



Analysis of Approximate Factorization in Iteration Methods

C. Eichler-Liebenow, P.J. van der Houwen, B.P. Sommeijer

Modelling, Analysis and Simulation (MAS)

**MAS-R9718 August 31, 1997**

Report MAS-R9718  
ISSN 1386-3703

CWI  
P.O. Box 94079  
1090 GB Amsterdam  
The Netherlands

CWI is the National Research Institute for Mathematics and Computer Science. CWI is part of the Stichting Mathematisch Centrum (SMC), the Dutch foundation for promotion of mathematics and computer science and their applications.

SMC is sponsored by the Netherlands Organization for Scientific Research (NWO). CWI is a member of ERCIM, the European Research Consortium for Informatics and Mathematics.

Copyright © Stichting Mathematisch Centrum  
P.O. Box 94079, 1090 GB Amsterdam (NL)  
Kruislaan 413, 1098 SJ Amsterdam (NL)  
Telephone +31 20 592 9333  
Telefax +31 20 592 4199

# Analysis of Approximate Factorization in Iteration Methods

C. Eichler-Liebenow  
University Halle-Wittenberg  
P.O. Box 8, 06099 Halle / Saale, Germany

P.J. van der Houwen & B.P. Sommeijer  
CWI  
P.O. Box 94079, 1090GB Amsterdam, The Netherlands

## ABSTRACT

We consider the systems of ordinary differential equations (ODEs) obtained by spatial discretization of multi-dimensional partial differential equations. In order to solve the initial value problem (IVP) for such ODE systems numerically, we need a stiff IVP solver, because the Lipschitz constant associated with the righthand side function  $\mathbf{f}$  becomes increasingly large as the spatial resolution is refined. Stiff IVP solvers are necessarily implicit, so that we are faced with the problem of solving large systems of implicit relations. In the solution process of the implicit relations one may exploit the fact that the righthand side function  $\mathbf{f}$  can often be split into functions  $\mathbf{f}_i$  which contain only the discretizations of derivatives with respect to one spatial dimension. In this paper, we analyse iterative solution methods based on approximate factorization which are suitable for implementation on parallel computer systems. In particular, we derive convergence and stability regions.

*1991 Mathematics Subject Classification:* 65L06

*Keywords and Phrases:* numerical analysis, partial differential equations, iteration methods, approximate factorization, parallelism.

*Notes.* The investigations reported in this paper were partly supported by the Dutch HPCN Program. Work carried out under project MAS 1.2 - 'Numerical Algorithms for Surface Water Quality Modelling'.

## 1. Introduction

The systems of ordinary differential equations (ODEs) obtained by spatial discretization of initial-boundary value problems in  $d$  spatial dimensions (the method of lines), are often of the form

$$(1.1) \quad \frac{dy(t)}{dt} = \mathbf{f}(t, \mathbf{y}(t)) = \sum_{i=1}^d \mathbf{f}_i(t, \mathbf{y}(t)), \quad \mathbf{y}, \mathbf{f} \in \mathbb{R}^G,$$

where  $G$  is a usually large integer depending on the number of spatial grid points used and where the splitting of the righthand side function  $\mathbf{f}$  is such that the function  $\mathbf{f}_i$  contains only the discretizations of derivatives with respect to the  $i$ th spatial dimension. In order to solve the initial value problem (IVP) for the system (1.1) numerically, we need a stiff IVP solver, because the Lipschitz constant with respect to  $\mathbf{y}$  associated with the righthand side function  $\mathbf{f}$  becomes increasingly large as the spatial

resolution is refined. Stiff IVP solvers are necessarily implicit, so that we are faced with the problem of solving large systems of implicit relations. In this paper, we construct and analyse iterative methods for solving these implicit relations which exploit the fact that the righthand side function  $\mathbf{f}$  can be split according to (1.1). Our analysis applies to IVP solvers that fit into the wide class of General Linear Methods introduced by Butcher in 1966 (see [2, p. 335] for a detailed discussion). These methods are of the form

$$(1.2) \quad \mathbf{Y}_{n+1} - \Delta t(\mathbf{A} \otimes \mathbf{I})\mathbf{F}(\mathbf{e}t_n + \mathbf{c}\Delta t, \mathbf{Y}_{n+1}) = (\mathbf{B} \otimes \mathbf{I})\mathbf{Y}_n + \Delta t(\mathbf{C} \otimes \mathbf{I})\mathbf{F}(\mathbf{e}t_{n-1} + \mathbf{c}\Delta t, \mathbf{Y}_n), \quad n = 0, 1, \dots$$

Here  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{C}$  denote  $s$ -by- $s$  matrices,  $\mathbf{I}$  is the identity matrix whose order equals that of the system (1.1),  $\mathbf{e}$  is an  $s$ -dimensional vector with unit entries,  $\mathbf{c} = (c_i)$  is an  $s$ -dimensional abscissae vector,  $\Delta t$  is the stepsize  $t_{n+1} - t_n$ , and  $\otimes$  denotes the Kronecker product, i.e. if  $\mathbf{A} = (a_{ij})$ , then  $\mathbf{A} \otimes \mathbf{I}$  denotes the matrix of matrices  $(a_{ij}\mathbf{I})$ . Furthermore, for any vector  $\mathbf{Y}_n = (\mathbf{y}_{ni})$ ,  $\mathbf{F}(\mathbf{e}t_{n-1} + \mathbf{c}\Delta t, \mathbf{Y}_n)$  contains the derivative values  $(\mathbf{f}(t_{n-1} + c_i\Delta t, \mathbf{y}_{ni}))$ . The  $s$  vector components  $\mathbf{y}_{n+1,i}$  of  $\mathbf{Y}_{n+1}$  represent numerical approximations to the  $s$  exact solution vectors  $\mathbf{y}(t_n + c_i\Delta t)$ . The quantities  $\mathbf{Y}_n$  are usually called the *stage vectors* and their components  $\mathbf{y}_{ni}$  the *stage values*. We assume that the step point value  $\mathbf{y}_n$  is defined by the last component of  $\mathbf{Y}_n$ , i.e.  $\mathbf{y}_n := (\mathbf{e}_s^T \otimes \mathbf{I})\mathbf{Y}_n$ , where  $\mathbf{e}_s$  is the  $s$ th unit vector.

Each step by the method (1.2) requires the solution of the nonlinear system  $\mathbf{R}_n(\mathbf{Y}) = \mathbf{0}$  with

$$(1.3) \quad \mathbf{R}_n(\mathbf{Y}) := \mathbf{Y} - \Delta t(\mathbf{A} \otimes \mathbf{I})\mathbf{F}(\mathbf{e}t_n + \mathbf{c}\Delta t, \mathbf{Y}) - (\mathbf{B} \otimes \mathbf{I})\mathbf{Y}_n - \Delta t(\mathbf{C} \otimes \mathbf{I})\mathbf{F}(\mathbf{e}t_{n-1} + \mathbf{c}\Delta t, \mathbf{Y}_n).$$

In order to solve this system, we consider the modified Newton iteration process:

$$(1.4) \quad \mathbf{M}(\mathbf{Y}^{(j)} - \mathbf{Y}^{(j-1)}) = -\mathbf{R}_n(\mathbf{Y}^{(j-1)}), \quad j = 1, 2, \dots,$$

where  $\mathbf{M}$  is an approximation to the Jacobian matrix of  $\mathbf{R}_n(\mathbf{Y})$ . Evidently,  $\mathbf{M}$  is given by

$$(1.5) \quad \mathbf{M} = \mathbf{I} - \mathbf{A} \otimes \Delta t \mathbf{J} = \frac{1}{d} \sum_{i=1}^d (\mathbf{I} - d\mathbf{A} \otimes \Delta t \mathbf{J}_i),$$

where  $\mathbf{J}$  and  $\mathbf{J}_i$  are approximations to the Jacobian matrices of  $\mathbf{f}$  and  $\mathbf{f}_i$  with respect to  $\mathbf{y}$ , respectively. This expression shows that solving the multi-dimensionally linear Newton systems (1.4) by a direct method is quite costly. It is the aim of this paper to reduce these costs by designing a parallel iterative linear system solver based on an approximate factorization of the matrix  $\mathbf{M}$ . This linear solver may be considered as the *inner* iteration process and the Newton process (1.4) as the *outer* iteration process. For the inner iteration process, a number of convergence results are derived and for a *finite* number of inner and outer iterations in the inner-outer iteration process, we derive the stability matrix and the order of accuracy. For two- and three-dimensional problems, we compute stability regions for the 2nd-order backward differentiation method, the 3rd-order Radau IIA method, and a 3rd-order diagonally implicit general linear method (that is,  $\mathbf{A}$  is diagonal in (1.2)). For numerical results obtained by the approach analysed in the present paper, we refer to the references [5] and [10].

## 2. Iteration methods based on approximate factorization

Consider the outer-inner iteration process

$$(2.1) \quad \begin{aligned} \Pi \left( \mathbf{Y}^{(j,v)} - \mathbf{Y}^{(j,v-1)} \right) &= -\mathbf{M}\mathbf{Y}^{(j,v-1)} + \mathbf{M}\mathbf{Y}^{(j-1,r)} - \mathbf{R}_n(\mathbf{Y}^{(j-1,r)}), \quad v = 1, 2, \dots, r, \\ \Pi &:= \prod_{i=d}^1 \left( \mathbf{I} - \mathbf{A}^* \otimes \Delta t \mathbf{J}_i \right), \end{aligned}$$

where  $\mathbf{A}^*$  is a 'convenient' matrix which 'approximates' the matrix  $\mathbf{A}$ . Evidently, if the iterates  $\mathbf{Y}^{(j,v)}$  converge, then they can only converge to the Newton iterate  $\mathbf{Y}^{(j)}$ , irrespective the choice of  $\mathbf{A}^*$ .

Each inner iteration in (2.1) requires the solution of  $d$  linear systems with system matrix  $\mathbf{I} - \mathbf{A}^* \otimes \Delta t \mathbf{J}_i$  of order  $sG$ . The  $d$  LU-decompositions of the system matrices  $\mathbf{I} - \mathbf{A}^* \otimes \Delta t \mathbf{J}_i$  can be done in parallel, irrespective the choice of  $\mathbf{A}^*$ . Moreover, the matrices  $\mathbf{J}_i$  each correspond with a one-dimensional differential operator, so that solving these linear systems is relatively cheap.

We consider two options for choosing the matrix  $\mathbf{A}^*$ , viz. (i)  $\mathbf{A}^* = \mathbf{A}$  and (ii)  $\mathbf{A}^*$  similar to a diagonal matrix with real positive diagonal entries. If  $\mathbf{A}^* = \mathbf{A}$ , then the matrix  $\Pi$  is called the *approximate factorization* of the matrix  $\mathbf{M}$  [3, p. 439]. If  $\mathbf{A}^*$  is similar to a diagonalizable matrix  $\mathbf{D}$ , then we can diagonalize the iteration method (2.1) by means of a transformation  $\mathbf{Y}^{(j,v)} = (\mathbf{Q} \otimes \mathbf{I}) \tilde{\mathbf{Y}}^{(j,v)}$ , where  $\mathbf{Q}$  is such that  $\mathbf{D} := \mathbf{Q}^{-1} \mathbf{A}^* \mathbf{Q}$  is diagonal. Thus,

$$(2.1') \quad \begin{aligned} \tilde{\Pi} \left( \tilde{\mathbf{Y}}^{(j,v)} - \tilde{\mathbf{Y}}^{(j,v-1)} \right) &= -(\mathbf{Q}^{-1} \otimes \mathbf{I}) \mathbf{M} (\mathbf{Q} \otimes \mathbf{I}) \tilde{\mathbf{Y}}^{(j,v-1)} + (\mathbf{Q}^{-1} \otimes \mathbf{I}) \left( \mathbf{M}\mathbf{Y}^{(j-1,r)} - \mathbf{R}_n(\mathbf{Y}^{(j-1,r)}) \right), \\ \tilde{\Pi} &:= (\mathbf{Q}^{-1} \otimes \mathbf{I}) \Pi (\mathbf{Q} \otimes \mathbf{I}) = \prod_{i=d}^1 \left( \mathbf{I} - \mathbf{D} \otimes \Delta t \mathbf{J}_i \right), \quad v = 1, 2, \dots, r. \end{aligned}$$

Evidently, the factor matrices  $\mathbf{I} - \mathbf{D} \otimes \Delta t \mathbf{J}_i$  of the system matrix  $\tilde{\Pi}$  are block-diagonal. Hence, the diagonalized iteration method (2.1') allows for a considerable amount of additional parallelism, because the diagonal structure of  $\mathbf{D}$  enables us to decouple each of the linear systems into  $s$  subsystems which can be solved concurrently.

**Remark 2.1.** If only one inner iteration is performed with initial iterate  $\mathbf{Y}^{(j,0)} = \mathbf{Y}^{(j-1,1)}$ , then the outer-inner iteration process  $\{(1.4), (2.1)\}$  reduces to

$$(2.2) \quad \Pi \left( \mathbf{Y}^{(j)} - \mathbf{Y}^{(j-1)} \right) = -\mathbf{R}_n(\mathbf{Y}^{(j-1)}), \quad j = 1, 2, \dots$$

This method is related to the PDIRK and PTIRK methods proposed in [4] and [6] which arise if we set  $\tilde{\Pi} = \mathbf{I} - \mathbf{A}^* \otimes \Delta t \mathbf{J}$ , where  $\mathbf{J}$  is the full Jacobian of the righthand side in (1.1) with  $\mathbf{A}^*$  respectively a diagonal and triangular matrix. ♦

## 2.1. Region of convergence

The iteration error  $\mathbf{Y}^{(j,v)} - \mathbf{Y}^{(j)}$  associated with (2.1) satisfies the recursion

$$(2.3) \quad \mathbf{Y}^{(j,v)} - \mathbf{Y}^{(j)} = \mathbf{Z} \left( \mathbf{Y}^{(j,v-1)} - \mathbf{Y}^{(j)} \right), \quad \mathbf{Z} := \mathbf{I} - \Pi^{-1}\mathbf{M}, \quad v = 1, 2, \dots$$

Insight into the convergence of the inner iteration process is obtained by applying a normal mode analysis by assuming that the Jacobian matrices  $\mathbf{J}_i$  share the same eigensystem [3]. For brevity of notation, we introduce the following definition.

**Definition 2.1.** Let  $\mathbf{E}(\Delta t \mathbf{J}_1, \dots, \Delta t \mathbf{J}_d)$  be a matrix depending on  $\Delta t \mathbf{J}_1, \dots, \Delta t \mathbf{J}_d$ . Then  $\mathbf{E}(\mathbf{z})$  is the  $s$ -by- $s$  matrix obtained by replacing the matrices  $\Delta t \mathbf{J}_i$  by the scalars  $z_i$  and  $\mathbf{z} = (z_1, \dots, z_d)$ . ♦

Thus, with the matrices  $\mathbf{M}$  defined in (1.5),  $\Pi$  defined in (2.1), and  $\mathbf{Z}$  defined in (2.3) we associate the matrices

$$(2.4) \quad \mathbf{M}(\mathbf{z}) = \mathbf{I} - (\mathbf{e}^T \mathbf{z}) \mathbf{A}, \quad \Pi(\mathbf{z}) = \prod_{i=1}^d (\mathbf{I} - z_i \mathbf{A}^*), \quad \mathbf{Z}(\mathbf{z}) := \mathbf{I} - \Pi^{-1}(\mathbf{z}) \mathbf{M}(\mathbf{z}).$$

Evidently, if we choose  $z_i := \lambda(\mathbf{J}_i) \Delta t$  where  $\lambda(\mathbf{J}_i)$  denotes an eigenvalue of  $\mathbf{J}_i$ , then the eigenvalues of the amplification matrix  $\mathbf{Z}$  in (2.3) are given by those of the matrix  $\mathbf{Z}(\mathbf{z})$ . The region of convergence is now defined by the region in the  $\mathbf{z}$ -plane where  $\mathbf{Z}(\mathbf{z})$  has its eigenvalues  $\lambda(\mathbf{Z}(\mathbf{z}))$  within the unit circle. Assuming that the eigenvalues of the 'partial' Jacobians  $\mathbf{J}_i$  are in the nonpositive halfplane, we shall call the iteration method (2.1) *A-convergent* if the region of convergence contains the region  $\{\mathbf{z}: \operatorname{Re}(z_i) \leq 0\}$  and *A( $\alpha$ )-convergent* if it contains the region  $\{\mathbf{z}: |\arg(-z_i)| \leq \alpha\}$ . Furthermore, if  $\mathbb{D}$  is a domain in the complex plane and if  $\mathbf{z} = (z_1, \dots, z_d)$  lies in the convergence region whenever  $z_i$  lies in  $\mathbb{D}$  for  $i = 1, \dots, d$ , then we shall call the iteration method (2.1) *A( $\mathbb{D}$ )-convergent*. The eigenvalues  $\lambda(\mathbf{Z}(\mathbf{z}))$  of  $\mathbf{Z}(\mathbf{z})$  will be called the *amplification factors* of the inner iteration method.

## 2.2. Iteration with $\mathbf{A}^* = \mathbf{A}$

For  $\mathbf{A}^* = \mathbf{A}$  the amplification factors  $\lambda(\mathbf{Z}(\mathbf{z}))$  are given by

$$(2.5) \quad \lambda(\mathbf{Z}(\mathbf{z})) = 1 - \pi^{-1}(\mathbf{z}) \mu(\mathbf{z}), \quad \mu(\mathbf{z}) := 1 - \lambda(\mathbf{A})(\mathbf{e}^T \mathbf{z}), \quad \pi(\mathbf{z}) := \prod_{i=1}^d (1 - \lambda(\mathbf{A}) z_i),$$

where  $\lambda(\mathbf{A})$  denotes an eigenvalue of  $\mathbf{A}$ .

Let us first consider the amplification factors for small values of  $|z_i|$ . It is easily verified that we may write  $\lambda(\mathbf{Z}(\mathbf{z})) = O((\Delta t)^2)$ , so that our first conclusion is:

**Theorem 2.1.** If  $\mathbf{A}^* = \mathbf{A}$ , then the amplification factors  $\lambda(\mathbf{Z}(\mathbf{z}))$  of the inner iteration method (2.1) are second-order in  $\Delta t$ . ♦

This result implies that the low frequencies in the iteration error are strongly damped (we remark that for  $A^* \neq A$ , the amplification factors are in general only of  $O(\Delta t)$ , as may be seen from (2.4)).

**2.2.1. The two-dimensional case.** We first consider the convergence region of (2.1) in the two-dimensional case. Then, the amplification factor can be factorized according to

$$(2.5') \quad \lambda(Z(\mathbf{z})) = \lambda(A)z_1(1 - \lambda(A)z_1)^{-1} \lambda(A)z_2(1 - \lambda(A)z_2)^{-1}.$$

This immediately leads to the result:

**Theorem 2.2.** Let the eigenvalues of  $A$  be written in the form  $\lambda(A) = \xi \pm i\eta$ . Then the inner iteration method  $\{(2.1), A^* = A\}$  is  $A(\mathbb{D})$ -convergent for  $d = 2$  with

$$\mathbb{D} := \bigcap_{\lambda(A)} W(\xi, \eta), \quad W(\xi, \eta) := \{z: \xi \operatorname{Re}(z) \pm |\eta| \operatorname{Im}(z) < \frac{1}{2}\}. \quad \blacklozenge$$

If  $A$  has its eigenvalues in the right halfplane (as is the usual situation), then it is easily verified that we have  $A(\alpha)$ -convergence with  $\alpha = \min(\arctan(\xi|\eta|^{-1}))$ . Thus, we have the corollary:

**Corollary 2.1.** If  $\operatorname{Re} \lambda(A) \geq 0$ , then the inner iteration method  $\{(2.1), A^* = A\}$  is  $A(\alpha)$ -convergent for  $d = 2$  with  $\alpha = \min(\arctan(\xi|\eta|^{-1}))$ .  $\blacklozenge$

For example, for the L-stable, third-order 2-stage Radau IIA method defined by

$$(2.6) \quad \mathbf{c} = \begin{pmatrix} 1/3 \\ 1 \end{pmatrix}, \quad A = \frac{1}{12} \begin{pmatrix} 5 & -1 \\ 9 & 3 \end{pmatrix}, \quad B = \begin{pmatrix} 0 & 1 \\ 0 & 1 \end{pmatrix}, \quad C = O,$$

we have  $\xi|\eta|^{-1} = \sqrt{2}$ , so that the generated iteration method is  $A(54.7^\circ)$ -convergent. A second example is a one-parameter family of 3-stage methods based on Lagrange interpolation formulas, defined by (cf. [4])

$$(2.7) \quad \mathbf{c} = \begin{pmatrix} 0 \\ c \\ 1 \end{pmatrix}, \quad A = \frac{1}{6(1-c)} \begin{pmatrix} 0 & 0 & 0 \\ 3c-4c^2+c^3 & c(3-2c) & -c^3 \\ -c^{-1}+4-3c & c^{-1} & 2-3c \end{pmatrix}, \quad B = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix}, \quad C = O,$$

where  $c \neq 0$  and  $c \neq 1$ . This method is (at least) third-order accurate and A-stable for  $c \geq 1/2$ . For  $c = 1/2$  it equals the fourth-order Lobatto IIIA method and for  $c > 1/2$  it is strongly A-stable. The eigenvalues of  $A$  are given by  $\lambda(A) = 0$  and  $\lambda(A) = \frac{1}{6}(1 + c \pm \sqrt{(1+c)^2 - 6c})$ . Hence, we have A-convergence for  $d = 2$  if  $c \leq 2 - \sqrt{3}$  or  $c \geq 2 + \sqrt{3}$ . In the Lobatto case, we have  $\xi|\eta|^{-1} = \sqrt{3}$ , yielding an  $A(60^\circ)$ -convergent method. We remark that the method (2.7) also could have been formulated as a *two-stage* method with  $\mathbf{c} = (c, 1)^T$  by exploiting the matrix  $C$  (see e.g. Section 402 in [2] for an equivalence of both formulations).

**2.2.2. The case  $\lambda(\mathbf{A}) \geq 0$ .** It often happens that all eigenvalues of  $\mathbf{A}$  are *real* and *nonnegative*. Examples are the L-stable, second-order backward differentiation formula (BDF) defined by

$$(2.8a) \quad \mathbf{c} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \quad \mathbf{A} = \frac{1}{3} \begin{pmatrix} 0 & 0 \\ 0 & 2 \end{pmatrix}, \quad \mathbf{B} = \frac{1}{3} \begin{pmatrix} 0 & 3 \\ -1 & 4 \end{pmatrix}, \quad \mathbf{C} = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix},$$

and the strongly A-stable, third-order accurate method [9]

$$(2.8b) \quad \mathbf{c} = \frac{1}{10} \begin{pmatrix} 21 \\ 10 \end{pmatrix}, \quad \mathbf{A} = \frac{1}{660} \begin{pmatrix} 462 & 0 \\ 0 & 1430 \end{pmatrix}, \quad \mathbf{B} = \begin{pmatrix} 0 & 1 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{C} = \frac{1}{660} \begin{pmatrix} 441 & 483 \\ -1000 & 230 \end{pmatrix}$$

(other examples where  $\mathbf{A}$  has only nonnegative eigenvalues are the methods of Orel [8] and Bendtsen [1]). In such cases, more general convergence results can be derived. If  $\lambda(\mathbf{A}) \geq 0$ , then it follows from (2.5) that A(0)-convergence is achieved if  $2\pi(\mathbf{z}) - \mu(\mathbf{z}) > 0$  for  $z_i \leq 0$ . Since we may write

$$\pi(\mathbf{z}) = \mu(\mathbf{z}) + p_2\lambda^2(\mathbf{A}) + p_3\lambda^3(\mathbf{A}) + \dots + p_d\lambda^d(\mathbf{A}),$$

where the coefficients  $p_i$  are nonnegative whenever  $z_i \leq 0$ , we see that

$$2\pi(\mathbf{z}) - \mu(\mathbf{z}) = \mu(\mathbf{z}) + 2(p_2\lambda^2(\mathbf{A}) + p_3\lambda^3(\mathbf{A}) + \dots + p_d\lambda^d(\mathbf{A}))$$

is positive for  $\lambda(\mathbf{A}) \geq 0$  and  $z_i \leq 0$ . Thus, we have:

**Theorem 2.3.** If  $\mathbf{A}$  has only eigenvalues  $\lambda(\mathbf{A}) \geq 0$ , then the inner iteration method  $\{(2.1), \mathbf{A}^* = \mathbf{A}\}$  is A(0)-convergent for all  $d$ . ♦

In order to draw conclusions on the A-convergence of the iteration method, we first observe that for  $\lambda(\mathbf{A}) \geq 0$  and  $\text{Re}(z_i) \leq 0$ ,  $\lambda(\mathbf{Z}(\mathbf{z}))$  is analytic in each of its arguments  $z_i$ . Hence, we may restrict  $z_i$  to purely imaginary values. In this way, the following result is straightforwardly verified:

**Theorem 2.4.** If  $\mathbf{A}$  has only eigenvalues  $\lambda(\mathbf{A}) \geq 0$ , then the inner iteration method  $\{(2.1), \mathbf{A}^* = \mathbf{A}\}$  is only A-convergent for  $d \leq 2$ . ♦

In a number of important applications, we do not need A-convergence with respect to all spatial dimensions. For example, in 3-dimensional hydrodynamical applications, the vertical mesh size is an order of magnitude smaller than in the horizontal dimensions. Hence, the "stiffness" of the linear Newton systems (1.4) comes from the vertical direction, so that we only need unconditional convergence with respect to this direction. Let us consider the 3-dimensional case where  $\mathbf{A}$  has real eigenvalues and where we require A( $\alpha$ )-convergence with respect to only one spatial direction, say with respect to  $z_3$ . Then, defining the region

$$(2.9) \quad \mathbb{D}(r_0, \alpha) := \{(z_1, z_2, z_3): |\arg(-z_i)| \leq \alpha \ (i = 1, 2, 3), |z_1| \leq r_0, |z_2| \leq r_0, |z_3| \leq \infty\}.$$



we want to compute the value of  $r_0$  such that the region of convergence in the  $(z_1, z_2, z_3)$ -plane contains the domain  $\mathbb{D}(r_0, \alpha)$ . Table 2.1 lists these values for a few values of  $\alpha$ .

**Table 2.1.** Values of  $r_0$  such that  $\{(2.1), A^* = A, \lambda(A) \geq 0\}$  is convergent in  $\mathbb{D}(r_0, \alpha)$ .

$\alpha$	$\frac{\pi}{4}$	$\frac{11\pi}{40}$	$\frac{12\pi}{40}$	$\frac{13\pi}{40}$	$\frac{15\pi}{40}$	$\frac{17\pi}{40}$	$\frac{\pi}{2}$
$r_0 \leq$	$\infty$	$\frac{8.5}{\rho(A)}$	$\frac{4.0}{\rho(A)}$	$\frac{2.5}{\rho(A)}$	$\frac{1.5}{\rho(A)}$	$\frac{1}{\rho(A)}$	$\frac{0.5}{\rho(A)}$

Thus, these numerical calculations indicate that we have  $A(\pi/4)$ -convergent for  $d = 3$ . In fact, this result can be proved analytically by means of the following lemma:

**Lemma 2.1.** Let  $\mathbf{z} := (z_1, \dots, z_d)$ ,  $m(\mathbf{z}) := 1 - \mathbf{e}^T \mathbf{z}$  and  $p(\mathbf{z}) := (1 - z_1)(1 - z_2) \dots (1 - z_d)$ . If  $d = 3$ , then  $1 - p^{-1}(\mathbf{z})m(\mathbf{z})$  assumes values within the unit circle in the region  $\{\mathbf{z}: |\arg(-z_k)| \leq \pi/4\}$ .

**Proof.** If  $|p(\mathbf{z})|^2 - |p(\mathbf{z}) - m(\mathbf{z})|^2 > 0$ , then the function  $1 - p^{-1}(\mathbf{z})m(\mathbf{z})$  assumes values within the unit circle. Let us write  $z_k = r_k \exp(i\alpha_k)$ . It can be verified that

$$|p(\mathbf{z})|^2 - |p(\mathbf{z}) - m(\mathbf{z})|^2 = 1 + r_1^2 + r_2^2 + r_3^2 + 2E_1 + 4E_2 + 2E_3 + 2E_4,$$

where

$$E_1 := -r_1 \cos(\alpha_1) - r_2 \cos(\alpha_2) - r_3 \cos(\alpha_3),$$

$$E_2 := r_1 r_2 \cos(\alpha_1) \cos(\alpha_2) + r_1 r_3 \cos(\alpha_1) \cos(\alpha_3) + r_2 r_3 \cos(\alpha_2) \cos(\alpha_3),$$

$$E_3 := -r_1(r_2^2 + r_3^2) \cos(\alpha_1) - r_2(r_1^2 + r_3^2) \cos(\alpha_2) - r_3(r_1^2 + r_2^2) \cos(\alpha_3) \\ - 4r_1 r_2 r_3 \cos(\alpha_1) \cos(\alpha_2) \cos(\alpha_3),$$

$$E_4 := r_1 r_2 r_3 (r_3 \cos(\alpha_1 + \alpha_2) + r_2 \cos(\alpha_1 + \alpha_3) + r_1 \cos(\alpha_2 + \alpha_3)).$$

If the expressions  $E_i$  are nonnegative for  $r_k \geq 0$  and for  $3\pi/4 \leq \alpha_k \leq 5\pi/4$ , then the assertion of the lemma is true. Because  $\cos(\alpha_k) < 0$  and  $\cos(\alpha_j + \alpha_k) \geq 0$  for  $j, k = 1, 2, 3$  provided  $3\pi/4 \leq \alpha_k \leq 5\pi/4$ , it is immediate that we do have  $E_i \geq 0$  for  $i = 1, \dots, 4$ . This proves the lemma.  $\blacklozenge$

From (2.5) it follows that we may express  $\lambda(Z(\mathbf{z}))$  as  $1 - p^{-1}(\mathbf{z})m(\mathbf{z})$  with  $p = \pi$ ,  $m = \mu$  and  $\mathbf{z}$  replaced by  $\lambda(A)\mathbf{z}$ . Hence, Lemma 2.1 can be applied showing that  $\lambda(Z(\mathbf{z}))$  assumes values within the unit circle in the region  $\{\mathbf{z}: |\arg(-z_k)| \leq \pi/4\}$  provided that  $\lambda(A) \geq 0$ . Thus, we have the result:

**Theorem 2.5.** If  $A$  has eigenvalues  $\lambda(A) \geq 0$ , then the inner iteration method  $\{(2.1), A^* = A\}$  is  $A(\pi/4)$ -convergent for  $d = 3$ .

Analogous results for  $d > 3$  can be obtained along the same lines, but the derivation becomes increasingly tedious. An alternative derivation might be the approach described by Hundsdorfer [7] who derives conditions such that the function  $1 + p^{-1}(\mathbf{z}/2)(\mathbf{e}^T \mathbf{z})$  assumes values within the unit circle for arbitrary values of  $d$ .

### 2.3. Iteration with $A^* \neq A$

It is only feasible to choose  $A^* = A$  either if the dimension  $sG$  of  $I - A^* \otimes \Delta t J_i$  is sufficiently small (thereby restricting the size of the problem) or if  $A$  is similar to a diagonal matrix with real eigenvalues. Hence, for large scale problems we should resort to methods with a real-spectrum matrix  $A$ . However, in order to achieve both  $A$ -stability and higher-order accuracy with respect to time, the dimension  $s$  of the matrix  $A$  should be sufficiently large, implying an increased number of processors (we cannot use higher-order BDFs requiring only one processor, because the third and higher-order BDFs are not  $A$ -stable). Hence, if we want a higher-order method and if we want to minimize the number of processors, then it is of interest to investigate what can be achieved by using matrices  $A^* \neq A$ . Moreover, it is of interest to know whether the convergence region can be improved by widening the set of matrices  $A^*$ .

A first consequence of using matrices  $A^* \neq A$  is that in general the amplification factors  $\lambda(Z(\mathbf{z}))$  for the nonstiff error components are not anymore  $O((\Delta t)^2)$  but  $O(\Delta t)$ . Secondly, at infinity the behaviour of  $\lambda(Z(\mathbf{z}))$  may also be quite different. For example, along the  $z_j$ -axis, we have  $Z(\mathbf{z}) = I - (I - z_j A^*)^{-1}(I - z_j A)$ , so that  $Z(\mathbf{z}) \approx I - (A^*)^{-1}A$  as  $z_j \rightarrow \infty$ . Hence, a *necessary* condition for  $A(\alpha)$ -convergence requires  $A^*$  to satisfy  $\rho(I - (A^*)^{-1}A) \leq 1$ . Evidently, this condition is trivially satisfied if  $A^* = A$ . Let us also look at  $Z(\mathbf{z})$  as all components  $z_j$  tend to infinity. Then, it is readily verified that

$$Z(\mathbf{z}) \approx I + (-1)^d \frac{z_1 + \dots + z_d}{z_1 \cdot \dots \cdot z_d} (A^*)^{-d} A.$$

In order to have  $A(\alpha)$ -convergence, it is necessary that this matrix has its eigenvalues on the unit disk for  $z_j = r_j \exp(i\phi)$  with  $|\phi - \pi| \leq \alpha$  and  $r_j \rightarrow \infty, j = 1, \dots, d$ . By writing

$$(2.10) \quad Z(\mathbf{z}) \approx I + x (-1)^d e^{i(1-d)\phi} (A^*)^{-d} A, \quad x := \frac{r_1 + \dots + r_d}{r_1 \cdot \dots \cdot r_d},$$

and defining  $\psi_d := (1-d)\phi$  and  $\lambda((A^*)^{-d} A) = \xi_d^* \pm i \eta_d^*$ , the eigenvalues of  $Z(\mathbf{z})$  are given by

$$\lambda(Z(\mathbf{z})) \approx 1 + x (-1)^d e^{i\psi_d} \lambda((A^*)^{-d} A) = 1 + x (-1)^d (\cos(\psi_d) + i \sin(\psi_d)) (\xi_d^* \pm i \eta_d^*).$$

Hence,

$$(2.11) \quad |\lambda(Z(\mathbf{z}))|^2 \approx (1 - x(-1)^d [-\xi_d^* \cos(\psi_d) \pm \eta_d^* \sin(\psi_d)])^2 + x^2 (\xi_d^* \sin(\psi_d) \pm \eta_d^* \cos(\psi_d))^2.$$

At  $x = 0$ , we have  $\lambda(Z(\mathbf{z})) = 1$ . In order to achieve that  $|\lambda(Z(\mathbf{z}))|$  does not immediately increase beyond 1 as  $x$  increases, we should require that  $(-1)^d [-\xi_d^* \cos(\psi_d) \pm |\eta_d^*| \sin(\psi_d)] \geq 0$ . On substitution of  $\psi_d = (1-d)(\pi \pm \alpha)$ , we find by an elementary manipulation that this inequality is satisfied if  $\pm \tan((d-1)\alpha) \leq \xi_d^* |\eta_d^*|^{-1}$  and  $0 \leq \alpha \leq \frac{\pi}{2(d-1)}$ . Thus, together with the earlier condition  $\rho(I - (A^*)^{-1}A) \leq 1$ , we arrive at the result:

**Theorem 2.6.** Let  $\xi_d^* \pm i \eta_d^* := \lambda((A^*)^{-d}A)$ . Then necessary conditions for  $A(\alpha)$ -convergence of the inner iteration method (2.1) are

- (i)  $-2\xi_1^* + [\xi_1^*]^2 + [\eta_1^*]^2 \leq 0$  for all eigenvalues  $\lambda((A^*)^{-1}A)$ ,
- (ii)  $\xi_d^* \geq 0$ ,  $0 \leq \alpha \leq \frac{1}{d-1} \arctan(\xi_d^* |\eta_d^*|^{-1})$  for all eigenvalues  $\lambda((A^*)^{-d}A)$ . ♦

From this theorem it follows that at best we may hope for  $A$ -convergence if  $d = 2$  and for  $A(45^\circ)$ -convergence if  $d = 3$ .

**2.3.1. The two-dimensional case.** Let us first look at the two-dimensional case. We have investigated the case where  $A$  is defined by the 2-stage Radau IIA matrix in (2.6) and  $A^*$  is the PDIRK matrix used in PDIRK methods (cf. [4], see also Remark 2.1), i.e.

$$(2.12a) \quad A^* = \frac{1}{30} \begin{pmatrix} 20-5\sqrt{6} & 0 \\ 0 & 12+3\sqrt{6} \end{pmatrix}.$$

The first condition of Theorem 2.6 is satisfied because this matrix was constructed such that  $(A^*)^{-1}A$  has eigenvalues 1. In order to check the second necessary condition, we should consider the matrix  $(A^*)^{-d}A$  for  $d = 2$ . We find

$$(A^*)^{-2}A \approx \begin{pmatrix} 6.24 & -1.25 \\ 1.80 & 0.60 \end{pmatrix}.$$

This matrix has two positive eigenvalues, so that the second condition of Theorem 2.6 is satisfied for  $0 \leq \alpha \leq \pi/2$ . By means of numerical verification, we found that we do have  $A$ -convergence.

**Result 2.1a.** If  $A$  is defined by the 2-stage Radau IIA matrix, then the inner iteration method  $\{(2.1), (2.12a)\}$  is  $A$ -convergent for  $d = 2$ . ♦

Note that the iteration method  $\{(2.1), (2.6)\}$  is not  $A$ -convergent but only  $A(54.7^\circ)$ -convergent (cf. Corollary 2.1). Hence, this is an example which shows that the convergence region can be improved by widening the set of matrices  $A^*$ .

An alternative option for  $A^*$  is such that the eigenvalues of  $Z(z_1, z_2)$  are  $O((\Delta t)^2)$ , resulting in a strong damping of the nonstiff error components. This is achieved by requiring  $\rho(A - A^*) = 0$  and yields two potential matrices  $A^*$ . However, the first necessary condition of Theorem 2.6 leaves us with the matrix

$$(2.12b) \quad A^* = \frac{1}{6} \begin{pmatrix} 1 & 0 \\ 0 & 3 \end{pmatrix}.$$

In this case, the iteration method is not anymore A-convergent, because the matrix

$$(A^*)^{-2}A = \begin{pmatrix} 15 & -3 \\ 3 & 1 \end{pmatrix}$$

has eigenvalues  $8 \pm \sqrt{40}$ . Hence, Theorem 2.6 implies that at best we have  $A(\alpha)$ -convergence with  $\alpha \leq \arctan(8/\sqrt{40}) \approx 51.6^\circ$ . In fact, a numerical calculation yields:

**Result 2.1b.** If  $A$  is defined by the 2-stage Radau IIA matrix, then the inner iteration method  $\{(2.1), (2.12b)\}$  is  $A(48^\circ)$ -convergent for  $d = 2$ .  $\blacklozenge$

**2.3.2. The three-dimensional case.** We conclude our convergence considerations with applying the approach described above to the 2-stage Radau IIA method in a three-dimensional problem. Omitting the details, we state the result:

**Result 2.2.** If  $A$  is defined by the 2-stage Radau IIA matrix, then the inner iteration method  $\{(2.1), (2.12a)\}$  is  $A(45^\circ)$ -convergent for  $d = 3$ .  $\blacklozenge$

## 2.4. Stability

In actual computation, we often do not iterate the outer and inner iteration process until convergence, particularly in the case of three-dimensional problems. Consequently, the stability properties of the resulting integration scheme will not be identical to those of the underlying integration method (1.2). In order to see the effect of the number of inner iterations on the stability, we consider the stability test equation  $\mathbf{y}' = \mathbf{J}\mathbf{y}$ . Let the predictor for the inner and outer iteration process be of the form

$$(2.13) \quad \mathbf{Y}^{(j,0)} = \mathbf{Y}^{(j-1,r)}, \quad \mathbf{Y}^{(0,r)} = \mathbf{P}\mathbf{Y}_n,$$

where  $\mathbf{P}$  is the predictor matrix. From (1.4) it follows that for the test equation  $\mathbf{Y}^{(j)} = \mathbf{M}^{-1}\mathbf{N}\mathbf{Y}_n$ , so that by virtue of (2.3)

$$(2.14) \quad \mathbf{Y}^{(j,r)} = (\mathbf{I} - \mathbf{Z}^r)\mathbf{M}^{-1}\mathbf{N}\mathbf{Y}_n + \mathbf{Z}^r\mathbf{Y}^{(j-1,r)}, \quad \mathbf{N} := \mathbf{B} \otimes \mathbf{I} + \mathbf{C} \otimes \Delta t \mathbf{J},$$

where  $r$  is the number of iterations in the inner iteration process. Suppose that  $m$  outer iterations are performed and that  $\mathbf{Y}_{n+1}$  is identified with  $\mathbf{Y}^{(m,r)}$ . Then, we can write  $\mathbf{Y}_{n+1}$  in the form

$$(2.15) \quad \mathbf{Y}_{n+1} = \mathbf{S}_{\text{mr}} \mathbf{Y}_n, \quad \mathbf{S}_{\text{mr}} := \mathbf{M}^{-1} \mathbf{N} + \mathbf{Z}^{\text{mr}} (\mathbf{P} - \mathbf{M}^{-1} \mathbf{N}).$$

We have stability if the stability matrix  $\mathbf{S}_{\text{mr}}$  has its eigenvalues on the unit disk. Assuming that the matrix  $\mathbf{P}$  is an expression in terms of  $\Delta t \mathbf{J}_i$ , we find that the eigenvalues of  $\mathbf{S}_{\text{mr}}$  are given by the eigenvalues of the  $s$ -by- $s$  matrix  $\mathbf{S}_{\text{mr}}(\mathbf{z})$ , where  $\mathbf{z} = (z_1, \dots, z_d)$  with  $z_i := \lambda(\mathbf{J}_i) \Delta t$  and where  $\mathbf{S}_{\text{mr}}(\mathbf{z})$  is defined according to Definition 2.1. Thus, we have the stability result:

**Theorem 2.7.** Let the inner and outer predictors be of the form (2.13) and let the matrices  $\mathbf{M}(\mathbf{z})$ ,  $\mathbf{\Pi}(\mathbf{z})$  and  $\mathbf{Z}(\mathbf{z})$  be defined by (2.4). Then, the outer-inner iteration process is stable at the point  $\mathbf{z}$  if the  $s$ -by- $s$  stability matrix

$$(2.16) \quad \mathbf{S}_{\text{mr}}(\mathbf{z}) := \mathbf{M}^{-1}(\mathbf{z}) \mathbf{N}(\mathbf{z}) + \mathbf{Z}^{\text{mr}}(\mathbf{z}) (\mathbf{P}(\mathbf{z}) - \mathbf{M}^{-1}(\mathbf{z}) \mathbf{N}(\mathbf{z})), \quad \mathbf{N}(\mathbf{z}) := \mathbf{B} + (\mathbf{e}^T \mathbf{z}) \mathbf{C}$$

has its eigenvalues on the unit disk.  $\blacklozenge$

Thus, the region of stability is given by the region in the  $\mathbf{z}$ -plane where  $\mathbf{S}_{\text{mr}}(\mathbf{z})$  has its eigenvalues on the unit disk. Assuming that the eigenvalues of the 'partial' Jacobians  $\mathbf{J}_i$  are in the nonpositive halfplane, we shall call the outer-inner iteration process *A-stable* if the region of stability contains the region  $\mathbb{S} := \{\mathbf{z}: \text{Re}(z_i) \leq 0\}$  and *A( $\alpha$ )-stable* if it contains the region  $\mathbb{S}(\alpha) := \{\mathbf{z}: |\arg(-z_i)| \leq \alpha\}$ .

In the case of Runge-Kutta methods (RK methods) using the last step value (LSV) predictor as outer predictor, the stability matrix  $\mathbf{S}_{\text{mr}}(\mathbf{z})$  can be reduced to a scalar stability function. Formulating such a method with a zero  $\mathbf{C}$ -matrix (see the discussion at the end of Subsection 2.2.1.), we obtain that  $\mathbf{N}(\mathbf{z})$  equals the matrix  $\mathbf{B}$  and is of the form  $\mathbf{e} \mathbf{e}_s^T$ . Furthermore, also the matrix  $\mathbf{P}(\mathbf{z})$  corresponding to the LSV predictor is of this form. As a result, the stability matrix  $\mathbf{S}_{\text{mr}}(\mathbf{z})$  in (2.16) will have zero vectors in its first  $s-1$  columns, as well. Recalling that  $\mathbf{y}_n = (\mathbf{e}_s^T \otimes \mathbf{I}) \mathbf{Y}_n$ , we see that we only need the last element in the last column of  $\mathbf{S}_{\text{mr}}(\mathbf{z})$ , resulting in a scalar stability function.

Using the above observations, (2.15) takes the form

$$(2.15') \quad \mathbf{y}_{n+1} = (\mathbf{e}_s^T \otimes \mathbf{I}) \mathbf{S}_{\text{mr}} \mathbf{Y}_n = (\mathbf{e}_s^T \otimes \mathbf{I}) (\mathbf{M}^{-1} + \mathbf{Z}^{\text{mr}} (\mathbf{I} - \mathbf{M}^{-1})) (\mathbf{e} \otimes \mathbf{I}) \mathbf{y}_n.$$

The eigenvalues of the matrix  $(\mathbf{e}_s^T \otimes \mathbf{I}) (\mathbf{M}^{-1} + \mathbf{Z}^{\text{mr}} (\mathbf{I} - \mathbf{M}^{-1})) (\mathbf{e} \otimes \mathbf{I})$  are given by the eigenvalues of  $\mathbf{e}_s^T [\mathbf{M}^{-1}(\mathbf{z}) + \mathbf{Z}^{\text{mr}}(\mathbf{z}) (\mathbf{I} - \mathbf{M}^{-1}(\mathbf{z}))] \mathbf{e}$ . Observing that  $\mathbf{e}_s^T \mathbf{M}^{-1}(\mathbf{z}) \mathbf{e}$  is the stability function of the underlying RK method, we obtain the result:

**Theorem 2.8.** Let the outer predictor be the LSV predictor  $\mathbf{P} = (\mathbf{e} \mathbf{e}_s^T \otimes \mathbf{I})$  and let  $\mathbf{R}$  be the stability function of the underlying RK method. Then, the iterated RK method is stable at the point  $\mathbf{z}$  if the stability function

$$(2.17) \quad \mathbf{R}_{\text{mr}}(\mathbf{z}) := \mathbf{R}(\mathbf{e}^T \mathbf{z}) + \mathbf{e}_s^T \mathbf{Z}^{\text{mr}}(\mathbf{z}) (\mathbf{I} - \mathbf{M}^{-1}(\mathbf{z})) \mathbf{e},$$

assumes values on the unit disk.  $\blacklozenge$

In the case of (2.17), more insight into the stability region can be obtained by considering the behaviour of the stability function at the origin and at infinity. At the origin, we obtain

$$(2.18) \quad \mathbf{R}_{\text{mr}}(\mathbf{z}) = \mathbf{R}(\mathbf{z}) - z^{\text{mr}+1} \mathbf{e}_s^T(\mathbf{A} - \mathbf{A}^* + \mathbf{O}(\mathbf{z}))^{\text{mr}}(\mathbf{A} + \mathbf{O}(\mathbf{z}))\mathbf{e}, \quad \mathbf{z} := \mathbf{e}^T\mathbf{z}.$$

To achieve that the stability region at least contains a 'left' neighbourhood of the origin in all complex  $z_j$ -planes, we write (2.18) as  $\mathbf{R}_{\text{mr}}(\mathbf{z}) = 1 + z + az^2 + bz^3 + cz^4 + \dots$ , and we require that this function has a nonzero imaginary stability boundary, that is there is a nonzero interval  $(0, i\beta)$  on the  $z_j$ -axes where  $|\mathbf{R}_{\text{mr}}(\mathbf{z})| < 1$ . This leads to the condition  $1 - 2a + \beta^2(a^2 - 2b + 2c) + \mathbf{O}(\beta^3) < 0$ , so that the first 5 terms of  $\mathbf{R}_{\text{mr}}(\mathbf{z})$  determine whether  $\beta$  is zero or not. Hence, if  $\mathbf{A}^* = \mathbf{A}$  and  $\text{mr} \geq 2$  or if  $\mathbf{A}^* \neq \mathbf{A}$  and  $\text{mr} \geq 4$ , then we obtain a positive imaginary stability boundary  $\beta$  whenever the stability function  $\mathbf{R}$  of the underlying RK method has a nonzero imaginary stability boundary.

In order to investigate the behaviour at infinity, we write  $z_j = r_j \exp(i\phi)$  with  $|\phi - \pi| \leq \alpha$  and  $r_j \rightarrow \infty$ ,  $j = 1, \dots, d$ . Writing the amplification matrix  $\mathbf{Z}(\mathbf{z})$  for large  $r_j$  as  $\mathbf{I} + x\mathbf{B}$ , where  $x$  and  $\mathbf{B}$  follow from (2.10) and using the expansion

$$(2.19) \quad \mathbf{M}^{-1}(\mathbf{z}) = (\mathbf{I} - (\mathbf{e}^T\mathbf{z})\mathbf{A})^{-1} = -(\mathbf{e}^T\mathbf{z})^{-1}\mathbf{A}^{-1} + \mathbf{O}(\mathbf{z}^{-2}),$$

we find

$$\mathbf{R}_{\text{mr}}(\mathbf{z}) \approx \mathbf{e}_s^T(\mathbf{M}^{-1}(\mathbf{z}) + (\mathbf{I} + x\mathbf{B})^{\text{mr}}(\mathbf{I} - \mathbf{M}^{-1}(\mathbf{z})))\mathbf{e} \approx \mathbf{e}_s^T(\mathbf{I} + \text{mr}x\mathbf{B} + \mathbf{O}(x\mathbf{z}^{-1}))\mathbf{e}.$$

Hence, using (2.10), we obtain

$$(2.20) \quad \mathbf{R}_{\text{mr}}(\mathbf{z}) \approx 1 + \text{mr}x (-1)^d e^{i(1-d)\phi} \mathbf{e}_s^T((\mathbf{A}^*)^{-d}\mathbf{A})\mathbf{e}, \quad x := \frac{r_1 + \dots + r_d}{r_1 \cdot \dots \cdot r_d},$$

$$|\mathbf{R}_{\text{mr}}(\mathbf{z})|^2 \approx \left(1 + \text{mr}x (-1)^d \mathbf{e}_s^T((\mathbf{A}^*)^{-d}\mathbf{A})\mathbf{e} \cos(\psi_d)\right)^2 + (\text{mr}x)^2 \left(\mathbf{e}_s^T((\mathbf{A}^*)^{-d}\mathbf{A})\mathbf{e} \sin(\psi_d)\right)^2,$$

where  $\psi_d := (1 - d)\phi$ . Similar to our discussion of (2.11), we conclude that a necessary condition for  $\mathbf{A}(\alpha)$ -stability is that  $(-1)^d \mathbf{e}_s^T((\mathbf{A}^*)^{-d}\mathbf{A})\mathbf{e} \cos(\psi_d) \leq 0$ . This leads to the result:

**Theorem 2.9.** Let the conditions of Theorem 2.8 be satisfied. Then, irrespective the value of  $\text{mr}$ , necessary conditions for  $\mathbf{A}(\alpha)$ -stability of the iterated RK method are

$$(2.21) \quad \mathbf{e}_s^T((\mathbf{A}^*)^{-d}\mathbf{A})\mathbf{e} \geq 0, \quad 0 \leq \alpha \leq \frac{\pi}{2(d-1)}. \quad \blacklozenge$$

It should be remarked that the condition on  $\alpha$  is identical with the  $\mathbf{A}(\alpha)$ -stability condition obtained by Hundsdorfer [7] for the Douglas splitting method.

Similar to the convergence behaviour (compare Theorem 2.6), we may only hope for  $\mathbf{A}$ -stability if  $d = 2$ , but not anymore if  $d > 2$ , whatever the value of  $\text{mr}$  is. The *actual* value of  $\alpha$  is expected to depend on  $\text{mr}$ . Furthermore, observing that  $\mathbf{e}_s^T((\mathbf{A}^*)^{-d}\mathbf{A})\mathbf{e}$  is the row sum of the last row in  $(\mathbf{A}^*)^{-d}\mathbf{A}$ ,

we see that the first condition of the theorem is always satisfied in the following two special situations:

- (i)  $A^*$  diagonal, last diagonal entry of  $A^*$  and last row sum of  $A$  positive
- (ii)  $A^* = A$ , last row sum of  $A^{1-d}$  positive.

As an example of the first situation, consider the 2-stage Radau IIA matrix in (2.6) with  $A^*$  defined by (2.12a) or (2.12b). In both cases, the first condition of Theorem 2.9 is satisfied, so that  $A(\pi/(2d-2))$ -stability may be possible. A numerical verification reveals that for  $d = 2$  the choice (2.12a) leads to  $A$ -stability for all values of  $mr$ , while (2.12b) leads to  $A(\alpha)$ -stability where  $\alpha$  *decreases* as  $mr$  increases (see Table 3.1 for details). For  $d = 3$ , the case (2.12a) yields  $A(45^\circ)$ -stability (see Table 3.2).

In order to illustrate the second situation, we consider the 2-stage Radau IIA matrix with  $A^* = A$  and  $d = 2$ , to obtain

$$(A^*)^{-2}A = A^{-1} = \frac{1}{2} \begin{pmatrix} 3 & 1 \\ -9 & 5 \end{pmatrix}.$$

Hence, it follows from Theorem 2.9 that there is no  $\alpha$ -value for which we have  $A(\alpha)$ -stability. However, if we take the 3-stage Radau IIA method with  $A^* = A$  and  $d = 2$ , then the last row sum of the matrix  $(A^*)^{-2}A = A^{-1}$  turns out to be 3 showing that the necessary  $A(\alpha)$ -stability condition is satisfied for  $\alpha \leq \pi/2$ . However, a numerical verification reveals that the iterated method is not  $A(\alpha)$ -stable, irrespective the value of  $mr$ . Using for  $A^*$  the matrix (2.12a) does yield an  $A(\alpha)$ -stable process with  $\alpha$  *increasing* as  $mr$  increases.

## 2.5. Order of accuracy

In order to derive the order of accuracy after a finite number of inner and outer iterations, we consider the iteration error  $\epsilon^{(j,v)} := \mathbf{Y}^{(j,v)} - \mathbf{Y}_{n+1}$ . From (1.2), (1.3) and (2.1) it follows that we may write

$$(2.22) \quad \begin{aligned} \epsilon^{(j,v)} &= Z \epsilon^{(j,v-1)} + \Delta t \Pi^{-1}(A \otimes I) \mathbf{G}_n(\epsilon^{(j-1,r)}), \\ \mathbf{G}_n(\epsilon) &:= \mathbf{F}(\mathbf{e}t_n + \mathbf{c}\Delta t, \mathbf{Y}_{n+1} + \epsilon) - \mathbf{F}(\mathbf{e}t_n + \mathbf{c}\Delta t, \mathbf{Y}_{n+1}) - (I \otimes J)\epsilon, \end{aligned}$$

where  $J$  is the same approximation to the Jacobian matrix as used in (1.5). After  $r$  inner iterations, this recursion yields

$$(2.23) \quad \epsilon^{(j,r)} = Z^r \epsilon^{(j-1,r)} + \Delta t (I - Z^r) M^{-1}(A \otimes I) \mathbf{G}_n(\epsilon^{(j-1,r)}),$$

where we assumed that  $\mathbf{Y}^{(j,0)} = \mathbf{Y}^{(j-1,r)}$ , i.e.  $\epsilon^{(j,0)} = \epsilon^{(j-1,r)}$ . Let  $\mathbf{G}_n$  possess a Lipschitz constant  $L_n(\Delta t)$  in the neighbourhood of the origin (with respect to the norm  $\|\cdot\|$ ) and let

$$(2.24) \quad Z^r = O((\Delta t)^{\theta_r}), \quad L_n(\Delta t) = O((\Delta t)^u).$$

The value of  $\theta$  follows from the expansion  $Z^r = ((A - A^*) \otimes \Delta t J + O(\Delta t^2))^r$ . Hence,  $\theta = 1$  if  $A^* \neq A$  and  $\theta = 2$  if  $A^* = A$ . Furthermore,  $u = 1$  if the Jacobian is updated every few integration steps, and  $u = 0$  if very crude approximations to the Jacobian are used. Let the method (1.2) have step point order  $p$ . Then, after  $m$  outer iterations and  $r$  inner iterations in each outer iteration, we have for the local error at the step point  $t_{n+1}$

$$\mathbf{e}_s^T \mathbf{Y}^{(m,r)} - \mathbf{y}(t_{n+1}) = \mathbf{e}_s^T (\mathbf{Y}^{(m,r)} - \mathbf{Y}_{n+1}) + \mathbf{e}_s^T \mathbf{Y}_{n+1} - \mathbf{y}(t_{n+1}) = \boldsymbol{\varepsilon}^{(m,r)} + O((\Delta t)^{p+1}).$$

From (2.23) it follows that

$$\|\boldsymbol{\varepsilon}^{(j,r)}\| \leq (O((\Delta t)^{\theta r}) + O((\Delta t)^{u+1})) \|\boldsymbol{\varepsilon}^{(j-1,r)}\|, \quad j \geq 1.$$

Hence,

$$\|\boldsymbol{\varepsilon}^{(m,r)}\| \leq (O((\Delta t)^{\theta r}) + O((\Delta t)^{u+1}))^m \|\mathbf{Y}^{(0,r)} - \mathbf{Y}_{n+1}\|.$$

so that

$$\mathbf{e}_s^T \mathbf{Y}^{(m,r)} - \mathbf{y}(t_{n+1}) = O((\Delta t)^{m\theta r + q + 1}) + O((\Delta t)^{m(u+1) + q + 1}) + O((\Delta t)^{p+1}),$$

where  $q$  denotes the order of the predictor. Thus, we have proved the result:

**Theorem 2.10.** Let the underlying integration method have step point order  $p$ , let the predictor have order  $q$ , and let (2.24) be satisfied. Then, the step point order  $p(m,r)$  of the iterated method is given by

$$p(m,r) = \min \{p, q + m \cdot \min \{\theta r, u + 1\}\}. \blacklozenge$$

This theorem implies that for a given number of inner iterations  $r$  the order of the corrector is obtained if the number of outer iterations satisfies

$$m \geq \frac{p - q}{\min \{\theta r, u + 1\}}.$$

Since an inner iteration is cheaper than an outer iteration, it is a good strategy to choose  $r \geq \theta^{-1}(u + 1)$ . Then, the order of the corrector is reached after  $\lceil (p - q)(u + 1)^{-1} \rceil$  outer iterations.

### 3. Summary of results

We conclude this paper with a summary of convergence and stability results in two and three dimensions. The stability results were obtained in the case of the LSV predictor. We did also compute stability regions in the case of extrapolation predictors, but the results were relatively poor. For



example, for two-dimensional problems, the two-step BDF was only A(0)-stable for all  $m\tau$ -values and only A(4<sup>o</sup>)-stable for large  $m\tau$ . Therefore, we restricted our considerations to the LSV predictor. Table 3.1 lists convergence and stability results for the two-dimensional case. Table 3.2 presents results for the three-dimensional case. As to the stability, we restricted our computations to the 3rd-order Radau method (2.6) with  $A^*$  defined by (2.12a) and the 2nd-order BDF method (2.8a) with  $A^* = A$ . For these methods, we computed the values  $(r_0, \alpha)$  for  $m\tau = 1, 2, \dots, 5$  such that the region  $\mathbb{D}(r_0, \alpha)$  defined in (2.9) is contained in the stability domain of the method.

**Table 3.1.** Convergence and stability properties for two-dimensional problems

Method	Matrix $A^*$	Convergence	Stability
3rd-order Radau (2.6)	A	A(54.7°)-convergent	A( $\alpha$ )-unstable for $mr \geq 1$ and $\alpha \geq 0$
	(2.12a)	A(90°)-convergent	A(90°)-stable for all $mr$
	(2.12b)	A(48°)-convergent	A(90°)-stable for $mr = 1$ A(80°)-stable for $mr \leq 2$ A(70°)-stable for $mr \leq 10$ A(60°)-stable for $mr \leq 30$ A(48°)-stable for all $mr$
2nd-order BDF (2.8a)	A	A(90°)-convergent	A(90°)-stable for all $mr$
3rd-order method (2.8b)	A	A(90°)-convergent	A(70°)-stable for $mr = 1$ A( $\alpha$ )-unstable for $mr > 1$ and $\alpha \geq 0$

**Table 3.2.** Stability domains for three-dimensional problems

Method	Matrix $A^*$	Convergence	Stability domain contains $\mathbb{D}(r_0, \alpha)$
3rd-order Radau (2.6)	A	A(90°)-convergent	A( $\alpha$ )-unstable for $mr > 1$ and $\alpha \geq 0$
	(2.12a)	A(45°)-convergent	$\alpha \leq 45^\circ$ , $r_0 = \{\infty, \dots, \infty\}$ $\alpha = 50^\circ$ , $r_0 = \{11.7, 13.5, 15.1, 16.4, 17.7\}$ $\alpha = 60^\circ$ , $r_0 = \{3.4, 4.5, 5.5, 6.3, 7.1\}$ $\alpha = 70^\circ$ , $r_0 = \{1.7, 2.5, 3.3, 2.3, 2.1\}$ $\alpha = 80^\circ$ , $r_0 = \{1.0, 1.7, 1.2, 1.2, 1.2\}$ $\alpha = 90^\circ$ , $r_0 = \{0.3, 0.4, 0.3, 0.4, 0.4\}$
2nd-order BDF (2.8a)	A	(2.12b) A(45°)-convergent	
		A(45°)-convergent see also Table 2.1	$\alpha \leq 45^\circ$ , $r_0 = \{\infty, \dots, \infty\}$ $\alpha = 50^\circ$ , $r_0 = \{11.6, \dots, 11.6\}$ $\alpha = 60^\circ$ , $r_0 = \{3.6, \dots, 3.6\}$ $\alpha = 70^\circ$ , $r_0 = \{2.0, \dots, 2.0\}$ $\alpha = 80^\circ$ , $r_0 = \{1.4, 1.4, 1.4, 1.1, 1.1\}$ $\alpha = 90^\circ$ , $r_0 = \{0.3, 0.3, 0.3, 0.3, 0.4\}$
3rd-order method (2.8b)	A	A(45°)-convergent see also Table 2.1	

**References**

- [1] Bendtsen, C. [1996]: Highly stable parallel Runge-Kutta methods, *Appl. Numer. Math.* 21, 1-8.
- [2] Butcher, J.C. [1987]: *The Numerical Analysis of Ordinary Differential Equations, Runge-Kutta and General Linear Methods*, Wiley.
- [3] Hirsch, C. [1988]: *Numerical Computation of Internal and External Flows, Vol. 1: Fundamentals of Numerical Discretization*, Wiley.
- [4] Houwen, P.J. van der & Sommeijer, B.P. [1991]: Iterated Runge-Kutta methods on parallel computers, *SIAM J. Sci. Stat. Comput.* 12, 1000-1028.
- [5] Houwen, P.J. van der, Sommeijer, B.P. & Kok, J. [1997]: The iterative solution of fully implicit discretizations of three-dimensional transport models, to appear in *Appl. Numer. Math.*
- [6] Houwen, P.J. van der, & Swart, J.J.B. de [1996]: Triangularly implicit iteration methods for ODE-IVP solvers, *SIAM J. Sci. Comput.*, 18, 41-55.
- [7] Hundsdorfer, W. [1996]: A note on the stability of the Douglas splitting method, Preprint NM-R9606, CWI Amsterdam, to appear in *Math. Comp.* 1998.
- [8] Orel, B. [1993]: Parallel Runge-Kutta methods with real eigenvalues, *Appl. Numer. Math.* 11, 241-250.
- [9] Sommeijer, B.P., Couzy, W. & van der Houwen, P.J. [1992]: A-stable parallel block methods for ordinary and integro-differential equations, *Appl. Numer. Math.* 9, 267-281.
- [10] Sommeijer, B.P. & Kok, J. [1997]: Domain decomposition for an implicit shallow-water transport solver. In: B.Hertzberger & P. Sloot (eds.), *Proceedings of the HPCN Europe 1997 Conference*, April 1997, Vienna, *Lect. Notes in Comp. Science* 1225, Springer, 379-388.