R}^{m}. \\ & \underbrace{\text{On the interval } [1 - T_{1}\varepsilon, 1]: \text{ In precise analogy to the above, define} \\ & \underbrace{\text{Will and } \underset{P\Delta}{\text{w}}_{P\Delta}^{III}, \text{ again omitting superscripts:} \\ & (3.38a) \qquad H_{\Delta}W_{\Delta} = 0, P_{-}U_{\Delta-}(1 - T_{1}\varepsilon) = I, P_{+}U_{\Delta-}(1) = 0, V_{\Delta-}(1 - T_{1}\varepsilon) = 0; \\ & (3.38b) \qquad H_{\Delta}W_{\Delta+} = 0, P_{-}U_{\Delta+}(1 - T_{1}\varepsilon) = 0, P_{+}U_{\Delta+}(1) = I, V_{\Delta+}(1 - T_{1}\varepsilon) = 0; \\ & (3.38c) \qquad H_{\Delta}W_{\Delta0} = 0, P_{-}U_{\Delta0}(1 - T_{1}\varepsilon) = P_{-}U_{0}(1 - T_{1}\varepsilon), P_{+}U_{\Delta0}(1) = S_{+}(\varepsilon), \\ & V_{\Delta0}(1 - T_{1}\varepsilon) = 0 \\ & (3.38d) \qquad H_{\Delta}w_{P\Delta} = \underbrace{g_{\Delta}}, P_{-}\underbrace{w_{P\Delta}(1 - T_{1}\varepsilon)} = P_{-}\underbrace{w_{P}(1 - T_{1}\varepsilon)}, P_{+}\underbrace{w_{P\Delta}(1)} = P_{+}\underbrace{w_{P}(1)}, \\ & \underbrace{w_{P\Delta}(1 - T_{1}\varepsilon)} = \underbrace{w_{P}(1 - T_{1}\varepsilon)} = \underbrace{w_{P}(1 - T_{1}\varepsilon)} \\ & \end{array}$$

The general solution of (3.34) on the right layer mesh is written formally as

(3.39) 
$$\begin{split} & w_{\Delta}^{III} = W_{\Delta-}^{III} \zeta_{\Delta}^{III} + W_{\Delta+}^{III} \zeta_{\pm}^{III} + W_{0}^{III} \zeta_{0}^{III} + w_{P\Delta}^{III} \\ & \text{with } \zeta_{-}^{III} \in \mathbb{R}^{n-}, \zeta_{\pm}^{III} \in \mathbb{R}^{n+}, \zeta_{0} \in \mathbb{R}^{m}. \\ & \underline{On \text{ the interval } [T_{0}\varepsilon, 1-T_{1}\varepsilon]: \text{ Define } W_{\Delta}^{II}, W_{\Delta0}^{II} \text{ and } w_{P\Delta}^{II} \text{ as follows: Let} \\ & Y_{0}(t), Z_{0}(t) \text{ be obtained from } U_{0}(t), V_{0}(t), \text{ via the (inverse) transformation} \\ & (3.22) \text{ from } X_{\Delta} = \begin{pmatrix} Y_{\Delta} \\ Z_{\Delta} \end{pmatrix} \text{ and } X_{\Delta0} = \begin{pmatrix} Y_{\Delta0} \\ Z_{\Delta0} \end{pmatrix}, \text{ respectively, which are defined as follows:} \\ & (3.40a) X_{\Delta} \text{ and } X_{\Delta0} \text{ satisfy the homogeneous equations (3.8), (3.9) with} \\ & \xi = 0 \end{split}$$

$$f = 0.$$

(3.40b) 
$$Y_{A}(T_{0}\varepsilon) = I, Z_{A}(T_{0}\varepsilon) = 0;$$

(3.40c)  $Y_{\Delta 0}(T_0 \varepsilon) = Y_0(T_0 \varepsilon), Z_{\Delta 0}(T_0 \varepsilon) = Z_0(T_0 \varepsilon).$ 

The particular solution  $\underset{\sim}{\mathtt{w}}_{P\Delta}^{I\,I}$  is defined, e.g., by

(3.40d) 
$$H_{\Delta \tilde{w} P \Delta} = \tilde{g}_{\Delta}, w_{P \Delta}(T_0 \varepsilon) = \tilde{w}_P(T_0 \varepsilon)$$

The general solution of (3.34) on the long interval away from the layers is written formally as

(3.41) 
$$\begin{split} & \underset{\Delta}{\text{w}_{\Delta}^{\text{II}}} = W_{\Delta}^{\text{II}} \zeta^{\text{II}} + W_{\Delta}^{\text{II}} \zeta^{\text{II}} + \underset{P\Delta}{\text{w}_{P\Delta}}^{\text{II}} \\ & \text{with } \zeta^{\text{II}} \in \mathbb{R}^{n}, \ \zeta_{0} \in \mathbb{R}^{m}. \end{split}$$

# 3.4 Convergence results

Below we state the various results regarding the convergence of the numerical methods, culminating in theorem 3.3 The proofs for those results which have not been proven elsewhere are contained in the next two sections. Denote by p the "regular" superconvergence order of the schemes under consideration, i.e.

(3.42) p = 2k for a k-stage Gauss scheme

= 2(k-1) for a k-stage Lobatto scheme.

Also, define the seminorms on collocation solutions,

(3.43a)  $||\psi_{\Lambda}||_{\Lambda} := \max\{||\psi_{i}||; 1 \le i \le N + 1\}$ 

(3.43b)  $||\psi_{\Lambda}||_{c} := \max\{|\psi_{i,i}||; 1 \le i \le N, 1 \le j \le k\}$ 

where the vector norms used are maximum norms. Thus

 $||\psi_{\Delta}^{C}|| = \max\{||\psi_{\Delta}||_{\Delta}, ||\psi_{\Delta}||_{C}\}$ . Also  $||\psi_{\Delta}||_{\Delta}^{M}$ ,  $||\psi_{\Delta}||_{C}^{M}$  will denote the seminorm where the range of i in (3.43) is restricted to  $M_{0} \le i \le M_{1}$ , M = I, II or III. For a matrix whose columns are collocation solutions,

a maximum on the column norms is taken.

For the "short" intervals [0,  $t_{\tilde{i}}$ ] and [ $t_{\tilde{i}}$ , 1] we have

Theorem 3.1

- (a) The solution representations (3.37) and (3.39) are valid (i.e. their components can be computed in a stable way)
- (b) With  $h_L$  the maximum step size in the layers ( $h_L \le \epsilon \max\{K_0, K_1\}$  by (3.32)), the "smooth" components satisfy

$$(3.44) \qquad ||W_{\Delta 0} - W_0||_{\Delta}^{I}, ||W_{\Delta 0} - W_0||_{\Delta}^{III}, ||w_{P\Delta} - w_{P}||_{\Delta}^{I}, ||w_{P\Delta} - w_{P}||_{\Delta}^{III} \le ch_{L}^{P}$$

(c) The auxiliary solution components in the layers (for which there are no counterparts in the exact solution decomposition) satisfy

$$(3.45) \qquad \mathsf{W}_{\Delta^+}^{\mathrm{I}}(\mathsf{T}_0\varepsilon) = \begin{pmatrix} \mathsf{0}_{\mathsf{n}_{-}}\mathsf{x}\mathsf{n}_{+} \\ \mathsf{I}_{\mathsf{n}_{+}}\mathsf{x}\mathsf{n}_{+} \\ \mathsf{0} \end{pmatrix} + \mathsf{O}(\varepsilon), \ \mathsf{W}_{\Delta^-}^{\mathrm{III}}(\mathsf{1} - \mathsf{T}_{1}\varepsilon) = \begin{pmatrix} \mathsf{I}_{\mathsf{n}_{-}}\mathsf{x}\mathsf{n}_{-} \\ \mathsf{0} \end{pmatrix} + \mathsf{O}(\varepsilon)$$

(d) For a given accuracy tolerance δ, δ ≥ cε, the layer meshes can be constructed as follows: With u:= max{|λ<sub>j</sub>(0)|, j=1,...,n\_}, v:= min{-re(λ<sub>j</sub>(0)), j=1,...,n\_} > 0, define

define

(3.46) 
$$h_1 := \frac{\varepsilon}{\mu} \left[ \frac{v}{\mu | c_{\gamma} |} \right]^{1/P} \delta^{1/P}$$

(3.47) 
$$h_i := h_{i-1} \exp\{\frac{1}{p} \frac{v}{\epsilon} h_{i-1}\}$$
 until  $t_{i+1} \ge T_0 \epsilon$   
where  $c_\gamma$  is a known constant (cf. Part I) and

(3.48) 
$$T_0 := v^{-1} |\ln \delta|$$
.

The right end layer is constructed in an analogous way, depending on  $\lambda_j(1)$ ,  $j = n-n_{+1}, \dots, n$ . Then  $(3.49a) \qquad W^{I}_{\Delta_{-}}(t_i) = \begin{pmatrix} \exp(\Lambda_{-}(0)t_i/\epsilon) \\ 0 \end{pmatrix} + O(\delta) \qquad 1 \le i \le \underline{i}$ 

$$(3.49b) \qquad \mathbb{W}_{\Delta^{+}}^{\text{III}}(t_{i}) = \begin{pmatrix} 0 \\ \exp(\Lambda_{+}(1)(t_{i}-1)/\epsilon) \end{pmatrix} + O(\delta) \qquad \overline{i} \le i \le N+1$$

The proof of this theorem is given in §4. We note that the assumption  $\delta \ge c_{\epsilon}$  is not essential, see §6.

For the "long" interval  $[t_i, t_{\overline{i}}]$  we have

<u>Theorem 3.2</u> Let (3.50)  $\kappa := \epsilon h^{-1}$ ,  $h := \min\{h_i; \quad i \le i \le i\}$  (a) The solution representation (3.41) is valid for  $\varepsilon$  sufficiently small such that  $\kappa < c$  for Lobatto points or  $\varepsilon \sum_{i=1}^{i} h_i^{-1} < c$  for Gauss points, where c is a constant of order 1.

(b) The first n fundamental solution components satisfy

$$(3.51a) \quad W_{\Delta}^{\text{II}}(t_{i}) = (-1)^{\binom{k}{i-i}\binom{E^{-1}(t_{i})}{0}} + 0(\kappa) \qquad \underline{i} \le i \le \overline{i}$$
  
for Gauss points and

$$(3.51b) \quad W_{\Delta}^{II}(t_{i}) = (-1)^{(k+1)(i-i)} \left( \Lambda^{-1}(t_{i})E^{-1}(t_{i})A_{11}(t_{i}) \right) + 0(\kappa)$$

$$0$$

$$i \le i \le \overline{i}$$

for Lobatto points.

(c) Define the error e as follows:

For a k-stage Gauss scheme, e:= h<sup>k</sup> and, if k is odd and the mesh is locally almost uniform, i.e.,

(3.52)  $h_{i+1} = h_i(1 + O(h_i))$  for all i odd <u>or</u> all i even, then  $e = h^{k+1}$ .

For a k-stage Lobatto scheme e:=  $Kh^{P} + \epsilon h^{k-1}$ , and, if k is even and the mesh is locally almost uniform then  $e = Kh^{P} + \epsilon h^{k}$ . Also, if the slow components z are absent from (1.1), (1.3) then K = 0.

Then we have

 $(3.53) \qquad ||w_{P\Delta}^{II} - w_{P}||_{\Delta}^{II}, ||W_{\Delta 0}^{II} - W_{0}||_{\Delta}^{II} \le ce.$ 

The proof of this theorem is given in §5.

The condition on  $\varepsilon$  and the mesh for Gauss points in (a) is slightly annoying. However, it can be argued that for  $O(\kappa)$  perturbations to be small, say O(h), this condition has to be satisfied anyway.

The central theorem, summarizing our efforts for linear problems with variable coefficients, follows:

<u>Theorem 3.3</u>: Assume that the boundary value problem (1.1) - (1.3) is well posed, uniformly for  $0 < \epsilon \le \epsilon_0$ , and denote the solution by  $x(t) = \begin{pmatrix} y(t) \\ z(t) \end{pmatrix}$ . Assume further that (3.54)  $det \begin{pmatrix} P_- E^{-1}(0) \\ P_- E^{-1}(1) \end{pmatrix} \neq 0$ .

Then for  $0 < \varepsilon \le \varepsilon_1$  there are positive constants  $\delta_0$ ,  $h_0$  and  $\kappa_0$  such that for any  $\delta$ ,  $0 < c_0 \varepsilon \le \delta \le \delta_0$  and any mesh  $\Delta = \Delta(\varepsilon)$  satisfying (3.46), (3.47), similar conditions at the right end layer,  $\kappa \le \kappa_0$  and  $h \le h_0$ , the k-stage collocation scheme based on Gauss or Lobatto points has a unique solution x (t) which satisfies

$$(3.55) \qquad ||\mathbf{x}_{\Delta} - \mathbf{x}||_{\Delta} \le c(e + \delta)$$

where e is defined in part (c) of theorem 3.2.

The proof of theorem 3.3 is essentially identical to that of theorem 5.3 in Weiss [9], so we only give an outline here.

### Proof outline for theorem 3.3

The basic task is to patch up the solution representations (3.37), (3.41) and (3.39) on the various segments of the mesh. The problem in transformed variables (3.23), (3.24), (or (3.34) for short) is considered, first under the boundary conditions

(3.56)  $P_{-1} = \eta_{-1} \in \mathbb{R}^{n_{-}}, P_{+1} = \eta_{+} \in \mathbb{R}^{n_{+}}, v_{1} = \eta_{0} \in \mathbb{R}^{m}.$ 

This corresponds to the differential problem (2.8), (2.10) which theorem 2.1 guarantees to be well-behaved for any given parameter vectors  $\eta_{-}, \eta_{+}$  and  $\eta_{0}$ .

Thus, the 3(n+m) components of the parametric representations (3.37), (3.41) and (3.39), i.e. of  $\zeta^{\Delta} := (\zeta_{-}^{I}, \zeta_{+}^{I}, \zeta_{0}^{I}, \zeta_{-}^{II}, \zeta_{0}^{III}, \zeta_{-}^{III}, \zeta_{+}^{III}, \zeta_{0}^{III}),$  are fixed by the 3(n+m) linear equations consisting of (3.56) plus the matching conditions

 $(3.57) \qquad \underset{\Delta}{\text{w}_{\Delta}^{I}}(t_{\underline{i}}) = \underset{\Delta}{\text{w}_{\Delta}^{II}}(t_{\underline{i}}) , \quad \underset{\Delta}{\text{w}_{\Delta}^{II}}(t_{\overline{i}}) = \underset{\Delta}{\text{w}_{\Delta}^{III}}(t_{\overline{i}}).$ 

In analogy to (2.11), the resulting  $3(n+m) \times 3(n+m)$  constraint matrix should have a uniformly bounded inverse for  $\delta$ , h and  $\kappa$  sufficiently small. Theorems 3.1 and 3.2 furnish us with information on the structure of key blocks of this matrix in (3.45), (3.49) and (3.51). Examining the resulting structure, it becomes apparent that the principal part of the constraint matrix is nonsingular iff the matrix

$$\begin{pmatrix} P_E^{-1}(t_{\underline{i}}) \\ P_{\underline{i}}E^{-1}(t_{\overline{i}}) \end{pmatrix}$$

is nonsingular. The condition for the latter to hold uniformly in  $\varepsilon$  is (3.54).

Now, theorem 2.3 guarantees that the exact solution w(t) has a similar decomposition and by (2.14) its parameter vector corresponding to  $\zeta^{\Delta}$ , which is determined so as to satisfy (2.10), can be written as

 $\zeta = (\zeta_{-}, 0, \zeta_{0}, 0, \zeta_{0}, 0, \tau_{+}, \zeta_{0}).$ 

The stability of the constraint matrix plus the convergence results of theorems 3.1 and 3.2 imply that

 $(3.58) \qquad ||\zeta - \zeta^{\Delta}|| \le c(e + \delta)$ 

Finally, the well-posedness of the problem guarantees, as described in §2, the safe transformation back to the x(t) variables with the choice of bounded  $n_+$ ,  $n_-$ ,  $n_0$  to satisfy the original boundary conditions. This completes the proof outline of theorem 3.3.

In practice, the transformation (3.22) or (2.5) is not needed and the difference equations are solved in the original variables. Thus, the matrices and vectors  $\Gamma_i$ ,  $g_i$ , of (3.14) are assembled according to (3.15) or (3.21) and the linear system of (n+m)(N+1) equations

(3.59)	Г - <sup>г</sup> 1	Ι					Т	[×1]	٢s	47
		- <sup>r</sup> 2	I					×2	=	2
								•		
				•	•					2
					236			•		3
						- 1 N	I	žΝ	2	I <sub>N</sub>
	Lв <sub>о</sub>						B	×N+1	E	

is solved for the superconvergent mesh values. The structure of the matrix in (3.59) is the familiar one for the trapezoidal or midpoint (box) scheme and is independent of k. The condition number of the matrix is a modest O(N) and in particular is independent of  $\varepsilon$  (cf. theorem 6.2 of Part I).

Indeed, it is a good practice in actual computation to roughly estimate the condition number of the above matrix for two values of  $\varepsilon$ , say. If that condition number seems to get large as  $\varepsilon$  decreases then something is "wrong": The mesh may be inadequate, or (3.54) does not hold or, perhaps most commonly, the differential problem is not well posed uniformly in  $\varepsilon$ . How to deal with the latter two cases will be discussed in a subsequent paper.

The meshes under consideration, which are described in fig. 1 and (3.32), (3.33), make sense from the simple point of view of approximating the solution profile independently of the differential equations. Thus, small subintervals are used where the solution varies fast and much larger subintervals are used elsewhere. However, the simple approximation point of view of the mesh may be somewhat misleading in the case where there is

no boundary layer, say at t = 0, even though n > 0, because the reduced solution happens to hit the boundary at the prescribed value(s) of  $\beta$ . In this case caution should be exercised, as the following discussion indicates.

Suppose for simplicity that the reduced solution, obtained by setting  $\varepsilon = 0$  in (1.1) - (1.3), satisfies (1.3), and that a uniform mesh with step size h >>  $\varepsilon$ , adequate to approximate the smooth solution well, is used in (3.5), (3.6). In §5 we show that setting  $\varepsilon = 0$  we get expressions (5.22) for Gauss points and (5.51) for Lobatto points. Here  $\underline{i} = 1$ ,  $\overline{i} = N+1$ , and we obtain

(3.60a)  $\bar{y}_{N+1} = (-1)^{kN} \bar{y}_1 + \phi_2$ 

for Gauss points and

(3.60b)  $\bar{y}_{N+1} = (-1)^{(k+1)N} A_{11}^{-1}(1)A_{11}(0)\bar{y}_1 + \phi$ 

for Lobatto points, where  $\phi$  is some inhomogeneity. Splitting the boundary condition matrices as

(3.61) 
$$B_0 = [B_0^y; B_0^z], B_1 = [B_1^y; B_1^z]$$

and substituting into the boundary conditions (3.16) gives the n+m equations

(3.62)  $B_y \bar{y}_1 + B_z \bar{z}_1 = \hat{\beta}$ where  $B_z$  and  $\hat{\beta}$  are appropriate quantities of no interest and

(3.63) 
$$B_{y} = \begin{cases} B_{0}^{y} + (-1)^{KN}B_{1}^{y} & \text{for Gauss points} \\ (k+1)N \\ B_{0}^{y} + (-1)B_{1}^{y}A_{11}^{-1}(1)A_{11}(0) & \text{for Lobatto points} \end{cases}$$

Now, clearly a necessary condition for (3.62) to have a unique solution is

(3.64) rank  $(B_y) = n$ 

This turns out not to be the case for some very simple examples, e.g., the

example in §6 of Part I.

If (3.64) holds then the uniform mesh is adequate for this problem with a smooth solution. The condition number of the matrix in (3.59) is then bounded independently of  $\varepsilon$ . However, if (3.64) does not hold then the reduced problem is singular. The condition number of the matrix in (3.59) is then, at best,  $O(h/\varepsilon)$ , and roundoff errors of this size are introduced in the course of computation.

There is a very simple remedy to this situation: Add to the mesh <u>one point</u>  $O(\varepsilon)$  away from each boundary where there is no layer (despite the existence of eigenvalues of the corresponding sign in  $A_{11}$ ). This one point layer mesh is sufficient for theorem 3.3 to hold: The more elaborate layer mesh of (3.46), (3.47) is constructed for accuracy reasons and therefore is not needed here.

# 4. Boundary layer regions

In this section we consider the linear problem (1.1), (1.2) on the subinterval [0,  $T_0 \varepsilon$ ], where a fine mesh satisfying (3.32) is assumed. Analogous results hold for the subinterval [1 -  $T_1 \varepsilon$ , 1]. Let

(4.1)  $h_i := \max\{h_i, 1 \le i \le N_0\}.$ 

Following Weiss [9] we consider the transformed system (2.6) - (2.7) for an easier analysis.

#### 4.1 Stability and results for smooth and auxiliary solution components

First, consider one equation

Solving the recurrence relation (4.4) we get

(4.6) 
$$y_{i+1} = \begin{bmatrix} i \\ \pi \\ \ell = 1 \end{bmatrix} \left( \frac{\lambda(t_{\ell})h_{\ell}}{\epsilon} \right) y_{1} + \frac{1}{\epsilon} \sum_{j=0}^{i-1} \prod_{\ell=i-j+1}^{i} \gamma(\frac{\lambda(t_{\ell})h_{\ell}}{\epsilon})h_{i-j \in i-j-i-j}$$

where  $\xi_{i-j}^{T} := \hat{b}^{T} \left[ \frac{\varepsilon}{\lambda(t_{i})h_{i}} - \hat{A} \right)^{-1} \hat{A} + I \right]$  is a bounded vector by (3.32). Now,

since the method is A-stable we have

 $(4.7a) |\gamma(\zeta)| \leq 1 \qquad re(\zeta) < 0.$ 

Furthermore, since  $\gamma'(0) = 1$  it follows that for any set S of the form

 $S \equiv S(\alpha_1, \alpha_2, \beta) = \{\zeta \mid 0 < |\zeta| \le \beta, \frac{\pi}{2} + \alpha_1 \le \arg \zeta \le \frac{3\pi}{2} - \alpha_2\}$ 

with  $\alpha_1, \alpha_2 > 0$ ,  $\alpha_1 + \alpha_2 \le \Pi$ ,  $\beta < \infty$ , there is a positive constant  $\mu = \mu(\alpha_1, \alpha_2, \beta)$  such that (4.7b)  $|\gamma(\zeta)| \le e^{\mu re(\zeta)}$ ,  $\zeta \in S$ 

By (3.32),

$$\left|\frac{\lambda(t_{\ell})h_{\ell}}{\varepsilon}\right| \leq |\lambda(t_{\ell})|K_{0} \leq \beta$$

for some well defined constant  $\beta$  of moderate size. Using (4.7b) we get

$$| \prod_{\substack{\ell=i-j+1}}^{i} \gamma(\frac{\lambda(t_{\ell})h_{\ell}}{\epsilon}) | \leq \prod_{\substack{\ell=i-j+1}}^{i} e^{-\mu \overline{\lambda}h_{\ell}/\epsilon} = e^{-\mu \overline{\lambda}(t_{i+1} - t_{i+1-j})/\epsilon}$$

and

$$\begin{split} \left| \frac{1}{\varepsilon} \sum_{j=0}^{i-1} [\prod_{\ell=i-j+1}^{i} \gamma(\frac{\lambda(t_{\ell})h_{\ell}}{\varepsilon})]h_{i-j} \right| &\leq \frac{1}{\varepsilon} \sum_{j=0}^{i-1} e^{-\mu\overline{\lambda}(t_{i+1}-t_{i+1-j})/\varepsilon} h_{i-j} \\ &\leq \frac{1}{\varepsilon} \int_{t_{1}}^{t_{i+1}} e^{-\mu\overline{\lambda}(t_{i+1}-s)/\varepsilon} ds &\leq \frac{1}{\mu\overline{\lambda}} \end{split}$$

So, substituting in (4.6) we get

(4.8) 
$$|y_{i+1}| \le |y_1| + c||f^c||$$

where

(4.9) 
$$c = (\bar{\lambda}_{\mu})^{-1} \max_{\substack{1 \le j \le i}} ||\xi_j||$$

It is straightforward to show that a similar result is obtained also for the Lobatto points.

This is the desired stability result for one equation. Next, consider the differential system (2.6) - (2.7) and its corresponding collocation discretization (3.23) - (3.24). <u>Theorem 4.1</u> The difference equations (3.23), (3.24) subject to the boundary conditions

 $\varepsilon_0$  is sufficiently small to enable a contraction argument below and depends on the bounds in lemma 3.1 and on max  $||\Lambda'(t)||$ . (To recall,  $0 \le t \le T_0 \varepsilon$  by  $u^C$  we mean the restriction of  $u_{\Delta}(t)$  to the mesh points plus the collocation points).

<u>Proof</u>: We consider the case for Gauss points; the case for Lobatto points is treated similarly. Our strategy is to consider first the simplified difference equations

 $(4.12) \quad \frac{\varepsilon}{h_{i}} \left( \underset{i \neq j}{u_{ij}} - \underset{i \neq l}{u_{i}} \right) = \sum_{l=1}^{k} \hat{a}_{jl} \{ \Lambda(t_{i}) \underset{i \neq l}{u_{il}} + f_{l}(t_{il}) \}$   $1 \le j \le k+1$ 

(4.13) 
$$\frac{1}{h_i} (v_{ij} - v_i) = \sum_{\ell=1}^{\kappa} \hat{a}_{j\ell} \{B_{22}(t_{i\ell}) v_{i\ell} + g_2(t_{i\ell})\}$$

where  $f_1(t_{il}) := g_1(t_{il}) + B_{12}(t_{il})v_{il}$ , and to treat the difference between (4.12) - (4.13) and (3.23) - (3.24) as a perturbation term of order  $h_1$ .

The components  $\{v_{ij}\}$  in (4.13), (4.10) are now completely separated from the components  $\{u_{ij}\}$ . For  $v_{\Delta}(t)$  the usual theory applies. This is a Runge-Kutta scheme for a non-stiff initial value problem, and certainly for  $\varepsilon$  small and  $h_L$  satisfying (3.32),  $v_{\Delta}(t)$  exists and satisfies (4.14)  $||v^{C}|| \le c\{||n_0|| + ||g_2||\}$ 

Now, for (4.12) note that since  $\Lambda(t)$  is diagonal, the vector system decouples into n scalar components. For each of the first n\_ components we

can apply the estimate (4.8) directly, since  $re(\lambda_j(t_i)) < 0$ ,  $1 \le j \le n_-$ . For the last  $n_+$  components,  $re(\lambda_j(t_i)) > 0$ , and we have to reverse the direction of integration, from right to left. Thus, for such a component, (4.8) is changed to read

(4.15)  $|y_i| \le |y_i| + c||f^c||$ 

which is compatible with the end conditions (4.10). We obtain that the difference equations (4.12) subject to (4.10) possess a solution  $u_{\tilde{\Delta}}$  satisfying

$$(4.16) \qquad ||u_{\Delta}||_{\Delta} \le ||n_{-}|| + ||n_{+}|| + c||f_{1}|| \le ||n_{-}|| + ||n_{+}|| + c_{1} \{||g_{\Delta}|| + ||n_{0}||\}$$

It is now easy to show a similar result for  $u_{ij}$  by expressing them in terms of  $u_i$  using (4.12).

This completes the treatment of the major part of the difference operator. Now, the equations (4.12) - (4.13) differ from (3.23) - (3.24) by terms of order  $h_i$  (or  $\varepsilon$ ) only, and a standard perturbation argument completes the proof

#### Q.E.D.

Now, with the stability result (4.11) and the linearity of the problem, part (a) of theorem 3.1 easily follows. Next, consider the "smooth" components  $W_{\Delta 0}(t)$  and  $w_{P\Delta}(t)$ . These correspond to the components in the exact solution decomposition which vary slowly across the boundary layer region. Substitution of  $W_{\Delta 0} - W_0$  and  $w_{P\Delta} - w_P$  into (4.11) immediately yields that

(4.17)  $||W_{\Delta 0}^{C} - W_{0}^{C}||, ||w_{P\Delta}^{C} - w_{P}^{C}|| \le ch_{L} = 0(\varepsilon)$ 

and this is really all we need. However, more can be obtained by applying the standard collocation analysis (Russell [7], Weiss [8]) to the original variables (i.e., analyzing the error in (3.5), (3.6)). After transforming

back to w, part (b) of theorem 3.1 is obtained.

Consider part (c) of theorem 3.1. We write

(4.18) 
$$W_{\Delta+}^{I} = F + G$$
,  $F = \begin{pmatrix} 0 \\ F_{+} \\ 0 \end{pmatrix}$ ,  $G = \begin{pmatrix} U_{+}^{I} \\ G_{+} \\ V_{+}^{I} \end{pmatrix}$ 

with F satisfying the homogeneous equations (4.12), (4.13), subject to

(4.19) 
$$F(t_{\underline{i}}) = \begin{pmatrix} 0 \\ I \\ 0 \end{pmatrix}$$

and G is the rest. The difference equations for F are again decoupled and so A-stability immediately implies that F is bounded. We now have to show that G is small. But comparing (3.36b) with what F satisfies, it is apparent that G satisfies the difference equations (3.23), (3.24), with inhomogeneous terms of size  $O(\varepsilon + h_2) = O(\varepsilon)$  and under homogeneous boundary conditions as in (3.36b). Using stability, part (c) of theorem 3.1 is proven.

# 4.2. Mesh selection in the layer regions

In §4.1 we have shown parts (a) - (c) of theorem 3.1. Here we treat the dominant components of the solution decomposition,  $W_{\Delta-}^{I}(t)$ . Analogous results for  $W_{\Delta+}^{III}(t)$  will be omitted.

First, consider one homogeneous equation with constant coefficients, (4.20)  $\varepsilon y' = \lambda y$ , y(0) = 1, with the solution  $y(t) = \exp(\frac{\lambda}{\varepsilon}t)$ , and denote  $\hat{\lambda} := -\operatorname{re}(\lambda) > 0$ . In Part I it was shown that, given a tolerance  $\delta < 1$ , the following mesh generates an approximation accurate to within this tolerance,

(4.21) 
$$h_1 := \frac{\varepsilon}{|\lambda|} c_p \delta^{1/p}, \quad c_p := \left[\frac{\hat{\lambda}}{|\lambda||c_{\gamma}|}\right]^{1/p}$$

 $(4.22) h_i := h_{i-1} \exp\{\frac{1}{p} \frac{\lambda}{\varepsilon} h_{i-1}\} = h_1 \exp\{\frac{1}{p} \frac{\lambda}{\varepsilon} t_i\}, i = 2, \dots, N_0.$ 

Here  $c_{\gamma}$  is a known constant depending only on p, and N<sub>0</sub> is determined so that  $t_{N_0+1} \ge T_0 \varepsilon > t_{N_0}$ . Since we would like the contribution of the fast decaying solution to be below  $\delta$  on the "long" interval  $[t_1, t_{\bar{1}}]$ , it is natural to set  $T_0$  so that

$$\exp(\frac{-\lambda}{\varepsilon} T_0 \varepsilon) = \delta$$

i.e.

(4.23) 
$$T_0 = \frac{1}{\hat{\lambda}}(-\ln\delta)$$

Note that the mesh defined by (4.21), (4.22) satisfies (3.32) because its steps are monotonically increasing and  $h_{\underline{i}} = c_{P/|\lambda|}\epsilon$ . Also, beyond  $t_{\underline{i}}$  the mesh becomes much sparser, depending only upon the accuracy needs for the reduced solution. Thus, the magnitude of  $|y(t_{\underline{i}})|$  is propagated essentially undamped by the numerical scheme outside the layer region.

Next we let  $\lambda$  in (4.20) vary as in (4.2), i.e.,  $\lambda(t) = \lambda(t_i)$ ,

 $t_i \le t < t_{i+1}$ . Then (cf. Part I, §4.2)

(4.24) 
$$y_{i+1} = \prod_{j=1}^{i} \gamma(\frac{\lambda(t_j)h_j}{\varepsilon}) = \prod_{j=1}^{i} \gamma(\frac{\lambda(0)h_j}{\varepsilon}) + R_i$$

where

(4.25) 
$$R_{i} = \prod_{j=1}^{i} \gamma(\frac{\lambda(0)h_{j}}{\varepsilon}) [\sum_{j=1}^{i} (1 + 0(\frac{t_{j}h_{j}}{\varepsilon})) - 1]$$

Lemma 4.1 The residue R<sub>i</sub> satisfies

 $(4.26) |R_i| \le c\varepsilon$ 

provided that  $\epsilon(\ln \delta)^2$  is bounded by a constant. Thus, the mesh (4.21), (4.22) with  $\lambda(0)$  replacing  $\lambda$  guarantees an approximation error of  $O(\delta) + O(\epsilon)$ .

<u>Proof</u>: Let  $t_{i+1} \leq -c_1 \epsilon \ln \delta$  with  $\delta = O(\epsilon)$ . As for (4.6) and the following arguments, we obtain

$$|\sum_{j=1}^{i} \gamma(\frac{\lambda(0)h_{j}}{\varepsilon})| \leq \exp(-\mu t_{i+1}/\varepsilon).$$

Writing

$$\exp(-\kappa t_j h_j / \epsilon) \le 1 + O(t_j h_j / \epsilon) \le \exp(\kappa t_j h_j / \epsilon)$$

we have

 $\exp(-\kappa t_{i+1}^2/2\epsilon) \leq \prod_{j=1}^{i} (1 + 0(t_j^h_j/\epsilon) \leq \exp(\kappa t_{i+1}^2/2\epsilon)$ So, by (4.25)

$$\begin{aligned} |R_{i}| &\leq \exp(-\mu t_{i+1}/\epsilon) [\exp(\kappa t_{i+1}^{2}/2\epsilon) - 1] \\ \text{Now, for } t_{i+1} &= -c_{1}\epsilon \ln\delta, \ \kappa t_{i+1}^{2}/2\epsilon &= 0(\epsilon(1\hbar\delta)^{2}) \text{ and so} \\ &\qquad \exp(ct_{i+1}^{2}/2\epsilon) - 1 \leq c_{2} \ t_{i+1}^{2}/2\epsilon \end{aligned}$$

Also

 $\exp(-\mu t_{i+1}/\epsilon) = c_3 \delta$ 

So

 $|R_i| \le C_4 \epsilon \delta (\ln \delta)^2 \le c \epsilon$ 

QED

Turning to the differential system (2.6), (2.7), we once again consider the difference equations (3.23), (3.24) as an O( $\varepsilon$ ) perturbation of (4.12), (4.13). The homogeneity and boundary conditions of (3.36a) plus the decoupling of (4.13) from (4.12) clearly imply that  $V_{\Delta_{-}}(t) \equiv 0$  and  $P_{+} \cup_{\Delta_{-}}(t) \equiv 0$ . Also, for each of the first n\_ components, the previous results for one equation apply provided that the mesh in (4.21), (4.22) is chosen according to it. Taking the most stringent of these choices will produce O( $\delta$ ) accuracy for all components. This is clearly achieved by the choice (3.46), (3.47). Part (d) of theorem 3.1 is then proven and hence, the proof of theorem 3.1 is complete.

The practical importance of using the mesh (3.46), (3.47) instead of, say, a uniform mesh has been demonstrated in table 4.2 of Part I. We now supplement this by some a priori estimates of  $N_0$ , the number of mesh points needed in the layer.

<u>Theorem 4.2</u> Asymptotically, for  $\varepsilon$  and  $\delta$  small, (4.27)  $N_0 = O(\delta^{-1/p}).$ 

Note that in (4.27) N<sub>0</sub> is independent of  $\epsilon$ . Also, a uniform mesh with T<sub>0</sub> given by (4.23) would yield N<sub>0</sub> =  $0(\delta^{-1/p}(-\ln\delta))$ .

Proof: It is sufficient to consider (4.21), (4.22). Then  

$$N_{0} = \sum_{i=1}^{N_{0}} h_{i}/h_{i} = \frac{1}{h_{1}} \sum_{i=1}^{N_{0}} h_{i} \exp\{-\frac{1}{p} \hat{\lambda}_{\varepsilon} t_{i}\} \approx \frac{1}{h_{1}} \int_{0}^{T_{0}\varepsilon} \exp\{-\frac{1}{p} \hat{\lambda}_{\varepsilon} t\} dt = -\frac{p\varepsilon}{\hat{\lambda}h_{1}} [1 - \exp\{-\frac{\hat{\lambda}T_{0}}{p}\}]$$

Substituting (4.23) for  $T_0$  and (4.21) for  $h_1$ ,

(4.28) 
$$N_0 \approx \frac{p}{c_p} \frac{|\lambda|}{\hat{\lambda}} (\delta^{-1/p} - 1) \approx \frac{|\lambda|}{\hat{\lambda}} pc_p^{-1} \delta^{-1/p}$$

This proves our claim.

QED

Further, the constants  $c_p$  of (4.21) can be shown to increase as p is increased (see Part I for  $|c_{\gamma}|$ ). Thus, the estimate (4.28) also indicates that for a given accuracy  $\delta$ ,  $N_0$  decreases as p (or k) is increased. Note that  $C_p$  also reflects a relative efficiency of higher order methods for problems where the eigenvalues have significant imaginary parts.

The mesh (3.46), (3.47) may be more demanding than necessary in case that eigenvalues of different magnitude are present in  $\Lambda_{(t)}$ . At a given t,  $0 \le t \le T_{0^{\varepsilon}}$ , the eigenvalue which imposes the smallest step size is the one for which the magnitude of the p-th derivative of the solution,  $(\frac{|\lambda|}{\varepsilon})^{p} \exp\{-\frac{\hat{\lambda}}{\varepsilon}t\}$ , is largest. Thus, if for instance,  $|\lambda_1| = \max\{|\lambda_j|, j = 1, ..., n_\}$ , then we can use (4.21) with  $\lambda := \lambda_1$  in place of (3.46) and then construct the mesh using (4.22) (with  $\lambda := \lambda_1$ ) until  $t_{i+1} \ge \hat{t}_1$ , where

(4.29) 
$$\hat{t}_{1} := \min{\{\hat{t}_{1j}; \hat{t}_{1j} > 0\}}, \quad \hat{t}_{ij} := \epsilon p \frac{re(\lambda_{j}) - re(\lambda_{1})}{\ln |\lambda_{1/\lambda_{j}}|}$$

Then, in case that  $\hat{t}_1 < T_0 \varepsilon$ , switch to  $\lambda := \lambda_k$  where k gives the minimum in (4.29) and continue with (4.22), etc. However, the overhead involved in constructing such a mesh is worth it only in special cases, as described above.

### 5. The interval away from boundaries

On the "long" interval  $[t_i, t_{\overline{i}}]$  we use the original problem variables and do not apply the transformation (2.5), because we can deal with the system (1.1), (1.2) directly in a simpler fashion. Thus our difference equations are (3.5), (3.6).

Consider first the "reduced" equations, where (3.5) is replaced by (5.1)  $0 = \sum_{\substack{k=1 \ l \neq 1}}^{k} \widehat{A_{11}(t_{i_k})} y_{i_k} + A_{12}(t_{i_k}) z_{i_k} + f_1(t_{i_k})$ }  $1 \le j \le k$ (For simplicity, assume that  $A_{rs}$  and  $f_r$  are independent of  $\varepsilon$ ,  $1 \le r$ ,  $s \le 2$ ). We will show that the solutions of (5.1), (3.6), denoted by  $\overline{y}_{\Delta}(t)$ ,  $\overline{z}_{\Delta}(t)$ , are bounded at the collocation points in terms of their data and approximate their analytic counterparts well at these points.

Let  $x(t) = \begin{pmatrix} y(t) \\ \tilde{z}(t) \end{pmatrix}$  be a smooth solution to the problem (1.1) - (1.2) on  $[t_{\underline{i}}, t_{\overline{i}}]$ . Recall that we can write

(5.2)  $y(t) = \overline{y}(t) + n(t)$ ,  $z(t) = \overline{z}(t) + z(t)$ where  $\overline{y}(t)$ ,  $\overline{z}(t)$  solve the "reduced" equations (5.3)  $0 = A_{11}\overline{y} + A_{12}\overline{z} + f_1$ 

$$t_i \leq t \leq t_i$$

(5.4)  $\bar{z}' = A_{21}\bar{y} + A_{22}\bar{z} + f_2$ (5.5)  $\bar{z}(t_1) = z(t_1)$ 

and, for any integer  $q \ge 0$ ,

(5.6) 
$$\begin{array}{c} q\\ j=0 \end{array} \left| \left| \frac{d^{j} n}{dt^{j}} \right| \right| = 0(\varepsilon) , \qquad \begin{array}{c} q\\ j=0 \end{array} \left| \left| \frac{d^{j} \zeta}{dt^{j}} \right| \right| = 0(\varepsilon) \\ j=0 \end{array} \right| \left| \frac{d^{j} \zeta}{dt^{j}} \right| = 0(\varepsilon) \qquad \begin{array}{c} t_{\underline{i}} \leq t \leq t_{\overline{i}} \\ \underline{i} \leq t \leq t_{\overline{i}} \end{array}$$

Now, for each j,  $1 \le j \le k$ , we can write by (5.1)

(5.7) 
$$\bar{y}_{ij} = -A_{11}(t_{ij})[A_{12}(t_{ij})\bar{z}_{ij} + f_1(t_{ij})].$$

Substituting into (3.6), this gives

$$(5.8) \qquad \frac{1}{h_{i}}(\bar{z}_{ij} - \bar{z}_{i}) = \sum_{\ell=1}^{k} \hat{a}_{j\ell} \left[A_{22}(t_{i\ell}) - A_{21}(t_{i\ell})A_{11}^{-1}(t_{i\ell})A_{12}(t_{i\ell})\right]\bar{z}_{i\ell} + \frac{f_{2}(t_{i\ell}) - A_{21}(t_{i\ell})A_{11}^{-1}(t_{i\ell})}{2} \int_{1}^{1} (t_{i\ell}) \int_{1}^{1} (t_{i\ell}) f_{1}(t_{i\ell}) dz_{i\ell} dz_{i\ell}$$

The solution  $\bar{z}_{\Delta}(t)$  is then the collocation solution of the non-stiff differential equations

(5.9)  $z' = [A_{22} - A_{21} A_{11}^{-1} A_{12}]z + [f_2 - A_{21} A_{11}^{-1} f_1]$   $t_1 \le t \le t_1$ which result from substituting (5.3) into (5.4), and the usual collocation theory implies not only that

(5.10) 
$$||\vec{z}^{c}|| \leq c_{1}(||\vec{z}_{\underline{i}}|| + ||f_{1}^{c}|| + ||f_{2}^{c}||)$$

but also that

(5.11)  $||\bar{z}_{\Delta} - \bar{z}||_{\Delta} = 0(h^p)$ ,  $||\bar{z}_{\Delta} - \bar{z}||_{c} = 0(h^{k+1}) + 0(h^p)$ if we take  $\bar{z}_{\underline{i}} := z(t_{\underline{i}})$ . Here p is defined in (3.42) and the seminorms are defined in (3.43). In the seminorms above and throughout this section, the superscript II is omitted.

Substituting into (5.7) we obtain similar results to (5.10), (5.11) for the fast components at the collocation points:

$$(5.12a) \qquad ||\bar{y}_{\Delta}||_{c} \leq c_{2}(||\bar{z}_{i}|| + ||f_{1}^{c}|| + ||f_{2}^{c}||)$$

(5.12b)  $||\bar{y}_{\Delta} - \bar{y}||_{c} = 0(\bar{h}^{k+1})$ 

In fact, if m = 0, i.e. there are no slow components z in (1.1), then from (5.7), (5.3),

(5.12c)  $\overline{y}_{ij} = \overline{y}(t_{ij})$   $1 \le i \le N, 1 \le j \le k.$ 

Now, for <u>Lobatto points</u>, (or any other set of collocation points with  $p_k = 1$ ), since the mesh points are also collocation points we get from (5.11), (5.3), (5.4),

(5.12d) 
$$||\bar{y}_{\Delta} - \bar{y}||_{\Delta} = O(h^p)$$

(5.12e) 
$$||\bar{y}_{A} - \bar{y}||_{A} = 0$$
 if  $m = 0$ 

This leads to the essential difference between Gauss and Lobatto points in part (c) of theorem 3.2

The above results, however, do not yield even general stability for

the fast components, without being supplemented. Thus we take a closer look at equations (3.5), (3.6) and their implementation when (5.13)  $\varepsilon \ll \underline{h} := \min \{\underline{h}_i; i \le i \le \overline{i}\}$ 

# 5.1 The case for Gauss points

Writing

(5.14)  $\hat{f}_{ij} := A_{12}(t_{ij})z_{ij} + f_1(t_{ij}) \hat{f}_i = (\hat{f}_{ij}, \dots, \hat{f}_{ik})^T$ we can consider (3.5) separately from (3.6), with  $\hat{f}_{i\ell}$  the inhomogeneous terms. As in (3.10) - (3.15) we obtain (5.15) $y_{i+1} = \Gamma_i y_i + g_i$ where  $\Gamma_i$  and  $g_i$  are given by (3.15), (3.11) - (3.13) with n replacing n+m,  $\varepsilon^{-1}A_{11}$  replacing A and  $\varepsilon^{-1}\hat{f}_i$  replacing  $f_i$ . We have (5.16)  $J_i^{-1} = (I - h_i \varepsilon^{-1} D_{A_{11}} (\hat{A} \otimes I))^{-1} = \varepsilon h_i^{-1} (\varepsilon h_i^{-1} I - \hat{G}_i)^{-1}, \ \hat{G}_i = D_{A_{11}} (\hat{A} \otimes I)$ As in equation (4.15) of Part I we write for the nonsingular matrix  $\hat{G}_{i}$  $(\varepsilon h_i^{-1}I - \hat{G}_i)^{-1} = (\varepsilon h_i^{-1}\hat{C}_i - I)\hat{G}_i^{-1}$ (5.17)provided that  $\epsilon < h_i ||\hat{G}_i^{-1}||^{-1}$ . Now, (5.18a)  $\hat{B}\hat{G}_{i}^{-1}C_{A_{11}} = \hat{B}(\hat{A}\otimes I)^{-1}D_{A_{11}}^{-1}C_{A_{11}} = (\hat{b}^{T}\hat{A}^{-1}I)I,$ where  $I = (1, ..., 1)^{T} \in \mathbb{R}^{k}$ . From (4.17) of Part I and (5.18a), (5.18b) I -  $\hat{B}\hat{G}_{1}^{-1}C_{A_{11}} = (1 - \hat{b}^{T}\hat{A}^{-1}\hat{1})I = (-1)^{k}I$ and this is the leading term of  $\Gamma_i$ . We get in (5.15)  $y_{i+1} = (-1)^{k}y_{i} + \varepsilon h_{i}^{-1}\hat{H}_{i}y_{i} + \hat{B}(\varepsilon h_{i}^{-1}\hat{C}_{i} - I)\hat{G}_{i}^{-1}\hat{f}_{i}$ (5.19)where the matrix  $\hat{H}_{i}$  is bounded and depends on  $\hat{b}$ ,  $\hat{A}$  and  $A_{11}(t_{ij})$ , j = 1, ..., k. Clearly, both  $\hat{H}_i$  and the matrix multiplying  $\hat{f}_i$  in (5.19) vary smoothly with i. Further, since (5.20)  $y_{ii} = A_{11}^{-1}(t_{ii})(-\hat{f}(t_{ii}) + \varepsilon F_{ii})$ (cf. (3.9)) we get for

(5.21a) 
$$\hat{y}_{i} := (y_{i1}, \dots, y_{ik})^{T}$$

the equations

(5.21b) 
$$\hat{y}_i = -D_{A_{11}}^{-1} \hat{f}_i - \varepsilon h_i^{-1} D_{A_{11}}^{-1} (\varepsilon h_i^{-1} \hat{C}_i - I) \hat{G}_i^{-1} (\hat{f}_i + C_{A_{11}} y_i).$$
  
Now, again set  $\varepsilon = 0$ . Then we get for the "reduced" solution by (5.19)

(5.22) 
$$\bar{y}_{i} = (-1)^{k(i-1)} \bar{y}_{i} - \sum_{j=i}^{i-1} (-1)^{k(i-1-j)} \hat{B} \hat{G}_{j-j}^{-1} \hat{f}_{j-j}$$

This general form of the "reduced" solution yields the following stability result.

<u>Lemma 5.1</u>. The difference equations (5.1), (3.6), using Gauss points for  $\bar{y}_{\Delta}$ ,  $\bar{z}_{\Delta}$  with  $\bar{y}_{\underline{i}}$  and  $\bar{z}_{\underline{i}}$  specified, have a unique solution provided that h is small enough. This solution satisfies (5.10) - (5.12) and

$$(5.23) \qquad ||\bar{y}_{\Delta}||_{\Delta} \leq ||\bar{y}_{i}|| + c\{||z_{i}|| + ||f_{1}^{c}|| + ||f_{2}^{c}|| + |\hat{f}_{2}^{c}|| + \hat{i} \\ \hat{i} \\ \Sigma \\ i = 0 \\ j = 1 \\ i = 0 \\ j = 1 \\ i = 1 \\ i = 0 \\ j = 1 \\ i = 1 \\$$

where i is the integral part of  $\frac{1}{2}(i - i)$ .

<u>Proof</u>: Firstly, note that by (5.14) and the stability result (5.10) for the slow components,

$$(5.24) \qquad ||\hat{f}_{j}|| \le c(||z_{\underline{i}}|| + ||f_{1}^{c}|| + ||f_{2}^{c}||$$

Next, we distinguish between two cases.

I. <u>k is odd</u>. We have (5.25)  $\hat{B}_{\Sigma}^{i-1}(-1)^{i-1-j}\hat{G}_{j}^{-1}\hat{f}_{j} = \hat{B}_{\Sigma}^{i-1}(\hat{G}_{j}^{-1}\hat{f}_{j} - \hat{G}_{j-1}^{-1}\hat{f}_{j-1})$ 

where j takes on only even or only odd values, with a possible additional end term. Each term in the right hand sum of (5.25) is written as (5.26)  $\hat{G}_{j}^{-1}\hat{f}_{j} - \hat{G}_{j-1}^{-1}\hat{f}_{j-1} = \hat{G}_{j}^{-1}(\hat{f}_{j} - \hat{f}_{j-1}) + (\hat{G}_{j}^{-1} - \hat{G}_{j-1}^{-1})\hat{f}_{j-1}$  The second term in (5.26) is bounded by  $O(h_j)||\hat{f}_{j-1}||$  and so, even after we sum these quantities up in (5.25), the bound of (5.24) holds. For the first term on the right hand side of (5.26) we have to examine expressions of the form

(5.27) 
$$[A(t_{j\ell})\bar{z}_{j\ell} - A(t_{j-1,\ell})\bar{z}_{j-1,\ell}] + [f_1(t_{j\ell}) - f_1(t_{j-1,\ell})]$$

The term in the first square brackets is again bounded in norm by  $O(h_j)||\bar{z}_{\Delta}^{c}||$ , in view of (5.11). The remaining term gives rise to the rightmost term in the bound (5.23).

II. <u>k is even</u>. Here we have to examine the sum  $\hat{B}_{\Sigma}^{i-1} \hat{G}_{j-1}^{-1} \hat{f}_{j}$  and make sure that it does not grow like  $h^{-1}$ , despite the fact that its terms do not alternate in sign. But by (5.16), (5.18b),

(5.28) 
$$\hat{B}\hat{G}_{j}^{-1} = \hat{B}(\hat{A}\otimes I)^{-1}D_{A_{11}}^{-1}, \quad \hat{b}^{T}\hat{A}^{-1}I = 0.$$

$$(5.29) \qquad \hat{B}\hat{G}_{j}^{-1}\hat{f}_{j} = \hat{B}(\hat{A}\otimes I)^{-1} \begin{pmatrix} A_{11}^{-1}(t_{j1})\hat{f}_{j1} \\ \vdots \\ A_{11}^{-1}(t_{jk})\hat{f}_{jk} \end{pmatrix} = \hat{B}(\hat{A}\otimes I)^{-1} \begin{bmatrix} A_{11}^{-1}(t_{j1})\hat{f}_{j1} \\ \vdots \\ A_{11}^{-1}(t_{jk})\hat{f}_{jk} \end{pmatrix} \\ - \begin{bmatrix} A_{11}^{-1}(t_{j1})\hat{f}_{jk} \\ A_{11}^{-1}(t_{j1})\hat{f}_{jk} \end{bmatrix} \\ A_{11}^{-1}(t_{j1})\hat{f}_{j1} \\ A_{11}^{-1}(t_{j1})\hat{f}_{j1} \end{bmatrix} \end{bmatrix}$$

This brings us to examine again terms like (5.27) and the remainder of the proof is, therefore, as for the case when k is odd

QED

Consider now the approximation error  $\bar{e}(t) = \bar{y}_{\Delta}(t) - \bar{y}(t)$ , where  $\bar{y}(t)$ is defined in (5.2) - (5.5). This error satisfies the difference equations (5.1), (3.6) (with  $\bar{z}_{\Delta} - \bar{z}$  replacing  $\bar{z}_{\Delta}$ ) with the local truncation error, appropriately represented, as the inhomogeneous term. Choosing  $\bar{y}_{\Delta}(t_{\underline{i}}) := \bar{y}(t_{\underline{i}})$  we obtain for the error  $\bar{e}_{i}$  at mesh points, by (5.22), (5.30)  $\bar{e}_{i} = \sum_{j=\underline{i}}^{i-1} (-1)^{k(i-1-j)} \kappa_{j} h_{j}^{k+1}$ with  $K_{j}$  varying smoothly with j, i.e. (5.31)  $K_{\underline{j}} - K_{\underline{j}-1} = 0(h_{\underline{j}})$ (This is well-known to hold for the principal term of the local truncation error). Thus, in general (5.32)  $||\bar{y}_{\Delta} - \bar{y}||_{\Delta} = 0(h^{k})$ and a sharper estimate is obtained when k is odd and (3.52) holds. For then we can arrange the sum in (5.30) in pairs of terms of different sign and equal value up to 0(h). Thus we get in this case

(5.33) 
$$||\bar{y}_{\Delta} - \bar{y}||_{\Delta} = 0(h^{k+1})$$

Note that at the collocation points we always have this better estimate.

Next, consider the "full" scheme (3.5), (3.6). We write the collocation solution for (1.1), (1.2) as

(5.34)  $y_{\Delta}(t) = \overline{y}_{\Delta}(t) + \eta_{\Delta}(t), \quad z_{\Delta}(t) = \overline{z}_{\Delta}(t) + \zeta_{\Delta}(t)$ choosing  $\overline{y}_{\underline{i}} = \overline{y}(t_{\underline{i}}), \quad \overline{z}_{\underline{i}} = \overline{z}(t_{\underline{i}}).$  Thus,  $\eta_{\underline{i}} = \eta(t_{\underline{i}}) = 0(\varepsilon), \quad \zeta_{\underline{i}} = \zeta(t_{\underline{i}}) = 0.$ Substituting in (3.5), (3.6) we get for  $\eta_{\Delta}$  and  $\zeta_{\Delta}$  (note that  $\overline{y}(t)$  is smooth and (5.12b), (5.32) hold),

$$(5.35) \qquad \sum_{\substack{k=1 \\ k=1}}^{k} \hat{a}_{jk} \{A_{11}(t_{ik})_{ik}^{n} + A_{12}(t_{ik})_{ik}^{n}\} = \varepsilon h_{i}^{-1}(n_{ij} - n_{i} + \bar{y}_{ij} - \bar{y}_{i}) = \varepsilon h_{i}^{-1}(n_{ij} - n_{i}) + 0(\varepsilon)$$

$$1 \le j \le k$$

(5.36) 
$$\sum_{\ell=1}^{k} \hat{a}_{j\ell} \{A_{21}(t_{i\ell}) = A_{22}(t_{i\ell}) \leq i\ell \} = h_{i}^{-1} (z_{ij} - z_{i})$$

(5.37) 
$$\eta_{\underline{i}} = 0(\varepsilon), \quad \zeta_{\underline{i}} = 0.$$

Now, the right hand side terms of (5.35) are considered as inhomogeneities (to conform to (5.1), (3.6) we multiply them as a vector by  $\hat{A} \hat{A}^{-1}$ ; this

does not change anything essential) and lemma 5.1 is applied. By (5.23), (5.10), (5.21),

- (5.38)  $||n_{\Delta}^{c}|| \leq c \in \{ \sum_{i=1}^{i} h_{i}^{-1} ||n_{\Delta}^{c}|| + h^{-1} ||\bar{y}_{\Delta}^{c}|| \}$
- (5.39)  $||\underline{\varsigma}^{c}_{\Delta}|| \leq c \varepsilon h^{-1} (||\underline{\eta}^{c}_{\Delta}|| + ||\underline{\tilde{y}}^{c}_{\Delta}||)$ Hence, if  $\overline{c} := c \varepsilon \sum_{\Sigma} h_{i}^{-1} < 1$  then, with  $\kappa := \varepsilon h^{-1}$ , i = i
- (5.40)  $\left|\left|\underset{\sim}{\eta}_{\Delta}^{\mathsf{C}}\right|\right| \leq \frac{\mathsf{C}\kappa}{1-\bar{\mathsf{C}}} \left|\left|\underset{\sim}{\bar{\mathsf{y}}}_{\mathsf{C}}^{\mathsf{C}}\right|\right| , \left|\left|\underset{\sim}{\zeta}_{\Delta}^{\mathsf{C}}\right|\right| \leq \frac{\mathsf{C}\kappa}{1-\bar{\mathsf{C}}} \left|\left|\underset{\sim}{\bar{\mathsf{y}}}_{\Delta}^{\mathsf{C}}\right|\right|.$

A combination of (5.40) with (5.23), (5.10) and (5.12a) now gives the desired stability result for the collocation solution  $y_{\Delta}(t)$ ,  $z_{\Delta}(t)$ . Part (a) of theorem 3.2 then follows for the Gaussian points.

Next, consider the first n fundamental solution components. The collocation approximations  $Y_{\Delta}(t)$ ,  $Z_{\Delta}(t)$  are defined by the homogeneous difference schemes (3.5), (3.6) and the initial conditions (3.40b). Thus, for the "reduced" problem with  $\varepsilon = 0$  we get by (5.8), (5.19), (5.41)  $\overline{Z}_{\Delta}(t) \equiv 0$ ,  $\overline{Y}_{\Delta}(t_i) = (-1)^{k(i-\underline{i})}I$   $\underline{i} < i \leq \overline{i}$ . Furthermore, by repeating the argument leading to (5.40) it becomes clear

that  $Y_{\Delta}(t_i)$ ,  $Z_{\Delta}(t_i)$  are only  $O(\kappa)$  away from  $\bar{Y}_{\Delta}(t_i)$  and  $\bar{Z}_{\Delta}(t_i)$ . Thus, applying the transformation (3.22), (3.51a) is obtained.

Finally, for a smooth solution  $x(t) = \begin{pmatrix} y(t) \\ \tilde{z}(t) \end{pmatrix}$  which may be a transformation of  $w_p(t)$  or of a column of  $W_0(t)$ , consider the approximation error. We write

(5.42) 
$$y_{\Delta}(t) - y(t) = (\bar{y}_{\Delta}(t) - \bar{y}(t)) + (n_{\Delta}(t) - n(t)), z_{\Delta}(t) - z(t) = (\bar{z}_{\Delta}(t) - \bar{z}(t)) + (\bar{z}_{\Delta}(t) - \bar{z}(t))$$

(see (5.2), (5.34)) and estimate the quantities on the right hand sides. We already have (5.11), (5.12b), (5.32) and (5.33). For the remaining terms we write the difference equations as in (5.35) - (5.37) with  $n_{il} - n(t_{il})$  replacing  $n_{il}$  etc. Additional local error terms of size  $O(\epsilon h_i^k)$  are easily incorporated into the right hand sides and the initial conditions corresponding to (5.37) are homogeneous. Thus, a result similar to (5.40) is obtained for the errors. This completes the proof of theorem 3.2 for Gauss points.

### 5.2 The case for Lobatto points

For Lobatto points we proceed in a similar way as above, establishing stability first. With  $\hat{f}_{ij}$  as in (5.14) we obtain (5.15) where  $\Gamma_i$  and  $g_i$ are given by (3.21), (3.20), with n replacing n+m,  $\epsilon^{-1}A_{11}$  replacing A and  $e^{-1}\hat{f}_i$  replacing  $f_i$ . Then,  $\overline{J}_i^{-1} = \varepsilon h_i^{-1} (\varepsilon h_i^{-1} I - \overline{G}_i)^{-1}$ ,  $\overline{G}_i = (\overline{A} \otimes I) \overline{D}_{A_{11}}$ (5.43) and we write, as in (5.17), provided that  $\epsilon$  <  $h_i^{}||\bar{G}_i^{-1}||^{-1}$  , (5.44)  $(\varepsilon h_i^{-1}I - \bar{G}_i)^{-1} = (\varepsilon h_i^{-1}\bar{C}_i - I)\bar{G}_i$ . The last rows of  $\bar{G}_i^{-1}$  are (5.45)  $A_{11}^{-1}(t_{i+1})(\hat{a}_{k2}I,\ldots,\hat{a}_{kk}I)$ where  $(\hat{a}_{k2}, \dots, \hat{a}_{kk})$  is the last row of  $\overline{A}^{-1}$  and hence  $-\sum_{\substack{k=2\\ k=2}}^{k} \hat{a}_{kk} \hat{a}_{kl} = \gamma(-\infty) = (-1)^{k+1}$  (cf. equations (4.19), (4.20) of Part I). Substitution in (3.21) yields  $y_{i+1} = (-1)^{k+1} A_{11}^{-1}(t_{i+1}) A_{11}(t_i) y_i + \varepsilon h_i^{-1} \bar{H}_i y_i + (\varepsilon h_i^{-1} \bar{C}_i - I) \bar{Q}_i \hat{f}_i$ (5.46)where  $\bar{H}^{}_i$  and  $\bar{Q}^{}_i$  are bounded matrices independently of  $\epsilon$ , which vary smoothly with i. For the other collocation points we obtain in precisely the same way (cf. (3.18))

(5.47) 
$$y_{ij} = c_{ij}A_{11}^{-1}(t_{ij})A_{11}(t_i)y_i + \varepsilon h_i^{-1}\bar{H}_{ij}y_i + (\varepsilon h_i^{-1}\bar{C}_i - I)\bar{Q}_{ij}f_i$$
  
 $2 \le j \le k - 1$ 

for appropriate constants  $c_{ij}$  and bounded matrices  $\bar{H}_{ij}$ ,  $\bar{Q}_{ij}$ .

Now, the equivalent of lemma 5.1 is covered already in (5.12a), and in a preferable way: Here for the "reduced" solution no  $\bar{y}_1$  values are specified (like for the reduced solution of the continuous problem) and no factor like the sum on  $f_1$  values in (5.23) appears. Then, writing (5.34) for the "full" solution, followed by (5.35), (5.36) and (5.48)  $z_1 = 0$ , we obtain

(5.49)  $\begin{aligned} \left| \left| \begin{array}{c} \left| \begin{array}{c} n_{\Delta}^{\mathbf{C}} \right| \right| &\leq c \ \varepsilon \ \underline{h}^{-1} \left( \ \left| \left| \begin{array}{c} n_{\Delta}^{\mathbf{C}} \right| \right| + \left| \left| \begin{array}{c} \overline{y}_{\Delta}^{\mathbf{C}} \right| \right| \right) \\ \\ \left| \left| \begin{array}{c} \zeta_{\Delta}^{\mathbf{C}} \right| \right| &\leq c \ \varepsilon \ h^{-1} \left( \left| \left| \begin{array}{c} n_{\Delta}^{\mathbf{C}} \right| \right| + \left| \left| \begin{array}{c} \overline{y}_{\Delta}^{\mathbf{C}} \right| \right| \right) \end{aligned} \right. \end{aligned}$ 

Hence if  $c_{\kappa} < 1$  then (5.40) is obtained. The distinction from the case for Gauss points is, however, that the condition relating  $\varepsilon$  and the mesh in order to enable the contraction argument is more pleasant here.

Next, consider the collocation approximation to  $Y_{\Delta}(t)$ ,  $Z_{\Delta}(t)$ . The solution is an  $O(\kappa)$  perturbation of the collocation solution to the "reduced" problem, where we consider the homogeneous equations (5.1), (3.6) subject to the initial conditions (3.40b). By (5.8) and (5.46),

(5.50)  $\bar{Z}_{\Delta}(t) \equiv 0$ ,  $\bar{Y}_{\Delta}(t_{i}) = (-1)^{(k+1)(i-1)}A_{11}^{-1}(t_{i})A_{11}(t_{i})$ Applying the transformation (3.22) and noting that  $E^{-1}A_{11}^{-1} = \Lambda^{-1}E^{-1}$ , we obtain (3.51b).

Finally, consider the approximation error for a smooth solution x(t). Writing the error as in (5.42), we need to consider  $v_{i\ell} := n_{i\ell} - n(t_i)$ . Again writing the difference equations as in (5.35), (5.36), (5.48) for the errors in  $n_{i\ell}$  and  $z_{i\ell}$ , we obtain the "reduced" form of the difference equations with inhomogeneous terms given by the local truncation error terms of size  $0(\epsilon h_i^k)$  with a smoothly varying principal error function. This is precisely as for the Gauss points, except that there the error in  $\eta$  is always dominated by the error in  $\tilde{y}$  (given that  $\kappa \ll 1$ ), while here the error in  $\eta$  may sometimes dominate, see (5.12d), (5.12e).

The error  $v_i = n_i - n(t_i)$  is obtained analogously to  $\overline{e}_i$  for Gauss points. We write the general solution to (5.46) for  $\epsilon = 0$  as

$$(5.51) \qquad \bar{y}_{i} = (-1)^{(i-\underline{i})(k+1)} A_{11}^{-1}(t_{i}) A_{11}(t_{\underline{i}}) \bar{y}_{\underline{i}} - \sum_{j=\underline{i}}^{i-1} (-1)^{(i-j-1)(k+1)} A_{11}^{-1}(t_{j+1}) \bar{Q}_{j} \hat{f}_{j}$$

and obtain for  $v_i$  upon substitution

(5.52) 
$$v_{i} = \sum_{j=i}^{i-1} (-1)^{(k+1)(i-j-1)} K_{j} \varepsilon h_{j}^{k}$$

with  $\underset{\sim}{K}_{j}$  varying smoothly with j as in (5.31). Thus (5.53)  $||n_{\Delta} - n||_{\Delta} \le c \varepsilon h^{k-1}$ 

and, in case that k is even and the mesh is locally almost uniform, (5.54)  $||_{\tilde{n}_{\Delta}} - \underline{n}||_{\Delta} \le c \varepsilon h^{k}$ .

The remaining error in the slow components clearly satisfies similar bounds. This completes the proof of theorem 3.2.

# 6. Numerical examples

The following numerical results were computed on an Amdahl 470-V/8 computer with a 14-hexadecimal-digits mantissa. The notation  $a-b \equiv ax10^{-b}$  is used throughout.

Example (Hemker [3]). Consider

(6.1a) 
$$\varepsilon u'' + (2 + \cos \pi t)u' - u = f(t)$$
  $0 \le t \le 1$ 

where

(6.1b) 
$$f(t) = -(1 + \varepsilon \pi^2) \cos \pi t - \pi (2 + \cos \pi t) \sin \pi t + (1 - \alpha + \frac{3}{2\varepsilon} \pi^2 t^2) e^{-3t/\varepsilon}$$

subject to

(6.1c) 
$$u(0) = \alpha$$
,  $u(1) = -1$ .

The solution is

(6.2) 
$$u(t) = \cos \pi t + (\alpha - 1)e^{-3t/\epsilon} + 0(\epsilon^2)$$

Thus, when  $\alpha \neq 1$  we have a boundary layer at t = 0 only.

When converting to a first order system note that if we use the usual variables u, u', then the problem does not have a bounded inverse (since  $u'(0) \sim '/\epsilon$ ). Instead we integrate once, as in Kreiss and Kreiss [4], Kreiss and Nichols [5], obtaining with y:= u the system

(6.3a) 
$$\varepsilon y' = -(2 + \cos \pi t)y + z$$

(6.3b)  $z' = (1 - \pi \sin \pi t)y + f(t)$ 

(6.3c) 
$$y(0) = \alpha, y(1) = -1.$$

So our matrix  $A_{11}$  is a negative scalar function of t here.

First we choose  $\alpha$  = 1 and use uniform meshes (clearly (3.64) holds). The results are listed in table 1, where under "E" we list the maximum error at mesh points and under "rate" the measured convergence rate in h. The results for Gauss and Lobatto points confirm part (c) of theorem 3.2. In addition, we list for comparison numerical results using collocation at Table 1: Example ! with a smooth solution throughout,  $\varepsilon = 10^{-10}$ 

		Gauss p	oints				<u>Radau</u> p	points				Lobatto	points	
<u>~</u>	z	cond	шţ	rate	<u>~</u>	Z	cond	ш	rate	<u>~</u> ]	2	cond	шļ	rate
-	10	.23+3	.64-1		-	10	.62+2	.20		2	10	.40+3	.65-1	
	20	.45+3	.16-1	2.0		20	.12+3	.10	6.0		20	.76+3	.17-1	2.0
	40	.88+3	.40-2	2.0		40	.25+3	.53-1	1.0		40	.15+4	.43-2	2.0
2	10	.87+2	.47-2		2	10	.66+2	.33-3		e	10	.16+3	.30-4	
	20	.16+3	.12-2	2.0		20	.13+3	.40-4	3.0		20	.29+3	.19-5	4.0
	40	.31+3	.29-3	2.0		40	.25+3	.49-5	3.0		40	.55+3	.12-6	4.0
e	10	.23+3	.16-3		ę	10	.66+2	.18-5		4	10	.40+3	.41-6	
	20	.45+3	.98-5	4.0		20	.13+3	.54-7	5.0		20	.76+3	.68-8	5.9
	40	.88+3	.61-6	4.0		40	.25+3	.17-8	5.0		40	.15+4	.11-9	6.0
4	10	.88+2	.88-5		4	10	.66+2	.21-8		2	10	.16+3	.70-10	
	20	.16+3	.55-6	4.0		20	.13+3	.17-10	7.0		20	.29+3	.28-12	8.0
	40	.31+3	.34-7	4.0		070	.25+3	.13-12	7.0		40	.55+3	.12-13	*

\* rate polluted by roundoff errors

the unsymmetric Radau points (see Part I). The usage of the latter schemes is possible here because all the eigenvalues of  $A_{11}$  have the same sign in their real part. For the examples discussed in Weiss [9] or in §6 of Part I, for instance, the Radau schemes are unstable unless the transformation 2.5 is explicitly applied (and this time not just for analysis) and the schemes are upwinded. Therefore, we stay with the symmetric schemes.

Next we set  $\alpha = 0$ , obtaining a steep boundary layer near t = 0. Results are listed in table 2. Here the meshes are constructed by taking the corresponding meshes of table 1 and adding a layer mesh according to (4.21), (4.22) with  $\hat{\lambda} = -\lambda = 3$ . The accuracy tolerance  $\delta$  is chosen to be just below the smooth solution error for the finest mesh in table 1, for each scheme. Here we list under "E" the maximum error on all mesh points with "rate" the convergence rate in the maximum mesh width h. Note the relatively small number of mesh points needed to achieve high accuracy with the higher order schemes, particularly of Lobatto type.

Also listed in table 2 are some results when  $\delta \ll \epsilon \ll 1$ . This is not covered by our analysis (see part (d) of theorem 3.1), because we are primarily concerned in this paper with what happens when  $\epsilon \rightarrow 0$ . However, as indicated by the numerical results, the analysis can be extended to cover this case as well. Indeed, when  $\delta \ll \epsilon$  then a denser mesh is constructed in the layer regions and this only makes the situation more regular.

Other examples have been tried as well. In particular, numerical solutions for the example in Weiss [9], which for some parameter values violates condition (3.54), have been computed. Their behaviour is similar to that reported in [9] for the midpoint and trapezoidal schemes and their discussion is therefore omitted.

Table 2: Example 1 with a boundary layer,  $\alpha = 0$ 

Gauss points

<u>Gauss points</u>						Lobatto points					
<u>a</u>	<u>k</u>	<u>ð</u>	N	E	<u>rate</u>	k	<u>ð</u>	N	E	rate	
10-10	1	13	32	.21-1		2	13	32	.13-1		
			42	.54-2	2.0			42	.32-2	2.0	
			62	.15-2	1.8			62	.80-3	2.0	
	2	14	20	.63-2		3	17	57	.22-4		
			30	.16-2	2.0			67	.13-5	4.0	
			50	.39-3	2.0			87	.82-7	4.0	
	3	17	26	.10-3		4	110	54	.75-7		
			36	.62-5	4.1			64	.11-8	6.0	
			56	.39-6	4.0			84	.10-9	3.5	
	4	18	22	.12-4		5	110	30	.11-9		
			32	.73-6	4.0			40	.70-10		
			52	.45-7	4.0			60	.70-10		
10-4	3	17	25	.10-3		3	17	56	.20-4		
		/////R	35	.62-5	4.1			66	.11-5	4.2	
			55	.38-6	4.0			86	.86-7	3.7	
	4	18	21	.12-4		4	110	53	.61-7		
			31	.66-6	4.1			63	.11-8	5.7	
			51	.26-7	4.6			83	.94-10	3.6	

#### REFERENCES

- U. Ascher, S. Pruess and R.D. Russell, "On spline basis selection for solving differential equations", (1981) To appear in SIAM J. Numer. Anal.
- U. Ascher and R. Weiss, "Collocation for singular perturbation problems I: First order systems with constant coefficients". (1981) To appear in SIAM J. Numer. Anal.
- P.W. Hemker, "A numerical study of stiff two-point boundary value problems", Math. Centrum, Amsterdam (1977).
- 4. B. Kreiss and H.O. Kreiss, "Numerical methods for singular perturbation problems", SIAM J. Numer. Anal. 18 (1981), 262-276.
- H.O. Kreiss and N. Nichols, "Numerical methods for singular perturbation problems", Uppsala University, Dept. of Computer Science Report #57 (1975).
- P.A. Markowich and C.A. Ringhofer, "Collocation methods for boundary value problems on 'long' intervals", (1981) To appear in Math. Comp.
- R.D. Russell, "Collocation for systems of boundary value problems", Numer. Math. <u>23</u> (1974), 119-133.
- R. Weiss, "The application of implicit Runge-Kutta and collocation methods to boundary-value problems", Math. Comp. 28 (1974), 449-464.
- 9. R. Weiss, "An analysis of the box and trapezoidal schemes for linear singularly perturbed boundary value problems", Manuscript 1982.