



Centrum voor Wiskunde en Informatica

**REPORTRAPPORT**

Approximate factorization for time-dependent partial differential equations

P.J. van der Houwen, B.P. Sommeijer

Modelling, Analysis and Simulation (MAS)

**MAS-R9915 June 16, 1999**

Report MAS-R9915  
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

# Approximate Factorization for Time-dependent Partial Differential Equations

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

## ABSTRACT

The first application of approximate factorization in the numerical solution of time-dependent partial differential equations (PDEs) can be traced back to the celebrated papers of Peaceman and Rachford and of Douglas of 1955. For linear problems, the Peaceman-Rachford-Douglas method can be derived from the Crank-Nicolson method by the approximate factorization of the system matrix in the linear system to be solved. This factorization is based on a splitting of the system matrix. In the numerical solution of time-dependent PDEs we often encounter linear systems whose system matrix has a complicated structure, but can be split into a sum of matrices with a simple structure. In such cases, it is attractive to replace the system matrix by an approximate factorization based on this splitting. This contribution surveys various possibilities for applying approximate factorization to PDEs and presents a number of new stability results for the resulting integration methods.

*1991 Mathematics Subject Classification:* 65L06, 65L20, 65M12, 65M20

*Keywords and Phrases:* numerical analysis, partial differential equations, approximate factorization, stability.

*Note:* Work carried out under project MAS 1.2 - 'Numerical Algorithms for Surface Water Quality Modelling'.

## 1. Introduction

The first application of approximate factorization in the numerical solution of time-dependent partial differential equations (PDEs) can be traced back to the celebrated papers of Peaceman and Rachford [21] and of Douglas [5] of 1955. More explicitly, approximate factorization was formulated by Beam and Warming [1] in 1976.

In order to illustrate the idea of approximate factorization, consider the initial-boundary value problem for the two-dimensional diffusion equation

$$\frac{\partial u(t,x,y)}{\partial t} = \frac{\partial^2 u(t,x,y)}{\partial x^2} + \frac{\partial^2 u(t,x,y)}{\partial y^2}$$

and let this problem be discretized in space by finite differences. Then, we obtain an initial-value problem (IVP) for a system of ordinary differential equations (ODEs)

$$(1.1) \quad \frac{d\mathbf{y}(t)}{dt} = J_1\mathbf{y} + J_2\mathbf{y},$$

where  $\mathbf{y}(t)$  contains approximations to  $u(t,x,y)$  at the grid points and  $J_1$  and  $J_2$  are matrices representing finite difference approximations to  $\partial^2/\partial x^2$  and  $\partial^2/\partial y^2$ . The system (1.1) can be integrated by e.g. the second-order trapezoidal rule, yielding the well known Crank-Nicolson method [3]

$$(1.2) \quad \left(I - \frac{1}{2}\Delta t (J_1 + J_2)\right)\mathbf{y}_{n+1} = \mathbf{y}_n + \frac{1}{2}\Delta t (J_1 + J_2)\mathbf{y}_n.$$

Here,  $I$  denotes the identity matrix,  $\Delta t$  is the timestep and  $\mathbf{y}_n$  represents a numerical approximation to  $\mathbf{y}(t_n)$ . Each step requires the solution of a linear system with system matrix  $I - \frac{1}{2}\Delta t(J_1 + J_2)$ . Due to the relatively large bandwidth, the solution of this system by a *direct* factorization of the system matrix is quite expensive. Following Beam and Warming [1], (1.2) is written in the equivalent form

$$(1.2') \quad \left(I - \frac{1}{2}\Delta t (J_1 + J_2)\right)(\mathbf{y}_{n+1} - \mathbf{y}_n) = \Delta t (J_1 + J_2)\mathbf{y}_n,$$

and the system matrix is replaced by an *approximate* factorization, to obtain

$$(1.3) \quad \left(I - \frac{1}{2}\Delta t J_1\right)\left(I - \frac{1}{2}\Delta t J_2\right)(\mathbf{y}_{n+1} - \mathbf{y}_n) = \Delta t (J_1 + J_2)\mathbf{y}_n.$$

This method is easily verified to be identical with the alternating direction implicit method (ADI method) of Peaceman-Rachford and Douglas, usually represented in the form

$$(1.3') \quad \mathbf{y}_{n+1/2} = \mathbf{y}_n + \frac{1}{2}\Delta t (J_1\mathbf{y}_{n+1/2} + J_2\mathbf{y}_n), \quad \mathbf{y}_{n+1} = \mathbf{y}_{n+1/2} + \frac{1}{2}\Delta t (J_1\mathbf{y}_{n+1/2} + J_2\mathbf{y}_{n+1}).$$

Although we now have to solve two linear systems, the small bandwidth of the matrices  $I - \frac{1}{2}\Delta t J_k$  causes that direct solution methods are not costly. Since the factorized system matrix in (1.3) is a second-order approximation to the system matrix in (1.2'), the ADI method is a third-order perturbation of (1.2'), and hence of (1.2), so that it is second-order accurate. Note that directly applying approximate factorization to the system matrix in (1.2) would yield a first-order accurate method. Hence, the intermediate step which replaces (1.2) by (1.2') is essential.

The application of approximate factorization is not restricted to schemes resulting from time discretizations by the trapezoidal rule. For example, one may replace the trapezoidal rule (1.2) by a second-order linear multistep method and proceed as described above. In fact, approximate factorization can be applied in many more cases where linear or nonlinear time-dependent PDEs are solved numerically. We mention (i) the linear multistep approach of Warming and Beam [28] described in Section 2.1, (ii) linearly implicit integration methods like Rosenbrock methods (see Section 2.2), (iii) linearization of a nonlinear method (Section 2.3), and (iv) iterative application of approximate factorization for solving linear systems (Section 3). In all these cases, we are faced with linear systems whose system matrix has the form  $I - \Delta t M$ , where the matrix  $M$  itself has a complicated structure, but can be split into a sum  $\sum M_k$  with matrices  $M_k$  possessing a simple structure. This leads us to replace  $I - \Delta t M$  by the approximate factorization  $\prod(I - \Delta t M_k)$ .

In this paper, we discuss the application of the approximate factorization technique to the four cases mentioned above and we present stability theorems for the resulting integration methods, many of which are new results. One of the results is that in the case of three-component splittings  $M = \sum M_k$ , where the  $M_k$  have purely imaginary eigenvalues, iterative approximate factorization leads to methods with substantial stability boundaries. Such methods are required in the numerical solution of 3-dimensional, convection-dominated transport problems.

## 2. Noniterative factorized methods

Consider an initial-boundary value problem for the PDE

$$(2.1) \quad \frac{\partial u(t, \mathbf{x})}{\partial t} = L(t, \mathbf{x}, u(t, \mathbf{x})),$$

where  $L$  is a differential operator in the  $d$ -dimensional space variable  $\mathbf{x} = (x_1, \dots, x_d)$ . Spatial discretization yields an IVP for a system of ODEs

$$(2.2) \quad \frac{d\mathbf{y}(t)}{dt} = \mathbf{f}(t, \mathbf{y}), \quad \mathbf{y}(t_0) = \mathbf{y}_0.$$

In order to simplify the notations, we shall assume that (2.2) is rewritten in autonomous form. Furthermore, it will be assumed that the Jacobian matrix  $J(\mathbf{y}) := \partial \mathbf{f}(\mathbf{y}) / \partial \mathbf{y}$  can be split into a sum of  $m$  matrices, i.e.  $J(\mathbf{y}) = \sum J_k$ , where the splitting is either according to the spatial dimensions (as in the early papers on splitting methods), or to the physical terms in the PDE (2.1), or according to any other partition leading to matrices  $J_k$  with a convenient structure. In this paper, we only use splittings of the Jacobian and not of the righthand side function  $\mathbf{f}(\mathbf{y})$ . This is often convenient in the case of nonlinear PDEs.

We discuss three options for applying noniterative approximate factorization techniques, viz. (i) the ADI method of Warming and Beam, (ii) approximate factorization of linearly implicit integration methods and (iii) approximate factorization in the linearization of nonlinear methods.

### 2.1. The method of Warming and Beam

Consider the linear multistep method (LM method)

$$(2.3) \quad \rho(E)\mathbf{y}_{n-\mu+1} = \Delta t \sigma(E)\mathbf{f}(\mathbf{y}_{n-\mu+1}), \quad \rho(z) := \sum_{i=0}^{\mu} a_i z^{\mu-i}, \quad \sigma(z) := \sum_{i=0}^{\mu} b_i z^{\mu-i}, \quad a_0 = 1,$$

where  $E$  is the forward shift operator and  $\mu \geq 1$ . Warming and Beam [28] rewrite (2.3) in the form

$$(2.3') \quad \rho(E) \left( \mathbf{y}_{n-\mu+1} - b_0 \Delta t \mathbf{f}(\mathbf{y}_{n-\mu+1}) \right) = \Delta t \left( \sigma(E) - b_0 \rho(E) \right) \mathbf{f}(\mathbf{y}_{n-\mu+1}).$$

Since the degree of  $\rho$  is larger than that of  $\sigma - b_0 \rho$ , the righthand side does not depend on  $\mathbf{y}_{n+1}$ . In [28] it is assumed that  $\mathbf{f}$  is linear, i.e.  $\mathbf{f}(\mathbf{y}) = J\mathbf{y}$ , so that (2.3') becomes a linear system for  $\rho(E)\mathbf{y}_{n-\mu+1}$ . However, by replacing (2.3') with

$$(2.3'') \quad \rho(E) \left( \mathbf{y}_{n-\mu+1} - b_0 \Delta t J \mathbf{y}_{n-\mu+1} \right) = \Delta t \left( \sigma(E) - b_0 \rho(E) \right) \mathbf{f}(\mathbf{y}_{n-\mu+1}),$$

we can also deal with ODE systems where  $\mathbf{f}$  is nonlinear (see [2]). Assuming that (2.3) is consistent, so that  $\rho(1) = 0$ , it can be shown that (2.3'') is an  $O((\Delta t)^3)$  perturbation of (2.3'), and hence of (2.3). The method (2.3'') is linearly implicit in the quantity  $\mathbf{q}_n := \rho(E)\mathbf{y}_{n-\mu+1}$  with system matrix  $I - b_0 \Delta t J = I - b_0 \Delta t \sum J_k$ , where  $I$  denotes the identity matrix (in the following, the identity matrix will always be denoted by  $I$  without specifying its order, which will be clear from the

context). Approximate factorization of this system matrix leads to the method of Warming and Beam:

$$(2.4) \quad \begin{aligned} \Pi \mathbf{q}_n &= \Delta t \left( \sigma(\mathbf{E}) - b_0 \rho(\mathbf{E}) \right) \mathbf{f}(\mathbf{y}_{n-\mu+1}), \quad \Pi := \prod_{k=1}^m (\mathbf{I} - b_0 \Delta t \mathbf{J}_k), \\ \mathbf{y}_{n+1} &= \mathbf{q}_n - (\rho(\mathbf{E}) - \mathbf{E}^\mu) \mathbf{y}_{n-\mu+1}. \end{aligned}$$

Since  $\mathbf{q}_n = O(\Delta t)$  it follows that (2.4) is an  $O((\Delta t)^3)$  perturbation of (2.3) which was itself an  $O((\Delta t)^3)$  perturbation of (2.3). Thus, if (2.3) is at least second-order accurate, then (2.4) is also second-order accurate. Since the LM method (2.3) cannot be A-stable if its order is higher than two and because A-stability of (2.3) will turn out to be a necessary condition for (2.4) to be A-stable (see Section 2.4), this order limitation is not restrictive.

If the PDE is linear and if (2.3) is defined by the trapezoidal rule, then (2.4) is identical with the Peaceman-Rachford method (1.3) for  $m = 2$ . Hence, (2.4) might be considered as an extension of the Peaceman-Rachford method (1.3) (or (1.3')) to nonlinear PDEs with multicomponent splittings. The computational efficiency of (2.4) depends on the structure of the successive system matrices  $\mathbf{I} - b_0 \Delta t \mathbf{J}_k$ . Let us consider the case of an  $m$ -dimensional convection-dominated problem where the convection terms are discretized by third-order upwind formulas. Using dimension splitting, the  $\mathbf{J}_k$  become block-diagonal whose blocks are penta-diagonal matrices. The LU-decomposition of  $\mathbf{I} - b_0 \Delta t \mathbf{J}_k$  and the forward/backward substitution each requires about  $8N$  flops for large  $N$ ,  $N$  denoting the dimension of  $\mathbf{J}_k$  (see e.g. [9, p. 150]). Hence, the total costs are only proportional to  $N$ , viz.  $8mN$  flops per step and an additional  $8mN$  flops if the LU-decompositions are recomputed. Moreover, there is scope for a lot of vectorization, so that on vector computers the solution of the linear systems in (2.4) is extremely fast. Furthermore, there is a lot of intrinsic parallelism, because of the block structure of  $\mathbf{J}_k$ . However, the crucial point is the magnitude of the stepsize for which the method is stable. This will be the subject of Section 2.4.

Finally, we remark that the Warming-Beam method (2.4) was originally designed as an ADI method based on dimension splitting, but it can of course be applied to any Jacobian splitting  $\mathbf{J} = \sum \mathbf{J}_k$ .

## 2.2. Factorized linearly implicit methods

In the literature, various families of linearly implicit methods have been proposed. The first methods of this type are the Rosenbrock methods, proposed in 1962 by Rosenbrock [20]. A more general family contains the linearly implicit Runge-Kutta methods developed by Strehmel and Weiner [26]. Here, we illustrate the factorization for Rosenbrock methods which are defined by (cf. [11, p. 111])

$$(2.5) \quad \begin{aligned} \mathbf{y}_{n+1} &= \mathbf{y}_n + (\mathbf{b}^T \otimes \mathbf{I}) \mathbf{K}, \quad \mathbf{K} := (\mathbf{k}_i), \\ (\mathbf{I} - \mathbf{T} \otimes \Delta t \mathbf{J}) \mathbf{K} &= \Delta t \mathbf{F}(\mathbf{e} \otimes \mathbf{y}_n + (\mathbf{L} \otimes \mathbf{I}) \mathbf{K}), \quad \mathbf{J} \approx \mathbf{J}(\mathbf{y}_n) := \frac{\partial \mathbf{f}(\mathbf{y}_n)}{\partial \mathbf{y}}, \end{aligned} \quad i = 1, \dots, s, \quad n \geq 0,$$

where  $\mathbf{b}$  and  $\mathbf{e}$  are  $s$ -dimensional vectors,  $\mathbf{e}$  has unit entries,  $\mathbf{T}$  is an  $s$ -by- $s$  diagonal or lower triangular matrix,  $\mathbf{L}$  is a strictly lower triangular  $s$ -by- $s$  matrix, and  $\otimes$  denotes the Kronecker or direct matrix product. Furthermore, for any vector  $\mathbf{V} = (\mathbf{v}_i)$ ,  $\mathbf{F}(\mathbf{V})$  is defined by  $(\mathbf{f}(\mathbf{v}_i))$ . If the order of the method (2.5) is independent of the choice of the Jacobian approximation  $\mathbf{J}$ , then (2.5) is called a Rosenbrock-W method [25]. Note that the steppoint formula in (2.5) is explicit, so that the main computational effort goes into the computation of the implicitly defined vector  $\mathbf{K}$ . Since  $\mathbf{T}$  is lower triangular and  $\mathbf{L}$  is strictly lower triangular, the  $s$  subsystems for  $\mathbf{k}_i$  can be solved successively. Moreover, although the system for  $\mathbf{K}$  is nonlinear, these subsystems are linear. Let us rewrite the system for  $\mathbf{K}$  in the equivalent form

$$(\mathbf{I} - \mathbf{D} \otimes \Delta t \mathbf{J}) \mathbf{K} = \Delta t \mathbf{F}(\mathbf{e} \otimes \mathbf{y}_n + (\mathbf{L} \otimes \mathbf{I}) \mathbf{K}) + ((\mathbf{T} - \mathbf{D}) \otimes \Delta t \mathbf{J}) \mathbf{K},$$

where  $\mathbf{D}$  is a diagonal matrix whose diagonal equals that of  $\mathbf{T}$ . Then, approximately factorizing the block-diagonal system matrix  $\mathbf{I} - \mathbf{D} \otimes \Delta t \mathbf{J} = \mathbf{I} - \mathbf{D} \otimes \Delta t \sum \mathbf{J}_k$  leads to the *factorized Rosenbrock method*

$$(2.6) \quad \mathbf{y}_{n+1} = \mathbf{y}_n + (\mathbf{b}^T \otimes \mathbf{I}) \mathbf{K},$$

$$\mathbf{\Pi} \mathbf{K} = \Delta t \mathbf{F}(\mathbf{e} \otimes \mathbf{y}_n + (\mathbf{L} \otimes \mathbf{I}) \mathbf{K}) + ((\mathbf{T} - \mathbf{D}) \otimes \Delta t \mathbf{J}) \mathbf{K}, \quad \mathbf{\Pi} := \prod_{k=1}^m (\mathbf{I} - \mathbf{D} \otimes \Delta t \mathbf{J}_k).$$

If the Rosenbrock method (2.5) is at least second-order accurate and if  $\mathbf{J} = \mathbf{J}(\mathbf{y}_n) + \mathcal{O}(\Delta t)$ , then (2.6) is also at least second-order accurate. However, as observed in [27], if (2.5) is a Rosenbrock-W method with a diagonal matrix  $\mathbf{T}$  with constant diagonal entries  $\kappa$ , then the approximate factorization does not affect the order of accuracy. This follows from the fact that for  $\mathbf{T} = \kappa \mathbf{I}$  we can write  $\mathbf{\Pi} = \mathbf{I} - \kappa \mathbf{I} \otimes \Delta t \mathbf{J}^*$ . Hence, we may consider the factorized Rosenbrock method (2.6) as the original Rosenbrock-W method with  $\mathbf{J} = \mathbf{J}^*$ . Since in Rosenbrock-W methods the Jacobian can be freely chosen, Rosenbrock-W methods and their factorized versions have the same order of accuracy.

As to the computational efficiency of factorized Rosenbrock methods, we observe that if in the underlying Rosenbrock method  $\mathbf{T} = \mathbf{D}$  and  $\mathbf{L} = \mathbf{O}$ , then the  $s$  subsystems for  $\mathbf{k}_i$  in (2.6) can be solved concurrently. These subsystems have the same structure as in the Warming-Beam method (2.4), so that the computational efficiency is comparable on a parallel computer system. As an example of such a parallel Rosenbrock method, we have

$$(2.7) \quad \mathbf{b} = \frac{1}{2(\kappa_2 - \kappa_1)} \begin{pmatrix} 2\kappa_2 - 1 \\ -2\kappa_1 + 1 \end{pmatrix}, \quad \mathbf{T} = \begin{pmatrix} \kappa_1 & 0 \\ 0 & \kappa_2 \end{pmatrix}, \quad \mathbf{L} = \mathbf{O}, \quad \kappa_1 \neq \kappa_2,$$

which is second-order accurate if  $\mathbf{J} = \mathbf{J}(\mathbf{y}_n) + \mathcal{O}(\Delta t)$ .

However, if either  $\mathbf{T} \neq \mathbf{D}$  or  $\mathbf{L} \neq \mathbf{O}$ , then the  $s$  subsystems in (2.6) have to be solved sequentially.

### 2.3. Approximate factorization of linearized methods

Instead of starting with a linearly implicit integration method, we also may linearize a nonlinear method. In fact, the Rosenbrock methods of the preceding section can be introduced by linearizing

diagonally implicit Runge-Kutta (DIRK) methods (cf. [11, p. 111]). In the literature, many other examples of linearization can be found. For instance, the linearization of the  $\theta$ -method applied to the porous media equation (in Richtmyer and Morton [22, p. 203]), the linearization of the Crank-Nicolson method for hyperbolic conservation laws (in Beam and Warming [1]) and the linearization of LM methods for the compressible Navier-Stokes equations (in Beam and Warming [2]). In this paper, we consider the linearization of a class of methods which contains most methods from the literature:

$$(2.8) \quad \begin{aligned} \mathbf{y}_{n+1} &= (\mathbf{a}^T \otimes \mathbf{I}) \mathbf{Y}_{n+1} + \mathbf{g}_n, & \mathbf{Y}_{n+1} &:= (\mathbf{y}_{n+c_i}), \\ \mathbf{Y}_{n+1} - \Delta t (\mathbf{T} \otimes \mathbf{I}) \mathbf{F}(\mathbf{Y}_{n+1}) &= \mathbf{G}_n, & \mathbf{F}(\mathbf{Y}_{n+1}) &:= (\mathbf{f}(\mathbf{y}_{n+c_i})), \end{aligned} \quad \begin{matrix} i = 1, \dots, s, \\ n \geq 0. \end{matrix}$$

Here,  $\mathbf{a}$  is an  $s$ -dimensional vector and  $\mathbf{T}$  is again an  $s$ -by- $s$  matrix. The steppoint value  $\mathbf{y}_{n+1}$  and the components  $\mathbf{y}_{n+c_i}$  of  $\mathbf{Y}_{n+1}$  represent numerical approximations to the exact solution values  $\mathbf{y}(t_n + \Delta t)$  and  $\mathbf{y}(t_n + c_i \Delta t)$ , where the  $c_i$  are given abscissae.  $\mathbf{Y}_{n+1}$  is called the *stage vector*, its components  $\mathbf{y}_{n+c_i}$  the *stage values*.  $\mathbf{G}_n$  and  $\mathbf{g}_n$  are assumed to be defined by preceding steppoint values  $\mathbf{y}_n, \mathbf{y}_{n-1}, \dots$  and by the preceding stage vectors  $\mathbf{Y}_n, \mathbf{Y}_{n-1}, \dots$  and their derivatives. Again, the steppoint formula is explicit, so that the main computational effort goes into the solution of the stage vector  $\mathbf{Y}_{n+1}$ .

Let us linearize the stage vector equation in (2.8) to obtain for  $\mathbf{Y}_{n+1}$  the linear system

$$(2.9) \quad \mathbf{Y}_{n+1} - \Delta t (\mathbf{T} \otimes \mathbf{I}) (\mathbf{F}(\mathbf{Y}^0) + (\mathbf{I} \otimes \mathbf{J})(\mathbf{Y}_{n+1} - \mathbf{Y}^0)) = \mathbf{G}_n, \quad \mathbf{J} \approx \mathbf{J}(\mathbf{y}_n) := \frac{\partial \mathbf{f}(\mathbf{y}_n)}{\partial \mathbf{y}}.$$

Here,  $\mathbf{Y}^0$  is an approximation to  $\mathbf{Y}_{n+1}$ , for example,  $\mathbf{Y}^0 = \mathbf{Y}_n$  or  $\mathbf{Y}^0 = \mathbf{e} \otimes \mathbf{y}_n$ . However, with this simple choice, the order of the linearized method is not necessarily the same as the original method (2.8). For instance, if (2.8) has order  $p \geq 2$  and if  $\mathbf{J} = \mathbf{J}(\mathbf{y}_n) + \mathcal{O}(\Delta t)$ , then the order of the linearized method is in general not higher than two. If (2.8) has order  $p \geq 3$ , then higher-order formulas for  $\mathbf{Y}^0$  should be used. Of course, if the ODE system (2.2) is already linear, i.e.  $\mathbf{y}' = \mathbf{J}\mathbf{y}$ , then  $\mathbf{Y}^0$  does not play a role, because (2.9) is identical with the stage vector equation in (2.8) for all  $\mathbf{Y}^0$ . Note that this also implies that the *linear* stability properties of (2.8) and its linearization are identical for all  $\mathbf{Y}^0$ .

It turns out that approximate factorization of linear systems of the type (2.9) is most effective if  $\mathbf{T}$  is either diagonal or (lower) triangular as in the case of the Rosenbrock method (2.5). Therefore, from now on, we impose this condition on  $\mathbf{T}$ . Furthermore, instead of directly applying approximate factorization to the linear system (2.9), we first rewrite it into the equivalent form (compare (1.2'))

$$(2.9') \quad (\mathbf{I} - \Delta t \mathbf{D} \otimes \mathbf{J})(\mathbf{Y}_{n+1} - \mathbf{Y}^0) = \mathbf{G}_n - \mathbf{Y}^0 + \Delta t (\mathbf{T} \otimes \mathbf{I}) \mathbf{F}(\mathbf{Y}^0) + \Delta t ((\mathbf{T} - \mathbf{D}) \otimes \mathbf{J})(\mathbf{Y}_{n+1} - \mathbf{Y}^0),$$

where again  $\mathbf{D} = \text{diag}(\mathbf{T})$ . Proceeding as in the preceding section leads to the factorized method

$$(2.10) \quad \begin{aligned} \mathbf{y}_{n+1} &= (\mathbf{a}^T \otimes \mathbf{I}) \mathbf{Y}_{n+1} + \mathbf{g}_n, \\ \Pi(\mathbf{Y}_{n+1} - \mathbf{Y}^0) &= \mathbf{G}_n - \mathbf{Y}^0 + \Delta t (\mathbf{T} \otimes \mathbf{I}) \mathbf{F}(\mathbf{Y}^0) + \Delta t ((\mathbf{T} - \mathbf{D}) \otimes \mathbf{J})(\mathbf{Y}_{n+1} - \mathbf{Y}^0) \end{aligned}$$



with  $\Pi$  defined as in (2.6). If  $\mathbf{Y}_{n+1} - \mathbf{Y}^0 = \mathcal{O}(\Delta t)$ , then (2.10) presents a third-order perturbation of (2.9). Hence, by setting  $\mathbf{Y}^0 = \mathbf{Y}_n$  or  $\mathbf{Y}^0 = \mathbf{e} \otimes \mathbf{y}_n$ , the resulting method is second-order accurate provided that (2.8) is also (at least) second-order accurate. We shall refer to the approximately factorized, linearized method (2.10) as the *AFL method*.

If  $\mathbf{T}$  is diagonal, then the subsystems for the components of  $\mathbf{Y}_{n+1} - \mathbf{Y}^0$  can be solved concurrently, and if  $\mathbf{T}$  is lower triangular, then these subsystems should be solved successively (note that  $\mathbf{T} - \mathbf{D}$  is strictly lower triangular). The computational efficiency of solving the linear systems in (2.10) is comparable with that of (2.6).

## 2.4. Stability

As already remarked, the crucial point is the stability of the factorized methods. We shall discuss stability with respect to the model problem  $\mathbf{y}' = \mathbf{J}\mathbf{y} = \sum \mathbf{J}_k \mathbf{y}$ , where the matrices  $\mathbf{J}_k$  commute. Application of the factorized methods to this model problem leads to linear recursions. The roots  $\zeta$  of the corresponding characteristic equations define the *amplification factors* of the method. These amplification factors are functions of the vector  $\mathbf{z} = (z_1, \dots, z_m)^T$ , where  $z_k$  runs through the eigenvalues of  $\Delta t \mathbf{J}_k$ . We call a method *stable* at the point  $\mathbf{z}$  if its amplification factor  $\zeta(\mathbf{z})$  is on the unit disk. Likewise, we shall call a function  $R(\mathbf{z})$  *stable* at  $\mathbf{z}$  if  $R(\mathbf{z})$  is on the unit disk. In the stability definitions and stability theorems given below, we shall use the notation

$$\begin{aligned} \mathbb{W}(\alpha) &:= \{w \in \mathbb{C} : |\arg(-w)| \leq \alpha\}, \\ \mathbb{R}(\beta) &:= (-\beta, 0], \\ \mathbb{I}(\beta) &:= \{w \in \mathbb{C} : \arg(w) = \pm \frac{\pi}{2}, |w| < \beta\}. \end{aligned}$$

**Definition 2.1.** A method or a function is called

- A*( $\alpha$ )-*stable* if it is stable for  $z_k \in \mathbb{W}(\alpha)$ ,  $k = 1, \dots, m$ ,
- A*-*stable* if it is stable for  $z_k \in \mathbb{W}(\pi/2)$ ,  $k = 1, \dots, m$ ,
- A<sub>r</sub>*( $\alpha$ )-*stable* if it is stable for  $z_1, \dots, z_r \in \mathbb{R}(\infty) \wedge z_{r+1}, \dots, z_m \in \mathbb{W}(\alpha)$ .  $\blacklozenge$

The first two definitions of stability are in analogy with the definitions in numerical ODE theory. The third type of stability was introduced by Hundsdorfer [17] and will be referred to as *A<sub>r</sub>*( $\alpha$ )-stability. This type of stability is relevant in the case of convection-diffusion-reaction equations. For example, for *systems* of two-dimensional convection-diffusion-reaction equations in which the Jacobian of the reaction terms has real, stiff eigenvalues, we would like *A<sub>1</sub>*( $\pi/2$ )-stability for  $m = 3$ , that is, stability in the region  $\mathbb{R}(\infty) \times \mathbb{W}(\pi/2) \times \mathbb{W}(\pi/2)$ . Then, by choosing the splitting such that  $\mathbf{J}_1$  corresponds with the reaction terms, and  $\mathbf{J}_2$  and  $\mathbf{J}_3$  with the convection-diffusion terms in the two spatial directions, we achieve unconditional stability. We remark that in the case of a *single* two-dimensional convection-diffusion-reaction equation, we need only *A*-stability for  $m = 2$ , because we can choose the splitting such that  $\mathbf{J}_1$  corresponds with the reaction term and the convection-diffusion in one spatial direction, and  $\mathbf{J}_2$  with convection-diffusion in the other spatial direction. Note that in this splitting, the matrices  $\mathbf{J}_1$  and  $\mathbf{J}_2$  both have a band structure with small band width.

In the following, we shall often encounter stability regions containing subregions of the form  $\mathbb{S}_1 \times \mathbb{S}_2 \times \mathbb{S}_3$ . In the case of approximate factorizations that are *symmetric* with respect to the Jacobians  $J_1, J_2$  and  $J_3$ , as in the methods (2.4), (2.6) and (2.10), this means that the stability region also contains the subregions  $\mathbb{S}_1 \times \mathbb{S}_3 \times \mathbb{S}_2, \mathbb{S}_2 \times \mathbb{S}_1 \times \mathbb{S}_3, \dots$ .

In the next sections, we give stability theorems for the method of Warming and Beam, and a few AFL and factorized Rosenbrock methods.

## 2.5. Method of Warming and Beam

Applying the Warming-Beam method (2.4) to the stability test problem yields the following characteristic equation for the amplification factor  $\zeta = \zeta(\mathbf{z})$  of the method:

$$(2.11) \quad \rho(\zeta) - \psi(\mathbf{z})\sigma(\zeta) = 0, \quad \psi(\mathbf{z}) := \mathbf{e}^T \mathbf{z} \left[ b_0 \mathbf{e}^T \mathbf{z} + \prod_{k=1}^m (1 - b_0 z_k) \right]^{-1}.$$

Stability properties of (2.4) can be derived by using the following lemma of Hundsdorfer [17]:

**Lemma 2.1.** Let  $H_m$  be the function defined by

$$H_m(\mathbf{w}) := 1 + \mathbf{e}^T \mathbf{w} \prod_{k=1}^m \left( 1 - \frac{1}{2} w_k \right)^{-1}, \quad \mathbf{w} = (w_1, \dots, w_m)^T, \quad m \geq 2.$$

$H_m$  is  $A(\alpha)$ -stable if and only if  $\alpha \leq \frac{1}{2} \pi(m-1)^{-1}$  and  $A_r(\alpha)$ -stable if and only if  $\alpha \leq \frac{1}{2} \pi(m-r)^{-1}$  ♦

**Theorem 2.1.** Let the LM method (2.3) be  $A$ -stable. Then the Warming-Beam method (2.4) is:

- (a)  $A(\alpha)$ -stable for  $m \geq 2$  if and only if  $\alpha \leq \frac{1}{2} \pi(m-1)^{-1}$ .
- (b)  $A_r(\alpha)$ -stable for  $m \geq 2$  and  $r \geq 1$  if and only if  $\alpha \leq \frac{1}{2} \pi(m-r)^{-1}$ .
- (c) Stable in the region  $\mathbb{I}(\beta_1) \times \mathbb{I}(\beta_1) \times \mathbb{R}(\beta_2)$  for  $m = 3$  if  $b_0^2 \beta_1^2 (b_0 \beta_2 - 3) = 1$ .

**Proof.** If the method (2.3) is  $A$ -stable, then (2.4) is stable at the point  $\mathbf{z}$  if  $\text{Re}(\psi(\mathbf{z})) \leq 0$ , or equivalently, if  $|(1 + c\psi(\mathbf{z}))(1 - c\psi(\mathbf{z}))^{-1}| \leq 1$  for some positive constant  $c$ . Let us choose  $c = b_0$  (the  $A$ -stability of (2.3) implies that  $b_0 > 0$ ). Then, it follows from (2.11) that

$$(2.12) \quad \frac{1 + b_0 \psi(\mathbf{z})}{1 - b_0 \psi(\mathbf{z})} = H_m(2b_0 \mathbf{z}),$$

where  $H_m$  is defined in Lemma 2.1. Applying this lemma with  $\mathbf{w} = 2b_0 \mathbf{z}$  proves part (a) and (b). Part (c) is proved by analysing the inequality  $|H(2b_0 \mathbf{z})| \leq 1$  for  $\mathbf{z} = (iy_1, iy_2, x_3)$ . For  $m = 3$  this leads to

$$\begin{aligned} & ((1 - b_0^2 y_1 y_2)(1 - b_0 x_3) + 2b_0 x_3)^2 + (b_0(y_1 + y_2)(1 + b_0 x_3))^2 \\ & \leq (1 + b_0^2 y_1^2)(1 + b_0^2 y_2^2)(1 - b_0 x_3)^2. \end{aligned}$$

The most critical situation is obtained if  $y_1$  and  $y_2$  assume their maximal value. Setting  $y_1 = y_2 = \beta_1$  and taking into account that  $x_3 \leq 0$ , the inequality reduces to  $x_3 \geq -(1 + 3b_0^2\beta_1^2) / b_0^3 \beta_1^2$  from which assertion (c) is immediate.  $\blacklozenge$

This theorem implies A-stability for  $m = 2$ , a result already obtained by Warming and Beam [28]. Furthermore, the theorem implies A(0)-stability for all  $m \geq 2$ , and A( $\alpha$ )-stability with  $\alpha \leq \pi/4$  for  $m \geq 3$ . Hence, we do not have unconditional stability in the case where all Jacobians  $J_k$  have eigenvalues close to the imaginary axis. We even do not have stability in regions of the form  $\mathbb{I}(\beta) \times \mathbb{I}(\beta) \times \mathbb{I}(\beta)$  or  $\mathbb{I}(\beta) \times \mathbb{I}(\beta) \times \mathbb{R}(\infty)$  with  $\beta > 0$  (see also [12]). However, part (c) of the theorem implies for  $m = 3$  stability in  $\mathbb{W}(\pi/2) \times \mathbb{W}(\pi/2) \times \mathbb{R}(3/b_0)$ . Such regions are suitable for *systems* of two-dimensional convection-diffusion-reaction equations with real, *nonstiff* eigenvalues in the reaction part. Note that it is advantageous to have a small  $b_0$ -value, whereas L-stability of the underlying LM method does not lead to better stability properties.

**Remark 2.1.** The amplification factors of the stabilizing corrections method of Douglas (cf. [6], [7]) are given by  $\zeta = H(\mathbf{z})$ , so that it has similar stability properties as the Warming-Beam method  $\blacklozenge$

**Remark 2.2.** As already remarked in Section 2.1, (2.4) can be seen as a generalization of the Peaceman-Rachford method (1.3) to nonlinear PDEs with multicomponent splittings. In the literature, a second, direct generalization of (1.3') is known, however its stability is less satisfactory. For the definition of this generalization, let  $\mathbf{F}$  be a splitting function with  $m$  arguments satisfying the relation  $\mathbf{F}(\mathbf{y}, \dots, \mathbf{y}) = \mathbf{f}(\mathbf{y})$  and define  $\mathbf{F}_k$  by setting the  $k$ th argument of  $\mathbf{F}$  equal to  $\mathbf{y}^k$  and all other arguments equal to  $\mathbf{y}^{k-1}$ . Then, the direct generalization of (1.3') reads

$$\mathbf{y}^0 = \mathbf{y}_n, \quad \mathbf{y}^k = \mathbf{y}^{k-1} + \frac{\Delta t}{m} \mathbf{F}_k, \quad \mathbf{y}_{n+1} = \mathbf{y}^m, \quad k = 1, \dots, m$$

(cf. e.g. [19, p. 278] and [16]). This scheme is second-order accurate for all  $\mathbf{F}$ . Evidently, it reduces to (1.3') for linear problems and  $m = 2$ . Its amplification factor is given by

$$\zeta(\mathbf{z}) = \prod_{k=1}^m \frac{m + \mathbf{e}^T \mathbf{z} - z_k}{m - z_k},$$

showing that unlike the Warming-Beam method, it is not even A(0)-stable for  $m \geq 3$  (e.g.  $\zeta(\mathbf{e}z_0) \approx (1 - m)^m$  as  $z_0 \rightarrow \infty$ ), so that it is only of use for  $m = 2$ .  $\blacklozenge$

## 2.6. AFL-LM methods

We start with AFL methods based on the class of LM methods (2.3). Writing (2.3) in the form (2.8) and applying the AFL method (2.10) to the stability test problem leads to the characteristic equation

$$(2.13) \quad \rho(\zeta) - \mathbf{e}^T \mathbf{z} \sigma(\zeta) = \left[ 1 - b_0 \mathbf{e}^T \mathbf{z} - \prod_{k=1}^m (1 - b_0 z_k) \right] (\zeta - 1) \zeta^{\mu-1}.$$

This equation does not allow such a general stability analysis as in the case (2.11). Therefore, we confine our considerations to two particular cases, viz. the AFL methods based on the trapezoidal rule and the BDF method.

**2.6.1. The trapezoidal rule.** The trapezoidal rule is defined by  $\rho(\zeta) = \zeta - 1$ ,  $\sigma(\zeta) = \frac{1}{2}(\zeta + 1)$ . This leads to the characteristic equation  $\zeta = H(\mathbf{z})$ , where  $H$  is defined in Lemma 2.1. Hence, according to the proof of Theorem 2.1 the AFL-trapezoidal rule and the Warming-Beam method with  $b_0 = 1/2$  possess the same stability region, so that Theorem 2.1 applies (with  $b_0 = 1/2$  in part (c)).

**2.6.2. The BDF.** For the BDF with  $\rho(\zeta) = \zeta^2 - \frac{4}{3}\zeta + \frac{1}{3}$ ,  $\sigma(\zeta) = \frac{2}{3}\zeta^2$  the characteristic equation (2.13) assumes the form

$$(2.14) \quad \zeta^2 - C_1\zeta + C_2 = 0,$$

$$C_1 = \frac{P(\mathbf{z})}{Q(\mathbf{z})}, \quad C_2 = \frac{1}{Q(\mathbf{z})}, \quad P(\mathbf{z}) := Q(\mathbf{z}) + 1 + 2\mathbf{e}^T\mathbf{z}, \quad Q(\mathbf{z}) := 3 \prod_{k=1}^m \left(1 - \frac{2}{3}z_k\right).$$

In order to find the stability region, we use Schur's criterion stating that the amplification factors are on the unit disk if  $|C_2|^2 + |C_1 - C_1^*C_2| \leq 1$ .

**Theorem 2.2.** The AFL-BDF method is A-stable for  $m = 2$  and  $A(\pi/4)$ -stable for  $m = 3$ .

**Proof.** Let  $P^*$  and  $Q^*$  denote the complex conjugates of  $P$  and  $Q$ . Then, in terms of  $P$  and  $Q$ , the Schur criterion requires that the polynomial  $E(\mathbf{z}) := (|Q(\mathbf{z})|^2 - 1)^2 - |PQ^* - P^*|^2$  is nonnegative in the stability region. Writing  $z_k = iy_k$  with  $y_k$  real, we straightforwardly find for  $m = 2$  that

$$E(iy_1, iy_2) = \frac{16}{9} (y_1 + y_2)^2 (9y_1^2 + 9y_2^2 + 4y_1^2y_2^2 + 6y_1y_2).$$

It is easily seen that  $E(iy_1, iy_2) \geq 0$  for all  $y_1$  and  $y_2$ , proving the A-stability for  $m = 2$ .

For  $m = 3$  we set  $z_k = x_k - ix_k$ ,  $k = 1, 2, 3$ , with  $x_k \leq 0$  and derived an expression for  $E(\mathbf{z})$  with the help of Maple. This expression has the form  $-\mathbf{e}^T\mathbf{x} s(\mathbf{x})$ , where  $\mathbf{x} = (x_1, x_2, x_3)^T$  and  $s(\mathbf{x})$  consists of a sum of terms each term being of the form  $x_1^p x_2^q x_3^r$ , where  $p, q$  and  $r$  are nonnegative integers. We verified that the coefficients of these terms are all positive if  $p+q+r$  is even and negative otherwise (the length of the formulas prevents us from presenting  $s(\mathbf{x})$  here). Hence,  $E(\mathbf{z}) \geq 0$  for all  $z_k = x_k - ix_k$  with  $x_k \leq 0$ . Likewise, it can be shown that  $E(\mathbf{z}) \geq 0$  for all  $z_k = x_k + ix_k$  with  $x_k \leq 0$ , proving the  $A(\pi/4)$ -stability for  $m = 3$ . ♦

In addition, we determined stability regions of the form  $\mathbb{I}(\beta_1) \times \mathbb{I}(\beta_1) \times \mathbb{R}(\beta_2)$  by analysing the stability boundary curve  $E(iy_1, iy_2, x_3) = 0$  with the help of Maple. In particular, we found that in the region  $\mathbb{W}(\pi/2) \times \mathbb{W}(\pi/2) \times \mathbb{R}(\beta)$  the value of  $\beta$  is determined by the equation  $E(i\infty, i\infty, \beta) \equiv 0$  and in the region  $\mathbb{I}(\beta) \times \mathbb{I}(\beta) \times \mathbb{R}(\infty)$  by the equation  $E(i\beta, i\beta, \infty) = 0$ . This leads to  $\beta = \frac{9 + 3\sqrt{17}}{4}$  and  $\beta = \frac{3}{4}\sqrt{2}$ , respectively. For the sake of easy comparison, we have listed a number of stability

results derived in this paper in Table 5.1. This table shows that the AFL-BDF regions  $\mathbb{I}(\beta) \times \mathbb{I}(\beta) \times \mathbb{R}(\infty)$  and  $\mathbb{W}(\pi/2) \times \mathbb{W}(\pi/2) \times \mathbb{R}(\beta)$  are larger than the corresponding stability regions of the Warming-Beam method generated by the BDF ( $b_0 = 2/3$ ). We now even have stability in  $\mathbb{I}(\beta) \times \mathbb{I}(\beta) \times \mathbb{I}(\beta)$  with nonzero imaginary stability boundary  $\beta$ , but these boundaries are quite small ( $\beta < 1/10$ ).

## 2.7. AFL-DIRK methods

If we define in (2.8)  $\mathbf{g}_n = (1 - \mathbf{a}^T \mathbf{e}) \mathbf{y}_n$  and  $\mathbf{G}_n = \mathbf{e} \otimes \mathbf{y}_n$ , then (2.8) becomes a diagonally implicit Runge-Kutta (DIRK) method. We define an AFL-DIRK method by approximating  $\mathbf{Y}_{n+1}$  by means of (2.10) with  $\mathbf{Y}^0 = \mathbf{e} \otimes \mathbf{y}_n$ . The amplification factor with respect to the stability test model becomes

$$(2.15) \quad \zeta(\mathbf{z}) = 1 + \mathbf{e}^T \mathbf{z} \mathbf{a}^T \left( \prod_{k=1}^m (\mathbf{I} - z_k \mathbf{D}) - \mathbf{e}^T \mathbf{z} (\mathbf{T} - \mathbf{D}) \right)^{-1} \mathbf{T} \mathbf{e}.$$

Let us consider the second-order, L-stable DIRK methods

$$(2.16a) \quad \mathbf{y}_{n+1} = (\mathbf{e}_2^T \otimes \mathbf{I}) \mathbf{Y}_{n+1}, \quad \mathbf{Y}_{n+1} = \begin{pmatrix} \mathbf{y}_{n+\kappa} \\ \mathbf{y}_{n+1} \end{pmatrix},$$

$$\mathbf{Y}_{n+1} - \Delta t (\mathbf{T} \otimes \mathbf{I}) \mathbf{F}(\mathbf{Y}_{n+1}) = \mathbf{e} \otimes \mathbf{y}_n, \quad \mathbf{T} = \begin{pmatrix} \kappa & 0 \\ 1-\kappa & \kappa \end{pmatrix},$$

and

$$(2.16b) \quad \mathbf{y}_{n+1} = (1 - \mathbf{a}^T \mathbf{e}) \mathbf{y}_n + \mathbf{a}^T \mathbf{Y}_{n+1}, \quad \mathbf{Y}_{n+1} = \begin{pmatrix} \mathbf{y}_{n+\kappa} \\ \mathbf{y}_{n+1-\kappa} \end{pmatrix}, \quad \mathbf{a} = \frac{1}{2\kappa^2} \begin{pmatrix} 3\kappa-1 \\ \kappa \end{pmatrix},$$

$$\mathbf{Y}_{n+1} - \Delta t (\mathbf{T} \otimes \mathbf{I}) \mathbf{F}(\mathbf{Y}_{n+1}) = \mathbf{e} \otimes \mathbf{y}_n, \quad \mathbf{T} = \begin{pmatrix} \kappa & 0 \\ 1-2\kappa & \kappa \end{pmatrix},$$

where  $\mathbf{e}_2 = (0, 1)^T$  and  $\kappa = 1 \pm \frac{1}{2}\sqrt{2}$ . The amplification factor (2.15) becomes in both cases

$$(2.17) \quad \zeta(\mathbf{z}) = 1 + \frac{\mathbf{e}^T \mathbf{z}}{\pi(\mathbf{z})} + \frac{\kappa(1-\kappa)(\mathbf{e}^T \mathbf{z})^2}{\pi^2(\mathbf{z})}, \quad \pi(\mathbf{z}) := \prod_{k=1}^m (1 - \kappa z_k), \quad \kappa = 1 \pm \frac{1}{2}\sqrt{2}.$$

**Theorem 2.3.** The AFL versions of (2.16) are A-stable for  $m = 2$  and  $A(\pi/4)$ -stable for  $m = 3$ .

**Proof.** Writing  $\zeta(\mathbf{z}) = \mathbf{P}(\mathbf{z})\mathbf{Q}^{-1}(\mathbf{z})$ , where  $\mathbf{P}$  and  $\mathbf{Q}$  are polynomials in  $z_1, z_2$  and  $z_3$ , it follows that we have A-stability if the E-polynomial  $\mathbf{E}(\mathbf{z}) := |\mathbf{Q}(\mathbf{z})|^2 - |\mathbf{P}(\mathbf{z})|^2$  is nonnegative for all purely imaginary  $z_k$ . Using Maple, we found for  $\kappa = 1 \pm \frac{1}{2}\sqrt{2}$

$$\mathbf{E}(iy_1, iy_2, 0) = \frac{1}{4} (17 \pm 12\sqrt{2}) (y_1 + y_2)^4,$$

which proves the A-stability for  $m = 2$ . Similarly, the  $A(\pi/4)$  stability can be shown for  $m = 3$ . ♦

Thus, the A-stability and  $A(\alpha)$ -stability properties of the Warming-Beam, AFL-trapezoidal, AFL-BDF, and the above AFL-DIRK methods are comparable for  $m \leq 3$ . However, for the AFL-DIRK methods we found (numerically) the stability regions  $\mathbb{I}(\beta_1) \times \mathbb{I}(\beta_1) \times \mathbb{R}(\infty)$  and  $\mathbb{W}(\pi/2) \times \mathbb{W}(\pi/2) \times \mathbb{R}(\beta_2)$  with  $\beta_1 \approx 1.26$  and  $\beta_2 \approx 10.2$  for  $\kappa = 1 - \frac{1}{2}\sqrt{2}$  and with  $\beta_1 \approx 0.28$  and

$\beta_2 \approx 1.75$  for  $\kappa = 1 + \frac{1}{2}\sqrt{2}$ . Hence, choosing  $\kappa = 1 - \frac{1}{2}\sqrt{2}$  we have larger stability regions than the corresponding stability regions of the other methods (see Table 5.1). We also have stability in regions of the type  $\mathbb{I}(\beta) \times \mathbb{I}(\beta) \times \mathbb{I}(\beta)$  with  $\beta > 0$ , but  $\beta$  is uselessly small.

## 2.8. Factorized Rosenbrock methods

Finally, we consider the factorized Rosenbrock method (2.6). With respect to the stability test model its amplification factor is given by

$$(2.18) \quad \zeta(\mathbf{z}) = 1 + \mathbf{e}^T \mathbf{z} \mathbf{b}^T \left( \prod_{k=1}^m (\mathbf{I} - z_k \mathbf{D}) - \mathbf{e}^T \mathbf{z} (\mathbf{L} + \mathbf{T} - \mathbf{D}) \right)^{-1} \mathbf{e}.$$

We consider the original, second-order, L-stable Rosenbrock method [20] defined by (2.5) with

$$(2.19a) \quad \mathbf{b} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \quad \mathbf{T} = \begin{pmatrix} \kappa & 0 \\ 0 & \kappa \end{pmatrix}, \quad \mathbf{L} = \frac{1}{2} \begin{pmatrix} 0 & 0 \\ 1 - 2\kappa & 0 \end{pmatrix}, \quad \kappa = 1 \pm \frac{1}{2}\sqrt{2},$$

and the second-order, L-stable Rosenbrock-W method (see Dekker and Verwer [4, p. 233]) with

$$(2.19b) \quad \mathbf{b} = \frac{1}{2} \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \mathbf{T} = \begin{pmatrix} \kappa & 0 \\ -2\kappa & \kappa \end{pmatrix}, \quad \mathbf{L} = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}, \quad \kappa = 1 \pm \frac{1}{2}\sqrt{2}.$$

The amplification factor (2.18) is for both methods (2.19) identical with the amplification factor (2.17) of the DIRK methods (2.16), so that all results of the preceding section apply to (2.19). The factorization of the Rosenbrock-W method (2.19b) has successfully been used by Sandu [23] and Verwer et al. [27] for the solution of large scale air pollution problems.

## 3. Factorized iteration

Except for the factorized Rosenbrock-W methods, the factorized methods discussed in the preceding section are at most second-order accurate. As already observed by Beam and Warming [2], a simple way to arrive at higher-order methods that are still computationally efficient, is factorized iteration of higher-order integration methods. Evidently, if the iteration method converges, then we retain the order of accuracy of the underlying integration method (to be referred to as the corrector). Likewise, if the convergence conditions are satisfied, then the stability properties of the iterated method are the same as those of the corrector. Hence, the stability region of the iterated method is the intersection of the convergence region of the iteration method and the stability region of the corrector. Thus, if we restrict our considerations to A-stable, preferably L-stable correctors, then the stability region of the iterated method is the same as the convergence region of the iteration method.

Perhaps even more important than the possibility of constructing higher-order methods is the increased robustness of the iterative approach. The reason is that the stability problem for the noniterative approach is replaced by a convergence problem for the iterative approach. However,

unlike stability, which concerns accumulation of perturbations through a large number of integration steps, convergence can be controlled in each single step.

In Section 3.1, we discuss (i) *AFN iteration*, that is, approximately factorized Newton iteration of the nonlinear stage vector equation in (2.8), and (ii) *AF iteration*, that is, approximately factorized iteration of the linearized stage vector equation (2.9). AFN and AF iteration enables us to achieve stability in regions of the form  $\mathbb{I}(\beta) \times \mathbb{I}(\beta) \times \mathbb{W}(\pi/2)$ .

The AFN and AF methods treat all terms in the ODE system implicitly. In the case where the ODE system contains terms that are nonstiff or mildly stiff with respect to the other terms, it may be advantageous to treat these terms explicitly. This will be illustrated in Section 3.2.

Finally, in Section 3.3 we show how A-stability for three-component Jacobian splittings can be obtained, albeit at the cost of an increase of the computational complexity.

### 3.1. The AFN and AF iteration methods

Applying Newton iteration to the stage vector equation in (2.8) yields the linear Newton systems

$$(3.1) \quad (\mathbf{I} - \Delta t \mathbf{T} \otimes \mathbf{J})(\mathbf{Y}^j - \mathbf{Y}^{j-1}) = \mathbf{G}_n - \mathbf{Y}^{j-1} + \Delta t(\mathbf{T} \otimes \mathbf{I})\mathbf{F}(\mathbf{Y}^{j-1}), \quad j \geq 1.$$

Next we apply approximate factorization to obtain the *AFN iteration method*

$$(3.2) \quad \Pi(\mathbf{Y}^j - \mathbf{Y}^{j-1}) = \mathbf{G}_n - \mathbf{Y}^{j-1} + \Delta t(\mathbf{T} \otimes \mathbf{I})\mathbf{F}(\mathbf{Y}^{j-1}) + \theta \Delta t ((\mathbf{T} - \mathbf{D}) \otimes \mathbf{J})(\mathbf{Y}^j - \mathbf{Y}^{j-1}), \quad j \geq 1,$$

where  $\Pi$  is defined as before,  $\mathbf{Y}^0$  is a suitable initial approximation to  $\mathbf{Y}_{n+1}$ , and where  $\theta$  is a free parameter to be explained later. Note that after one iteration the AFN process is identical with (2.10) if we set  $\theta = 1$  and if (2.10) and (3.2) use the same approximation  $\mathbf{Y}^0$ .

In the case of the linear system (2.9), we apply the *AF iteration method*

$$(3.3) \quad \Pi(\mathbf{Y}^j - \mathbf{Y}^{j-1}) = \mathbf{G}_n - \mathbf{Y}^{j-1} + \Delta t(\mathbf{T} \otimes \mathbf{I})\mathbf{F}(\mathbf{Y}^0) + \Delta t (\mathbf{T} \otimes \mathbf{J})(\mathbf{Y}^{j-1} - \mathbf{Y}^0) \\ + \theta \Delta t ((\mathbf{T} - \mathbf{D}) \otimes \mathbf{J})(\mathbf{Y}^j - \mathbf{Y}^{j-1}), \quad j \geq 1,$$

which is of course just the linearization of (3.2). The AFN and AF processes are consistent for all  $\theta$ , that is, if the iterates  $\mathbf{Y}^j$  converge, then they converge to the solutions  $\mathbf{Y}_{n+1}$  of (2.8) and (2.9), respectively. Since the formulas (3.2) and (3.3) have the same structure as the AFL method (2.10), we conclude that, given the LU-decompositions of the factor matrices in  $\Pi$ , the costs of performing one iteration are comparable with those of applying (2.10). Hence, the efficiency of the AFN and AF processes is largely determined by the number of iterations needed to more or less solve the implicit system. The large scale 3D shallow water transport experiments reported in [15], [24] and [12] indicate that two or three iterations suffice. Also note that for  $\theta = 0$  the subsystems in the resulting iteration processes can be solved in parallel, even if  $\mathbf{T}$  is a triangular matrix.

AFN iteration can also be applied for solving *simultaneously* the subsystems for the components  $\mathbf{k}_i$  of  $\mathbf{K}$  from the Rosenbrock method (2.5). Similarly, AF iteration can be applied *successively* to these (linear) subsystems. Here, we shall concentrate on the iteration of (2.8) and (2.9). For details on the AFN and AF iteration of Rosenbrock methods we refer to [13].

**3.1.1. The iteration error.** Let us consider the recursions for the error  $\epsilon^j := \mathbf{Y}^j - \mathbf{Y}_{n+1}$ . From (2.8) and (3.2) it follows that the AFN error satisfies the *nonlinear* recursion

$$(3.4) \quad \begin{aligned} \epsilon^j &= \mathbf{Z} \epsilon^{j-1} + \Delta t \Phi(\epsilon^{j-1}), \quad j \geq 1, \\ \mathbf{Z} &:= \mathbf{I} - (\Pi - \theta \Delta t ((\mathbf{T} - \mathbf{D}) \otimes \mathbf{J}))^{-1} (\mathbf{I} - \Delta t \mathbf{T} \otimes \mathbf{J}), \\ \Phi(\epsilon) &:= (\Pi - \theta \Delta t ((\mathbf{T} - \mathbf{D}) \otimes \mathbf{J}))^{-1} (\mathbf{T} \otimes \mathbf{I}) (\mathbf{F}(\mathbf{Y}_{n+1} + \epsilon) - \mathbf{F}(\mathbf{Y}_{n+1}) - (\mathbf{I} \otimes \mathbf{J}) \epsilon). \end{aligned}$$

Similarly, we deduce from (2.9) and (3.3) for the AF error the *linear* recursion

$$(3.5) \quad \epsilon^j = \mathbf{Z} \epsilon^{j-1}, \quad j \geq 1.$$

It is difficult to decide which of the two iteration processes has a better rate of convergence. However, in a first approximation, the rates of convergence are comparable, because in the neighbourhood of the origin the Lipschitz constant of the function  $\Phi$  is quite small, provided that  $\mathbf{J}$  is a close approximation to  $\mathbf{J}(\mathbf{y}_n)$ . Therefore, we will concentrate on the amplification matrix  $\mathbf{Z}$ .

First of all, we consider the convergence for small  $\Delta t$ . Since  $\mathbf{Z} = (1 - \theta) \Delta t (\mathbf{T} - \mathbf{D}) \otimes \mathbf{J} + O((\Delta t)^2)$ , the following theorem is easily proved (cf. [13]):

**Theorem 3.1.** The iteration errors of the AFN and AF iteration processes (3.2) and (3.3) satisfy

$$\begin{aligned} \epsilon^j &= O((\Delta t)^{2j}) \epsilon^0, \quad j \geq 1 && \text{if } \mathbf{T} \text{ is diagonal or if } \theta = 1, \\ \epsilon^j &= \begin{cases} O((\Delta t)^j) \epsilon^0 & \text{for } 1 \leq j \leq s-1 \\ O((\Delta t)^{2j+1-s}) \epsilon^0 & \text{for } j \geq s \end{cases} && \text{if } \mathbf{T} \text{ is lower triangular and } \theta \neq 1. \quad \blacklozenge \end{aligned}$$

This theorem shows that we always have convergence if  $\Delta t$  is sufficiently small. It also indicates that the nonstiff error components (corresponding with eigenvalues of  $\mathbf{J}_k$  of modest magnitude) are rapidly removed from the iteration error. Furthermore, we now see the price to be paid if we set  $\theta = 0$ , while  $\mathbf{T}$  is lower triangular (and not diagonal). In such cases, the subsystems in (3.2) and (3.3) can still be solved in parallel, however, at the cost of a lower order of convergence.

**3.1.2. Convergence and stability regions.** The eigenvalues  $\lambda(\mathbf{Z})$  of the amplification matrix  $\mathbf{Z}$  will be called the *amplification factors* in the iteration process. As in the stability analysis, we consider the test equation where the Jacobian matrices  $\mathbf{J}_k$  commute. For this model problem, they are given by the eigenvalues of the matrix  $\mathbf{I} - \Pi^{-1}(\mathbf{I} - \Delta t \mathbf{T} \otimes \mathbf{J})$ , so that

$$\lambda(\mathbf{Z}) = 1 - (1 - \lambda(\mathbf{T}) \mathbf{e}^T \mathbf{z}) \prod_{k=1}^m (1 - \lambda(\mathbf{T}) z_k)^{-1}.$$

Note that  $\lambda(\mathbf{Z})$  does not depend on the parameter  $\theta$ . We shall call a method convergent at  $\mathbf{z}$  if  $\lambda(\mathbf{Z})$  is within the unit circle at  $\mathbf{z}$ . This leads us to the following analogue of Definition 2.1:



**Definition 3.1.** The iteration method is called

- $A(\alpha)$ -convergent if it is convergent for  $z_k \in \mathbb{W}(\alpha)$ ,  $k = 1, \dots, m$ ,  
 $A$ -convergent if it is convergent for  $z_k \in \mathbb{W}(\pi/2)$ ,  $k = 1, \dots, m$ ,  
 $A_r(\alpha)$ -convergent if it is convergent for  $z_1, \dots, z_r \in \mathbb{R}(\infty) \wedge z_{r+1}, \dots, z_m \in \mathbb{W}(\alpha)$ . ♦

From now on, we shall explicitly assume that

- (3.6) the corrector method is A-stable or L-stable,  
the matrix T has nonnegative eigenvalues,  
the iteration process is performed until convergence.

These assumptions imply that the region of stability equals the region of convergence. The following theorem provides information on the  $A(\alpha)$ -stability characteristics [8].

**Theorem 3.2.** Let the conditions (3.6) be satisfied. Then, AFN and AF iteration is  $A(0)$ -stable for  $m \geq 2$ , A-stable for  $m = 2$ , and  $A(\pi/4)$ -stable for  $m = 3$ . ♦

A comparison with the Theorems 2.1, 2.2 and 2.3 reveals that for  $m \leq 3$  AFN and AF iteration have the same  $A(\alpha)$ -stability characteristics as obtained for the noniterative methods discussed in this paper. However, the stability results of Theorem 3.2 apply to any A-stable or L-stable integration method of the form (2.8) or (2.9) with  $\lambda(T) \geq 0$ , so that the order of accuracy can be raised beyond 2.

Furthermore, we found the stability region  $\mathbb{W}(\pi/2) \times \mathbb{W}(\pi/2) \times \mathbb{R}(\beta)$  with  $\beta = (1 + \sqrt{2}) \rho^{-1}(T)$ , where  $\rho(T)$  denotes spectral radius of T (do not confuse  $\rho(T)$  with the Dahlquist polynomial  $\rho(\zeta)$  used in (2.3)). Hence, this region can be made greater than the corresponding stability regions of all preceding noniterative methods (see Table 5.1) by choosing corrector methods such that  $\rho(T)$  is sufficiently small. In Section 4, we give methods with  $\rho(T)$  in the range [0.13, 0.5].

An even greater advantage is that factorized iteration leads to stability in regions of the form  $\mathbb{I}(\beta_1) \times \mathbb{I}(\beta_1) \times \mathbb{I}(\beta_2)$  with substantial values of  $\beta_1$  and  $\beta_2$ . In [13] it was shown that

$$\beta_1 = \min_{\lambda \in \Lambda(T)} \min_{0 \leq x \leq \lambda \beta_2} \frac{g(x)}{\lambda},$$

where  $\Lambda(T)$  denotes the spectrum of T and  $g$  is defined by  $4xg^3 + 2(x^2 - 1)g^2 - x^2 - 1 = 0$ . Thus, if we choose  $\beta_1$  not larger than the minimal value of  $g(x)\rho^{-1}(T)$  in the interval  $[0, \infty]$ , then we have stability in the region  $\mathbb{I}(\beta_1) \times \mathbb{I}(\beta_1) \times \mathbb{W}(\pi/2)$ . This optimal value of  $\beta_1$  is given by

$$(3.7) \quad \beta_1 = \frac{1}{6\rho(T)} \left( 2 + (26 + 6\sqrt{33})^{1/3} - 8(26 + 6\sqrt{33})^{-1/3} \right) \approx \frac{0.65}{\rho(T)}.$$

Since usually  $\rho(T)$  is less than 1, we obtain quite substantial values for  $\beta_1$ . This makes the iterative approach superior to the noniterative approach, where we found stability regions of the form  $\mathbb{I}(\beta) \times \mathbb{I}(\beta) \times \mathbb{I}(\beta)$  with at best quite small  $\beta$ . The stability region  $\mathbb{I}(\beta_1) \times \mathbb{I}(\beta_1) \times \mathbb{W}(\pi/2)$  enables us to integrate shallow water problems where we need unconditional stability in the vertical direction

(because of the usually fine vertical resolutions) and substantial imaginary stability boundaries in the horizontal directions (because of the convection terms). The AFN-BDF method was successfully used in [15] and [24] for the solution of large scale, three-dimensional shallow water transport problems.

### 3.2. Partially implicit iteration methods

The AFN and AF iteration methods (3.2) and (3.3) are implicit with respect to all Jacobians  $J_k$  in the splitting  $J(\mathbf{y}) = \sum J_k$ . However, Table 5.1 frequently shows finite values for the stability boundaries. This raises the question whether it is necessary to treat all terms in the corresponding splitting implicitly. Afterall, when applying the standard, explicit, fourth-order Runge-Kutta method, we have real and imaginary stability boundaries of comparable size, viz.  $\beta \approx 2.8$  and  $\beta = 2\sqrt{2}$ , respectively.

In [13] this question is addressed and preliminary results are reported for iteration methods where  $\Pi$  does not contain all Jacobians  $J_k$ . In this approach, the iteration method can be fully tuned to the problem at hand. In this paper, we illustrate the partially implicit approach for transport problems in three-dimensional air pollution, where the horizontal spatial derivatives are often treated explicitly. In such problems, the Jacobian matrix  $J(\mathbf{y})$  can be split into three matrices where  $J_1$  corresponds with the convection terms and the two horizontal diffusion terms,  $J_2$  corresponds with the vertical diffusion term, and  $J_3$  corresponds with the chemical reaction terms. It is typical for air pollution terms that  $J_2$  and  $J_3$  are extremely stiff (that is, possess eigenvalues of large magnitude), and that  $J_1$  is moderately stiff in comparison with  $J_2$  and  $J_3$  (see e.g. [23] and [27]). This leads us to apply (3.2) or (3.3) with  $\Pi$  replaced with  $\Pi_1 := (I - \Delta t D \otimes J_2)(I - \Delta t D \otimes J_3)$ . Thus, only the vertical diffusion and the chemical interactions are treated implicitly. In the error recursions (3.4) and (3.5), the amplification matrix  $Z$  should be replaced by

$$Z_1 := I - (\Pi_1 - \theta \Delta t ((T-D) \otimes J))^{-1} (I - \Delta t T \otimes J), \quad \Pi_1 := (I - \Delta t D \otimes J_2)(I - \Delta t D \otimes J_3).$$

Since  $Z_1 = O(\Delta t)$ , the nonstiff components in the iteration error are less strongly damped than by the AFN and AF processes (see Theorem 3.1). This is partly compensated by the lower iteration costs when using  $\Pi_1$  instead of  $\Pi$ .

Let us assume that the eigenvalues  $z_2$  of  $\Delta t J_2$  are negative (vertical diffusion) and the eigenvalues  $z_3$  of  $J_3$  are in the left halfplane (chemical reactions). We are now interested to what region we should restrict the eigenvalues  $z_1$  of  $\Delta t J_1$  in order to have convergence. For the model problem this region is determined by the intersection of the domains bounded by the curve

$$(3.8) \quad |\lambda z_1|^2 + 2\lambda^3 z_2 \operatorname{Im}(z_3) \operatorname{Im}(z_1) = (1 + \lambda^2 |z_3|^2)(1 - \lambda z_2)^2 - \lambda^4 z_2^2 |z_3|^2,$$

where  $\lambda \in \Lambda(T)$ ,  $z_2 \in \mathbb{R}(\infty)$  and  $z_3 \in \mathbb{W}(\pi/2)$ . It can be verified that this intersection is given by the points  $|z_1| < \rho^{-1}(T)$ . Thus, we have proved:

**Theorem 3.3.** Let the conditions (3.6) be satisfied, let  $m = 3$ , let  $\Pi$  be replaced by  $\Pi_1$  in the AFN and AF iteration methods, and define the disk  $\mathbb{D}(\beta) := \{w \in \mathbb{C}: |w| < \beta\}$ . Then, the stability region contains the region  $\mathbb{D}(\rho^{-1}(T)) \times \mathbb{R}(\infty) \times \mathbb{W}(\pi/2)$ .  $\blacklozenge$

We remark that the approximate factorization operator  $\Pi_1$  is not symmetric with respect to all three Jacobians. This means that the stability region of the methods of Theorem 3.3 also contain the region  $\mathbb{D}(\rho^{-1}(T)) \times \mathbb{W}(\pi/2) \times \mathbb{R}(\infty)$ , but not e.g. the region  $\mathbb{R}(\infty) \times \mathbb{D}(\rho^{-1}(T)) \times \mathbb{W}(\pi/2)$ .

### 3.3. A-stability for three-component Jacobian splitting

So far, the approximate factorization methods constructed in this paper are not A-stable for three-component splittings. However, all these methods can be modified such that they become A-stable for  $m = 3$ . The idea is to start with e.g. a two-component splitting  $J = J_1 + J^*$ , where  $J_1$  has the desired simple structure, but  $J^*$  has not, and to solve the linear system containing  $J^*$  iteratively with approximate factorization iteration. We illustrate this for the AFN process (3.2). Consider the process

$$(3.9a) \quad (\mathbf{I} - \Delta t \mathbf{D} \otimes \mathbf{J}_1) \tilde{\Delta}^j = \mathbf{G}_n - \mathbf{Y}^{j-1} + \Delta t (\mathbf{T} \otimes \mathbf{I}) \mathbf{F}(\mathbf{Y}^{j-1}) + \theta \Delta t ((\mathbf{T} - \mathbf{D}) \otimes \mathbf{J}) (\mathbf{Y}^j - \mathbf{Y}^{j-1}),$$

$$(3.9b) \quad (\mathbf{I} - \Delta t \mathbf{D} \otimes (\mathbf{J}_2 + \dots + \mathbf{J}_m)) \Delta^j = \tilde{\Delta}^j, \quad \mathbf{Y}^j - \mathbf{Y}^{j-1} = \Delta^j,$$

where  $j \geq 1$ . This method can be interpreted as the AFN method (3.2) with  $m$  replaced by 2 and  $J_2$  replaced by  $J_2 + \dots + J_m$ . Hence, Theorem 3.2 implies that we have A-stability with respect to the eigenvalues of  $J_1$  and  $J_2 + \dots + J_m$  (assuming that the corrector is A-stable). If the matrix  $J_2 + \dots + J_m$  does not have a 'convenient' structure (e.g. a small band width), then the system for  $\Delta^j$  cannot be solved efficiently. In such cases, we may solve this system by an AFN (inner) iteration process:

$$(3.9c) \quad \begin{aligned} \Delta^{j,0} &= \tilde{\Delta}^j, \\ \Pi_1(\Delta^{j,i} - \Delta^{j,i-1}) &= \tilde{\Delta}^j - (\mathbf{I} - \Delta t \mathbf{D} \otimes (\mathbf{J}_2 + \dots + \mathbf{J}_m)) \Delta^{j,i-1}, \quad i = 1, \dots, r, \\ \mathbf{Y}^j &= \mathbf{Y}^{j-1} + \Delta^{j,r}, \quad \Pi_1 := \prod_{k=2}^m (\mathbf{I} - \Delta t \mathbf{D} \otimes \mathbf{J}_k). \end{aligned}$$

Since the stability theory of Section 3.1 applies to this process, we can apply the stability results from this section. The following theorem summarizes the main results for  $\{(3.9a), (3.9c)\}$ :

**Theorem 3.4.** Let (3.6) be satisfied. Then, the inner-outer AFN process  $\{(3.9a), (3.9c)\}$  is

- (a) A(0)-stable for  $m \geq 2$ , A-stable for  $m = 2$  and  $m = 3$ , and A( $\pi/4$ )-stable for  $m = 4$ ,
- (d) Stable in  $\mathbb{W}(\pi/2) \times \mathbb{I}(\beta) \times \mathbb{I}(\beta) \times \mathbb{W}(\pi/2)$  with  $\beta \approx 0.65 \rho^{-1}(T)$  for  $m = 4$ ,
- (e) Stable in  $\mathbb{W}(\pi/2) \times \mathbb{W}(\pi/2) \times \mathbb{W}(\pi/2) \times \mathbb{R}(\beta)$  with  $\beta = (1 + \sqrt{2}) \rho^{-1}(T)$  for  $m = 4$ .  $\blacklozenge$

Thus the process  $\{(3.9a), (3.9c)\}$  has excellent stability properties, but its computational complexity is considerably larger than that of (3.2). In order to compare this, let the number of outer iterations

be denoted by  $q$ . Then (3.2) requires the solution of  $qm$  linear systems, whereas  $\{(3.9a),(3.9c)\}$  requires  $q(rm + 1 - r)$  linear system solutions. For example, if  $m = 3$ , then we need  $3q$  and  $(2r+1)q$  linear system solutions, respectively, so that for  $r \geq 2$  the nested approach is more expensive.

Evidently, the above approach can also be applied to the AF method (3.3), but also to the noniterative methods (2.4), (2.6), and (2.10). In the noniterative methods, the computational complexity increases from  $m$  to  $rm + 1 - r$  linear system solutions, i.e. by the same factor as in the iterative case.

It should be remarked that there exist several splitting methods, not based on approximate factorization, that are also A-stable for three-component Jacobian splittings. We mention the ADI method of Gourlay and Mitchell [10] and the trapezoidal and midpoint splitting methods of Hundsdorfer [18]. These methods are second-order accurate and possess the same amplification factor  $\zeta(\mathbf{z}) = \prod_{k=1}^3 (1 + \frac{1}{2}z_k)(1 - \frac{1}{2}z_k)^{-1}$  from which the A-stability is immediate. However, the internal stages of these methods are not consistent, that is, in a steady state the internal stage values are not stationary points of the method. This leads to loss of accuracy (cf. [17]).

#### 4. Methods with minimal $\rho(\mathbf{T})$

The stability regions in Table 5.1 and the Theorems 3.3 and 3.4 indicate that small values of  $\rho(\mathbf{T})$  increase the stability regions of the iterated methods. Similarly, Theorem 2.1 shows that small values of  $b_0$  increases the stability region  $\mathbb{I}(\beta_1) \times \mathbb{I}(\beta_1) \times \mathbb{R}(\beta_2)$  of the Warming-Beam method (note that for LM methods  $b_0 = \rho(\mathbf{T})$ ). Therefore, it is relevant to look for methods with small  $\rho(\mathbf{T})$ . Let us first consider the two-parameter family of all second-order, A-stable linear two-step methods (cf. [28, Figure 2]). Taking  $b_0$  and  $a_2$  as the free parameters (see (2.3)), this family is defined by

$$(4.1) \quad \rho(\zeta) = \zeta^2 - (a_2 + 1)\zeta + a_2, \quad \sigma(\zeta) = b_0\zeta^2 + \frac{1}{2}(3 - a_2 - 4b_0)\zeta + b_0 - \frac{1}{2}(1 + a_2),$$

where  $-1 \leq a_2 < 1$  and  $b_0 \geq \frac{1}{2}$ . Hence, the smallest value of  $b_0$  is  $\frac{1}{2}$ . Moreover, from an implementation point of view, the trapezoidal rule choice  $a_2 = 0$  is attractive.

Next, we consider the family of DIRK methods. We recall that they are defined by (2.8) with  $\mathbf{g}_n = (1 - \mathbf{a}^T \mathbf{e})\mathbf{y}_n$  and  $\mathbf{G}_n = \mathbf{Y}^0 = \mathbf{e} \otimes \mathbf{y}_n$ . In [14] a number of methods with minimal  $\rho(\mathbf{T})$ , relative to the number of stages, have been derived. Here, we confine ourselves to presenting a few second-order, L-stable methods by specifying the matrix  $\mathbf{T}$  (in all cases  $\mathbf{a} = \mathbf{e}_s$  and  $\rho(\mathbf{T}) = \kappa$ ).

$$(4.2) \quad \mathbf{T} = \begin{pmatrix} \kappa & 0 \\ 1 - \kappa & \kappa \end{pmatrix}, \quad \kappa = 1 - \frac{1}{2}\sqrt{2} \approx 0.29,$$

$$(4.3) \quad \mathbf{T} = \begin{pmatrix} \kappa & 0 & 0 \\ \frac{1 - 4\kappa + 2\kappa^2}{2(1 - \kappa)} & \kappa & 0 \\ 0 & 1 - \kappa & \kappa \end{pmatrix}, \quad \kappa = \frac{1}{12} \left( 9 + 3\sqrt{3} - \sqrt{72 + 42\sqrt{3}} \right) \approx 0.18,$$

$$(4.4) \quad T = \frac{1}{4} \begin{pmatrix} 4\kappa & 0 & 0 & 0 \\ \frac{1-8\kappa+16\kappa^2+8\kappa^3}{1-4\kappa+2\kappa^2} & 4\kappa & 0 & 0 \\ 0 & \frac{2-8\kappa+4\kappa^2}{1-\kappa} & 4\kappa & 0 \\ 0 & 0 & 4(1-\kappa) & 4\kappa \end{pmatrix} \quad \kappa = \frac{4 + 2\sqrt{2} - \sqrt{20 + 14\sqrt{2}}}{4} \approx 0.13.$$

## 5. Summary of stability results

We conclude this paper with Table 5.1 which compares a number of stability results for various factorized methods based on three-component splittings.

**Table 5.1.** Stability regions of approximate factorization methods with 3-component splittings.

Methods	Stability region	Stability boundaries
Warming-Beam (2.4) and AFL-trapezoidal ( $b_0 = \frac{1}{2}$ )	$\mathbb{W}(\pi/4) \times \mathbb{W}(\pi/4) \times \mathbb{W}(\pi/4)$ $\mathbb{W}(\pi/2) \times \mathbb{W}(\pi/2) \times \mathbb{R}(\beta)$ $\mathbb{I}(\beta) \times \mathbb{I}(\beta) \times \mathbb{R}(\infty),$	$\beta = \frac{3}{b_0}$ $\beta = 0$
AFL-BDF	$\mathbb{W}(\pi/4) \times \mathbb{W}(\pi/4) \times \mathbb{W}(\pi/4)$ $\mathbb{W}(\pi/2) \times \mathbb{W}(\pi/2) \times \mathbb{R}(\beta),$ $\mathbb{I}(\beta) \times \mathbb{I}(\beta) \times \mathbb{R}(\infty),$	$\beta = \frac{9 + 3\sqrt{17}}{4} \approx 5.34$ $\beta = \frac{3}{4}\sqrt{2} \approx 1.06$
AFL-DIRK (2.16) and Factorized Rosenbrock (2.19) with $\kappa = 1 - \frac{1}{2}\sqrt{2}$	$\mathbb{W}(\pi/4) \times \mathbb{W}(\pi/4) \times \mathbb{W}(\pi/4)$ $\mathbb{W}(\pi/2) \times \mathbb{W}(\pi/2) \times \mathbb{R}(\beta)$ $\mathbb{I}(\beta) \times \mathbb{I}(\beta) \times \mathbb{R}(\infty),$	$\beta \approx 10.2$ $\beta \approx 1.26$
AFN/AF iteration (3.2) / (3.3)	$\mathbb{W}(\pi/4) \times \mathbb{W}(\pi/4) \times \mathbb{W}(\pi/4)$ $\mathbb{W}(\pi/2) \times \mathbb{W}(\pi/2) \times \mathbb{R}(\beta)$ $\mathbb{I}(\beta) \times \mathbb{I}(\beta) \times \mathbb{W}(\pi/2)$	$\beta = (1 + \sqrt{2}) \rho^{-1}(T) \approx 2.41 \rho^{-1}(T)$ $\beta \approx 0.65 \rho^{-1}(T)$
AFN/AF iteration (3.2) / (3.3) with $\Pi_1$ (Section 3.2)	$\mathbb{D}(\beta) \times \mathbb{R}(\infty) \times \mathbb{W}(\pi/2)$	$\beta = \rho^{-1}(T)$
Nested AFN {(3.9a),(3.9c)}	$\mathbb{W}(\pi/2) \times \mathbb{W}(\pi/2) \times \mathbb{W}(\pi/2)$	

## References

- [1] Beam, M. & Warming, R.F. [1976]: An implicit finite-difference algorithm for hyperbolic systems in conservation-law form, *J. Comput. Physics* 22, 87-110.
- [2] Beam, M. & Warming, R.F. [1979]: An implicit factored scheme for the compressible Navier-Stokes equations II, The numerical ODE connection, paper no. 79-1446, Proc. AIAA 4th Computational Fluid Dynamics Conference, Williamsburg, Virginia, 1979.
- [3] Crank, J. & Nicolson, P. [1947]: A practical method for numerical integration of solutions of partial differential equations of heat-conduction type, *Proc. Cambridge Philos. Soc.* 43, 50-67.
- [4] Dekker, K. & Verwer, J.G. [1984]: *Stability of Runge-Kutta methods for stiff nonlinear differential equations*, North-Holland.
- [5] Douglas, J. Jnr. [1955]: On the numerical integration of  $u_{xx} + u_{yy} = u_t$ , *J. Soc. Indust. App. Math.* 3, 42-65.
- [6] Douglas, J. Jnr. [1962]: Alternating direction methods for three space variables, *Num. Math.* 4, 41-63.
- [7] Douglas, J. Jnr. & Gunn, J.E. [1964]: A general formulation of alternating direction methods. Part I. Parabolic and hyperbolic problems, *Num. Math.* 6, 428-453.
- [8] Eichler-Liebenow, C., Houwen, P.J. van der & Sommeijer, B.P. [1998]: Analysis of approximate factorization in iteration methods, *Appl. Numer. Math.* 28, 245-258.
- [9] Golub, G. H. & Van Loan, C.F. [1989]: *Matrix computations*, Second edition, The Johns Hopkins University Press, Baltimore, Maryland.
- [10] Gourlay, A.R. & Mitchell, A.R. [1972]: On the structure of alternating direction implicit (A.D.I.) and locally one dimensional (L.O.D.) difference methods, *J. Inst. Math. Applics.* 9, 80-90.
- [11] Hairer, E. & Wanner, G. [1991]: *Solving ordinary differential equations II. Stiff and differential-algebraic problems*, Springer-Verlag, Berlin.
- [12] Houwen, P.J. van der & Sommeijer, B.P. [1998]: Approximate factorization in shallow water applications, Report MAS R9835, CWI, Amsterdam, submitted for publication.
- [13] Houwen, P.J. van der & Sommeijer, B.P. [1999]: Factorization in block-triangularly implicit methods for shallow water applications, Report MAS R9906, CWI, Amsterdam, submitted for publication.
- [14] Houwen, P.J. van der & Sommeijer, B.P. [1999]: Diagonally implicit Runge-Kutta methods for 3D shallow water applications, Report MAS R9907, CWI, Amsterdam, submitted for publication.
- [15] Houwen, P.J. van der, Sommeijer, B.P. & Kok, J. [1997]: The iterative solution of fully implicit discretizations of three-dimensional transport models, *Appl. Numer. Math.* 25, 243-256.
- [16] Houwen, P.J. van der & Verwer, J.G. [1979]: One-step splitting methods for semi-discrete parabolic equations, *Computing* 22, 291-309.

- [17] Hundsdorfer, W. [1998]: A note on the stability of the Douglas splitting method, *Math. Comp.* 67, 183-190.
- [18] Hundsdorfer, W. [1998]: Trapezoidal and midpoint splittings for initial-boundary value problems, *Math. Comp.* 67, 1047-1062.
- [19] Marchuk, G.I. [1990], Splitting methods and alternating directions methods, in: *Handbook of Numerical Analysis* (eds. P.G. Ciarlet and J.L. Lions), Vol. I, Part I, 197-459.
- [20] Rosenbrock, H.H. [1962-1963]: Some general implicit processes for the numerical solution of differential equations, *Computer J.* 5, 329-330.
- [21] Peaceman, D.W. & Rachford, H.H. Jnr. [1955]: The numerical solution of parabolic and elliptic differential equations, *J. Soc. Indust. App. Math.* 3, 28-41.
- [22] Richtmyer, R.D. & Morton, K.W. [1967]: *Difference methods for initial-value problems*, Wiley.
- [23] Sandu, A. [1997]: *Numerical aspects of air quality modelling*, PhD thesis, University of Iowa.
- [24] Sommeijer, B.P. [1999]: The iterative solution of fully implicit discretizations of three-dimensional transport models, in: *Parallel Computational Fluid Dynamics - Implementation and Results Using Parallel Technology* (eds. C.A. Lin, A. Ecer, J. Periaux, N. Satofuka), *Proceedings of the 10th Int. Conf. on Parallel CFD*, May 1998, Hsinchu, Taiwan, Elsevier.
- [25] Steihaug, T. & Wolfbrandt, A. [1979]: An attempt to avoid exact Jacobian and nonlinear equations in the numerical solution of stiff differential equations, *Math. Comp.* 33, 521-534.
- [26] Strehmel, K. & Weiner, R. [1992]: *Linearly implicit Runge-Kutta methods and their application* (German), Teubner.
- [27] Verwer, J.G., Hundsdorfer, W. & Blom, J.G. [1998]: *Numerical time integration for air pollution models*, Report MAS R9825, CWI, Amsterdam, submitted for publication.
- [28] Warming, R.F. & Beam, M. [1979]: An extension of A-stability to alternating direction implicit methods, *BIT* 19, 395-417.