



Centrum voor Wiskunde en Informatica

**REPORTRAPPORT**

Factorization in block-triangularly implicit methods for shallow water applications

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

Modelling, Analysis and Simulation (MAS)

**MAS-R9906 March 1999**

Report MAS-R9906  
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

# Factorization in Block-Triangularly Implicit Methods for Shallow Water Applications

P.J. van der Houwen & B.P. Sommeijer  
CWI

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

## ABSTRACT

The systems of first-order ordinary differential equations obtained by spatial discretization of the initial-boundary value problems modelling phenomena in shallow water in 3 spatial dimensions have righthand sides of the form  $\mathbf{f}(t, \mathbf{y}) := \mathbf{f}_1(t, \mathbf{y}) + \mathbf{f}_2(t, \mathbf{y}) + \mathbf{f}_3(t, \mathbf{y}) + \mathbf{f}_4(t, \mathbf{y})$ , where  $\mathbf{f}_1$ ,  $\mathbf{f}_2$  and  $\mathbf{f}_3$  contain the spatial derivative terms with respect to the x, y and z directions, respectively, and  $\mathbf{f}_4$  represents the forcing terms and/or reaction terms. The number N of components of  $\mathbf{f}$  is usually extremely large. It is typical for shallow water applications that the function  $\mathbf{f}_4$  is nonstiff and that the function  $\mathbf{f}_3$  corresponding with the vertical spatial direction is much more stiff than the functions  $\mathbf{f}_1$  and  $\mathbf{f}_2$  corresponding with the horizontal spatial directions. The reason is that in shallow seas the gridsize in the vertical direction is several orders of magnitude smaller than in the horizontal directions. In order to solve the initial value problem for the system of ordinary differential equations numerically, we need a stiff solver. Stiff IVP solvers are necessarily implicit, requiring the solution of large systems of implicit relations. In a few earlier papers, we considered implicit Runge-Kutta methods leading to *fully coupled*, implicit systems whose dimension is a multiple of N, and *block-diagonally* implicit methods in which the implicit relations can be decoupled into subsystems of dimension N. In the present paper, we analyse Rosenbrock type methods and the related DIRK methods (diagonally implicit Runge-Kutta methods) leading to *block-triangularly* implicit relations. In particular, we shall present a convergence analysis of various iterative methods based on approximate factorization for solving the triangularly implicit relations.

1991 Mathematics Subject Classification: 65L06

Keywords and Phrases: numerical analysis, shallow water applications, iteration methods, approximate factorization, parallelism.

Note. The investigations reported in this paper were partly supported by the Dutch HPCN Program.

## 1. Introduction

We consider initial-boundary value problems modelling phenomena in shallow water in 3 spatial dimensions. The systems of ordinary differential equations (ODEs) obtained by spatial discretization (method of lines) of the governing partial differential equations can be written in the form

$$(1.1) \quad \frac{d\mathbf{y}(t)}{dt} = \mathbf{f}(t, \mathbf{y}(t)), \quad \mathbf{f}(t, \mathbf{y}) := \mathbf{f}_1(t, \mathbf{y}) + \mathbf{f}_2(t, \mathbf{y}) + \mathbf{f}_3(t, \mathbf{y}) + \mathbf{f}_4(t, \mathbf{y}), \quad \mathbf{y}, \mathbf{f}_k \in \mathbb{R}^N,$$

where  $\mathbf{f}_1$ ,  $\mathbf{f}_2$  and  $\mathbf{f}_3$  contain the spatial derivative terms with respect to the x, y and z directions, respectively,  $\mathbf{f}_4$  represents the forcing terms and/or reaction terms, and N is a large integer proportional to the number of spatial grid points used for the spatial discretization. It is typical for shallow water applications that the function  $\mathbf{f}_4$  is nonstiff and that the function  $\mathbf{f}_3$  corresponding with the vertical spatial direction is much more stiff than the functions  $\mathbf{f}_1$  and  $\mathbf{f}_2$  corresponding with the horizontal spatial directions. As a consequence, the spectral radius of the Jacobian matrix  $\partial \mathbf{f}_3 / \partial \mathbf{y}$  is much larger than the spectral radius of  $\partial \mathbf{f}_1 / \partial \mathbf{y}$  and  $\partial \mathbf{f}_2 / \partial \mathbf{y}$ . The reason is that in shallow seas the

gridsize in the vertical direction is several orders of magnitude smaller than in the horizontal directions.

In order to solve the initial value problem (IVP) for the system (1.1) numerically, we need a stiff IVP solver, because the Lipschitz constants with respect to  $\mathbf{y}$  associated with the functions  $\mathbf{f}_1$ ,  $\mathbf{f}_2$  and  $\mathbf{f}_3$  become increasingly large as the spatial resolution is refined. Stiff IVP solvers are necessarily implicit, requiring the solution of large systems of implicit relations. In a few earlier papers, we considered implicit Runge-Kutta methods leading to *fully coupled*, implicit systems whose dimension is a multiple of  $N$  (cf. [3], [5] and [10]), and *block-diagonally* implicit methods in which the implicit relations can be decoupled into subsystems of dimension  $N$  (cf. [6]). In the present paper, we analyse Rosenbrock type methods and the related DIRK methods (diagonally implicit Runge-Kutta methods) leading to *block-triangularly* implicit relations (this is also the case for the DIRK methods, in spite of the terminology 'diagonally implicit'). Rosenbrock type methods, and in particular factorized versions of these methods, are quite popular in air pollution simulations (see e.g. [9], [12], and [13]). This motivated us to look whether Rosenbrock and the related DIRK methods can also be useful in shallow water modelling. First we show that in shallow water applications, factorized Rosenbrock methods are less suitable. However, iteration of Rosenbrock and DIRK methods using approximate factorization looks quite promising. This paper will focus on the convergence analysis of approximate factorization iteration of the triangularly implicit Rosenbrock and DIRK relations.

## 2. Rosenbrock methods and their factorization

We start with an example of a family of two-stage Rosenbrock methods:

$$\begin{aligned}
 \mathbf{y}_{n+1} &= \mathbf{y}_n + b\mathbf{k}_1 + (1-b)\mathbf{k}_2, \\
 (2.1) \quad (\mathbf{I} - \kappa_1 \Delta t \mathbf{J})\mathbf{k}_1 &= \Delta t \mathbf{f}(\mathbf{y}_n), \\
 (\mathbf{I} - \kappa_2 \Delta t \mathbf{J})\mathbf{k}_2 &= \Delta t \mathbf{f}(\mathbf{y}_n + \mu \mathbf{k}_1) + v \Delta t \mathbf{J} \mathbf{k}_1, \quad \kappa_i > 0, \quad \mu := \frac{\frac{1}{2} - b\kappa_1 + (b-1)\kappa_2}{1 - b} - v.
 \end{aligned}$$

Here,  $b$ ,  $\kappa_1$ ,  $\kappa_2$  and  $v$  are free parameters and  $\mathbf{J}$  is an approximation to the Jacobian matrix  $\partial \mathbf{f} / \partial \mathbf{y}$  at  $t_n$ . For simplicity of notation, we assumed the ODE of autonomous form. The nonautonomous version can be obtained by applying (2.1) to the augmented system  $\{\mathbf{y}' = \mathbf{f}(y_0, \mathbf{y}), y_0' = 1\}$ . The method (2.1) is triangularly implicit, that is,  $\mathbf{k}_1$  and  $\mathbf{k}_2$  can be computed by successively solving 2 linear systems of dimension  $N$ .

If  $\mathbf{J} = \partial \mathbf{f} / \partial \mathbf{y}(t_n) + O(\Delta t)$ , then the formulas (2.1) are all second-order accurate Rosenbrock methods. The stability function for (2.1) is given by

$$(2.2) \quad R(z) = \frac{1 + (1 - \kappa_1 - \kappa_2)z + \frac{1}{2}(1 - 2\kappa_1 - 2\kappa_2 + 2\kappa_1\kappa_2)z^2}{(1 - \kappa_1 z)(1 - \kappa_2 z)}.$$

From this expression it follows that the methods (2.1) are A-stable if  $\frac{1}{2} \leq \kappa_1 + \kappa_2 \leq 2\kappa_1\kappa_2 + \frac{1}{2}$  and L-stable if  $\kappa_1 + \kappa_2 = \kappa_1\kappa_2 + \frac{1}{2}$ .

The first examples of Rosenbrock methods were given by Rosenbrock [8] in 1962 and are obtained by choosing in (2.1)

$$(2.3) \quad \mathbf{b} = 0, \quad \kappa_1 = \kappa_2 = \kappa := 1 \pm \frac{1}{2}\sqrt{2}, \quad \mathbf{v} = 0.$$

Of particular interest are the methods which remain second-order accurate if we choose an arbitrary matrix for  $\mathbf{J}$ . Such methods are called Rosenbrock-W methods and were proposed by Steihaug and Wolfbrandt [11]. If we choose in (2.1)  $\kappa_1 = \kappa_2 = \kappa$  and  $\mathbf{v} = -\kappa(1-\mathbf{b})^{-1}$ , then (2.1) becomes a W-method (see Dekker and Verwer [2, p. 233]). The special case

$$(2.4) \quad \mathbf{b} = \frac{1}{2}, \quad \kappa_1 = \kappa_2 = \kappa := 1 \pm \frac{1}{2}\sqrt{2}, \quad \mathbf{v} = -2\kappa$$

was used by Verwer et al. [12] for solving atmospheric transport problems. Note, however, that for stability reasons,  $\mathbf{J}$  should be a reasonably close approximation to the true Jacobian  $\partial \mathbf{f} / \partial \mathbf{y}$  at  $\mathbf{t}_n$ .

## 2.1. General Rosenbrock methods

More generally, we consider Rosenbrock methods of the form (cf. [4, p. 111])

$$(2.5) \quad \mathbf{y}_{n+1} = \mathbf{y}_n + (\mathbf{b}^T \otimes \mathbf{I}) \mathbf{K}, \quad (\mathbf{I} - \mathbf{T} \otimes \Delta t \mathbf{J}) \mathbf{K} = \Delta t \mathbf{F}(\mathbf{e} \otimes \mathbf{y}_n + (\mathbf{L} \otimes \mathbf{I}) \mathbf{K}),$$

where  $\mathbf{b}$  is an  $s$ -dimensional vector,  $\mathbf{K} := (\mathbf{k}_1^T, \dots, \mathbf{k}_s^T)^T$ , and  $\mathbf{T}$  and  $\mathbf{L}$  are lower and strictly lower triangular  $s$ -by- $s$  matrices, respectively. This property of  $\mathbf{T}$  and  $\mathbf{L}$  implies that (2.5) is triangularly implicit, so that the components  $\mathbf{k}_i$  of  $\mathbf{K}$  can be computed by successively solving  $s$  linear systems of dimension  $N$  with system matrices  $\mathbf{I} - \kappa_i \Delta t \mathbf{J}$ , where the  $\kappa_i$  denote the diagonal entries of  $\mathbf{T}$ . 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.

If  $\mathbf{T}$  is not diagonal (as in (2.4)), then for an actual implementation one often transforms the linear system for  $\mathbf{K}$  by a Butcher similarity transformation  $\mathbf{U} = (\mathbf{T} \otimes \mathbf{I}) \mathbf{K}$ , where  $\mathbf{T}$  is assumed invertible (cf. [4, p. 120]). Writing  $\mathbf{T}^{-1} = \mathbf{S} + \mathbf{D}^{-1}$  with  $\mathbf{S}$  strictly lower triangular and  $\mathbf{D} = \text{diag}(\mathbf{T})$ , (2.5) becomes

$$(2.6) \quad \begin{aligned} \mathbf{y}_{n+1} &= \mathbf{y}_n + (\mathbf{b}^T \mathbf{T}^{-1} \otimes \mathbf{I}) \mathbf{U}, \\ (\mathbf{I} - \mathbf{D} \otimes \Delta t \mathbf{J}) \mathbf{U} &= \Delta t (\mathbf{D} \otimes \mathbf{I}) \mathbf{F}(\mathbf{e} \otimes \mathbf{y}_n + (\mathbf{L} \mathbf{T}^{-1} \otimes \mathbf{I}) \mathbf{U}) - (\mathbf{D} \mathbf{S} \otimes \mathbf{I}) \mathbf{U}. \end{aligned}$$

As in (2.5) the components  $\mathbf{u}_i$  of  $\mathbf{U}$  can be computed by again successively solving  $s$  linear systems of dimension  $N$ . As an example of a transformed Rosenbrock method, we give the transformation of the method (2.4):

$$(2.4') \quad \begin{aligned} \mathbf{y}_{n+1} &= \mathbf{y}_n + \frac{1}{2} \kappa^{-1} (3\mathbf{u}_1 + \mathbf{u}_2), \\ (\mathbf{I} - \kappa \Delta t \mathbf{J}) \mathbf{u}_1 &= \kappa \Delta t \mathbf{f}(\mathbf{y}_n), \quad \kappa = 1 \pm \frac{1}{2}\sqrt{2}, \\ (\mathbf{I} - \kappa \Delta t \mathbf{J}) \mathbf{u}_2 &= \kappa \Delta t \mathbf{f}(\mathbf{y}_n + \kappa^{-1} \mathbf{u}_1) - 2\mathbf{u}_1. \end{aligned}$$

Note that unlike (2.5), no Jacobian multiplications are involved in transformed Rosenbrock methods. In general, this is considered as an advantage because such Jacobian multiplications can be quite expensive. However, it should be remarked that in the case of shallow water applications the matrix  $J$  is extremely sparse, so that Jacobian multiplications are not so costly.

## 2.2. Factorized Rosenbrock methods

In order to further reduce the linear algebra costs in the method (2.4), Sandu [9] and Verwer et al. [13] applied to the system matrix  $I - \kappa\Delta tJ$  the technique of *approximate factorization* based on some splitting  $\Sigma J_k$  of the Jacobian  $J$ . This leads to the *factorized Rosenbrock method*.

This technique goes back to Peaceman and Rachford [7] who used it for approximately solving the linear systems originating from a finite difference discretization of two-dimensional parabolic problems. In such problems, the system matrix is of the form  $I - \frac{1}{2}\Delta tJ$ , where  $J$  is the discretization of the Laplace operator  $\partial^2/\partial x^2 + \partial^2/\partial y^2$ . By writing  $J = J_1 + J_2$ , where  $J_1$  and  $J_2$  correspond with  $\partial^2/\partial x^2$  and  $\partial^2/\partial y^2$ , respectively, Peaceman and Rachford replaced  $I - \frac{1}{2}\Delta tJ$  by the approximate factorization  $(I - \frac{1}{2}\Delta tJ_1)(I - \frac{1}{2}\Delta tJ_2)$ .

The same approximate factorization technique can be applied to the matrix  $I - T\otimes\Delta tJ$  in (2.5) or to the matrix  $I - D\otimes\Delta tJ$  in (2.6). We shall illustrate this for the case (2.6). Since we are concerned with shallow water applications, we use the splitting  $J = J_1 + J_2 + J_3$ , where the matrices  $J_k$  denote the Jacobian matrices of the terms  $\mathbf{f}_k$  at  $t_n$  occurring in the righthand side function  $\mathbf{f}$  in (1.1) and where the nonstiff interaction terms are ignored. This leads to the factorized method

$$(2.7) \quad \begin{aligned} \mathbf{y}_{n+1} &= \mathbf{y}_n + (\mathbf{b}^T T^{-1} \otimes I) \mathbf{V}, \\ \Pi \mathbf{V} &= \Delta t (D \otimes I) \mathbf{F}(\mathbf{e} \otimes \mathbf{y}_n + (L T^{-1} \otimes I) \mathbf{V}) - (D S \otimes I) \mathbf{V}, \end{aligned}$$

where  $\Pi$  is defined by

$$(2.8) \quad \Pi := (I - D \otimes \Delta t J_1)(I - D \otimes \Delta t J_2)(I - D \otimes \Delta t J_3), \quad D = \text{diag}(T).$$

Each step of the factorized Rosenbrock method (2.7) requires the solution of  $3s$  one-dimensional, linear systems. All LU-decompositions can be computed in parallel, but the  $3s$  forward-backward substitutions have to be done sequentially.

Since  $\Pi = I - D \otimes \Delta t J + O((\Delta t)^2)$ , we can interpret the factorized method as the original Rosenbrock method with an  $O(\Delta t)$ -perturbed matrix  $J$ . Hence, factorization will not affect the order of Rosenbrock-W methods. Furthermore, any factorized Rosenbrock method has at least order two if the original Rosenbrock method has at least order two.

## 2.3. Stability

Next, we define the stability region  $\mathbb{S}$  for the factorized versions of the methods (2.5) and (2.6). We first define the stability function by applying them to the test equation  $\mathbf{y}' = (J_1 + J_2 + J_3)\mathbf{y}$ . Assuming that the matrices  $J_k$  commute and ignoring the interaction terms in  $\mathbf{F}$ , the factorized versions of the methods (2.5) and (2.6) will reduce to recursions of the form

$$\mathbf{y}_{n+1} = \mathbf{R}(\Delta t \mathbf{J}_1, \Delta t \mathbf{J}_2, \Delta t \mathbf{J}_3) \mathbf{y}_n,$$

where  $\mathbf{R}(z_1, z_2, z_3)$  is a rational function of its arguments. Using the identity

$$1 + \mathbf{p}^T \mathbf{M}^{-1} \mathbf{q} = \frac{\det(\mathbf{M} + \mathbf{q} \mathbf{p}^T)}{\det(\mathbf{M})},$$

which holds for any  $m$ -by- $m$  matrix  $\mathbf{M}$  and any two  $m$ -dimensional vectors  $\mathbf{p}$  and  $\mathbf{q}$  (cf. [1, p. 475]), we find that the stability functions corresponding to the factorized versions of (2.5) and (2.6) can be respectively expressed as

$$(2.9) \quad \mathbf{R}(z_1, z_2, z_3) = \frac{\det(\mathbf{P} + z(\mathbf{e} \mathbf{b}^T - \mathbf{L}))}{\det(\mathbf{P})}, \quad \mathbf{P} := (\mathbf{I} - z_1 \mathbf{T})(\mathbf{I} - z_2 \mathbf{T})(\mathbf{I} - z_3 \mathbf{T}),$$

$$(2.10) \quad \mathbf{R}(z_1, z_2, z_3) = \frac{\det(\mathbf{P} + \mathbf{D} \mathbf{S} + z \mathbf{D}(\mathbf{e} \mathbf{b}^T - \mathbf{L}) \mathbf{T}^{-1})}{\det(\mathbf{P})}, \quad \mathbf{P} := (\mathbf{I} - z_1 \mathbf{D})(\mathbf{I} - z_2 \mathbf{D})(\mathbf{I} - z_3 \mathbf{D}).$$

where  $z := z_1 + z_2 + z_3$ .

The *stability region* is defined by the region  $\mathbb{S}$  in the  $(z_1, z_2, z_3)$ -space where  $|\mathbf{R}(z_1, z_2, z_3)| \leq 1$ . The method (2.7) is called *stable* if all eigenvalue triples  $(\Delta t \lambda(\mathbf{J}_1), \Delta t \lambda(\mathbf{J}_2), \Delta t \lambda(\mathbf{J}_3))$  are in  $\mathbb{S}$ . Since in shallow water applications, many of the eigenvalues of  $\mathbf{J}_k$ ,  $k = 1, 2, 3$ , are close to the imaginary axis, we are particularly interested in the most critical case where the eigenvalues of  $\mathbf{J}_k$  are purely imaginary, i.e.  $z_k = iy_k$  with  $y_k$  real-valued. Let us introduce for a given value of  $y_3$  the stability boundary  $\beta(y_3)$  which is such that the method is stable in a region of the form

$$(2.11) \quad \mathbb{S}(y_3) := \{(y_1, y_2): |y_k| \leq \beta(y_3), \quad k = 1, 2\},$$

where the *stability boundary*  $\beta(y_3)$  is not too small. Since the spectral radius of  $\Delta t \mathbf{J}_1$  and  $\Delta t \mathbf{J}_2$  is much smaller than that of  $\Delta t \mathbf{J}_3$ , we would like stability in all regions  $\mathbb{S}(y_3)$ ,  $|y_3| \leq \infty$ . The corresponding timestep condition is given by

$$(2.12) \quad \Delta t \leq \frac{\beta}{\max\{\rho(\mathbf{J}_1), \rho(\mathbf{J}_2)\}}, \quad \beta := \min_{y_3} \beta(y_3).$$

Let us consider the stability of the factorized versions of (2.3) and (2.4'). It is easily verified that their stability functions respectively take the form

$$(2.13) \quad \mathbf{R}_1(z_1, z_2, z_3) := 1 + \frac{z}{(1 - \kappa z_1)(1 - \kappa z_2)(1 - \kappa z_3)} + \frac{\frac{1}{2}(1 - 2\kappa)z^2}{(1 - \kappa z_1)^2(1 - \kappa z_2)^2(1 - \kappa z_3)^2},$$

$$(2.14) \quad \mathbf{R}_2(z_1, z_2, z_3) := 1 + \frac{2z}{(1 - \kappa z_1)(1 - \kappa z_2)(1 - \kappa z_3)} + \frac{\frac{1}{2}z^2 - z}{(1 - \kappa z_1)^2(1 - \kappa z_2)^2(1 - \kappa z_3)^2}$$

and that  $|R_1(0,0,iy_3)| < 1$ ,  $|R_2(0,0,iy_3)| < 1$  for  $y_3 \neq 0$ . Hence, we have a nonzero stability boundary  $\beta$ . However, a numerical calculation reveals that  $\beta$  is quite small (less than 1/10). Hence, the factorized versions of (2.3), (2.4) and (2.4') are of no use in shallow water applications.

### 3. Approximate factorization iteration

The quite poor stability properties of the factorized Rosenbrock methods can be explained by observing that the vector  $(T^{-1} \otimes I)\mathbf{V}$  defined by the factorized-Rosenbrock method (2.7) is too far away from the vector  $\mathbf{K} = (T^{-1} \otimes I)\mathbf{U}$  defined by the Rosenbrock method (2.5). In this section, we improve the stability by really solving the implicit relations in the underlying Rosenbrock method by an iteration process. We shall also study the iterative solution of the related implicit methods

$$(3.1) \quad \mathbf{y}_{n+1} = \mathbf{y}_n + \Delta t(\mathbf{b}^T \otimes I)\mathbf{F}(\mathbf{X}), \quad \mathbf{X} - \Delta t(\mathbf{A} \otimes I)\mathbf{F}(\mathbf{X}) = \mathbf{e} \otimes \mathbf{y}_n,$$

where  $\mathbf{A}$  is a lower triangular matrix. In [4, p.97] these methods are called DIRK methods (diagonally implicit Runge-Kutta methods). Like Rosenbrock methods, DIRK methods are triangularly implicit (in spite of the terminology 'diagonally implicit' now commonly accepted in the literature).

An advantage of the iterative approach is that we can rely on the stability of the underlying integration method. Thus, by choosing an A-stable integration method, we only have to deal with the region of *convergence* of the iteration method. The iteration processes considered below are based on the approximate factorization technique used in the preceding section and lead to acceptably large convergence regions.

#### 3.1. Iterative solution of the Rosenbrock equations

We consider two iterative approximate factorization approaches for actually solving the implicit Rosenbrock relations. The first approach solves the components  $\mathbf{u}_i$  from (2.6) one by one by *repeated* application of a linear system solver, the second approach solves all components  $\mathbf{k}_i$  from (2.5) *simultaneously* by a nonlinear system solver. We shall refer to these iteration methods as repeated and simultaneous approximate factorization iteration of the Rosenbrock method, briefly, the RAF-Rosenbrock and SAF-Rosenbrock processes, respectively.

**3.1.1. The RAF-Rosenbrock process.** The  $s$  linear systems in (2.6) have the form

$$(3.2) \quad (\mathbf{I} - \kappa_i \Delta t \mathbf{J})\mathbf{u}_i = \mathbf{g}_i, \quad i = 1, \dots, s,$$

$$\mathbf{g}_i := (\mathbf{e}_i^T \otimes I) \left( \Delta t (\mathbf{D} \otimes I)\mathbf{F}(\mathbf{e} \otimes \mathbf{y}_n + (\mathbf{L}T^{-1} \otimes I)\mathbf{U}) - (\mathbf{D}\mathbf{S} \otimes I)\mathbf{U} \right),$$

where  $\kappa_i$  is the  $i$ th diagonal entry of  $\mathbf{T}$ . Since  $\mathbf{L}$  and  $\mathbf{S}$  are strictly lower triangular, these  $s$  systems can be solved successively. We solve the  $i$ th linear system by the linear solver

$$(3.3) \quad (\mathbf{e}_i^T \otimes I)\Pi \left( \mathbf{u}_i^{(j)} - \mathbf{u}_i^{(j-1)} \right) = \mathbf{g}_i - (\mathbf{I} - \kappa_i \Delta t \mathbf{J})\mathbf{u}_i^{(j-1)}, \quad j = 1, 2, \dots, m, \quad i = 1, \dots, s,$$



where  $\Pi$  is defined in (2.8). In this RAF-Rosenbrock process the initial iterate  $\mathbf{u}_i^{(0)}$  should be provided by some predictor formula and the number of iterations  $m$  is assumed to be determined by some iteration strategy such that  $\mathbf{u}_i^{(m)}$  may be considered as the solution  $\mathbf{u}_i$  of (3.2).

If the iterates  $\mathbf{u}_i^{(j)}$  converge, then they can only converge to this solution  $\mathbf{u}_i$ . Each iteration in (3.3) requires the solution of 3 linear systems with system matrices  $I - \kappa_i \Delta t J_k$ ,  $k = 1, 2, 3$ , each of order  $N$ . Note that the three LU-decompositions of these system matrices can be done in parallel. These LU-decompositions and the corresponding forward-backward substitutions are relatively cheap, because the matrices  $J_k$  each correspond with a one-dimensional differential operator.

The convergence is determined by the error recursion satisfied by the iteration error  $\boldsymbol{\varepsilon}^{(j)}$ :

$$(3.4) \quad \boldsymbol{\varepsilon}^{(j)} := \mathbf{u}_i^{(j)} - \mathbf{u}_i$$

$$(3.5) \quad \boldsymbol{\varepsilon}^{(j)} = Z_1 \boldsymbol{\varepsilon}^{(j-1)}, \quad Z_1 := I - \Pi^{-1}(I - D \otimes \Delta t J), \quad j = 1, 2, \dots, m.$$

Before analysing the matrix  $Z_1$ , we first derive the error recursion for the other iterative approaches.

**3.1.2. The SAF-Rosenbrock process.** Instead of solving the linear systems in the Rosenbrock method (2.6) successively for the components  $\mathbf{u}_i$  of  $\mathbf{U}$ , we may iterate them simultaneously. Since in such an approach it is more convenient to go back to the untransformed method (2.5), we shall solve the components  $\mathbf{k}_i$  of  $\mathbf{K}$  simultaneously from (2.5). Consider the SAF-Rosenbrock process

$$(3.6) \quad \Pi (\mathbf{K}^{(j)} - \mathbf{K}^{(j-1)}) = - \left( (I - T \otimes \Delta t J) \mathbf{K}^{(j-1)} - \Delta t \mathbf{F}(\mathbf{e} \otimes \mathbf{y}_n + (L \otimes I) \mathbf{K}^{(j-1)}) \right), \quad j = 1, 2, \dots, m.$$

Note that this method is a *nonlinear* system solver.

Evidently, if the iterates  $\mathbf{K}^{(j)}$  converge and if (2.5) has a unique solution  $\mathbf{K}$ , then they can only converge to this solution  $\mathbf{K}$ . Each SAF-Rosenbrock iteration requires the solution of 3 linear systems with system matrices  $I - D \otimes \Delta t J_k$ ,  $k = 1, 2, 3$ , each of order  $sN$ . The  $3s$  LU-decompositions and the  $s$  forward-backward substitutions corresponding with each matrix  $I - D \otimes \Delta t J_k$  can be done in parallel. Again, the LU-decompositions and the forward-backward substitutions are relatively cheap, because  $J_k$  corresponds with a one-dimensional differential operator. A drawback is the matrix-vector multiplication in the righthand side of (3.6). Note that applying the SAF-Rosenbrock iteration process to (2.6) instead of (2.5) does not avoid such a matrix-vector multiplication.

Let us consider the iteration error  $\boldsymbol{\varepsilon}^{(j)} := \mathbf{K}^{(j)} - \mathbf{K}$ . From (2.5) and (3.6) it follows that

$$(3.7) \quad \boldsymbol{\varepsilon}^{(j)} = Z_2 \boldsymbol{\varepsilon}^{(j-1)} + \Delta t \Pi^{-1} \mathbf{G}(\boldsymbol{\varepsilon}^{(j-1)}), \quad Z_2 := I - \Pi^{-1}(I - (T+L) \otimes \Delta t J), \quad j = 1, 2, \dots, m,$$

$$\mathbf{G}(\boldsymbol{\varepsilon}) := \mathbf{F}(\mathbf{e} \otimes \mathbf{y}_n + (L \otimes I)(\mathbf{K} + \boldsymbol{\varepsilon})) - \mathbf{F}(\mathbf{e} \otimes \mathbf{y}_n + (L \otimes I)\mathbf{K}) - (L \otimes J)\boldsymbol{\varepsilon}.$$

Since  $\mathbf{G}(\boldsymbol{\varepsilon})$  has a small Lipschitz constant in the neighbourhood of the origin, the error recursion (3.7) essentially behaves as the linearized recursion

$$(3.8) \quad \boldsymbol{\varepsilon}^{(j)} \approx Z_2 \boldsymbol{\varepsilon}^{(j-1)}, \quad Z_2 := I - \Pi^{-1}(I - (T+L) \otimes \Delta t J), \quad j = 1, 2, \dots, m.$$

### 3.2. Iterative solution of DIRK equations

As for Rosenbrock methods, we may consider repeated and simultaneous approximate factorization iteration of the DIRK method (3.1). These processes are respectively given by

$$(3.9) \quad (\mathbf{e}_i^T \otimes \mathbf{I}) \Pi (\mathbf{x}_i^{(j)} - \mathbf{x}_i^{(j-1)}) = \mathbf{g}_i, \quad j = 1, 2, \dots, m, \quad i = 1, \dots, s,$$

$$\mathbf{g}_i := (\mathbf{e}_i^T \otimes \mathbf{I}) \left( (\mathbf{e} \otimes \mathbf{I}) \mathbf{y}_n - \mathbf{X}^{(j-1)} + \Delta t (\mathbf{A} \otimes \mathbf{I}) \mathbf{F}(\mathbf{X}^{(j-1)}) \right),$$

and

$$(3.10) \quad \Pi (\mathbf{X}^{(j)} - \mathbf{X}^{(j-1)}) = - \left( \mathbf{X}^{(j-1)} - \Delta t (\mathbf{A} \otimes \mathbf{I}) \mathbf{F}(\mathbf{X}^{(j-1)}) - (\mathbf{e} \otimes \mathbf{I}) \mathbf{y}_n \right), \quad j = 1, 2, \dots, m,$$

where  $\Pi$  is again defined by (2.8) with  $\mathbf{D} := \text{diag}(\mathbf{A})$ . They will be referred to as the RAF-DIRK and SAF-DIRK processes. A comparison with (3.3) and (3.6) shows that we have the same iteration costs except for the Jacobian multiplication.

Defining the iteration error  $\boldsymbol{\varepsilon}^{(j)} := \mathbf{X}^{(j)} - \mathbf{X}$ , we can write down the linearized error recursions. We find that the linearized error recursions associated with the RAF-DIRK method (3.9) and the SAF-DIRK method (3.10) are respectively given by

$$(3.11) \quad \boldsymbol{\varepsilon}^{(j)} = \mathbf{Z}_3 \boldsymbol{\varepsilon}^{(j-1)}, \quad \mathbf{Z}_3 := \mathbf{Z}_1, \quad j = 1, 2, \dots, m,$$

$$(3.12) \quad \boldsymbol{\varepsilon}^{(j)} \approx \mathbf{Z}_4 \boldsymbol{\varepsilon}^{(j-1)}, \quad \mathbf{Z}_4 := \mathbf{I} - \Pi^{-1}(\mathbf{I} - \mathbf{A} \otimes \Delta t \mathbf{J}), \quad j = 1, 2, \dots, m.$$

### 3.3. Convergence

The convergence of iterated Rosenbrock methods (3.2) and (3.6), and of the iterated DIRK methods (3.9) and (3.10) is determined by the amplification matrices

$$\begin{aligned} \mathbf{Z}_1 &:= \mathbf{I} - \Pi^{-1}(\mathbf{I} - \mathbf{D} \otimes \Delta t \mathbf{J}), & \mathbf{Z}_2 &:= \mathbf{I} - \Pi^{-1}(\mathbf{I} - (\mathbf{T} + \mathbf{L}) \otimes \Delta t \mathbf{J}), \\ \mathbf{Z}_3 &= \mathbf{Z}_1, & \mathbf{Z}_4 &:= \mathbf{I} - \Pi^{-1}(\mathbf{I} - \mathbf{A} \otimes \Delta t \mathbf{J}), \end{aligned}$$

occurring in the error recursions (3.5), (3.8), (3.11) and (3.12), respectively. They only differ by the matrix in front of  $\Delta t \mathbf{J}$  (we recall that  $\mathbf{D} = \text{diag}(\mathbf{T}) = \text{diag}(\mathbf{A})$ ). In the following subsections we respectively discuss the region of convergence where  $\rho(\mathbf{Z}_r) < 1$ , the rate of convergence of the nonstiff iteration error components, and the stability of the iterated methods.

**3.3.1. The region of convergence.** The matrices  $\mathbf{Z}_r$  are lower triangular block matrices with the same diagonal blocks

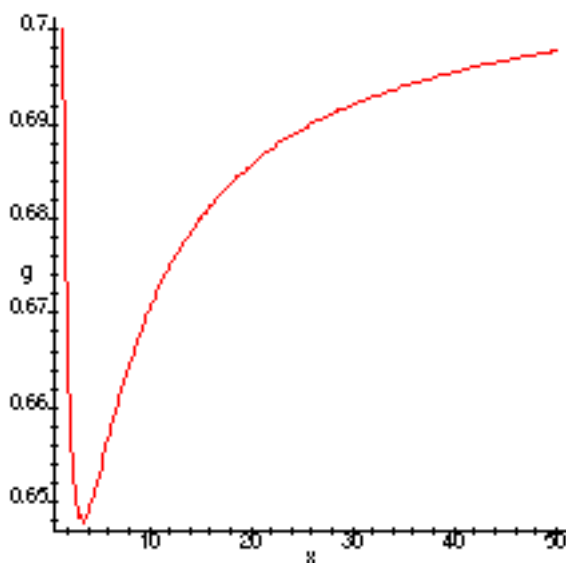
$$\mathbf{I} - (\mathbf{I} - \kappa_j \Delta t \mathbf{J}_1)^{-1} (\mathbf{I} - \kappa_j \Delta t \mathbf{J}_2)^{-1} (\mathbf{I} - \kappa_j \Delta t \mathbf{J}_3)^{-1} (\mathbf{I} - \kappa_j \Delta t \mathbf{J}), \quad j = 1, \dots, s,$$

for all  $r$ . Here, the  $\kappa_j$  denote the diagonal entries of  $\mathbf{D}$ . Hence, the eigenvalues of the matrices  $\mathbf{Z}_r$  are identical. They act as amplification factors for the eigenvalue components of the iteration error and are given by

$$(3.13) \quad \alpha_j = C(\kappa_j z_1, \kappa_j z_2, \kappa_j z_3), \quad C(x_1, x_2, x_3) := 1 - \frac{1 - x_1 - x_2 - x_3}{(1 - x_1)(1 - x_2)(1 - x_3)},$$

where  $j = 1, \dots, s$  and where  $z_k$  runs through the eigenvalues of  $\Delta t J_k$ . Evidently, we have convergence if  $|\alpha_j| < 1, j = 1, \dots, s$ . We consider the most critical case where the eigenvalues of  $J_k$  are purely imaginary, that is we consider the values of  $|\alpha_j| = |C(i\kappa_j y_1, i\kappa_j y_2, i\kappa_j y_3)|$ . Recalling that the spectral radius of  $\Delta t J_1$  and  $\Delta t J_2$  is much smaller than that of  $\Delta t J_3$ , we are interested in convergence regions of the form (cf. (2.11))

$$(3.14) \quad \mathbb{C}(y_3) := \{(y_1, y_2): |y_k| \leq \frac{\gamma(y_3)}{\rho(D)}, k = 1, 2\}, \quad |y_3| \leq \infty.$$



**Figure 3.1.** The function  $g(x)$  defined by (3.15).

**Theorem 3.1.** Let the function  $g(x)$  be defined by the relation

$$(3.15) \quad 4xg^3 + 2(x^2 - 1)g^2 - x^2 - 1 = 0.$$

Then, the convergence boundary  $\gamma(y_3)$  in (3.14) is given by

$$(3.16) \quad \gamma(y_3) = \rho(D) \min_j \frac{g(\kappa_j |y_3|)}{\kappa_j}$$

and the minimal value of  $\gamma(y_3)$  is given by the positive root of the equation  $4\gamma^4(\gamma^2 + 1) = 1$ .

**Proof.** We verified that for given values of  $y_3$ ,  $|C(i\kappa_j y_1, i\kappa_j y_2, i\kappa_j y_3)|$  increases most rapidly along the line  $y_1 = y_2$ , so that we may restrict our considerations to the values of

$$|\alpha_j| = |C(i\kappa_j y_1, i\kappa_j y_1, i\kappa_j y_3)| = \kappa_j^2 |y_1| \left( \frac{y_1^2(\kappa_j^2 y_3^2 + 1) + 4y_1 y_3 + 4y_3^2}{(1 + \kappa_j^2 y_1^2)^2 (1 + \kappa_j^2 y_3^2)} \right)^{1/2}.$$

If we set  $|\alpha_j| = 1$ ,  $x = \kappa_j y_3$  and  $y = \kappa_j y_1$ , then we find the relation

$$4xy^3 + 2(x^2 - 1)y^2 - x^2 - 1 = 0.$$

This relation determines a real-valued function  $y = g(x)$ . Hence, for given values of  $\kappa_j$  and  $y_3$ , i.e. of  $x$ , we have  $|\alpha_j| \leq 1$  provided that both  $\kappa_j |y_1|$  and  $\kappa_j |y_2|$  are bounded by  $g(\kappa_j |y_3|)$ . This proves (3.16).

In order to find the minimal value of  $\gamma(y_3)$ , we look at the plot of the function  $g(x)$  (see Figure 3.1). Let  $x_1(y)$  and  $x_2(y)$  denote the two solutions of the equation  $4xy^3 + 2(x^2 - 1)y^2 - x^2 - 1 = 0$ . Then, the minimal value of  $g(x)$  is determined by the relation  $x_1(y) = x_2(y)$ . This leads to the equation  $4y^4(y^2 + 1) = 1$  whose only positive root determines the minimal value of  $\gamma(y_3)$ . ♦

Since the positive root of the equation  $4\gamma^4(\gamma^2 + 1) = 1$  is given by  $\gamma = 0.647\dots$  we derive from this theorem the following convergence condition:

**Theorem 3.2.** Let  $\lambda(J_k)$ ,  $k = 1, 2, 3$ , be purely imaginary. Then, a sufficient condition for convergence of the iterated Rosenbrock methods (3.3) and (3.6), and of the iterated DIRK methods (3.9) and (3.10) is given by

$$\Delta t \leq \frac{\gamma}{\rho(D) \max\{\rho(J_1), \rho(J_2)\}}, \quad \gamma = 0.647\dots \quad \blacklozenge$$

**3.3.2. The rate of convergence of the nonstiff error components.** The rate of convergence of the *nonstiff* error components can be studied by the behaviour of the nonstiff amplification factors, that is, the eigenvalues of  $Z_r$  corresponding with small values of  $\Delta t \lambda(J_k)$ . From (3.13) it can be deduced that

$$\alpha_j = \kappa_j^2(z_1 z_2 + z_1 z_3 + z_2 z_3) + O((\Delta t)^3), \quad z_k = \Delta t \lambda(J_k), \quad j = 1, \dots, s.$$

Hence, after  $m$  iterations the amplification factors behave as  $O((\Delta t)^{2m})$  for all  $m$  and irrespective the value of  $r$ . However, this is not true for the amplification matrices  $Z_r^m$ .

**Theorem 3.3.** The amplification matrices  $Z_r$  satisfy the relations

$$\begin{aligned} r = 1, 3: \quad & Z_r^m = O((\Delta t)^{2m}) \quad \text{for all } m, \\ r = 2, 4: \quad & \begin{cases} Z_r^m = O((\Delta t)^m) & \text{for } m \leq s-1 \\ Z_r^m = O((\Delta t)^{2m+1-s}) & \text{for } m \geq s \end{cases} \end{aligned}$$

**Proof.** The relation  $Z_1 = Z_3 = O((\Delta t)^2)$  immediately follows from the definition of  $Z_1$  and  $Z_3$  in (3.5) and (3.11). For  $Z_2$  and  $Z_4$  it follows from (3.8) and (3.12) that

$$\begin{aligned} Z_2 &= I - (I + D \otimes \Delta t J)(I - (T+L) \otimes \Delta t J) + O((\Delta t)^2) = (T + L - D) \otimes \Delta t J + O((\Delta t)^2). \\ Z_4 &= I - (I + D \otimes \Delta t J)(I - A \otimes \Delta t J) + O((\Delta t)^2) = (A - D) \otimes \Delta t J + O((\Delta t)^2). \end{aligned}$$

Hence, we certainly have  $Z_r^m = O((\Delta t)^m)$  for  $r = 2, 4$ . However, writing  $Z_r = A_r + B_r$  with  $A_2 := (T + L - D) \otimes \Delta t J$  and  $A_4 := (A - D) \otimes \Delta t J$ , and observing that  $A_2$  and  $A_4$  are strictly lower block triangular, so that  $A_2^j$  and  $A_4^j$  vanish for  $j \geq s$ , we obtain for  $m \geq s$

$$Z_r^m = \binom{m}{m-s+1} A_r^{s-1} B_r^{m-s+1} + \dots + \binom{m}{m} B_r^m, \quad r = 2, 4.$$

Since  $A_r = O(\Delta t)$  and  $B_r = O((\Delta t)^2)$ , we find that

$$Z_r^m = O((\Delta t)^{2m-s+1}), \quad r = 2, 4. \quad \blacklozenge$$

From this theorem it follows that in all four approaches the nonstiff error components are rapidly removed from the iteration error. However, we may expect that the RAF processes (3.5) and (3.9) damp these nonstiff components stronger than the SAF processes (3.8) and (3.12).

**3.3.3. The region of stability.** Evidently, if the iteration process converges, then the stability of the iterated method is determined by the stability of the underlying integration method. Hence, with respect to the stability test equation, the stability region of the iterated method converges to the intersection of the convergence region and the stability region of the integration method, that is, to

$$\mathbb{S} := \mathbb{S}_0 \cap \mathbb{C}, \quad \mathbb{C} := \bigcap_{y_3} \mathbb{C}(y_3), \quad |y_3| \leq \infty,$$

where  $\mathbb{S}_0$  is the stability region of the integration method and  $\mathbb{C}(y_3)$  is defined by (3.14). For A-stable integration methods, the stability region  $\mathbb{S}$  equals the convergence region  $\mathbb{C}$ , so that the stability condition is given by the stepsize condition in Theorem 3.2. Thus, for iterated, A-stable integration methods we may define the stability boundary  $\beta := \gamma \rho^{-1}(D)$ .

For example, if the Rosenbrock methods (2.3) and (2.4) are iterated using the iteration matrix  $\Pi$ , then we find in both cases the stability boundary  $\beta \approx 2.20$ . If we choose  $\kappa_1 = \kappa_2 = \frac{1}{4}$  in (2.1), then (2.1) is still A-stable with a slightly greater stability boundary  $\beta \approx 2.59$ .

## 4. Explicit treatment of the horizontal terms

The modest values of the stability boundary  $\beta$  raises the question whether it is necessary to treat the horizontal terms implicitly. Afterall, when applying the standard, explicit, fourth-order Runge-Kutta method, we have an imaginary stability boundary of comparable size, viz.  $\beta = 2\sqrt{2}$ .

### 4.1. Fully explicit treatment of the horizontal terms

We once again consider the iteration methods (3.3), (3.6), (3.9) and (3.10), but we replace the iteration matrix  $\Pi$  by the matrix  $\Pi_3 := I - D \otimes \Delta t J_3$ . As a consequence, the amplification factors are now given by (cf. (3.13))

$$(4.1) \quad \alpha_j = C_3(\kappa_j z_1, \kappa_j z_2, \kappa_j z_3), \quad C_3(x_1, x_2, x_3) := 1 - \frac{1 - x_1 - x_2 - x_3}{1 - x_3}.$$

Hence,  $|C_3(i\kappa_j y_1, i\kappa_j y_2, i\kappa_j y_3)|^2 = \kappa_j^2 (y_1 + y_2)^2 (1 + \kappa_j^2 y_3^2)^{-1}$ , so that we have convergence in regions of the form (3.14) with

$$(4.2) \quad \gamma(y_3) = \frac{1}{2} \sqrt{1 + y_3^2 \rho^2(\mathbf{D})}, \quad |y_3| \leq \infty.$$

Thus, the convergence boundary  $\gamma$  in the timestep condition in Theorem 3.2 changes from  $\gamma \approx 0.647$  to  $\gamma = \frac{1}{2}$ . Since the iterations in (3.3), (3.6), (3.9) and (3.10) with  $\Pi$  replaced by  $\Pi_3$  are cheaper, the modest reduction of the convergence boundary seems to be a small price. Moreover, the convergence boundary (4.2) quickly increases with  $|y_3|$ , whereas (3.16) approaches a constant value  $\gamma(\infty) \approx 0.7$  (see Figure 3.1). Hence, the *stiff* error components will be more strongly damped when using the matrix  $\Pi_3$ . On the other hand, if we look at the behaviour of the amplification factor  $\alpha_j$  as  $\Delta t \rightarrow 0$ , then we find  $\alpha_j^m = O((\Delta t)^m)$ , so that the *nonstiff* iteration error components are expected to require more iterations to be removed from the iteration error (see Theorem 3.3).

Thus, the iteration processes with  $\Pi$  and  $\Pi_3$  both have advantages. This suggests a combination of the two iteration methods, for example, by iterating successively with iteration matrices  $\Pi_3, \Pi, \Pi_3, \Pi, \dots$  (the  $\Pi_3\Pi$  process) or with  $\Pi_3, \Pi, \Pi, \Pi_3, \Pi, \Pi, \dots$  (the  $\Pi_3\Pi^2$  process). Evidently, in such combined processes, the averaged iteration costs are still slightly higher than in the  $\Pi_3$  process and slightly lower than in the  $\Pi$  process.

**4.1.1. Convergence boundaries of the  $\Pi_3\Pi$  and  $\Pi_3\Pi^2$  processes.** Let us first consider the  $\Pi_3\Pi$  process more closely. After each two iterations, the corresponding amplification matrix  $Z_r$  in the linearized error recursion has diagonal blocks of the form

$$\left( \mathbf{I} - ((\mathbf{e}_j^T \otimes \mathbf{I}) \Pi_3)^{-1} (\mathbf{I} - \kappa_j \Delta t \mathbf{J}) \right) \left( \mathbf{I} - ((\mathbf{e}_j^T \otimes \mathbf{I}) \Pi)^{-1} (\mathbf{I} - \kappa_j \Delta t \mathbf{J}) \right), \quad j = 1, 2, \dots, s,$$

irrespective the value of  $r$ . Hence, the amplification factors for two iterations are given by

$$(4.3) \quad \alpha_j^2 = C_3(\kappa_j z_1, \kappa_j z_2, \kappa_j z_3) C(\kappa_j z_1, \kappa_j z_2, \kappa_j z_3),$$

where  $C$  and  $C_3$  are defined in (3.13) and (4.1). Again, we restrict  $z_k$  to imaginary values  $iy_k$  and again it turns out that for given values of  $y_3$  the amplification factor  $|\alpha_j^2|$  increases most rapidly along the line  $y_1 = y_2$ . Along the lines of the proof of Theorem 3.1 it can be shown that the convergence boundary  $\gamma(y_3)$  in (3.14) is given by (3.16) where the function  $g(x)$  is defined by

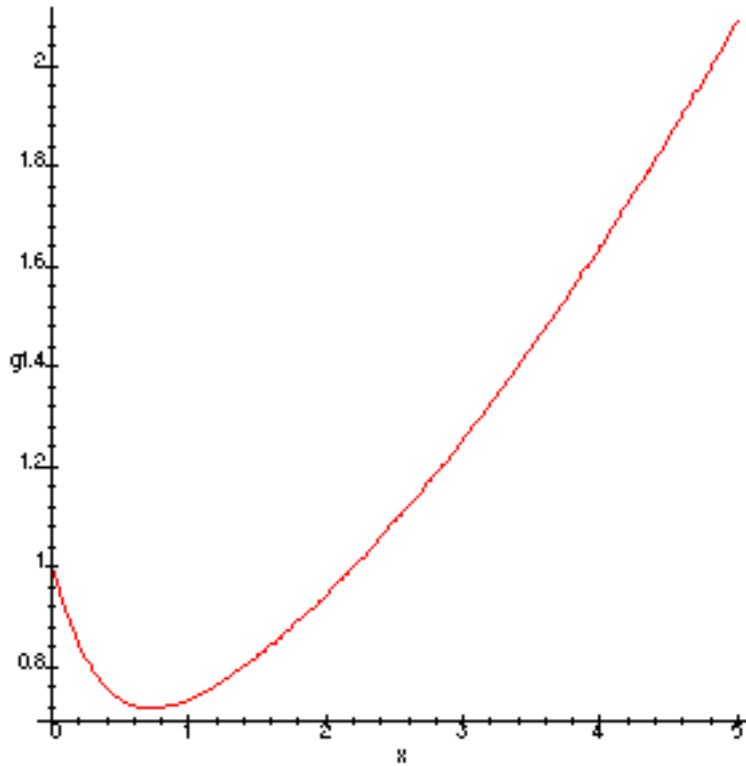
$$(4.4) \quad 4(x^2 + 1)g^6 + 16xg^5 - (x^4 - 14x^2 + 1)g^4 - 2(x^2 + 1)^2g^2 - (x^2 + 1)^2 = 0.$$

Figure 4.1 presents a plot of  $g(x)$ . It is easily seen that  $g(x) \approx \frac{1}{2}x$  as  $x \rightarrow \infty$ , so that the combined iteration process has the property that the convergence boundary  $\gamma(y_3)$  increases with  $y_3$  to infinity (compare (4.2)). In order to find the minimal value  $\gamma$  of  $g(x)$ , we write (4.4) as  $F(g, x) = 0$ . Then,  $\gamma$  is determined by (4.4) and by the equation

$$(4.5) \quad \frac{\partial F(g,x)}{\partial x} = 4(2xg^6 + 4g^5 - x(x^2 - 7)g^4 - 2x(x^2 + 1)g^2 - x(x^2 + 1)) = 0.$$

Solving (4.4) and (4.5) yields the minimal value  $\gamma \approx 0.72$ . This is even slightly greater than the value  $\gamma \approx 0.65$  given in Theorem 3.2 for the  $\Pi$  process.

Similarly, we find for the  $\Pi_3\Pi^2$  process a still greater value  $\gamma \approx 0.75$ .

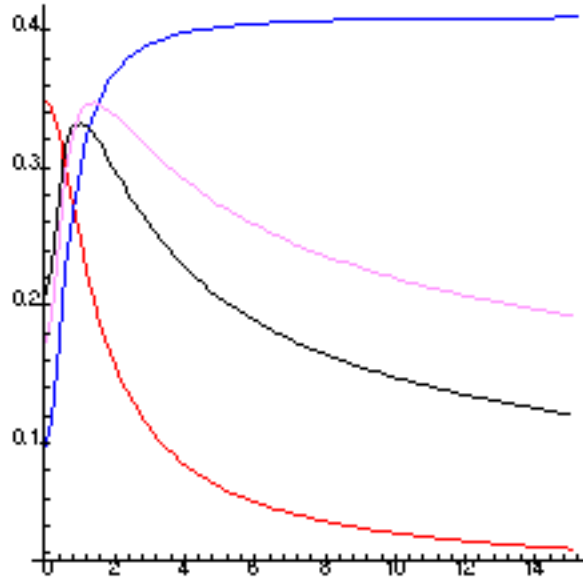


**Figure 4.1.** The function  $g(x)$  defined by (4.4).

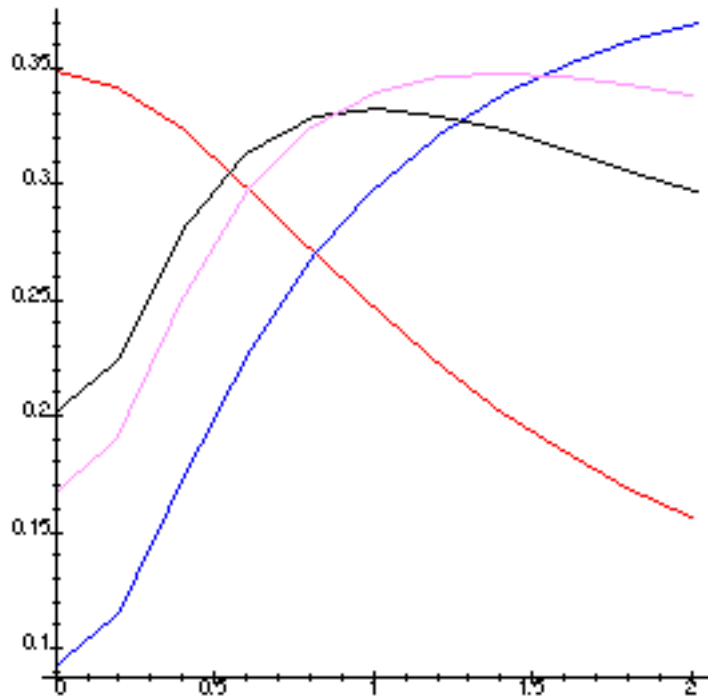
**4.1.2. The rate of convergence.** The size of the convergence region is not the only issue to be considered. The *rate* of convergence, that is, the magnitude of the amplification factors, is equally important. We consider the behaviour of the amplification factors for small values of  $\Delta t$   $\lambda(J_k)$  and their overall behaviour. For the  $\Pi_3\Pi$  process it can be deduced from (4.3) that for  $\Delta t \rightarrow 0$

$$\alpha_j^2 = \kappa_j^3(z_1z_2 + z_1z_3 + z_2z_3)(z_1 + z_2) + O((\Delta t)^4), \quad z_k = \Delta t \lambda(J_k), \quad j = 1, \dots, s,$$

so that after an even number of  $m$  iterations the amplification factors behave as  $O((\Delta t)^{3m/2})$ , irrespective the value of  $r$ . Similarly, we have for the  $\Pi_3\Pi^2$  process,  $\alpha_j^m = O((\Delta t)^{5m/3})$ , where  $m$  is assumed to be a multiple of 3. Recalling that the  $\Pi$  and  $\Pi_3$  processes yield after  $m$  iterations amplification by  $\alpha_j^m = O((\Delta t)^{2m})$  and  $\alpha_j^m = O((\Delta t)^m)$ , respectively, we see that with respect to damping of the nonstiff error components, the order of preference is  $\Pi$ ,  $\Pi_3\Pi^2$ ,  $\Pi_3\Pi$ , and  $\Pi_3$ .



**Figure 4.2a.** The functions  $\alpha_{\text{aver}}(x)$  for the  $\Pi_3$ ,  $\Pi_3\Pi$ ,  $\Pi_3\Pi^2$  and  $\Pi$  processes.



**Figure 4.2b.** Details near the origin of  $\alpha_{\text{aver}}(x)$  for the  $\Pi_3$ ,  $\Pi_3\Pi$ ,  $\Pi_3\Pi^2$  and  $\Pi$  processes.

As to the overall convergence, we compare the averaged amplification per iteration in the region  $\mathbb{C}(y_3)$  defined in (3.14), that is, the averaged values of the functions  $|C_3|$ ,  $|C|$ ,  $|C_3C|^{1/2}$  and  $|C_3C^2|^{1/3}$  in  $\mathbb{C}(y_3)$ , with  $C$  and  $C_3$  defined by (3.13) and (4.1). Writing  $x := \kappa_j y_3$  and denoting these averaged values by  $\alpha_{\text{aver}}(x)$ , we compute  $\alpha_{\text{aver}}(x)$  in the square  $\{|\kappa_j y_1| \leq \gamma, |\kappa_j y_2| \leq \gamma\}$ , where  $\gamma \approx 0.50, 0.65, 0.72, 0.75$  in the  $\Pi_3, \Pi, \Pi_3\Pi$  and  $\Pi_3\Pi^2$  processes, respectively. In Figure 4.2a, the function



$\alpha_{\text{aver}}(x)$  is plotted for  $0 \leq x \leq 15$ . Details in the interval  $0 \leq x \leq 2$  are given by Figure 4.2b. In these plots, the left and right end point values of the graphs respectively decrease and increase for the  $\Pi_3$ ,  $\Pi_3\Pi$ ,  $\Pi_3\Pi^2$  and  $\Pi$  processes. Thus, the  $\Pi_3$  and  $\Pi$  processes clearly are extreme cases, the  $\Pi_3$  process being superior for larger values of  $x$ , the  $\Pi$  process for small values of  $x$ . The  $\Pi_3\Pi$  and  $\Pi_3\Pi^2$  processes are good compromises. Moreover, they possess substantially larger convergence boundaries respectively allowing about 45% and 10% larger stepsizes than the  $\Pi_3$  and  $\Pi$  processes.

## 4.2. Partially explicit treatment of the horizontal terms

Instead of inserting iterations that are fully explicit with respect to the horizontal terms, one may consider the use of iterations with iteration matrices  $\Pi_{13} := (I - D \otimes \Delta t J_1)(I - D \otimes \Delta t J_3)$  and  $\Pi_{23} := (I - D \otimes \Delta t J_2)(I - D \otimes \Delta t J_3)$ . Since in shallow water applications, the two horizontal directions introduce a comparable degree of stiffness, we shall alternate these two iteration matrices. It should be remarked that we already used the iteration matrices  $\Pi_{13}$  and  $\Pi_{23}$  in [6] in another type of iteration methods for use in shallow water computations.

**4.2.1. Convergence boundary of the  $\Pi_{13}\Pi_{23}$  process.** After each two iterations, the corresponding amplification matrix  $Z_r$  in the linearized error recursion has diagonal blocks of the form

$$(I - ((e_j^T \otimes I)\Pi_{13})^{-1}(I - \kappa_j \Delta t J)) (I - ((e_j^T \otimes I)\Pi_{23})^{-1}(I - \kappa_j \Delta t J)), \quad j = 1, 2, \dots, s,$$

irrespective the value of  $r$ . Defining the functions

$$(4.6) \quad C_{k3}(x_1, x_2, x_3) := 1 - \frac{1 - x_1 - x_2 - x_3}{(1 - x_k)(1 - x_3)}, \quad k = 1, 2,$$

the amplification factors for two iterations are given by

$$(4.7) \quad \alpha_j^2 = C_{13}(\kappa_j z_1, \kappa_j z_2, \kappa_j z_3) C_{23}(\kappa_j z_1, \kappa_j z_2, \kappa_j z_3).$$

Restricting  $z_k$  to imaginary values  $iy_k$ , we obtain

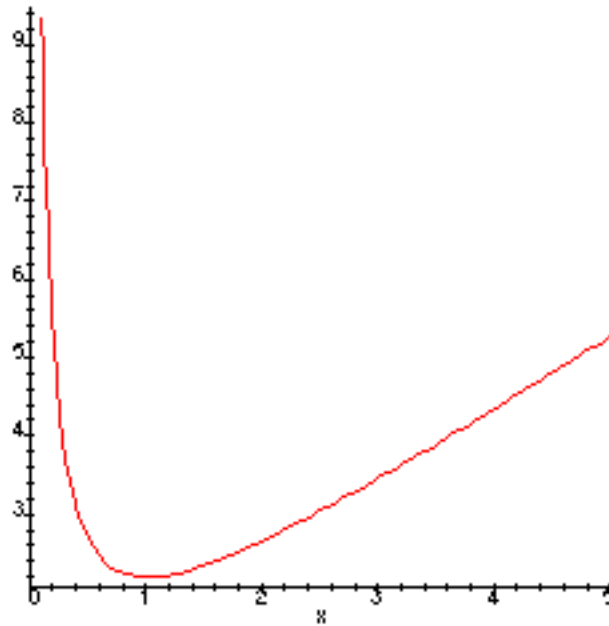
$$(4.8) \quad |\alpha_j^2| = \left( \frac{(\kappa_j^2 y_1^2 + \kappa_j^4 y_3^2 y_2^2)(\kappa_j^2 y_2^2 + \kappa_j^4 y_3^2 y_1^2)}{(1 + \kappa_j^2 y_1^2)(1 + \kappa_j^2 y_2^2)(1 + \kappa_j^2 y_3^2)^2} \right)^{1/2}.$$

We verified that for given values of  $\kappa_j^2 y_3^2 \geq 0.01$  the amplification factor  $|\alpha_j^2|$  increases most rapidly along the lines  $y_1 = 0$  and  $y_2 = 0$ . For  $\kappa_j^2 y_3^2 < 0.01$ , we found numerically that  $\gamma(y_3) \geq 10$ . Hence, proceeding as in the proof of Theorem 3.1, it can be shown that for  $\kappa_j |y_3| \geq 0.1$  the convergence boundary  $\gamma(y_3)$  in (3.14) is given by (3.16) where the function  $g(x)$  is defined by

$$(4.9) \quad x^2 g^4 - (x^2 + 1)^2 g^2 - (x^2 + 1)^2 = 0, \quad x \geq 0.1.$$

A comparison of the plot of this function in Figure 4.3 with Figure 4.1 shows that the convergence boundary  $\gamma(y_3)$  is considerably larger than that of the  $\Pi_3\Pi$  process. This also results in a much larger value for  $\gamma$ . To see this, we again write the implicit equation for  $g$  as  $F(g, x) = 0$  and solve  $\gamma$  from the

equations  $F(\gamma, x) = 0$  and  $\partial F(\gamma, x)/\partial x = 0$  to obtain  $\gamma = \sqrt{2 + 2\sqrt{2}} \approx 2.19$  which compares favourably with the value  $\gamma \approx 0.65$  of Theorem 3.2.

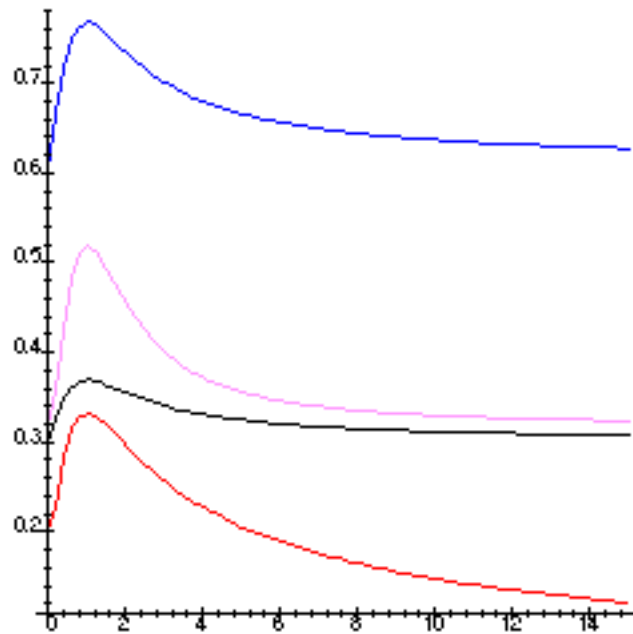


**Figure 4.3.** The function  $g(x)$  defined by (4.9).

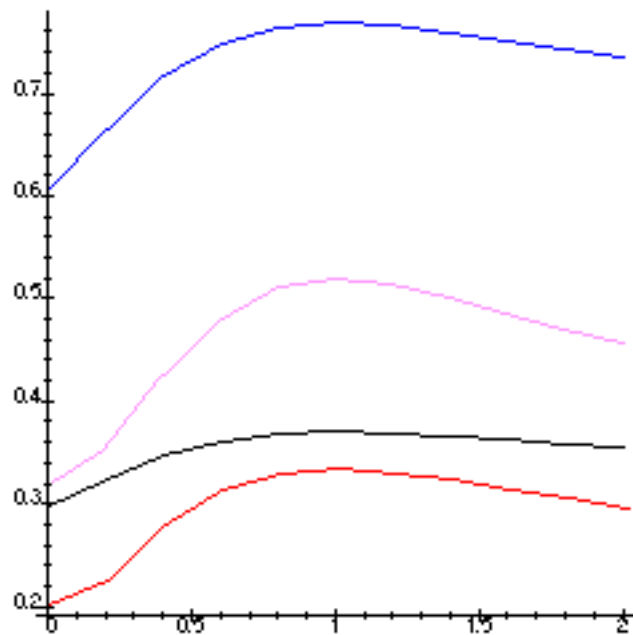
**4.2.2. The rate of convergence.** As to the behaviour of the amplification factors for small values of  $\Delta t \lambda(J_k)$ , it is immediate from (4.6) and (4.7) that  $\alpha_j^m = O((\Delta t)^m)$ . Thus, the  $\Pi_{13}\Pi_{23}$  process damps the nonstiff error components with the same order in  $\Delta t$  as  $\Pi_3$ , but less strongly than  $\Pi$ ,  $\Pi_3\Pi$  and  $\Pi_3\Pi^2$ . Next, we look for a given value of  $x := \kappa_j y_3$  at the averaged value  $\alpha_{\text{aver}}(x)$  of  $|C_{13}C_{23}|^{1/2}$  in the square  $\{|\kappa_j y_1| \leq \gamma, |\kappa_j y_2| \leq \gamma\}$  with  $\gamma = 2.19$ , and we compare this function with that obtained for the  $\Pi_3\Pi$  with  $\gamma = 0.72$ . The lowest and highest graph in the Figures 4.4a and 4.4b presents  $\Pi_3\Pi$  and  $\Pi_{13}\Pi_{23}$ , respectively. Hence, we may conclude that the three times larger convergence boundary of the process  $\Pi_{13}\Pi_{23}$  is paid by a considerably worse convergence rate. Nevertheless, the question remains which iteration method is most efficient. To answer this question, we may either count three  $\Pi_{13}\Pi_{23}$  iterations for one iteration, that is, we define  $\alpha_{\text{aver}}(x)$  by means of  $|C_{13}C_{23}|^{3/2}$ , or we apply  $\Pi_{13}\Pi_{23}$  with the same convergence boundary as used in the  $\Pi_3\Pi$  process, i.e.  $\gamma = 0.72$ . The middle graphs in the Figures 4.4a and 4.4b present these cases, where the lowest middle one corresponds with the  $\gamma = 0.72$  approach. Thus, our final conclusion is that on the basis of the  $\alpha_{\text{aver}}(x)$  profiles and the damping of the nonstiff error components, the  $\Pi_3\Pi$  is the most efficient iteration strategy.

## 5. Conclusions

In this paper we analysed the convergence of iterative solution methods for the linear systems arising in Rosenbrock integration methods and the nonlinear systems arising in DIRK integration methods. In



**Figure 4.4a.** The functions  $\alpha_{\text{aver}}(x)$  for the  $\Pi_3\Pi$  and  $\Pi_{13}\Pi_{23}$  processes.



**Figure 4.4b.** Details near the origin of  $\alpha_{\text{aver}}(x)$  for the  $\Pi_3\Pi$  and  $\Pi_{13}\Pi_{23}$  processes.

particular, we focused on the integration of ODEs originating from shallow water applications, where the ODE system contains a highly stiff part corresponding with the vertical derivative terms and a moderately stiff part corresponding with the horizontal derivative terms. We considered iteration methods based on the approximate factorizations of the system matrix associated with the application of modified Newton. The resulting convergence boundaries  $\gamma$  are given by

Process	$\Pi$	$\Pi_3$	$\Pi_3\Pi$	$\Pi_3\Pi^2$	$\Pi_{13}\Pi_{23}$
$\gamma \approx$	0.65	0.5	0.72	0.75	2.19

The best convergence characteristics were obtained for the  $\Pi_3\Pi$  process where alternately only the vertical direction and all coordinate directions are treated implicitly.

## References

- [1] Brunner, H. & Houwen, P.J. van der [1986]: The numerical solution of Volterra equations, North-Holland.
- [2] Dekker, K. & Verwer, J.G. [1984]: Stability of Runge-Kutta methods for stiff nonlinear differential equations, North-Holland.
- [3] 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.
- [4] Hairer, E. & Wanner, G. [1991]: Solving ordinary differential equations II. Stiff and differential-algebraic problems, Springer.
- [5] 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.
- [6] Houwen, P.J. van der & Sommeijer, B.P. [1998]: Approximate factorization in shallow water applications, Report MAS R9835, CWI, Amsterdam, submitted for publication.
- [7] 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.
- [8] Rosenbrock, H.H. [1962-1963]: Some general implicit processes for the numerical solution of differential equations, Computer J. 5, 329-330
- [9] Sandu, A. [1997]: Numerical aspects of air quality modelling, PhD thesis, University of Iowa.
- [10] 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.
- [11] 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.
- [12] Verwer, J.G., Spee, E.J., Blom, J.G. & Hundsdorfer, W. [1999]: A second-order Rosenbrock method applied to photochemical dispersion problems, to appear in SIAM J. Sci. Comput.
- [13] Verwer, J.G., Hundsdorfer, W. & Blom, J.G. [1998]: Numerical time integration for air pollution models, Report MAS R9825, CWI, Amsterdam, submitted for publication.