

On Error Structures and Extrapolation for Stiff Systems, with Application in the Method of Lines

W. Auzinger, Wien

Received November 7, 1989

Abstract — Zusammenfassung

On Error Structures and Extrapolation for Stiff Systems, with Application in the Method of Lines. In this paper, which carries on the considerations in [1], the structure of the global error is studied for some time discretization schemes, applied to a class of stiff initial value problems as they typically arise from the semi-discretization of parabolic initial/boundary value problems (method of lines). The implicit Euler and trapezoidal schemes and a locally one-dimensional splitting method are considered, and 'perturbed' asymptotic error expansions are derived which are valid independent of the stiffness (independent of the meshwidth in space). The key point within the analysis is a careful, quantitative description of the remainder term in such an expansion. The results are applicable in the method of lines setting and enable the prediction of the behavior of extrapolation algorithms for the class of problems under consideration. These theoretical considerations are illustrated by numerical examples.

AMS Subject Classifications: 65L05, 65M20

Key words: Stiff systems, method of lines, error structures, extrapolation

Über Fehlerstrukturen und Extrapolation bei steifen Systemen, mit Anwendung bei der Linienmethode. In dieser Arbeit wird, als Fortführung der Betrachtungen in [1], die Struktur des globalen Fehlers einiger Zeit-Diskretisierungsschemata bei Anwendung auf eine Klasse steifer Anfangswertprobleme studiert, wie sie typischerweise bei der Semi-Diskretisierung parabolischer Anfangs/Randwertprobleme auftreten (Linienmethode). Im Mittelpunkt der Betrachtungen stehen das implizite Eulerverfahren, die implizite Trapezregel und ein lokal eindimensionales Zwischenschritt-Schema, und es werden 'gestörte' asymptotische Entwicklungen hergeleitet, deren Gültigkeit unabhängig von der Steifheit (unabhängig von der Gitterfeinheit der Raum-Diskretisierung) besteht. Der entscheidende Punkt in der Analyse besteht in einer sorgfältigen, quantitativen Beschreibung des Restgliedes einer solchen Entwicklung. Die Resultate sind im Kontext der Linienmethode anwendbar und erlauben eine Vorhersage über das Verhalten von Extrapolationsverfahren bei der betrachteten Problemklasse. Die theoretischen Betrachtungen werden durch numerische Beispiele illustriert.

1. Introduction

In this paper we investigate the structure of the global error for some time discretization schemes applied a class of *stiff* systems of differential equations. A typical case we are having in mind are stiff systems that arise from the semi-discretization of certain PDE problems (method of lines). Some material concerning the implicit Euler method has already been presented in [1]. Here we also consider symmetric schemes (in particular, the implicit trapezoidal rule) and a locally one-dimensional

splitting method. It is helpful (but not necessary) that the reader is familiar with [1].

The main purpose of our study of global error structures is to provide a sound theoretical insight to the behavior of extrapolation methods (or other acceleration techniques). It is usually believed that, in those cases where the given problem (given stiff system or given parabolic PDE) admits a smooth solution (existence of a certain number of derivatives), the existence of an asymptotic error expansion¹ in the classical sense (cf. [7], [18]),

$$\varepsilon_v = \sum_{j=p}^q \tau^j e_j(t_v) + O(\tau^{q+1}) \quad \text{for } \tau \rightarrow 0, \quad (1.1)$$

ensures the full efficiency of extrapolation, with the full conventional order. The situation is, however, much more complicated: It turns out that even in the case where the solution of the given stiff problem has moderate-sized derivatives, it must be expected that the asymptotic error expansion (in particular, the O -constant in (1.1)) is influenced in a critical way by problem parameters which are very large in size, e.g., the conventional Lipschitz constant L .²

In the method lines context, for instance, L is proportional to some negative power of the spatial meshwidth h . Thus our goal must be to derive asymptotic error expansions that are valid uniformly for $h \rightarrow 0$ —a requirement that is by no means trivial. Note, however, that not only the case $h \rightarrow 0$ (arbitrarily refined space discretization) is of interest; also for *fixed*, small h the parameter L is large, and it has to be analyzed in a rigorous, *quantitative* way how this affects the asymptotic error expansion.³ Generally speaking: A careful, quantitative analysis of the behavior of the different terms in an asymptotic error expansion is necessary in the stiff case. In this paper, the use of the O -symbol is therefore always to be understood in the *quantitative sense*, i.e., with a moderate O -constant unaffected by critical, large problem parameters.

Recent results about error structures in the stiff case can be found in [1]–[3] and in [8]. The analysis in [1], in particular, is based on the important concept of *frequency domain arguments* (cf. e.g. [13], [19]). The present paper carries on the considerations of [1]. We consider linear, inhomogeneous, constant coefficient stiff problems and discuss the selfadjoint as well as the nonselfadjoint case. In contrast to [2], [3] and [8], we make no assumption about the distribution of the stiff eigenvalues in the negative complex half-plane (except a sectorial condition which is necessary in the non-selfadjoint case). Thus our results are applicable in the method of lines setting.

The paper is organized as follows: The rest of section 1 contains some preliminary considerations. In section 2 we recall the results of [1] and discuss their usefulness

¹ ε_v denotes the global error of a discretization method of order p , with time step τ , at $t_v = v\tau$.

² For a discussion of this point, cf. also [1], [2]. Useful remarks can also be found in [14], but no further analysis is done there.

³ Simultaneous extrapolation in space is of course possible and used in practice. However, since the essential difficulties caused by the stiffness appear also on a fixed space mesh, we confine ourselves with the case of a fixed space discretization.

in the method of lines context. In section 3 we investigate the structure of the global error for the implicit trapezoidal rule. The results of sections 2 and 3, which are formulated for the selfadjoint case, are extended to the nonselfadjoint case in section 4. Furthermore, section 5 contains a model problem analysis for a locally one-dimensional splitting method (LOD) applied to 2(space)D-problems. In section 6 we present some numerical experiments with extrapolation based on all these methods.

The essence of our theoretical results is the following: Asymptotic error expansions exist in a somewhat restricted sense. For the implicit Euler scheme, there occur 'irregular' error components that are, however, damped out away from the start. Symmetric schemes are in a sense less robust because they have no strong damping properties; however, it turns out that for the implicit trapezoidal rule the error structure is sufficiently regular to ensure the satisfactory performance of (at least) one extrapolation step. (In this respect, the trapezoidal rule is superior to the implicit midpoint rule.) For the locally one-dimensional splitting method based on implicit Euler steps w.r.t. the different space directions, the discussion of section 5 shows that 'irregular' error components play a dominant role and lead to the non-optimal performance of even the first extrapolation step. Despite the strong damping properties of the LOD scheme, not all of these perturbations are damped out away from the start (as it is the case for the fully implicit Euler scheme). All these results can be confirmed by numerical examples.

1.1. The Stiff ODE System

We consider a class of linear stiff initial value problems (of a finite dimension n)

$$\begin{aligned} u'(t) &= Au(t) + f(t), & t \in [0, T], \\ u(0) &= u_0 \end{aligned} \quad (1.2)$$

with a strictly dissipative matrix A , i.e., we assume⁴

$$m := \sup_{v \in \mathbb{C}^n, v \neq 0} \frac{\operatorname{Re} \langle Av, v \rangle}{\langle v, v \rangle} < 0. \quad (1.3)$$

In the case where A is not selfadjoint, we will need an additional 'sectorial condition' (see section 4).

The true solution of (1.2) will be denoted by $u = u(t)$. Furthermore, let H denote the solution operator of (1.2), i.e., $u(t) = e^{tA}u_0 + (Hf)(t)$ with

$$(Hf)(t) := \int_0^t e^{(t-s)A}f(s)ds. \quad (1.4)$$

The fact that the system (1.3) is stiff means that A has eigenvalues in the negative

⁴ $\langle \cdot, \cdot \rangle$ denotes the Euclidean inner product on \mathbb{C}^n , and $\|\cdot\|$ denotes the corresponding L_2 -norm resp. its associated operator norm.

complex half-plane which are large in size, and thus, $\|A\|$ is very large. A typical case we are having in mind is the semi-discretization in space of a parabolic initial/boundary value problem (cf. subsection 1.2 below). However, our considerations are of relevance for any stiff system (1.2) with a stiff spectrum that is distributed over a large part of the negative complex half-plane.

The inhomogeneity $f(t)$ of (1.3) is not required to be moderate-sized. We only assume that $A^{-1}f(t)$ is a moderate, smooth function; this is compatible with the existence of smooth solutions to (1.3). The case that $\|A^{-1}f\|$ is moderate-sized but not $\|f\|$ itself is, in particular, typical for stiff systems arising in the method of lines (cf. (1.6) below). $\|f\| = O(1)$ is a special case.

1.2. The PDE and Its Space Discretization, a Stiff ODE System

In the method of lines (MOL) context, our 'original problem' is an inhomogeneous, d -space dimensional parabolic initial/boundary value problem

$$\begin{aligned} \mathbf{u}'(t, x) &= \mathbf{A}\mathbf{u}(t, x) + \mathbf{g}(t, x) & (t, x) \in [0, T] \times \Omega, \quad \Omega \in \mathbb{R}^d, \\ \mathbf{u}(0, x) &= \mathbf{u}_0(x), & x \in \Omega, \\ \mathbf{u}(t, x) &= \mathbf{b}(t, x), & t \in [0, T], \quad x \in \partial\Omega \end{aligned} \quad (1.5)$$

with a strictly elliptic differential operator \mathbf{A} differentiating w.r.t. the space variable x (' denotes differentiation w.r.t. the time variable t). We assume that the problem data satisfy a sufficient number of *compatibility conditions* at $t = 0$ (cf. for instance [16]), such that $\mathbf{u} = \mathbf{u}(t, x)$ is a *smooth* solution of (1.5). The case of incompatible initial data is briefly discussed at the end of section 6.

Now we consider the stiff system (1.2) as a semi-discretization in space of (1.5), i.e., A represents an (e.g. finite difference or finite element) approximation of \mathbf{A} on a certain 'mesh' or 'triangulation' Ω_h (with 'meshwidth' h) over Ω . Thus, $\|A\|$ is proportional to some negative power of h and is unbounded for decreasing h .

In the sequel, Δ denotes an appropriate restriction operator (e.g., straight injection or local weighting) that maps functions $\mathbf{v} = \mathbf{v}(t, x)$ into their semi-discrete analoga $v(t) = \Delta\mathbf{v}(t)$.

The inhomogeneity f in (1.2) will now represent not only the space-discretization of the inhomogeneity $\mathbf{g}(t, x)$ from (1.5) but also contain a term representing, in a well-known way, the (a-priori eliminated) boundary condition. We will write

$$f(t) = g(t) + b(t) = \Delta\mathbf{g}(t) + b(t) \quad (1.6)$$

where $g = \Delta\mathbf{g}$ is the space discretization of $\mathbf{g}(t, x)$, and b corresponds to the boundary condition. By definition of b , $\|f\|$ and $\|b\|$ are large but, typically, $\|A^{-1}b\|$, $\|A^{-1}f\|$ are moderate-sized.

Global error and spatial truncation error. Let the *spatial truncation error* $s = s(t)$ be defined as

$$s := (A\Delta - \Delta\mathbf{A})\mathbf{u} + b. \quad (1.7)$$

For smooth $\mathbf{u}(t, x)$ and under reasonable assumptions about A , the function $s(t)$ will also be smooth (with moderate-sized derivatives up to a certain order). With definition (1.7), and with $u_0 = \Delta\mathbf{u}_0$ (start on true solution of the PDE), it is obvious that the global discretization error in space $e_s := u - \Delta\mathbf{u}$ is a solution of the stiff initial value problem

$$\begin{aligned} e_s'(t) &= Ae_s(t) + s(t), & t \in [0, T], \\ e_s(0) &= 0. \end{aligned} \quad (1.8)$$

It is well known that the smoothness of the solution of the semi-discrete problem (1.2) is not directly related to that of the original PDE (1.5), because the necessary compatibility conditions are not equivalent. This can be seen by successive differentiation of (1.8) which yields, for $t = 0$,⁵

$$u_0^{(j)} = \Delta\mathbf{u}_0^{(j)} + A^{j-1}s_0 + A^{j-2}s_0' + \dots + As_0^{(j-2)} + s_0^{(j-1)}. \quad (1.9)$$

Despite the smoothness of the spatial truncation error $s(t)$ with respect to t , it cannot be expected in general that the $\|A^k s_0^{(j-1-k)}\|$ ($k \geq 1$) are moderate sized, i.e., not affected by $\|A\|$. It is therefore preferable to center the analysis about the solution $\mathbf{u}(t, x)$ of the original PDE, avoiding reference to the solution $u(t)$ of the semi-discrete system (cf. also [17] for a discussion of this point). We will describe the structure of the 'full global error' $\eta_v - \Delta\mathbf{u}(t_v)$ in this spirit.

2. The Implicit Euler Scheme

The results in this section are formulated for the case where A is selfadjoint. The extension to the non-selfadjoint case is discussed in section 4.

The implicit (or backward) Euler scheme applied to a stiff system (1.2),

$$\frac{1}{\tau}(\eta_v - \eta_{v-1}) = A\eta_v + f_v, \quad t_v = v\tau, \quad f_v := f(t_v), \quad (2.1)$$

has already been considered in [1]. We recall the results:

2.1. Following a Smooth Solution of the Stiff ODE System

The following has been shown in [1]: Assume A is selfadjoint and negative definite, and assume that $u(t)$ is a smooth solution of the stiff system (1.2), i.e., $u(t)$ and its derivatives up to a certain order are moderate-sized. (I.e., we assume that certain compatibility conditions between the initial value u_0 and the inhomogeneity $f(t)$ are satisfied at $t = 0$.) Then the global error of the implicit Euler scheme (2.1) applied to (1.2) admits an asymptotic expansion

$$\eta_v - u(t_v) = \tau e_1(t_v) + \tau^2 e_2(t_v) + \rho, \quad (2.2)$$

⁵ For any function $v(t)$ we denote $v_0 := v(0)$.

where $e_1(t)$ and $e_2(t)$ are moderate-sized, τ -independent solutions of certain stiff differential equations of type (1.2) (the so-called 'variational equations', cf. [1]), (3.2) ff.); the remainder term ρ_v can be estimated by

$$\|\rho_v\| \leq \left(C_0 + \frac{C_1}{t_v}\right) \tau^3 \quad (2.3)$$

with certain τ -independent, moderate-sized bounds C_0 and C_1 . (See [1], Theorem 4.3.)

Remarks.

- The estimate (2.3) says that the remainder term ρ_v of expansion (2.2) shows a *reduced order*, namely $O(\tau^2)$, at the first grid points after $t = 0$; but these order reductions are *algebraically damped* (i.e., damped like $1/v = \tau/t_v$) with increasing v . The full, quantitative order $O(\tau^3)$ reappears 'away from' $t = 0$.
- The proof of Theorem 4.3 in [1] is heavily based on the estimate (see also [13])

$$\|(I - \tau A)^{-v} - e^{t_v A}\| \leq \frac{1}{v} \quad (2.4)$$

and another estimate of a similar type (cf. [1], Lemma 4.2 and equation (4.40)).

- The above assertion about the global error structure of the implicit Euler scheme is valid for any stiff system (1.2) and uniformly for arbitrary stepsizes τ . From the proof given in [1] (cf. [1], equation (4.34)) it can be seen that for a *fixed* stepsize τ the damping behavior of ρ_v can be characterized as *exponential*, i.e., like $(1 - \tau m)^{-v}$ (instead of $1/v$), $m < 0$ from (1.3), which leads to sharper bounds for $\tau|m|$ not too small.

2.2. Following a Smooth Solution of the Parabolic PDE

In the spirit of subsection 1.2, we now consider the stiff system (1.2) as a semi-discretization in space of the PDE (1.5), which is assumed to admit a sufficiently smooth solution $\mathbf{u} = \mathbf{u}(t, \mathbf{x})$. As indicated in section 1, it is now natural to formulate the asymptotic error expansion in terms of data of the original PDE (not of the auxiliary, semi-discrete problem). The necessary modifications are simple: Instead of (2.2) we make an analogous ansatz for the full global error,

$$\eta_v - \Delta \mathbf{u}(t_v) = \tau e_1(t_v) + \tau^2 e_2(t_v) + \rho_v, \quad (2.5)$$

and proceed in the same way as in [1], using Taylor expansions about $t = t_v$ of $\mathbf{u}(t, \mathbf{x})$ (instead of $u(t)$) and of the $e_i(t)$. Inserting into (2.5) and equating coefficients of powers of τ results in

$$e_1'(t) = A e_1(t) + f_1(t), \quad f_1(t) := \frac{1}{2} \Delta \mathbf{u}''(t), \quad (2.6a)$$

$$e_2'(t) = A e_2(t) + f_2(t), \quad f_2(t) := \frac{1}{2} e_1''(t) - \frac{1}{6} \Delta \mathbf{u}'''(t) \quad (2.6b)$$

and

$$\frac{1}{\tau}(\rho_v - \rho_{v-1}) = A \rho_v + \gamma_v + s(t_v). \quad (2.7)$$

Here, γ_v is defined in a way completely analogous as in [1], section 3, (3.5) ff., with \mathbf{u} replaced by the finite space-dimensional projection $\Delta \mathbf{u}$ of \mathbf{u} . $s(t)$ is the spatial truncation error (1.7), which is considered here as a fixed quantity limiting the overall accuracy achievable on the space mesh chosen.

Besides the occurrence of $s(t)$, the formal expansion (2.6a), (2.6b), (2.7) is completely analogous to that for the case considered in subsection 2.1, where $\Delta \mathbf{u}(t)$ now plays the role of $u(t)$. Under the present smoothness assumptions w.r.t. \mathbf{u} it is therefore obvious that the same conclusions can be drawn. We obtain:

The full global error $\eta_v - \Delta \mathbf{u}(t_v)$ of the implicit Euler scheme (2.1) applied to the semidiscretization (1.2) of (1.5) admits an asymptotic expansion

$$\eta_v - \Delta \mathbf{u}(t_v) = \tau e_1(t_v) + \tau^2 e_2(t_v) + \rho_v \quad (2.8)$$

where $e_1(t)$ and $e_2(t)$ are moderate-sized, τ -independent solutions of the 'variational equations' (2.6a), (2.6b), and where the remainder term ρ_v can be estimated by⁶

$$\|\rho_v\| \leq \left(C_0 + \frac{C_1}{t_v}\right) \tau^3 + C_s \max_{0 \leq t \leq t_v} \|s(t)\| \quad (2.9)$$

with certain τ -independent, moderate-sized bounds C_0 , C_1 and C_s .

3. The Implicit Trapezoidal Rule

The results in this section (Lemma 3.1, Theorem 3.2) are formulated for the case where A is selfadjoint. The extension to the non-selfadjoint case is discussed in section 4.

For the implicit trapezoidal rule applied to a stiff system (1.2),

$$\frac{1}{\tau}(\eta_v - \eta_{v-1}) = A \frac{1}{2}(\eta_{v-1} + \eta_v) + \frac{1}{2}(f_{v-1} + f_v), \quad t_v = v\tau, \quad f_v := f(t_v), \quad (3.1)$$

we will now investigate whether an asymptotic expansion

$$\eta_v - u(t_v) = \tau^2 e_2(t_v) + \rho_v \quad (3.2)$$

exists, with a remainder term satisfying $\rho_v = O(\tau^4)$.

3.1. The Formal Expansion

Following the classical work of GRAGG [7] (see also [18]), the variational equation defining $e_2(t)$ and the (discrete) remainder equation defining ρ_v in (3.2) can be derived

⁶ The influence of $s(t)$ on ρ_v can be estimated by a conventional B-stability argument.

by Taylor expansion of $u(t)$ and $e_2(t)$ about $t = \hat{t}_v := \frac{1}{2}(t_{v-1} + t_v)$ and equating coefficients of powers of τ^2 (cf. e.g. [3]). This results in

$$e_2'(t) = Ae_2(t) + f_2(t), \quad f_2(t) := \frac{1}{12}u'''(t) \tag{3.3}$$

and

$$\frac{1}{\tau}(\rho_v - \rho_{v-1}) = A \frac{1}{2}(\rho_{v-1} + \rho_v) + \gamma_v; \tag{3.4}$$

the inhomogeneity γ_v in (3.4) is given by

$$\gamma_v := \frac{1}{2}(j_{0,v}^- + j_{0,v}^+) + \frac{\tau^2}{2}(j_{2,v}^- + j_{2,v}^+) - \frac{1}{\tau}(i_{0,v}^+ - i_{0,v}^-) - \tau^2 \frac{1}{\tau}(i_{2,v}^+ - i_{2,v}^-), \tag{3.5}$$

where the $i_{k,v}^\pm$ and $j_{k,v}^\pm$ are the following Taylor remainder terms:

$$i_{0,v}^\pm := \frac{\tau^4}{96} \int_0^1 (1-\sigma)^3 u^{IV} \left(\hat{t}_v \pm \sigma \frac{\tau}{2} \right) d\sigma, \tag{3.6a}$$

$$i_{2,v}^\pm := \frac{\tau^2}{4} \int_0^1 (1-\sigma) e_2'' \left(\hat{t}_v \pm \sigma \frac{\tau}{2} \right) d\sigma,$$

$$j_{0,v}^\pm := \frac{\tau^4}{96} \int_0^1 (1-\sigma)^3 u^V \left(\hat{t}_v \pm \sigma \frac{\tau}{2} \right) d\sigma, \tag{3.6b}$$

$$j_{2,v}^\pm := \frac{\tau^2}{4} \int_0^1 (1-\sigma) \left[e_2''' - \frac{1}{12}u^V \right] \left(\hat{t}_v \pm \sigma \frac{\tau}{2} \right) d\sigma.$$

For a given initial value $e_2(0) := e_{2,0}$, the solution $e_2(t)$ of the variational equation (3.3) and its derivatives read

$$\begin{aligned} e_2(t) &= e^{tA}e_{2,0} + (Hf_2)(t), \\ e_2'(t) &= e^{tA} \underbrace{[Ae_{2,0} + f_{2,0}]}_{e_{2,0}'} + (Hf_2')(t), \\ e_2''(t) &= e^{tA} \underbrace{[A^2e_{2,0} + Af_{2,0} + f_{2,0}']}_{e_{2,0}''} + (Hf_2'')(t), \\ e_2'''(t) &= e^{tA} \underbrace{[A^3e_{2,0} + A^2f_{2,0} + Af_{2,0}']}_{e_{2,0}'''} + (Hf_2''')(t), \\ &\vdots \end{aligned} \tag{3.7}$$

with H from (1.4) and $f_2(t) = \frac{1}{12}u'''(t)$. Using $u' = Au + f$ at $t = 0$, the initial values in (3.7) can be expressed as

$$\begin{aligned} e_{2,0}' &= Ae_{2,0} + \frac{1}{12}u_0''', \\ e_{2,0}'' &= A^2e_{2,0} + \frac{1}{6}u_0^{IV} - \frac{1}{12}f_0''', \\ e_{2,0}''' &= A^3e_{2,0} + \frac{1}{4}u_0^V - \frac{1}{6}f_0^{IV} - \frac{1}{12}Af_0''', \\ &\vdots \end{aligned} \tag{3.8}$$

3.2. Following a Smooth Solution of the Stiff ODE System

In this subsection we consider the case where $u(t)$ is a smooth solution of the stiff system (1.2).

Let the initial value for $e_2(t)$ be fixed by $e_2(0) := 0$. Then

$$e_2(t) = (Hf_2)(t), \quad f_2(t) = \frac{1}{12}u'''(t) \tag{3.9}$$

is a moderate-sized function (due to our smoothness assumptions w.r.t. $u(t)$). The higher derivatives $e_2^{(j)}(t)$ do not remain uniformly moderate-sized because, even for $e_{2,0} = 0$, the initial values $e_{2,0}^{(j)}$ are influenced by higher and higher powers⁷ of A . Within the equation for ρ_v , there occur derivatives of $e_2(t)$ up to the third order (cf. (3.5), (3.6a), (3.6b)); now the crucial point is to study in what way this affects the remainder term ρ_v of expansion (3.2).

The inhomogeneity of the remainder equation. The inhomogeneity γ_v (cf. (3.5)) of the remainder equation (3.4) splits into

$$\gamma_v = \gamma_{v,0} + \gamma_{v,2} \tag{3.10}$$

where $\gamma_{v,0}$ and $\gamma_{v,2}$ are defined in terms of derivatives of $u(t)$ and $e_2(t)$, respectively (cf. (3.6a), (3.6b)). Due to the smoothness of $u(t)$ it is obvious that there holds

$$\gamma_{v,0} = O(\tau^4) \tag{3.11}$$

with a moderate-sized O -constant (unaffected by $\|A\|$). To obtain a more concise representation for the other inhomogeneous term $\gamma_{v,2}$, we rewrite the integral in (3.6b) (definition of $j_{2,v}^\pm$) using partial integration,

$$\int_0^1 (1-\sigma) e_2'' \left(\hat{t}_v \pm \sigma \frac{\tau}{2} \right) d\sigma = \pm \left[\frac{2}{\tau} \int_0^1 e_2'' \left(\hat{t}_v \pm \sigma \frac{\tau}{2} \right) d\sigma - e_2''(\hat{t}_v) \right], \tag{3.12}$$

which easily leads to

$$\gamma_{v,2} = \frac{\tau}{4} \int_0^1 \sigma \left[e_2'' \left(\hat{t}_v + \sigma \frac{\tau}{2} \right) - e_2'' \left(\hat{t}_v - \sigma \frac{\tau}{2} \right) \right] d\sigma. \tag{3.13}$$

Now, according to

$$e_{2,0}'' = e^{tA}e_{2,0}'' + (Hf_2'')(t), \quad e_{2,0}'' = \frac{1}{6}u_0^{IV} - \frac{1}{12}f_0''' = \frac{1}{6}u_0^{IV} - \frac{1}{12}A(A^{-1}f_0''') \tag{3.14}$$

(cf. (3.8)), $\gamma_{v,2}$ consists of

- a ‘genuine’ $O(\tau^4)$ -term, originating from the moderate function $Hf_2'' = \frac{1}{12}Hu^V$, and
- a ‘critical’ term, originating from $e^{tA}e_{2,0}''$, which we denote by $\gamma_{v,2}^{\text{crit}}$.

⁷ Cf. for instance the occurrence of a factor A within the expression for $e_{2,0}''$ in (3.8). Note further that the $A^{-1}f^{(j)}$ are assumed to be moderate-sized but not the $f^{(j)}$ (cf. section 1). Thus, $e_{2,0}''$ is even affected by a factor A^2 .

A simple calculation yields⁸

$$\gamma_{v,2}^{\text{crit}} = \frac{\tau^3}{4} e^{\tau A} \left(\frac{\tau}{2} A\right)^{-2} \left(I - \frac{\tau}{2} A\right) [Ge^{-\tau A} - I] \cdot e''_{2,0}; \quad (3.15)$$

here we have introduced the denotation

$$G := \left(I - \frac{\tau}{2} A\right)^{-1} \left(I + \frac{\tau}{2} A\right). \quad (3.16)$$

Summarizing, we have

$$\rho_v = \gamma_{v,2}^{\text{crit}} + O(\tau^4), \quad \gamma_{v,2}^{\text{crit}} \text{ from (3.15)}. \quad (3.17)$$

Estimation of the remainder term ρ_v . For $\rho_0 = 0$ we obtain from (3.4) and (3.17)

$$\rho_v = \left(I - \frac{\tau}{2} A\right)^{-1} \tau \sum_{i=1}^v G^{v-i} [\gamma_{i,2}^{\text{crit}} + O(\tau^4)] = \rho_v^{\text{crit}} + O(\tau^4). \quad (3.18)$$

Here, the $O(\tau^4)$ -term corresponds to the $O(\tau^4)$ -term in (3.17) and is the result of a conventional B-stability estimate based on $\|G\| < 1$. The ‘critical’ part ρ_v^{crit} , originating from $\gamma_{v,2}^{\text{crit}}$, remains to be studied. From (3.15),

$$\rho_v^{\text{crit}} = \left(I - \frac{\tau}{2} A\right)^{-1} \frac{\tau^3}{4} \left(\frac{\tau}{2} A\right)^{-2} \left(I - \frac{\tau}{2} A\right) [Ge^{-\tau A} - I] \tau \sum_{i=1}^v G^{v-i} (e^{\tau A})^i \cdot e''_{2,0}. \quad (3.19)$$

Due to

$$\sum_{i=1}^v G^{v-i} (e^{\tau A})^i = [G^v - e^{\tau A}] [Ge^{-\tau A} - I]^{-1} \quad (3.20)$$

and with the representation (3.14) for $e''_{2,0}$, (3.19) can be rewritten as

$$\rho_v^{\text{crit}} = \tau^4 (\tau A)^{-2} [G^v - e^{\tau A}] \cdot \left[\frac{1}{8} u_0^{IV} - \frac{1}{12} A(A^{-1} f_0''')\right]. \quad (3.21)$$

Estimating ρ_v^{crit} by a B-stability argument yields only $\rho_v^{\text{crit}} = O(\tau^2)$, which is far too pessimistic and, in fact, useless. The following sharper estimates are now essential.

Lemma 3.1. For G from (3.16) and $v \geq 1$ the estimates

$$\|(\tau A)^{-2} [G^v - e^{\tau A}]\| \leq C, \quad (3.22a)$$

$$\|(\tau A)^{-1} [G^v - e^{\tau A}]\| \leq \frac{C}{v} \quad (3.22b)$$

hold independent of the stepsize τ , with certain moderate-sized constants C .

⁸ Here we use the identities $\int_0^1 \sigma e^{\pm \sigma(\tau/2)A} d\sigma = \left(\frac{\tau}{2} A\right)^{-2} \left[I - \left(I \mp \frac{\tau}{2} A\right) e^{\pm(\tau/2)A}\right]$.

*Proof:*⁹ At first we prove (3.22a). Since A is assumed selfadjoint and negative definite, (3.22a) is equivalent to

$$\sup_{\xi > 0} |\alpha_v(\xi)| \leq C \quad (3.23)$$

where

$$\alpha_v(\xi) := \frac{1}{\xi^2} \left[\left(\frac{1 - \frac{\xi}{2}}{1 + \frac{\xi}{2}}\right)^v - e^{-v\xi} \right] = \frac{1}{\xi^2} \left[\frac{1 - \frac{\xi}{2}}{1 + \frac{\xi}{2}} - e^{-\xi} \right] \sum_{i=1}^{v-1} \left(\frac{1 - \frac{\xi}{2}}{1 + \frac{\xi}{2}}\right)^i e^{-(v-1-i)\xi}. \quad (3.24)$$

To prove (3.23), we estimate $\alpha_v(\xi)$ for $0 < \xi \leq 2$ and $\xi \geq 2$ separately.

- For $0 < \xi \leq 2$, we make use of the fact that $\left(1 - \frac{\xi}{2}\right) / \left(1 + \frac{\xi}{2}\right)$ is an approximation for $e^{-\xi}$:

$$\left| \frac{1}{\xi^2} \left[\frac{1 - \frac{\xi}{2}}{1 + \frac{\xi}{2}} - e^{-\xi} \right] \right| \leq C|\xi|; \quad (3.25)$$

furthermore,

$$\begin{aligned} \sum_{i=1}^{v-1} \left(\frac{1 - \frac{\xi}{2}}{1 + \frac{\xi}{2}}\right)^i e^{-(v-1-i)\xi} &= e^{-(v-1)\xi} \sum_{i=1}^{v-1} \underbrace{\left[\frac{1 - \frac{\xi}{2}}{1 + \frac{\xi}{2}} e^{\xi}\right]^i}_{\leq 1, \xi \in [0, 2]} \leq v e^{-(v-1)\xi} \\ &= \frac{1}{\xi} (v\xi) e^{-v\xi} e^{\xi} \leq \frac{C}{\xi}. \end{aligned} \quad (3.26)$$

From (3.25) and (3.26), the estimate $|\alpha_v(\xi)| \leq C$ follows for $0 < \xi \leq 2$.

- For $\xi \geq 2$, the desired estimate

$$|\alpha_v(\xi)| = \frac{1}{\xi^2} \left[\left(\frac{1 - \frac{\xi}{2}}{1 + \frac{\xi}{2}}\right)^v - e^{-v\xi} \right] \leq C \quad (3.27)$$

is trivial.

This proves (3.22a).

In a similar way we now prove (3.22b), which is equivalent to

$$\sup_{\xi > 0} |\beta_v(\xi)| \leq \frac{C}{v}, \quad v \geq 1, \quad (3.28)$$

where $\beta_v(\xi) := \xi \alpha_v(\xi)$ (with $\alpha_v(\xi)$ from (3.24)).

- For $0 < \xi \leq 2$, we conclude from (3.25) and (3.26)

$$|\beta_v(\xi)| \leq C \xi^2 v e^{(v-1)\xi} = C \frac{1}{v} (v\xi)^2 e^{-v\xi} e^{\xi} \leq \frac{C}{v}. \quad (3.29)$$

⁹ Throughout, C denotes a generic, moderate-sized constant.

- For $\xi \geq 2$, we have $\xi^{-1}e^{-v\xi} \leq C/v$, and¹⁰

$$\begin{aligned} \left| \frac{1}{\xi} \left(\frac{1 - \frac{\xi}{2}}{1 + \frac{\xi}{2}} \right)^v \right| &= \frac{1}{\xi} \frac{1}{v \left| \ln \left(\frac{\xi - 2}{\xi + 2} \right) \right|} \left| \ln \left(\frac{\xi - 2}{\xi + 2} \right) \right| e^{-v \left| \ln \left(\frac{\xi - 2}{\xi + 2} \right) \right|} \\ &\leq \frac{1}{\xi} \frac{1}{v \left| \ln \left(\frac{\xi - 2}{\xi + 2} \right) \right|} e^{-1}, \end{aligned} \tag{3.30}$$

where $1/(\xi |\ln((\xi - 2)/(\xi + 2))|)$ is uniformly bounded for $2 \leq \xi < \infty$. This proves $|\beta_v(\xi)| \leq C/v$ for $\xi \geq 2$. \square

Lemma 3.1 enables the estimation of the ‘critical’ component ρ_v^{crit} of the remainder term ρ_v (cf. (3.21)): Due to (3.22a) and (3.22b),

$$\begin{aligned} \|\rho_v^{\text{crit}}\| &\leq \tau^4 \|(\tau A)^{-2} [G^v - e^{t_v A}]\| \cdot \left\| \frac{1}{6} u_0^{IV} \right\| + \tau^3 \|(\tau A)^{-1} [G^v - e^{t_v A}]\| \cdot \left\| \frac{1}{12} A^{-1} f_0''' \right\| \\ &\leq C\tau^4 + \frac{C}{v} \tau^3 = \left(C + \frac{C}{t_v} \right) \tau^4. \end{aligned} \tag{3.31}$$

Summarizing all that, we end up with

Theorem 3.2. *The global error of the implicit trapezoidal rule (3.1) applied to (1.2) admits an asymptotic expansion*

$$\eta_v - u(t_v) = \tau^2 e_2(t_v) + \rho_v \tag{3.32}$$

where $e_2(t)$ is a moderate, τ -independent solution of the variational equation (3.3); the remainder term ρ_v can be estimated by

$$\|\rho_v\| \leq \left(C_0 + \frac{C_1}{t_v} \right) \tau^4 \tag{3.33}$$

with moderate bounds C_i .

Remarks.

- Theorem 3.2 shows that, similarly as for the implicit Euler scheme (cf. section 2), the remainder term ρ_v suffers from an *order reduction* at the first grid points; however, the full order $O(\tau^4)$ reappears away from $t = 0$. Thus, the implicit trapezoidal rule exhibits certain damping properties (despite the fact that it is not strongly stable).
- It is worth mentioning that these order reductions are mainly caused by the presence of ‘mildly stiff’ eigenvalues: Due to the occurrence of a common factor A^{-1} within the expression for ρ_v^{crit} (cf. (3.21)), the components of ρ_v

¹⁰ The case $\xi = 2$ is trivial.

corresponding to the large eigenvalues λ of A ($1/|\lambda| \ll 1$) are actually not significant.¹¹

- Concerning longer expansions, similar results with a remainder term of a higher order *cannot* be expected; this is due to the lack of a strong damping within the difference equation for ρ_v . This point will be illustrated in section 6 by means of a numerical experiment (τ^2 -extrapolation).

3.3. Following a Smooth Solution of the Parabolic PDE

As has been explained in subsection 2.2 for the implicit Euler scheme, no further analysis needs to be done for this case; only minor technical modifications are necessary. Thus we obtain, in a completely analogous way as in subsection 2.2:

Corollary 3.3. *The full global error $\eta_v - \Delta u(t_v)$ of the implicit trapezoidal rule (3.1) applied to the semi-discretization (1.2) of (1.5) admits an asymptotic expansion*

$$\eta_v - \Delta u(t_v) = \tau^2 e_2(t_v) + \rho_v \tag{3.34}$$

where $e_2(t)$ is a moderate, τ -independent solution of the variational equation

$$e_2'(t) = A e_2(t) + \frac{1}{12} \Delta u'''(t); \tag{3.35}$$

the remainder term ρ_v can be estimated by

$$\|\rho_v\| \leq \left(C_0 + \frac{C_1}{t_v} \right) \tau^4 + C_3 \max_{0 \leq t \leq t_v} \|s(t)\| \tag{3.36}$$

with certain τ -independent, moderate-sized bounds C_0, C_1 and C_3 .

4. The Non-selfadjoint Case

We now show how the results of sections 2 and 3 can be extended to the case where A is not selfadjoint. We briefly recall the necessary theoretical background but do not present all the technical details, which are quite tedious.

Instead of (1.3), we impose the usual *sectorial condition*

$$\langle Av, v \rangle \in \mathcal{S}_\theta \quad \text{for all } v \in \mathbb{C}^n, \tag{4.1}$$

¹¹ For problems where, besides non-stiff eigenvalues, there occur only stiff eigenvalues which are very large in size, this leads to a ‘pure’ asymptotic expansion, with a remainder term that shows the full conventional order over a large range of stepsizes (stepsize restricted from below). A detailed analysis of this case (including certain classes of non-autonomous and nonlinear problems) has been given in [3]. In the present context, we are of course not interested in problems with such a spectral structure.

where S_θ is the sector $\{z \in \mathbb{C}: |\text{Arg } z - \pi| \leq \theta\}$ ($0 < \theta < \frac{\pi}{2}$) in the negative complex half-plane.

A theoretical framework for the treatment of the non-selfadjoint case which is frequently used in the context of parabolic problems—and which turns out to be very suitable for our analysis of asymptotic error expansions—is based on the Dunford-Taylor integral representation for holomorphic functions φ acting on A (cf. [11]):

$$\varphi(A) = \frac{1}{2\pi i} \int_\Gamma \varphi(z)(zI - A)^{-1} dz, \tag{4.2}$$

where $\Gamma \in \mathbb{C}$ is a simple closed positively oriented curve containing all eigenvalues of A in its interior. For $(zI - A)^{-1}$, the so-called *resolvent* of A , the following estimate is well known:

Lemma 4.1 (Resolvent inequality). *If A satisfies (4.1) then, for all $z \in \mathbb{C}$ with $|\text{Arg } z - \pi| = \tilde{\theta}$ ($0 < \theta < \tilde{\theta} < \frac{\pi}{2}$),*

$$\|(zI - A)^{-1}\| \leq \frac{C}{|\tilde{\theta} - \theta| |z|} \tag{4.3}$$

with a moderate-sized constant C independent of θ and $\tilde{\theta}$.

Proof: Can be based on a simple duality argument and is not presented here (cf. e.g. [5]). □

On the basis of (4.2) and (4.3), estimates for $\|\varphi(A)\|$ can easily be reduced to certain estimates of $\varphi(z)$ for $z \in S_\theta$, as e.g. described in [13]. In particular, we shall use Lemma 1.2 of [13], which is heavily based on (4.2) and (4.3), and which can be reformulated as follows:¹²

Lemma 4.2 (Le Roux [13]). *Let A satisfy (4.1), and consider a function $\varphi(z)$ ($\varphi(-\infty) = 0$) which is continuous on the sector S_θ ($0 < \theta < \frac{\pi}{2}$) and holomorphic in the interior of S_θ , and which satisfies for some constant $R > 0$ and two functions φ_1 and φ_2 from \mathbb{R}_+ to \mathbb{R}_+ the following estimates:*

$$\forall z \in S_\theta, |z| \leq R: |\varphi(z)| \leq \varphi_1(|z|), \tag{4.4a}$$

$$\forall z \in S_\theta, |z| \geq R: |\varphi(z)| \leq \varphi_2(|z|). \tag{4.4b}$$

Then there is a moderate-sized constant C such that, for $0 < \theta < \tilde{\theta} < \frac{\pi}{2}$,

$$\|\varphi(A)\| \leq \frac{C}{|\tilde{\theta} - \theta|} \left[\int_0^R \varphi_1(r) \frac{dr}{r} + \int_R^\infty \varphi_2(r) \frac{dr}{r} \right]. \tag{4.5}$$

¹² We formulate Lemma 1.2 of [13] for the special case $\varphi(-\infty) = 0$. Note that A corresponds to $-A$ in [13].

Proof: See [13]. □

Lemma 4.2 provides the tool to prove

Proposition 4.3. *For A satisfying a sectorial condition (4.1), the results of sections 2 and 3 remain valid; the constants C appearing in the various estimates are moderate-sized if $\left(\frac{\pi}{2} - \theta\right)^{-1}$ is moderate-sized, and are proportional to a certain negative power of $\left(\frac{\pi}{2} - \theta\right)$ for $\theta \rightarrow \frac{\pi}{2}$.*

Proof (sketch):

- For the implicit Euler scheme, a look at the arguments in [1] shows that, apart from some minor technical details, it is sufficient to verify that the assertions of Lemma 4.2 in [1] (estimates ([1], (4.24)) and ([1], (4.25))) remain valid, with constants $C = C(\theta)$ depending on θ in the indicated way. The estimate ([1], (4.24)), which is equivalent to

$$\|(I - \tau A)^{-\nu} - (e^{\tau A})^\nu\| \leq \frac{C}{\nu} \tag{4.6}$$

(cf. (2.4) above), is covered by the results in [13], where inequalities of the type (4.6) are derived, with $C = C(\theta)$, for a general class of rational approximations to $e^{\tau A}$. The proof of the second estimate ([1], (4.25)) for the non-selfadjoint case is not covered by the results of [13] but can be worked out along similar lines on the basis of Lemma 4.2. The technical details are more complicated than for the proof of (4.6); however, we refrain from discussing these details here.

- For the implicit trapezoidal rule, it has to be verified that Lemma 3.1 ((3.22a), (3.22b)) carries over to the present situation. This can be done on the basis of Lemma 4.2:

—*Proof of (3.22a).* To apply Lemma 4.2, we derive estimates (4.4a) and (4.4b) for $\varphi(z) := \alpha_\nu(-z)$ (α_ν from (3.24)). To this end, we note that (3.25) remains valid for complex arguments ζ with $\text{Re } \zeta \geq 0$; moreover, it can be shown that

$$\left| \frac{1 - \frac{\zeta}{2}}{1 + \frac{\zeta}{2}} \right| \leq e^{-(1/2) \text{Re } \zeta} \quad \text{for } \zeta \in -S_\theta, |\zeta| \leq 2. \tag{4.7}$$

This easily leads to

$$|\varphi(z)| = |\alpha_\nu(-z)| \leq \varphi_1(|z|), \quad z \in S_\theta, |z| \leq 2, \quad \text{with } \varphi_1(r) = C\nu r e^{-\nu(r/2) \cos \theta}. \tag{4.8}$$

Furthermore, the estimate

$$|\varphi(z)| = |\alpha_\nu(-z)| \leq \varphi_2(|z|), \quad z \in S_\theta, |z| \geq 2, \quad \text{with } \varphi_2(r) = \frac{C}{r^2}, \tag{4.9}$$

is obvious. Now, calculation of the integrals appearing in Lemma 4.2, (4.5), with φ_1 and φ_2 given by (4.8) resp. (4.9), easily leads to the desired estimate (cf. (3.22a))

$$\|(\tau A)^{-2}[G^v - e^{tA}]\| \leq C(\theta); \tag{4.10}$$

it turns out that $C(\theta)$ is moderate-sized for θ away from $\frac{\pi}{2}$ and $O\left(\left(\frac{\pi}{2} - \theta\right)^{-2}\right)$ for $\theta \rightarrow \frac{\pi}{2}$.

—Proof of (3.22b). Similarly as above, estimates (4.4a) and (4.4b) are required for $\varphi(z) := \beta_v(-z)$ (where $\beta_v(\zeta) := \zeta\alpha_v(\zeta)$ with α_v from (3.24)). First of all, the estimate

$$|\varphi(z)| = |\beta_v(-z)| \leq \varphi_1(|z|), \quad z \in S_\theta, |z| \leq 2, \quad \text{with } \varphi_1(r) = C\nu r^2 e^{-\nu(r/2) \cos \theta} \tag{4.11}$$

can be derived in the very same way as (4.8). Furthermore, using the inequality

$$\left| \frac{1 - \frac{\zeta}{2}}{1 + \frac{\zeta}{2}} \right| \leq e^{-2 \operatorname{Re} 1/\zeta} \quad \text{for } \zeta \in -S_\theta, |\zeta| \geq 2, \tag{4.12}$$

we easily obtain

$$|\varphi(z)| = |\beta_v(-z)| \leq \varphi_2(|z|), \quad z \in S_\theta, |z| \geq 2, \quad \text{with } \varphi_2(r) = \frac{1}{r} e^{-\nu r \cos \theta} + \frac{1}{r} e^{-\nu(2/r) \cos \theta}. \tag{4.13}$$

Now, calculation of the integrals appearing in Lemma 4.2, (4.5), with φ_1 and φ_2 given by (4.11) resp. (4.13), easily leads to the desired estimate (cf. (3.22b))

$$\|(\tau A)^{-1}[G^v - e^{tA}]\| \leq \frac{C(\theta)}{\nu}, \quad \nu \geq 1; \tag{4.14}$$

it turns out that $C(\theta)$ is moderate-sized for θ away from $\frac{\pi}{2}$ and $O\left(\left(\frac{\pi}{2} - \theta\right)^{-3}\right)$ for $\theta \rightarrow \frac{\pi}{2}$. □

5. A Locally One-dimensional Splitting Method

Splitting schemes are very important for the efficient practical solution of multi-space dimensional PDEs (cf. e.g. [21]). Here we consider a locally one-dimensional splitting scheme (LOD) for 2-space dimensional problems. We assume that A in (1.2) splits into

$$A = A_1 + A_2, \tag{5.1}$$

where A_1 and A_2 ‘differentiate’ only w.r.t. one of the space directions but involve no coupling w.r.t. the other direction. Also the inhomogeneity f in (1.2) is split into $f = f_1 + f_2$ in a suitable way. The following LOD scheme is based on alternating

steps of backward Euler type:

$$\frac{1}{\tau}(\tilde{\eta}_v - \eta_{v-1}) = A_1 \tilde{\eta}_v + f_{1,v}, \tag{5.2a}$$

$$\frac{1}{\tau}(\eta_v - \tilde{\eta}_v) = A_2 \eta_v + f_{2,v}. \tag{5.2b}$$

Elimination of the intermediate quantity $\tilde{\eta}_v$ yields

$$\frac{1}{\tau}(\eta_v - \eta_{v-1}) = A \eta_v + f_v - \tau A_1(A_2 \eta_v + f_{2,v}). \tag{5.3}$$

The computational cost for performing one step (5.2a), (5.2b) of the LOD scheme is significantly below that for one ‘fully implicit’ backward Euler step. Moreover, the LOD scheme has excellent stability properties. On the other hand, its consistency properties are rather poor. The full local error can be written as the sum of the local error of the fully implicit Euler scheme (2.1) plus the spatial truncation error $s(t)$, plus the term

$$-\tau A_1 w(t_v), \quad w(t) := A_2 \Delta u(t) + f_2(t). \tag{5.4}$$

Assume $u(t, x)$ is a smooth solution of (1.5). Then, for reasonable splittings, the function $w(t)$ defined in (5.4) will be moderate-sized (despite the factor A_2); but $A_1 w(t)$ cannot be expected to be moderate-sized because $w(t)$ will not, in general, assume zero boundary values. This *reduced local accuracy* does not, however, necessarily reduce the global accuracy of the LOD scheme: For the case where A_1 and A_2 are negative definite and commuting, for instance, it is not difficult to show that the LOD scheme is unconditionally convergent of first order in time, i.e., its full global error satisfies

$$\|\eta_v - \Delta u(t_v)\| \leq C\tau + C_s \max_{0 \leq t \leq t_v} \|s(t)\| \tag{5.5}$$

with moderate bounds C and C_s that are independent of the spatial meshwidth. (Cf. also [9] and [10], where similar results are derived for other splitting schemes.)

Here we are interested in the *structure* of the global error. The main question is in what way the ‘parasitic’ local error term (5.4) affects the asymptotic error expansion. This point is studied in the following for a model situation.

5.1. The Formal Expansion

Let us consider a short expansion

$$\eta_v - \Delta u(t_v) = \tau e_1(t_v) + \rho_v; \tag{5.6}$$

the usual procedure leads to the first variational equation

$$e_1'(t) = A e_1(t) + \frac{1}{2} \Delta u''(t) - A_1 w(t), \quad w(t) \text{ from (5.4)}, \tag{5.7}$$

and the remainder equation

$$\frac{1}{\tau}(\rho_v - \rho_{v-1}) = (A - \tau A_1 A_2)\rho_v + \gamma_v \tag{5.8}$$

with the inhomogeneity

$$\gamma_v = \frac{1}{\tau}i_{0,v} + i_{1,v} - \tau^2 A_1 A_2 e_1(t_v) + s(t_v). \tag{5.9}$$

Here, $i_{0,v}$ and $i_{1,v}$ are defined in exactly the same way as for the implicit Euler scheme (cf. [1]):

$$i_{0,v} = -\frac{\tau^3}{2} \int_0^1 \sigma^2 \Delta u'''(t_{v-1} + \sigma\tau) d\sigma, \quad i_{1,v} = \tau^2 \int_0^1 \sigma e_1''(t_{v-1} + \sigma\tau) d\sigma. \tag{5.10}$$

Note the occurrence of the 'large' factor $A_1 A_2$ within (5.9); a similar term does not occur for the Euler scheme.

The solution of (5.7) with initial value $e_{1,0} = 0$ is

$$e_1(t) = \frac{1}{2}(H\Delta u''(t) - (HA_1 w)(t)) \tag{5.11}$$

(H from (1.4)); its second derivative can be written as

$$e_1''(t) = e^{tA} e_{1,0}'' + \frac{1}{2}(H\Delta u^{IV}(t) - (HA_1 w'')(t)), \tag{5.12}$$

with initial value

$$e_{1,0}'' = \Delta u_0'' - \frac{1}{2}AA^{-1}f_0'' - A_1 w_0' - AA_1 w_0. \tag{5.13}$$

It is reasonable to assume that not only $w(t)$ but also its first and second derivatives are moderate-sized.

5.2. Model Problem Analysis

In the following we restrict our considerations to a simple model problem, namely the inhomogeneous initial/boundary value problem for the 2D heat equation on a rectangular domain Ω , discretized by central second order difference quotients on a regular spatial grid Ω_h with meshwidth h . Then, A_1 and A_2 are the discrete versions of $\partial_{x_1 x_1}$ and $\partial_{x_2 x_2}$. Note that A_1 and A_2 are symmetric negative definite, and commute: $A_1 A_2 = A_2 A_1$.

We now have $\|A^{-1}A_i\| \leq 1$, and

$$\|(HA_i)(t)\| \leq \|(HA)(t)\| \|A^{-1}A_i\| \leq \|e^{tA} - I\| \leq 1 \tag{5.14}$$

($i = 1, 2$), and thus $e_1(t)$ from (5.11) is moderate-sized.

The behavior of the remainder term ρ_v is not a priori obvious. With

$$G := [I - \tau(A - \tau A_1 A_2)]^{-1} = (I - \tau A_2)^{-1}(I - \tau A_1)^{-1}, \quad \|G\| \leq 1, \tag{5.15}$$

the solution of (5.8) with initial value $\rho_0 = 0$ reads

$$\rho_v = \tau \sum_{l=1}^v G^{v+1-l} \gamma_l. \tag{5.16}$$

Now we discuss the influence of the various terms in (5.9) on ρ_v . For smooth $u(t, x)$, the influence of $\frac{1}{\tau}i_{0,v}$ (cf. (5.10)) is $O(\tau^2)$ due to stability ($\|G\| \leq 1$). Furthermore, the influence of the spatial truncation error $s(t)$ is unavoidable and can also be estimated by a conventional stability argument. The influence of the other terms within γ_v remains to be studied:

Influence of the inhomogeneous term $i_{1,v}$ on ρ_v (cf. (5.10), (5.12), (5.13)):

- By our smoothness assumptions and due to (5.14), the contribution of $\frac{1}{2}(H\Delta u^{IV}(t) - (HA_1 w'')(t))$ to $i_{1,v}$ and thus, by stability, its influence on ρ_v , is $O(\tau^2)$.
- Furthermore, the contribution of $e^{tA} e_{1,0}''$ to $i_{1,v}$ is

$$\tau^2 e^{t_{v-1}A} \underbrace{\int_0^1 \sigma e^{\sigma A} d\sigma}_{=: I_1} \cdot e_{1,0}'', \tag{5.17}$$

where $\|e^{t_{v-1}A} I_1\| \leq C$ and where $e_{1,0}''$ (cf. (5.13)) contains 'critical' terms which are affected by the factors A, A_1 and AA_1 .

—To estimate the influence of those terms in (5.13) which contain a factor A or A_1 , we note that

$$\tau \sum_{l=1}^v G^{v+1-l} = (G^v - I)(A - \tau A_1 A_2)^{-1}, \tag{5.18}$$

and use

$$\|(A - \tau A_1 A_2)^{-1} A_{[i]}\| \leq 1. \tag{5.19}$$

This leads to a sharpened stability estimate, from which the influence of $-\frac{1}{2}AA^{-1}f_0'' - A_1 w_0'$ on ρ_v is easily seen to be $O(\tau^2)$.

—To estimate the influence of $AA_1 w_0$ (cf. (5.13)) on ρ_v , stability arguments are too weak; we need a direct estimate for the norm of the positive definite operator

$$\tau \sum_{l=1}^v G^{v+1-l} e^{t_{l-1}A} I_1 A A_1. \tag{5.20}$$

Here, the situation is very similar as in [1] for the implicit Euler scheme. In fact, the desired estimate for (5.20) can be derived from the results of [1]: Using the fact that (5.20) is positive definite, and the using the obvious inequalities $\|G(I - \tau A)^{-1}\| \leq 1$ and $\|A^{-1}A_1\| \leq 1$, we obtain

$$\begin{aligned} \left\| \tau \sum_{i=1}^v G^{v+1-i} e^{t_i A} I_1 A A_1 \right\| &\leq \left\| \tau \sum_{i=1}^v (I - \tau A)^{-(v+1-i)} e^{t_i A} I_1 A^2 \right\| \\ &= \|\tau^{-1} [(I - \tau A)^{-v} - (e^{\tau A})^v]\| \leq \frac{1}{t_v} \end{aligned} \quad (5.21)$$

due to [1], Lemma 4.2. Using (5.21), the influence of $AA_1 w_0$ on ρ_v can now be estimated by $O(\tau^2/t_v) = O(\tau/v)$. Thus there occurs an *order reduction* in ρ_v at the first grid points; these order reductions are, however, damped out like $1/v$.

Influence of the inhomogeneous term $-\tau^2 A_1 A_2 e_1(t_v)$ on ρ_v : Due to (5.11), (5.16) and since A_2 and $H(t) = \int_0^t e^{(t-s)A} ds$ commute, the corresponding contribution to ρ_v can be written as

$$-\tau \sum_{i=1}^v G^{v+1-i} \tau^2 A_1 \frac{1}{2} (HA_2 \Delta u^v)(t_i) + \tau \sum_{i=1}^v G^{v+1-i} \tau^2 A_1 A_2 (HA_1 w)(t_i). \quad (5.22)$$

Due to (5.18), (5.19) and due to $\|(HA_2)(t)\| \leq 1$ (cf. (5.14)), the norm of the first term in (5.22) can easily be estimated by $O(\tau^2)$.

Concerning the second term in (5.22), the estimate

$$\begin{aligned} \left\| \tau \sum_{i=1}^v G^{v+1-i} \tau^2 A_1 A_2 (HA_1 w)(t_i) \right\| \\ \leq \tau \|G^v - I\| \|(A - \tau A_1 A_2)^{-1} \tau A_1 A_2\| \max_{0 \leq t \leq t_v} \|(HA_1 w)(t)\| \end{aligned} \quad (5.23)$$

leads only to a $O(\tau)$ result, because $\|(A - \tau A_1 A_2)^{-1} \tau A_1 A_2\|$ is only $O(1)$ but not $O(\tau)$. Moreover, a damping effect cannot be expected here.

At first sight this leads to the negative result that ρ_v contains a term which can only be estimated by $O(\tau)$ (and not $O(\tau^2)$) and which is not subject to damping. For our model problem, however, an improved estimate can be derived; the basis for this estimate is given by Lemma 5.1, which makes use of a result by HUNDSORFER and VERWER ([10]; see also BRENNER, CROUZEIX and THOMÉE [4]).

Lemma 5.1. For a 'sufficiently smooth'¹³ spatial grid function v on Ω_h , the estimate

$$\|(A - \tau A_1 A_2)^{-1} \tau A_1 A_2 v\| \leq C_\gamma \tau^\gamma \quad (5.24)$$

holds for any $\gamma \in [0, 0.5)$, with a moderate-sized, h -independent constant C_γ .

Proof: It was shown in [10] that for 'sufficiently smooth'¹³ v , the estimate

$$\left\| \left(I - \frac{\tau}{2} A_i \right)^{-1} \tau A_i v \right\| \leq C_\gamma \tau^{\gamma/2} \quad (5.25)$$

¹³ By 'sufficiently smooth' we mean that v can be interpreted as the restriction $v = \Delta v$ of a smooth function $v(x)$, $x \in \Omega$, with a certain number of continuous, moderate-sized derivatives. To be more precise: $v \in C^3(\bar{\Omega})$ is required (see [10]). The crucial point is that v is not required to assume zero boundary values (in which case a $O(\tau)$ -estimate would be trivial).

is valid for arbitrary $\gamma \in [0, 0.5)$, with a moderate-sized, h -independent constant C_γ .¹⁴ Next we note that $(A - \tau A_1 A_2)^{-1} \tau A_1 A_2$ can be rewritten as

$$\begin{aligned} (A - \tau A_1 A_2)^{-1} \tau A_1 A_2 &= \left[\left(I - \frac{\tau}{2} A_1 \right)^{-1} \tau A_1 + \left(I - \frac{\tau}{2} A_2 \right)^{-1} \tau A_2 \right]^{-1} \\ &\quad \times \left(I - \frac{\tau}{2} A_1 \right)^{-1} \tau A_1 \left(I - \frac{\tau}{2} A_2 \right)^{-1} \tau A_2. \end{aligned} \quad (5.26)$$

Now we apply (5.25) for $i = 2$. Furthermore, we note that the application of $\left(I - \frac{\tau}{2} A_2 \right)^{-1} \tau A_2$ to the grid function v leaves the smoothness properties of v w.r.t. the space variable x_1 unaltered, because A_2 operates only w.r.t. the x_2 -direction. Thus we may once more apply (5.25), for $i = 1$ and with $\left(I - \frac{\tau}{2} A_2 \right)^{-1} \tau A_2 v$ instead of v . Moreover, $\left\| \left[\left(I - \frac{\tau}{2} A_1 \right)^{-1} \tau A_1 + \left(I - \frac{\tau}{2} A_2 \right)^{-1} \tau A_2 \right]^{-1} \right\|$ is uniformly bounded due to the fact that A_1 and A_2 are strictly negative definite. Thus, (5.24) is proved. \square

Application of Lemma 5.2 to the second term in (5.22), with the smooth grid function $v = w(t)$, yields the estimate

$$\begin{aligned} \left\| \tau \sum_{i=1}^v G^{v+1-i} \tau^2 A_1 A_2 (HA_1 w)(t_i) \right\| \\ \leq \tau \|G^v - I\| \max_{0 \leq t \leq t_v} \{ \|(HA_1)(t)\| \cdot \|(A - \tau A_1 A_2)^{-1} \tau A_1 A_2 w(t)\| \} \\ \leq C_\gamma \tau^{1+\gamma} \end{aligned} \quad (5.27)$$

with $\gamma \approx 0.5$.

Summarizing all that, we end up with

Theorem 5.2. The full global error $\eta_v - \Delta u(t_v)$ of the LOD scheme (5.2a), (5.2b) applied to the model problem (2D heat equation in a rectangle) admits an asymptotic expansion

$$\eta_v - \Delta u(t_v) = \tau e_1(t_v) + \rho_v \quad (5.28)$$

where $e_1(t)$ is a moderate, τ -independent solution of the variational equation (5.7); the remainder term ρ_v can be estimated by

$$\|\rho_v\| \leq C_\gamma \tau^{1+\gamma} + \left(C_0 + \frac{C_1}{t_v} \right) \tau^2 + C_s \max_{0 \leq t \leq t_v} \|s(t)\| \quad (5.29)$$

with $\gamma \approx 0.5$ and with certain h - and τ -independent, moderate-sized bounds C_γ , C_0 , C_1 and C_s .

¹⁴ Naturally, (5.25) is based on the corresponding estimate for the discrete 1D heat operator. Note that for $\gamma = 0.5$, C_γ would be affected by an additional factor $|\log h|$ (see [10]).

For a numerical illustration of this somewhat imperfect error structure of the LOD scheme, see section 6.

6. Numerical Examples; Discussion

In this final section our intention is to demonstrate the sharpness of our results about global error structures by means of simple but instructive numerical experiments¹⁵ (polynomial τ - resp. τ^2 -extrapolation). To make the error structure w.r.t. the time discretization clearly visible, we integrated along a smooth solution of the (finite dimensional) stiff system, avoiding space discretization errors.

Let us consider the 'Prothero-Robinson type' model system

$$u'(t) = A(u(t) - g(t)) + g'(t), \tag{6.1}$$

$$u(0) = g(0)$$

where the stiff matrix A arises from the standard space discretization (using central 2nd order difference quotients) of the 1D resp. 2D heat equation on a uniform mesh with spatial meshwidth $h = 1/64$. Except for special choices of g , the inhomogeneity of (6.1) is large in size (i.e., affected by a negative power of h). The true solution of (6.1) is $u(t) = g(t)$.

Table 6.1 displays the L_2 -norm of the global error $\eta_v - u(t_v)$ for the extrapolated Euler scheme, applied to the 1D test problem,¹⁶ at $t = 0.1$. Tables 6.2(a) and 6.2(b) show the results at $t = 2.0$ after performing global resp. local extrapolation. A similar result was already presented and discussed in [1]; global extrapolation clearly benefits from the fact that those error components which show a reduced order at the first grid points are damped out as the integration proceeds.

global error at $t = 0.1$		$h = 1/64$		Table 6.1	
τ	Euler	1st EX	2nd EX	3rd EX	
1/20	1.024E-03				
1/40	5.450E-04	6.629E-05			
1/80	2.820E-04	1.895E-05	3.192E-06		
1/160	1.435E-04	5.097E-06	4.807E-07	9.459E-08	
observed order					
	0.91				
	0.95	1.81			
	0.97	1.89	2.73		

¹⁵ All experiments were performed in ANSI double precision arithmetic.

¹⁶ $g(t)$ was chosen as the restriction of $e^{-t} \cos x$ to the spatial mesh.

global error at $t = 2.0$		$h = 1/64$		Table 6.2(a)	
τ	Euler	1st EX	2nd EX	3rd EX	
1/20	3.028E-04				
1/40	1.499E-04	2.963E-06			
1/80	7.460E-05	7.318E-07	1.176E-08		
1/160	3.721E-05	1.819E-07	1.455E-09	1.774E-11	
observed order					
	1.01				
	1.01	2.02			
	1.00	2.01	3.02		

global error at $t = 2.0$		$h = 1/64$		Table 6.2(b)	
τ	Euler	1st EX	2nd EX	3rd EX	
1/20	1.531E-04				
1/40	8.152E-05	9.925E-06			
1/80	4.218E-05	2.845E-06	4.871E-07		
1/160	2.148E-05	7.720E-07	8.158E-08	2.375E-08	
observed order					
	0.91				
	0.95	1.80			
	0.97	1.88	2.58		

The analogous results for τ^2 -extrapolation based on the implicit trapezoidal rule are displayed in Tables 6.3 and 6.4(a), (b). The very satisfactory performance of one extrapolation step, predicted by Theorem 3.2, is clearly visible; but the higher extrapolation steps yield no further increase in accuracy (at the τ^4 -level the global error is affected by 'irregular', oscillating terms). Furthermore, global extrapolation is again superior to local extrapolation; however, the difference is not so striking as for the Euler scheme (damping is only relevant in the range of mildly stiff eigenvalues).

global error at $t = 0.1$		$h = 1/64$		Table 6.3	
τ	ITR	1st EX	2nd EX	3rd EX	
1/20	9.880E-06				
1/40	2.445E-06	4.858E-08			
1/80	6.097E-07	2.592E-09	2.134E-09		
1/160	1.523E-07	1.466E-10	9.098E-11	1.012E-10	
observed order					
	2.01				
	2.00	4.23			
	2.00	4.14	4.55		

global error at $t = 2.0$ $h = 1/64$ Table 6.4(a)				
τ	GLOBAL EXTRAPOLATION			
	ITR	1st EX	2nd EX	3rd EX
1/20	2.474E-06			
1/40	6.186E-07	7.965E-10		
1/80	1.547E-07	2.779E-11	4.486E-11	
1/160	3.867E-08	6.832E-13	1.604E-12	1.780E-12
observed order				
	2.00			
	2.00	4.84		
	2.00	5.35	4.81	

global error at $t = 2.0$ $h = 1/64$ Table 6.4(b)				
τ	LOCAL EXTRAPOLATION			
	ITR	1st EX	2nd EX	3rd EX
1/20	1.478E-06			
1/40	3.657E-07	7.266E-09		
1/80	9.119E-08	3.883E-10	3.198E-10	
1/160	2.278E-08	2.182E-11	1.344E-11	1.501E-11
observed order				
	2.01			
	2.00	4.23		
	2.00	4.15	4.57	

It should also be mentioned that for the implicit midpoint rule (also called the Crank-Nicolson scheme), even the first extrapolation step is only partially successful (the typical oscillations occur already at the $O(\tau)^2$ -level).

Tables 6.5 and 6.6(a), (b) display the corresponding results for the locally one-dimensional splitting scheme applied to the 2D test problem.¹⁷ The results are in perfect agreement with our theoretical considerations (see section 5): It pays to apply one extrapolation step, but the achieved accuracy is only $\approx O(\tau^{1.5})$. For the higher extrapolation steps only a marginal increase in accuracy is observed. Furthermore, global extrapolation is not superior to local extrapolation, an obvious consequence of the fact that the $O(\tau^{1.5})$ -error terms are not subject to damping (see Theorem 5.2).

global error at $t = 0.1$ $h = 1/64$ Table 6.5				
τ	GLOBAL EXTRAPOLATION			
	LOD	1st EX	2nd EX	3rd EX
1/20	4.724E-03			
1/40	2.813E-03	9.210E-04		
1/80	1.554E-03	3.150E-04	1.329E-04	
1/160	8.219E-04	1.040E-04	4.410E-05	3.381E-05
observed order				
	0.75			
	0.86	1.55		
	0.92	1.60	1.59	

¹⁷ $g(t)$ was chosen as the restriction of $e^{-t} \cos(x_1 + x_2)$ to the spatial mesh.

global error at $t = 2.0$ $h = 1/64$ Table 6.6(a)				
τ	GLOBAL EXTRAPOLATION			
	LOD	1st EX	2nd EX	3rd EX
1/20	7.929E-04			
1/40	4.586E-04	1.301E-04		
1/80	2.499E-04	4.525E-05	1.984E-05	
1/160	1.312E-04	1.516E-05	6.597E-06	5.058E-06
observed order				
	0.79			
	0.88	1.52		
	0.93	1.58	1.59	

global error at $t = 2.0$ $h = 1/64$ Table 6.6(b)				
τ	LOCAL EXTRAPOLATION			
	LOD	1st EX	2nd EX	3rd EX
1/20	7.066E-04			
1/40	4.208E-04	1.378E-04		
1/80	2.325E-04	4.712E-05	1.989E-05	
1/160	1.230E-04	1.557E-05	6.599E-06	5.058E-06
observed order				
	0.75			
	0.86	1.55		
	0.92	1.60	1.59	

Remarks on the case of incompatible initial data. In several existing papers, numerical experience with extrapolation techniques is reported for the case of nonsmooth (incompatible) initial data; cf. for instance [6], [12], [20]. Some remarks on this case are in order:

The implicit Euler scheme damps high-frequency error components very rapidly and can therefore also be successfully be applied in the case of nonsmooth initial data (see also [13]). Furthermore, numerical experience shows that it pays to apply extrapolation (naturally, global extrapolation is the method of our choice here); a rigorous analysis has not been done so far but will not differ fundamentally from our analysis in [1]. (Anyhow, also in the smooth case we had to deal with incompatible initial data for the variational equations.)

Symmetric schemes suffer from the drawback that highly oscillatory error components are only damped under severe restrictions on the time step τ (cf. e.g. [15]). This lack of damping can, however, be overcome by algorithmic measures: Combining, for instance, the Crank-Nicolson scheme with four backward Euler steps at the beginning leads to a satisfactory damping of high frequency components (see [15]). But such an algorithmic modification is not compatible with a formal expansion in even powers of τ , which will further impair the efficiency of extrapolation.

Concerning the LOD scheme, it is worth mentioning that damping effects, which are of minor significance in the smooth case, seem to become relevant in the case of

nonsmooth data (cf. the numerical results reported in [20] for the globally extrapolated LOD scheme).

References

- [1] W. Auzinger, On the error structure of the implicit Euler scheme applied to stiff systems of differential equations, *Computing* 43, 115–131 (1989).
- [2] W. Auzinger, R. Frank, F. Macsek, Asymptotic error expansions for stiff equations: The implicit Euler scheme, *SIAM J. Numer. Anal.* 27, 1990.
- [3] W. Auzinger, R. Frank, Asymptotic error expansions for stiff equations: An analysis for the implicit midpoint and trapezoidal rules in the strongly stiff case, *Numer. Math.* 56, 469–499 (1989).
- [4] P. Brenner, M. Crouzeix, V. Thomée, Single step methods for inhomogeneous linear differential equations in Banach space, *RAIRO Numer. Anal.* 16, 5–26 (1982).
- [5] M. Crouzeix, P. A. Raviart, Approximation d'équations d'évolution linéaires par des méthodes multi-pas, *Etude numérique des grands systèmes*, Rencontre INRIA, Novosibirsk, Dunod, Paris, 1978.
- [6] A. R. Gourlay, J. Ll. Morris, The extrapolation of first order methods for parabolic partial differential equations II, *SIAM J. Numer. Anal.* 17, 641–655 (1980).
- [7] W. B. Gragg, Repeated extrapolation to the limit in the numerical solution of ordinary differential equations, Ph.D. Thesis, UCLA 1963.
- [8] E. Hairer, Ch. Lubich, Extrapolation at stiff differential equations, *Numer. Math.* 52, 377–400 (1988).
- [9] W. H. Hundsdorfer, Local and global order reduction for some LOD schemes, Report NM-R8914, Dept. of Numerical Mathematics, Centre for Mathematics and Computer Science, Amsterdam, 1989.
- [10] W. H. Hundsdorfer, J. G. Verwer, Stability and convergence of the Peaceman-Rachford ADI method for initial-boundary value problems, *Math. Comp.* 53, 81–101 (1989).
- [11] T. Kato, *Perturbation theory for linear operators*, Springer, Berlin-Heidelberg-New York, 1966.
- [12] J. D. Lawson, J. Ll. Morris, The extrapolation of first order methods for parabolic partial differential equations I, *SIAM J. Numer. Anal.* 15, 1212–1224 (1978).
- [13] M.-N. Le Roux, Semidiscretization in time for parabolic problems, *Math. Comp.* 33, 919–931 (1979).
- [14] G. I. Marchuk, V. V. Shaidurov, *Difference methods and their extrapolation*, Springer Applications of Mathematics 19, 1983.
- [15] R. Rannacher, Discretization of the heat equation with singular initial data, *ZAMM* 62, T346–T348, 1982.
- [16] R. D. Richtmeyer, K. W. Morton, *Difference methods for initial value problems*, Interscience, 1967.
- [17] J. M. Sanz-Serna, J. G. Verwer, Stability and convergence at the PDE/stiff ODE interface, *Applied Numerical Mathematics* 5, 117–132 (1989).
- [18] H. J. Stetter, Asymptotic expansions for the error of discretization algorithms for non-linear functional equations, *Numer. Math.* 7, 18–31 (1965).
- [19] V. Thomée, Galerkin finite element methods for parabolic problems, *Lecture Notes in Mathematics* 1054, Springer, 1984.
- [20] J. G. Verwer, H. B. De Vries, Global extrapolation of a first order splitting method, *SIAM J. Sci. Stat. Comput.* 6, 771–780 (1985).
- [21] N. N. Yanenko, *The method of fractional steps*, Springer, Berlin-Heidelberg-New York, 1971.

W. Auzinger
 Institut für Angewandte und Numerische Mathematik
 Technische Universität Wien
 Wiedner Hauptstrasse 6–10
 A-1040 Wien
 Austria