HOUWEN, P. J. VAN DER; SOMMEIJER, B. P.

# Iteration of Runge-Kutta Methods with Block Triangular Jacobians

*Wir betrachten Iterationsprozesse zur Lösung der impliziten Relationen, die impliziten Runge-Kutta-Verfahren (RK-Verfahren) beigeordnet sind, wenn sie auf steife Anfangswertprobleme (ARP) angewendet werden. Der konventionelle Ansatz zur Lösung der RK-Gleichungen verwendet die Newton-Iteration, wobei die volle Jacobi-Matrix der rechten Seite ausgenutzt wird. Für ARP großer Dimension ist dieser Ansatz wegen der hohen Kosten, die bei der LU-Zerlegung der Jacobi-Matrix der RK-Gleichungen auftreten, nicht attraktiv. Es wurden verschiedene Vorschläge zur Reduzierung dieser hohen Kosten gemacht. Am weitesten bekannte Gegenmittel ist die Verwendung von Ähnlichkeitstransformationen, durch die die RK-Jacobi-Matrix in eine block-diagonale Matrix transformiert wird, deren Blöcke die ARP-Dimension haben. In dieser Arbeit untersuchen wir einen alternativen Ansatz, der die RK-Jacobi-Matrix direkt durch eine block-diagonale oder durch eine block-triangulare Matrix ersetzt, deren Blöcke selbst block-triangulare Matrizen sind. So ein äußerst ,vereinfachtes' Newton-Iterations-Verfahren gestattet ein Beträchtliches an Parallelität. Einen bedeutenden Beitrag stellt hier allerdings die Antwort auf die Frage dar, ob der block-triangulare Ansatz konvergiert. Ziel der Arbeit ist es, Einsicht in den Effekt auf die Konvergenz block-triangularer Jacobi-Matrix-Approximationen zu gewinnen.*

We shall consider iteration processes for solving the implicit relations associated with implicit Runge-Kutta (RK) methods applied to stiff initial value problems (IVPs). The conventional approach for solving the RK equations uses Newton iteration employing the full righthand side Jacobian. For IVPs of large dimension, this approach is not attractive because of the high costs involved in the LU-decomposition of the Jacobian of the RK equations. Several proposals have been made to reduce these high costs. The most well-known remedy is the use of similarity transformations by which the RK Jacobian is transformed to a block-diagonal matrix the blocks of which have the IVP dimension. In this paper we study an alternative approach which directly replaces the RK Jacobian by a block-diagonal or block-triangular matrix the blocks of which themselves are block-triangular matrices. Such a grossly 'simplified' Newton iteration process allows for a considerable amount of parallelism. However, the important issue is whether this block-triangular approach does converge. It is the aim of this paper to get insight into the effect on the convergence of block-triangular Jacobian approximations.

MSC (1991): 65L06, 65L05, 34A50

## 1. Introduction

We shall consider iteration processes for solving the implicit relations associated with implicit Runge-Kutta (RK) methods applied to the stiff initial value problem (IVP)

$$\mathbf{y}'(t) = \mathbf{f}(\mathbf{y}(t)), \qquad \mathbf{y}(t_0) = \mathbf{y}_0, \qquad \mathbf{y}, \mathbf{f} \in \mathbb{R}^d. \tag{1.1}$$

Let the (s-stage) RK method be given by

$$\mathbf{R}(\mathbf{Y}) = \mathbf{0}, \qquad \mathbf{R}(\mathbf{Y}) := \mathbf{Y} - h(A \otimes I)\,\mathbf{F}(\mathbf{Y}) - (\mathbf{e} \otimes \mathbf{I})\,\mathbf{y_n}, \qquad \mathbf{y}_{n+1} = (\mathbf{e}_s^{\mathsf{T}} \otimes I)\,\mathbf{Y}, \tag{1.2}$$

where $h$ is the integration step, $\mathbf{y}_n$ and $\mathbf{y}_{n+1}$ represent approximations to the exact solution vector $\mathbf{y}(t)$ at $t = t_n$ and $t = t_{n+1}$, $A$ is the $s$-by-$s$ RK matrix, $\otimes$ denotes the Kronecker product, the $s$-dimensional vectors $\mathbf{e}$ and $\mathbf{e}_i$, respectively, are the vector with unit entries and the $i$th unit vector, and $I$ is the $d$-by-$d$ identity matrix (in the following, the identity matrix will always be denoted by $I$ and its dimension will be clear from the context in which it appears). The $s$ components $\mathbf{Y}_i$ of $\mathbf{Y}$ represent intermediate approximations to the exact solution values and $\mathbf{F}(\mathbf{Y})$ contains the derivative values $(\mathbf{f}(\mathbf{Y}_i))$. The classical RK methods of this type, like the Radau IIA and Lobatto IIIA methods, are highly accurate and highly stable, and therefore reliable methods for solving the IVP (1.1).

The conventional approach for solving the system $\mathbf{R}(\mathbf{Y}) = \mathbf{0}$ uses Newton iteration which requires the Jacobian matrix $I - A \otimes hJ$ of the RK equations (1.2). Here, $J$ denotes the Jacobian $\partial\mathbf{f}/\partial\mathbf{y}$ of the righthand side function $\mathbf{f}$. For large $d$, this approach is not attractive because of the high costs involved in the LU-decomposition of the $sd$-by-$sd$ RK Jacobian $I - A \otimes hJ$. To be more precise, the LU costs are given by $2s^3d^3/3 + O(s^2d^2)$ flops. In the following, we shall ignore the last term in this expression. Several proposals have been made to reduce these high costs. The most well-known remedy is the use of similarity transformations by which $I - A \otimes hJ$ is transformed to a block-diagonal matrix $I - D \otimes hJ$ the blocks of which have dimension $d$ (cf. BUTCHER [1]). For the classical implicit RK methods that we want to use, the diagonal entries of $D$ are complex, so that further modifications are needed involving complex arithmetic (cf. HAIRER and WANNER [5]). The resulting iteration method is highly efficient and forms the basis for the by now famous RADAU5 code given in [5]. Moreover, this iteration method has intrinsic parallelism, so that it is suitable for implementation on a parallel system. In fact, by the similarity transformation approach, the sequential (or effective) LU costs associated with $s$-stage RK methods can be reduced to $8d^3/3$.

An alternative approach directly replaces the RK Jacobian $I - A \otimes hJ$ by a block-diagonal or block-triangular matrix $I - B \otimes hJ$, where $B$ is diagonal or triangular with real diagonal entries $B_{ii}$. This approach was analysed in [6] and [7]. The main costs involved in this method consist of the evaluation of the righthand side Jacobian $J$, the LU-

decompositions of the $s$ matrices $I - hB_{ii}J$, $ms$ forward/backward substitutions, and $ms$ righthand side evaluations. Here, $m$ denotes the number of iterations. It turns out that, except for the forward/backward substitutions, these costs reduce by a factor $s$ when a parallel system with $s$ processors is used. We shall be particularly interested in high-dimensional problems, i.e. $d \gg 1$. Therefore, only $O(md^2)$ and $O(d^3)$ operations will be taken into account. Furthermore, we assume that the RK Jacobian needs an up-date at the beginning of each RK step (which is quite realistic because of the relatively large steps allowed by implicit RK methods).

Denoting the computational complexity per step on $p$ processors by $E(p)$ flops, we have

$$E(p) \approx p^{-1}cd^2 + \tfrac{2}{3}\lceil p^{-1}s \rceil d^3 + 2msd^2 \,, \tag{1.3}$$

where $cd^2$ represents the computational complexity for computing $J$ and $\lceil x \rceil$ denotes the smallest integer greater than or equal to $x$. For large $d$, the expression (1.3) shows that on one processor, the costs of the block-triangular approach are $s/4$ times the costs required by the similarity transformation approach. However, on $s$ processors, this fraction becomes $1/4$, so that for large $d$ the block-triangular method should become 4 times faster than RADAU5.

In this paper, we want to reduce the computational complexity of the block-triangular method by tuning the iteration process to the problem at hand. For example, it often happens that the system (1.1) can be split into weakly coupled subsystems. In such cases, it may be advantageous to adapt the RK Jacobian to these coupling properties. Suppose that the righthand side Jacobian matrix $J$ is approximated by a $\sigma$-by-$\sigma$ block-triangular matrix $\tilde{J}$ ($\sigma$ is assumed to be greater than 1) where the blocks $\tilde{J}_{ik}$ are $d_i$-by-$d_k$ matrices with $i$, $k = 1, \ldots, \sigma$, and let the RK Jacobian be replaced by an $s$-by-$s$ block-triangular matrix of which each diagonal block is itself a $\sigma$-by-$\sigma$ block-triangular matrix the diagonal blocks of which are $d_i$-by-$d_i$ matrices. The block-triangular structure of the simplified RK Jacobian implies that the $sd$-dimensional linear system falls apart into $s\sigma$ linear subsystems, $s$ of which have dimensions $d_1, d_2, \ldots, d_\sigma$, respectively. The vector $\mathbf{d} := (d_1, d_2, \ldots, d_\sigma)^{\mathsf{T}}$ characterizes the partitioning into blocks of the matrix $\tilde{J}$ and will therefore be called the *partitioning vector*. For large $d$ and $\sigma$, the block-triangular approach reduces the computational work considerably, provided that the number of iterations, $\tilde{m}$, does not increase excessively. Such an increase of the number of iterations can be avoided by decreasing the stepsize. Let $h$ and $\tilde{h}$ denote the stepsizes taken by the full Jacobian and block-triangular versions, and let $\tilde{h}$ be such that for $\tilde{m} = m$, the block-triangular version produces the same accuracy as the full Jacobian version. Assuming that the block-triangular version up-dates its Jacobian and corresponding LU-decomposition with the same frequency as the full Jacobian version, the sequential computational complexity $\tilde{E}(p)$ of the block-triangular version over a step $h$ is given by

$$\tilde{E}(p) \approx p^{-1}\tilde{c}d^2 + \tfrac{2}{3}\lceil p^{-1}\sigma s \rceil \tilde{d}^3 + 2h\tilde{h}^{-1}ms\|\mathbf{d}\|_2^2 \,, \tag{1.4}$$

where $\|\mathbf{d}\|_2$ denotes the Euclidean norm of $\mathbf{d}$, $\tilde{c}d^2$ represents the computational complexity for computing $\tilde{J}$, and $\tilde{d}$ is the maximal value of the dimensions $d_i$. Furthermore, assuming that sufficiently many processors are available, the speed-up factor for the block-triangular approach on $p = \sigma s$ processors is given by

$$S := \frac{E(\sigma s)}{\tilde{E}(\sigma s)} \approx \frac{3c + 2\sigma s(d + 3sm)}{3\tilde{c} + 2\sigma s d^{-2}(\tilde{d}^3 + 3smh\tilde{h}^{-1}\|\mathbf{d}\|_2^2)} \,. \tag{1.5}$$

If the righthand side Jacobian $J$ is *expensive*, i.e., $c$ and $\tilde{c}$ are large, then we have speed-up by a factor $S \approx c\tilde{c}^{-1}$. Consequently, for expensive righthand side Jacobians, it is recommendable to choose $\tilde{J}$ as sparse as possible (e.g. block-diagonal).

In the case of *cheap* righthand side Jacobians ($c$ and $\tilde{c}$ can be ignored), it follows from (1.5) that

$$S \approx \frac{d^2}{\|\mathbf{d}\|_2^2} \frac{3sm + d}{3smh\tilde{h}^{-1} + \tilde{d}^3\|\mathbf{d}\|_2^{-2}} \,, \tag{1.6}$$

showing that $S = S(m)$ is a monotonically decreasing function of $m$ satisfying the inequality

$$\frac{d^2}{h\tilde{h}^{-1}\|\mathbf{d}\|_2^2} \leq S(m) \leq \frac{d^2}{\|\mathbf{d}\|_2^2} \frac{3s + d}{3sh\tilde{h}^{-1} + \tilde{d}^3\|\mathbf{d}\|_2^{-2}} \,. \tag{1.7}$$

We remark that these expressions for $S$ are related to the theoretical speed-up factor. Hence, an actual implementation on a parallel architecture will show speed-up factors that are bounded above by the theoretical ones, due to (machine dependent) communication costs and synchronization overhead. The important issue is whether the block-triangular iteration method does converge as $m \to \infty$. It is the aim of this paper to get insight into the effect on the convergence of block-triangular Jacobian approximations.

## 2. Iteration of RK methods

Consider the following Newton-type iteration scheme for approximating $\mathbf{y}_{n+1}$:

$$\left.\begin{aligned} (I - B \otimes h\tilde{J})\,(\mathbf{Y}^{(j)} - \mathbf{Y}^{(j-1)}) &= -\mathbf{R}(\mathbf{Y}^{(j-1)}) + h\Gamma(\mathbf{Y}^{(j)}, \mathbf{Y}^{(j-1)}) \,, \qquad j = 1, \ldots, m \,, \\ \mathbf{y}_{n+1} &= (\mathbf{e}_s^{\mathsf{T}} \otimes I)\,\mathbf{Y}^{(m)} \,, \end{aligned}\right\} \tag{2.1}$$

where $I$ is the $sd$-by-$sd$ identity matrix, $B$ is diagonal or (lower) triangular with positive diagonal entries, $\tilde{J}$ is an approximation to the true Jacobian $J$ at $\mathbf{y}_n$, and where $\Gamma$ is an appropriately chosen function based on the structure of $\tilde{J}$. It will be assumed that $\Gamma(\mathbf{U}, \mathbf{U})$ vanishes for any $\mathbf{U}$. Hence, if (2.1) converges, then it converges to the solution of $\mathbf{R}(\mathbf{Y}) = \mathbf{0}$. We remark that the case $\tilde{J} = J$ and $\Gamma(\mathbf{Y}^{(j)}, \mathbf{Y}^{(j-1)}) = \mathbf{0}$ has been analysed in [6] for $B$ diagonal and in [7] for $B$ triangular.

Let us define the function $\Gamma$ by

$$\Gamma(\mathbf{Y}^{(j)}, \mathbf{Y}^{(j-1)}) := (L \otimes I)\,\mathbf{F}(\mathbf{Y}^{(j)}) + (C \otimes I)\,\mathbf{G}(\mathbf{Y}^{(j)}, \mathbf{Y}^{(j-1)}) - ((L+C) \otimes I)\,\mathbf{F}(\mathbf{Y}^{(j-1)}), \qquad (2.2)$$

where $C$ is diagonal, $L$ is strictly lower triangular, and where for any $\mathbf{U}$ the function $\mathbf{G}$ satisfies the relation $\mathbf{G}(\mathbf{U}, \mathbf{U}) = \mathbf{F}(\mathbf{U})$. In fact, $\mathbf{G}(\mathbf{Y}^{(j)}, \mathbf{Y}^{(j-1)})$ is an approximation to $\mathbf{F}(\mathbf{Y})$ using the most recent iteration values available.

To motivate the introduction of the function $\Gamma$, we will discuss in some detail the various terms occurring in the righthand side of (2.2). In the first term we encounter the matrix $L$, acting on the c u r r e n t iterate $\mathbf{Y}^{(j)}$. As a consequence, the $s$ systems of dimension $d$ cannot be solved in parallel as it was the case in [6], even if $B$ is diagonal. Hence, owing to the strictly lower triangular form of $L$, these systems are solved sequentially. In this way we introduced a 'Gauss-Seidel type' approach, since the stage component vectors $\mathbf{Y}_1^{(j)}, \ldots, \mathbf{Y}_{k-1}^{(j)}$ at the new iteration level are used in solving for $\mathbf{Y}_k^{(j)}$. Next, we consider the last term in (2.2). This term does not complicate the algorithm, since here only the known, previous iterate $\mathbf{Y}^{(j-1)}$ is involved. Finally, we comment on the role of the function $\mathbf{G}$, occurring in the second term of (2.2). The major aim for introducing this function is to be able to use the most recent information available w i t h i n the solution of each of the $s$ linear systems. As already observed in the Introduction, if the matrix $\tilde{J}$ is a $\sigma$-by-$\sigma$ b l o c k - t r i a n g u l a r approximation $(\tilde{J}_{ik})$ to $J$ where the blocks $\tilde{J}_{ik}$ are $d_i$-by-$d_k$ matrices, then each of the $s$ linear systems in (2.1) falls apart into a sequence of $\sigma$ linear subsystems, respectively of dimensions $d_1, d_2, \ldots, d_\sigma$. The block-triangular structure of $\tilde{J}$ enables us to 'up-date' the components of $\mathbf{G}(\mathbf{Y}^{(j)}, \mathbf{Y}^{(j-1)})$ during the computation of each of the stage value approximations $\mathbf{Y}_k^{(j)}$, $k = 1, \ldots, s$. In the next subsection, an explicit formula for $\mathbf{G}$ is given in case of a linear problem. In conclusion, we might say that two complementary forms of Gauss-Seidel iteration have been introduced: one by means of the matrix $L$ (to use new information from one stage to the next), and the other by the function $\mathbf{G}$ (acting within each stage) by exploiting the block-triangular structure of $\tilde{J}$.

## 2.1. The error equation

In this section, we discuss the convergence for the linear case $\mathbf{y}' = J\mathbf{y}$. Let the righthand side Jacobian $J$ be partitioned according to $J = (J_{ik})$ where the blocks $J_{ik}$ are $d_i$-by-$d_k$ matrices, and let $J$ be split according to $J = J_{\mathrm{L}} + J_{\mathrm{D}} + J_{\mathrm{U}}$, where $J_{\mathrm{L}}$, $J_{\mathrm{D}}$, and $J_{\mathrm{U}}$ are (with respect to the block partitioning $J_{ik}$) strictly lower triangular, diagonal and strictly upper triangular block matrices. For this model equation, $\mathbf{G}(\mathbf{Y}^{(j)}, \mathbf{Y}^{(j-1)})$ can be expressed in the form

$$\mathbf{G}(\mathbf{Y}^{(j)}, \mathbf{Y}^{(j-1)}) = (I \otimes J_{\mathrm{L}})\,\mathbf{Y}^{(j)} + (I \otimes (J_{\mathrm{D}} + J_{\mathrm{U}}))\,\mathbf{Y}^{(j-1)},$$

so that

$$h\Gamma(\mathbf{Y}^{(j)}, \mathbf{Y}^{(j-1)}) = (L \otimes hJ + C \otimes hJ_{\mathrm{L}})\,\mathbf{Y}^{(j)} + (C \otimes h(J_{\mathrm{D}} + J_{\mathrm{U}}) - (L+C) \otimes hJ)\,\mathbf{Y}^{(j-1)}. \qquad (2.3)$$

The recursion for $\mathbf{Y}^{(j)}$ takes the form

$$(I - B \otimes h\tilde{J})(\mathbf{Y}^{(j)} - \mathbf{Y}^{(j-1)}) = (\mathbf{e} \otimes I)\,\mathbf{y}_n - \mathbf{Y}^{(j-1)} + ((L \otimes hJ) + (C \otimes hJ_{\mathrm{L}}))\,\mathbf{Y}^{(j)}$$
$$+ ((C \otimes h(J_{\mathrm{D}} + J_{\mathrm{U}})) + ((A - L - C) \otimes hJ))\,\mathbf{Y}^{(j-1)}. \qquad (2.4)$$

For the exact corrector solution we have

$$(I - B \otimes h\tilde{J})(\mathbf{Y} - \mathbf{Y}) = (\mathbf{e} \otimes I)\,\mathbf{y}_n - \mathbf{Y} + ((L \otimes hJ) + (C \otimes hJ_{\mathrm{L}}))\,\mathbf{Y}$$
$$+ ((C \otimes h(J_{\mathrm{D}} + J_{\mathrm{U}})) + ((A - L - C) \otimes hJ))\,\mathbf{Y}. \qquad (2.5)$$

From (2.4) and (2.5) it follows that the error recursion is given by

$$\left.\begin{aligned}
\mathbf{Y}^{(j)} - \mathbf{Y} &= M(\mathbf{Y}^{(j-1)} - \mathbf{Y}), \\
M &= h(I - hW)^{-1}(A \otimes J - W), \qquad W := B \otimes \tilde{J} + L \otimes J + C \otimes J_{\mathrm{L}}.
\end{aligned}\right\} \qquad (2.6)$$

The error amplification matrix $M$ is completely determined by the RK matrix $A$ and the lower block-triangular matrix $W$. In this paper, we shall restrict our analysis to the two special cases

$$B = C = D, \qquad \tilde{J} = J_{\mathrm{D}}, \qquad (2.7\,\mathrm{a})$$

$$B = D, \qquad C = O, \qquad \tilde{J} = J_{\mathrm{D}} + J_{\mathrm{L}}, \qquad (2.7\,\mathrm{b})$$

where $D$ denotes a diagonal matrix with nonnegative diagonal entries. The methods generated by (2.7a) and (2.7b) both lead to the same matrix $W$:

$$W = T \otimes J - D \otimes J_{\mathrm{U}}, \qquad T := L + D. \tag{2.8}$$

Hence, they possess identical error recursions, but will produce different solutions when applied to nonlinear problems.

A necessary and sufficient condition for convergence of the error recursion (2.6) requires the spectral radius $\varrho(M)$ to be less than 1. In the case $\tilde{J} = J$, $L = C = O$, analysed in [6], this spectral radius condition leads to a condition in terms of the eigenvalues of $hJ$. For the more general family of methods generated by (2.7), this is not possible and the condition $\varrho(M) < 1$ can only be verified by a direct numerical computation. However, if all diagonal entries of $hW$ are sufficiently large, then the condition $\varrho(M) < 1$ can be transformed into a spectrum condition for $J^{-1}J_{\mathrm{U}}$. In the case where not all diagonal entries are large, it is possible to derive bounds for the amplification factor $\mu$ occurring in the relation

$$\|M(\mathbf{a} \otimes \mathbf{v})\| = \mu\|\mathbf{a} \otimes \mathbf{v}\|, \tag{2.9}$$

where $\mathbf{v}$ is in the eigenspace of $J$ and $\mathbf{a}$ is in the eigenspace of the matrix

$$Z(z) = z(I - zT)^{-1}(A - T), \qquad z := h\lambda, \tag{2.10}$$

$\lambda$ denoting the eigenvalue of $J$ corresponding to $\mathbf{v}$. If $\mathbf{a} \otimes \mathbf{v}$ happens to be an eigenvector of $M$, then the amplification factor $\mu = \mu(h, z)$ equals the corresponding eigenvalue of $M$, so that convergence requires that all $\mu$ are less than 1. Hence, $\mu(h, z) < 1$ is a necessary condition for convergence.

The derivation of amplification factor bounds and the convergence condition for the large-diagonal-entries case will be the subjects of the following two sections.

## 2.2. Derivation of amplification factor bounds

The following theorem holds.

Theorem 2.1: *Let $W$ be of the form* (2.8), *let $Z(z)$ be defined by* (2.10), *and let $\mathbf{v}$ and $\mathbf{a}$ be eigenvectors of $J$, and $Z(z)$ with eigenvalues $\lambda$ and $\zeta(z)$, respectively. If*

$$h < \frac{1}{\gamma\|J_{\mathrm{U}}\|}, \qquad \gamma := \|D\| \max_{h} \|(I - T \otimes hJ)^{-1}\|, \tag{2.11a}$$

*then the amplification factor $\mu$ defined in* (2.9) *satisfies*

$$\mu \leq \frac{|\zeta(z)| + \gamma h\|J_{\mathrm{U}}\|}{1 - \gamma h\|J_{\mathrm{U}}\|}, \tag{2.11b}$$

*and the corresponding convergence region is given by*

$$\text{Spectrum } hJ \in \mathbb{C} := \{z : \varrho(Z(z)) < 1 - 2\gamma h\|J_{\mathrm{U}}\|\}. \tag{2.11c}$$

Proof: From (2.8) it follows that $M$ can be represented in the form

$$\left.\begin{aligned} M &= (I + Q)^{-1}(Q + V), \\ Q &:= (I - T \otimes hJ)^{-1}(D \otimes hJ_{\mathrm{U}}), \qquad V := (I - T \otimes hJ)^{-1}((A - T) \otimes hJ). \end{aligned}\right\} \tag{2.12}$$

By means of the conditions of the theorem, it is easily verified that

$$V(\mathbf{a} \otimes \mathbf{v}) = (h\lambda(I - h\lambda T)^{-1}(A - T) \otimes I)(\mathbf{a} \otimes \mathbf{v}) = \zeta(z)(\mathbf{a} \otimes \mathbf{v})$$

so that

$$\|V(\mathbf{a} \otimes \mathbf{v})\| \leq |\zeta(z)|\|\mathbf{a} \otimes \mathbf{v}\|.$$

Furthermore, assuming that $\|Q\| < 1$, we have

$$\|(I + Q)^{-1}\| \leq \frac{1}{1 - \|Q\|}.$$

Hence,

$$\|M(\mathbf{a} \otimes \mathbf{v})\| = \|(I + Q)^{-1}(Q + V)(\mathbf{a} \otimes \mathbf{v})\| \leq \frac{\|Q\| + |\zeta(z)|}{1 - \|Q\|}\|\mathbf{a} \otimes \mathbf{v}\|.$$

Since $\|Q\| \leq \gamma h\|J_{\mathrm{U}}\|$, where $\gamma$ is defined in (2.11a), we obtain the bound (2.11b) and the convergence region (2.11c). ■

The bound (2.11 b) on $\mu$ is sharp for $J_U = O$, i.e. $\mu = |\zeta(z)|$, but will be rather conservative as $\|J_U\|$ increases. If the spectrum of $J$ is assumed to cover the whole left halfplane, then the conditions (2.11) lead to the stepsize condition

$$h < [1 - \max_{\text{Re } z \le 0} \varrho(Z(z))]/2\gamma\|J_U\| .$$                                                                     (2.13 a)

Similarly, if the spectrum of $J$ is known to be negative, then we obtain

$$h < [1 - \max_{z \le 0} \varrho(Z(z))]/2\gamma\|J_U\| .$$                                                                                (2.13 b)

Given the IVP, the two crucial quantities determining the stepsize conditions (2.13) are the values of $\gamma$ and $\max \varrho(Z(z))$. In [7] matrices $T$ have been constructed such that $\varrho(Z(z))$ is small in the whole left halfplane. In order to get some idea of the magnitude of the coefficient $\gamma$, we consider the case where $J$ is a normal matrix, so that

$$\gamma = \|D\| \max_{\text{Re } z \le 0} \|(I - zT)^{-1}\| .$$

The following two examples compute the corresponding stepsize conditions (2.13).

Example 2.1: For the two-point Radau IIA corrector, the Butcher matrix $A$ and the matrix $T$ as constructed in [7] are given by

$$A = \begin{pmatrix} 5/12 & -1/12 \\ 3/4 & 1/4 \end{pmatrix}, \qquad T = \begin{pmatrix} 5/12 & 0 \\ 3/4 & 2/5 \end{pmatrix}, \qquad \varrho(Z(z)) = \left| \frac{9z}{(12 - 5z)(5 - 2z)} \right| .$$

From this we find $\gamma \approx 0.71$, $\max_{\text{Re } z \le 0} \varrho(Z(z)) \approx 0.18$, and $\max_{z \le 0} \varrho(Z(z)) \approx 0.09$. Hence, the convergence conditions (2.13) become $h < 0.58\|J_U\|^{-1}$ and $h < 0.64\|J_U\|^{-1}$, respectively. ■

Example 2.2: Similarly, the four-point Radau IIA corrector is defined by the Butcher matrix

$$A = \begin{pmatrix} .112\,999\,479\,323\,16 & -.040\,309\,220\,723\,52 & .025\,802\,377\,420\,34 & -.009\,904\,676\,507\,3 \\ .234\,383\,995\,747\,40 & .206\,892\,573\,935\,36 & -.047\,857\,128\,048\,54 & .016\,047\,422\,806\,52 \\ .216\,681\,784\,623\,25 & .406\,123\,263\,867\,37 & .189\,036\,518\,170\,06 & -.024\,182\,104\,899\,83 \\ .220\,462\,211\,176\,77 & .388\,193\,468\,843\,17 & .328\,844\,319\,980\,06 & .062\,500\,000\,000\,00 \end{pmatrix}$$

for which [7] derived the matrix

$$T = \begin{pmatrix} .1130 & 0 & 0 & 0 \\ .2344 & .2905 & 0 & 0 \\ .2167 & .4834 & .3083 & 0 \\ .2205 & .4668 & .4414 & .1176 \end{pmatrix} .$$

Numerically, we found $\gamma \approx 0.96$, $\max_{\text{Re } z \le 0} \varrho(Z(z)) \approx 0.51$, and $\max_{z \le 0} \varrho(Z(z)) \approx 0.16$, so that the conditions (2.13) become $h < 0.25\|J_U\|^{-1}$ and $h < 0.43\|J_U\|^{-1}$, respectively. ■

## 2.3. Large diagonal entries in the Jacobian

It sometimes happens that $hW$ has large diagonal entries (i.e. $hW - I \approx hW$), or equivalently,

$$\min_k |J_{kk}| \gg h^{-1}\left(\min_i D_{ii}\right)^{-1}, \qquad i = 1, \ldots, s, \qquad k = 1, \ldots, d,$$                                     (2.14 a)

where $J$ is assumed to be nonsingular. Outside the transient phase, where usually relatively large stepsizes $h$ are taken, condition (2.14a) is often satisfied. From (2.14a) it then follows that the error amplification matrix $M$ can be approximated by

$$M \approx I - W^{-1}(A \otimes J) = I - (T \otimes I - D \otimes J^{-1}J_U)^{-1}(A \otimes I) .$$

The eigenvalues of $M$ are given by those of the matrix $\tilde{M}(z) := I - (T \otimes I - zD \otimes I)^{-1}(A \otimes I)$, where $z$ runs through the spectrum of $J^{-1}J_U$. Hence, we have convergence if

$$\text{Spectrum of } J^{-1}J_U \in \mathbb{C} := \{z : \varrho(\tilde{M}(z)) < 1\} .$$                                                       (2.14 b)

Example 2.3: We derive the region of convergence for the two-point and four-point Radau IIA correctors of the Examples 2.1 and 2.2. The characteristic equation for the eigenvalues $\tilde{\mu}(z)$ of the matrix $\tilde{M}(z)$ takes the form $\det(A - T + zD + \tilde{\mu}(z)(T - zD)) = 0$. Inspection of the region where $\tilde{\mu}(z)$ is bounded by 1 reveals that for both correctors the region of convergence for the eigenvalues of $J^{-1}J_U$ contains a disk of radius $r$ which is centered at the origin and an infinite wedge in the left halfplane with half angle $\alpha$. For the two-point and four-point Radau IIA correctors, we obtain $\{r = 0.27, \alpha = 54°\}$ and $\{r = 0.11, \alpha = 18°\}$, respectively. ■

Remark 2.1: It often happens that the system of ODEs (1.1) contains nonstiff equations (an equation $y_i'(t) = f_i(\mathbf{y}(t))$ is called *nonstiff* if all derivative values $\partial f_i(\mathbf{y}(t))/\partial y_j$, $j = 1, \ldots, d$, are of moderate size, say bounded by 1). Such nonstiff equations do not need implicit treatment. Therefore, in applying the convergence conditions (2.13) and (2.14), we may delete all rows and all columns in $J$ and $J_U$ which correspond to nonstiff equations.

## 2.4. Permutation, transformation and scaling of the ODE system

It is often possible that the ordering of the equations in the system of ODEs (1.1) can be changed in such a way that entries of large magnitude in $J$ move to the lower left corner of the matrix. This may help to reduce the norm of the matrix $J_U$ in condition (2.13) or to relax the condition (2.14b), so that an attractive partitioning vector $\mathbf{d}$ can be obtained (i.e., $\mathbf{d}$ has small entries $d_i$). Let us write $\mathbf{z}(t) = P\mathbf{y}(t)$ where $P$ is a permutation matrix, which is such that the Jacobian $PJP^\mathsf{T}$ of the permuted system $\mathbf{z}'(t) = Pf(P^\mathsf{T}\mathbf{z}(t))$ has a dominant, lower block-triangular structure. We shall define a reordering by the permutation vector $\mathbf{p} = (p_1, p_2, \ldots, p_d)^\mathsf{T}$, where $p_i$ denotes the index of the $y$-component in the original system ($\mathbf{p} = (1, \ldots, d)^\mathsf{T}$ implies no reordering). Evidently, the permutation matrix $P$ associated with $\mathbf{p}$ is defined by $P := (\mathbf{e}_{p_1}, \mathbf{e}_{p_2}, \ldots, \mathbf{e}_{p_d})^\mathsf{T}$ and the entries of $PJP^\mathsf{T}$ are given by $J^*_{ij} = J_{p_i p_j}$, where $J_{rk}$ denote the entries of $J$.

It may happen that the solution vector $\mathbf{y}$ possesses components of large and small magnitude. In such cases, it is not clear when the permuted Jacobian has a 'dominant, lower block-triangular' structure, and it may be useful to scale the ODE system by writing $\tilde{\mathbf{y}}(t) = D\mathbf{y}(t)$, where $D = \mathrm{diag}\,(1/\mathbf{y}(t_0))$. Then, $\tilde{\mathbf{y}}'(t) = Df(D^{-1}\tilde{\mathbf{y}}(t))$ has the scaled Jacobian $DJD^{-1}$, and rather than choosing $P$ such that $PJP^\mathsf{T}$ is dominant, lower block-triangular, $P$ is chosen such that $PDJD^{-1}P^\mathsf{T}$ is dominant, lower block-triangular.

It should be remarked that it is possible to achieve a complete lower block-triangular structure by the real-Schur-decomposition of $J$. Writing $\mathbf{z}(t) = Q\mathbf{y}(t)$ and $\mathbf{z}'(t) = Qf(Q^\mathsf{T}\mathbf{z}(t))$, the (orthogonal) matrix $Q$ can be chosen such that $QJQ^\mathsf{T}$ has a lower block-triangular structure with blocks of at most dimension 2. However, the computation of $Q$ (by the QR-algorithm) is rather expensive and requires $15d^3$ (Moler) flops (cf. [4, p. 235]).

Finally, we remark that in actual computation, the reordering, the real-Schur-decomposition, and the scaling approach are most effective if the righthand side Jacobian is slowly changing over a large number of steps and if the transformed righthand sides $Pf(P^\mathsf{T}\mathbf{z}(t))$ and $Qf(Q^\mathsf{T}\mathbf{z}(t))$ can be provided in 'written out' form (otherwise the many additional matrix vector multiplications will reduce the efficiency considerably).

# 3. Numerical experiments

The crucial aspect of the block-triangular Jacobian approximations discussed in this paper is the convergence behaviour for $\sigma > 1$. In this section, we illustrate the performance and speed-up factors for a few test problems. Given the partitioning vector $\mathbf{d}$ and the iterated RK method $\{(2.1), (2.2)\}$, we shall apply the following three modes (see also (2.7)):

| | | | | |
|---|---|---|---|---|
| Full Jacobian: | $\tilde{J} = J$, | $B = D$, | $C = O$, | (3.1) |
| Triangular Jacobian: | $\tilde{J} = J_\mathrm{D} + J_\mathrm{L}$, | $B = D$, | $C = O$, | (3.2 a) |
| Diagonal Jacobian: | $\tilde{J} = J_\mathrm{D}$, | $B = C = D$. | | (3.2 b) |

We used the four-stage Radau IIA corrector and we define the matrices $A$ and $T = L + D$ as in Example 2.2. We shall refer to the methods generated by (3.1), (3.2a), and (3.2b) as the Full $J$, the Trian $J$, and the Diag $J$ version (for a discussion of the Full $J$ version we refer to [7]).

## 3.1. Convergence conditions

The Trian $J$ and Diag $J$ versions both lead to $W = T \otimes J - D \otimes J_U$ as defined in (2.8), so that the matrix $W$ is of the form as presupposed in the conditions (2.13) and (2.14). For the four-stage Radau IIA corrector, we have the stepsize condition

$$h < 0.25 \, \|J_U\|^{-1} \tag{3.3}$$

(see Example 2.2), so that there is no severe stepsize restriction, provided that the partitioning vector $\mathbf{d}$ is such that $J_U$ is nonstiff. Alternatively, we may check whether the conditions (2.14) are satisfied. For the four-stage Radau IIA corrector, these conditions read:

$$\min_k |J_{kk}| \gg 8.85 \, h^{-1}, \qquad k = 1, \ldots, d, \qquad \text{spectrum } J^{-1}J_U \in \mathbb{C}, \tag{3.4}$$

where $\mathbb{C}$ is specified in Example 2.3.

## 3.2. Test problems

In all experiments, constant stepsizes have been used. If needed, we adapted the initial condition such that the integration starts outside the transient phase. For a given number of iterations, $m$, the tables of results present the minimal number of correct digits $cd$ of the components of $\mathbf{y}$ at the end point $t = t_{\text{end}}$ of the integration interval (i.e., the absolute errors are written as $10^{-cd}$). Furthermore, we compute the corresponding speed-up factors. Since in all examples, the costs for computing the Jacobian are negligible, we shall use formula (1.6).

### 3.2.1. Problem of Davison

In ENRIGHT [3], the following 80-dimensional system of ODEs with a strongly dominant Jacobian matrix is advocated as a test problem for stiff solvers:

$$\mathbf{y}'(t) = A\mathbf{y}(t) + \frac{4}{\pi}\,\mathbf{e}_d \sum_{k=0}^{4} \frac{\sin\left((2k+1)\,\pi t\right)}{2k+1}\,, \qquad d = 80,\ \mathbf{y}(0) = \mathbf{0},\ t_{\text{end}} = 5\,. \tag{3.5}$$

Here, the entries of $A = (a_{ij})$ are 0.01, except for the diagonal entries, the lower and upper off-diagonal entries that are respectively given by $a_{ii} = -(1.5)^{80-i}$, $a_{i,i-1} = a_{i,i+1} = 0.1$. This problem originates from DAVISON [2]. It is an ideal example for applying a fully diagonal approximation to the Jacobian. Keeping the original ordering $\mathbf{p} = (1, \ldots, 80)^{\mathsf{T}}$, and using the maximum norm we have $\|J_{\mathrm{U}}\| = 0.88$, so that condition (3.3) becomes $h < 0.32$.

Since (3.5) is linear, the Diag $J$ and Trian $J$ modes are identical. We applied the method with $\{\sigma = 1,\ \mathbf{d} = (80)\}$, i.e. the Full $J$ version, and with $\{\sigma = 80,\ \mathbf{d} = (1, \ldots, 1)^{\mathsf{T}}\}$, where $\mathbf{d}$ is the partitioning vector. Table 3.1 presents the $cd$-values obtained. Not surprisingly, the accuracies are the same for $h\tilde{h}^{-1} = 1$, so that (1.7) shows that the speed-up factor is at least $S = 80$. Note that convergence is also obtained for $h > 0.32$, indicating that the convergence condition (3.3) is rather conservative.

Table 3.1. Davison problem (3.5)

| Version | $h$ | $\mathbf{d}^{\mathsf{T}}$ | $m = 1$ | $m = 2$ | $m = 3$ | $m = 4$ | ... | $m = 10$ | $m = \infty$ |
|---|---|---|---|---|---|---|---|---|---|
| Full $J$ | 0.5 | (80) | 1.6 | 2.2 | 2.1 | 2.1 | ... | 2.0 | 2.0 |
| Diag $J$ | 0.5 | $(1, \ldots, 1)$ | 1.6 | 2.2 | 2.1 | 2.1 | ... | 2.0 | |
| Full $J$ | 0.2 | (80) | 1.9 | 3.3 | 4.1 | 4.2 | ... | 4.2 | 4.2 |
| Diag $J$ | 0.2 | $(1, \ldots, 1)$ | 1.9 | 3.3 | 4.1 | 4.2 | ... | 4.2 | |
| Full $J$ | 0.1 | (80) | 2.2 | 4.0 | 5.7 | 7.0 | ... | 7.2 | 7.2 |
| Diag $J$ | 0.1 | $(1, \ldots, 1)$ | 2.2 | 4.0 | 5.7 | 7.0 | ... | 7.2 | |

### 3.2.2. HIRES problem of Schäfer

A second example is provided by the HIRES problem given in [5, p. 157] which originates from SCHÄFER [8] for explaining the 'High Irradiance Responses' of photomorphogenesis:

$$
\begin{aligned}
y_1' &= -1.71 y_1 + 0.43 y_2 + 8.32 y_3 + 0.0007\,, \\
y_2' &= +1.71 y_1 - 8.75 y_2\,, \\
y_3' &= -10.03 y_3 + 0.43 y_4 + 0.035 y_5\,, \\
y_4' &= +8.32 y_2 + 1.71 y_3 - 1.12 y_4\,, \\
y_5' &= -1.745 y_5 + 0.43 y_7 + 0.43 y_6\,, \\
y_6' &= -280 y_6 y_8 + 0.69 y_4 + 1.71 y_5 - 0.43 y_6 + 0.69 y_7\,, \\
y_7' &= +280 y_6 y_8 - 1.81 y_7\,, \\
y_8' &= -280 y_6 y_8 + 1.81 y_7\,,
\end{aligned}
\qquad
\begin{aligned}
y_1(5) &= 0.316\,516\,757\,046 \times 10^{-1}\,, \\
y_2(5) &= 0.648\,154\,953\,106 \times 10^{-2}\,, \\
y_3(5) &= 0.458\,345\,106\,475 \times 10^{-2}\,, \\
y_4(5) &= 0.897\,432\,327\,352 \times 10^{-1}\,, \\
y_5(5) &= 0.162\,451\,453\,753\,, \\
y_6(5) &= 0.685\,043\,896\,144\,, \\
y_7(5) &= 0.564\,670\,034\,192 \times 10^{-2}\,, \\
y_8(5) &= 0.532\,996\,580\,805 \times 10^{-4}\,,
\end{aligned}
\tag{3.6}
$$

with $t_{\text{end}} = 305$. Only the last three equations of the system (3.6) are relatively stiff, so that we can keep the original ordering. It is easily seen that setting $\sigma = 2$ and $\mathbf{d} = (4, 4)^{\mathsf{T}}$ yields a matrix $J_{\mathrm{U}}$ that contains only one non-zero entry, i.e. $(J_{\mathrm{U}})_{3,5} = 0.035$. Hence, in view of condition (3.3), we may expect amplification factors less than 1 without severe restrictions on the stepsize $h$.

The Diag $J$ and Trian $J$ modes produce almost the same results. Therefore, we listed results only for the Diag $J$ mode. The figures in Table 3.2 show that from the second iteration on, the Full $J$ and Diag $J$ version yield comparable accuracies for $h\tilde{h}^{-1} = 1$. The speed-up factor is given by $S \approx 2 + m^{-1}$.

Table 3.2. HIRES problem (3.6) of Schäfer

| Version | $h$ | $\mathbf{d}^\mathsf{T}$ | $m=1$ | $m=2$ | $m=3$ | $m=4$ | ... | $m=10$ | $m=\infty$ |
|---|---|---|---|---|---|---|---|---|---|
| Full $J$ | 15 | (8) | 3.1 | 4.0 | 3.9 | 4.1 | ... | 5.6 | 7.9 |
| Diag $J$ | 15 | (4, 4) | 2.2 | 3.8 | 4.0 | 4.1 | ... | 5.6 | |
| Full $J$ | 7.5 | (8) | 3.3 | 4.4 | 4.7 | 5.3 | ... | 7.0 | 9.0 |
| Diag $J$ | 7.5 | (4, 4) | 2.5 | 4.5 | 4.8 | 5.5 | ... | 7.0 | |

### 3.2.3. NUCREAC problem of Strehmel-Weiner

In STREHMEL and WEINER [9, p. 310], we find a simplified model of a nuclear reactor:

$$
\left.
\begin{aligned}
y_1' &= -\tfrac{1}{3}\left(500 y_2 - 374280\right) y_1 + \tfrac{1}{3} \sum_{i=3}^{8} \beta_i y_i \,, \\
y_2' &= -\frac{1}{1.67}\left(330 y_2 - 136000 y_1 - 9900\right), \\
y_i' &= -\gamma_i(y_i - y_1)\,, \qquad 3 \le i \le 8 \,,
\end{aligned}
\right\}
\tag{3.7}
$$

where

$$
\mathbf{y}(0.5) = \begin{pmatrix}
1.745\,794\,025\,602\,1 \\
749.478\,029\,221\,95 \\
1.579\,316\,355\,556\,2 \\
1.321\,865\,374\,099\,7 \\
1.104\,186\,334\,140\,0 \\
1.040\,256\,901\,940\,0 \\
1.011\,285\,091\,275\,3 \\
1.004\,608\,805\,868\,6
\end{pmatrix}, \quad
\beta = (\beta_i) = \begin{pmatrix}
30.2 \\
82.8 \\
284.4 \\
141.1 \\
157.7 \\
23.8
\end{pmatrix}, \quad
\gamma = (\gamma_i) = \begin{pmatrix}
3 \\
1.13 \\
0.301 \\
0.111 \\
0.0305 \\
0.0124
\end{pmatrix}.
$$

Only the first two equations are stiff, so that in the Diag $J$ and Trian $J$ modes we may set $\mathbf{d} = (2,2,2,2)^\mathsf{T}$ with $\sigma = 4$. Since the stiff subsystem is iterated with a full Jacobian, convergence is expected without stepsize restriction (see Remark 2.1). The results at $t_\text{end} = 15$ listed in Table 3.3 show that the Full $J$ and Diag $J$ versions produce comparable accuracies for $m > 1$ and $h\tilde{h}^{-1} = 1$ (again, the Trian $J$ and Diag $J$ modes yield almost identical results). The speed-up factor (1.6) is $S \approx 4 + 2.5 m^{-1}$. Here, and in the following examples, $N$ denotes the number of time steps, i.e., $h = (t_\text{end} - t_0)/N$.

Table 3.3. NUCREAC problem (3.7) of Strehmel-Weiner

| Version | $N$ | $\mathbf{d}^\mathsf{T}$ | $m=1$ | $m=2$ | $m=3$ | $m=4$ | ... | $m=10$ | $m=\infty$ |
|---|---|---|---|---|---|---|---|---|---|
| Full $J$ | 2 | (8) | 1.5 | 2.5 | 3.3 | 3.5 | ... | 3.5 | 3.5 |
| Diag $J$ | 2 | (2, 2, 2, 2) | 1.0 | 2.0 | 2.9 | 3.5 | ... | 3.5 | |
| Full $J$ | 5 | (8) | 1.9 | 3.2 | 4.2 | 5.2 | ... | 8.1 | 8.1 |
| Diag $J$ | 5 | (2, 2, 2, 2) | 1.6 | 2.9 | 4.1 | 5.2 | ... | 8.1 | |
| Full $J$ | 10 | (8) | 2.2 | 3.8 | 5.0 | 6.2 | ... | 10.1 | 10.1 |
| Diag $J$ | 10 | (2, 2, 2, 2) | 2.0 | 3.6 | 5.0 | 6.2 | ... | 10.1 | |

### 3.2.4. ATMOS20 problem of Verwer

The ATMOS20 problem is a stiff, nonlinear system of 20 ODEs originating from an air pollution model (see VERWER [10], we note that this paper contains a misprint: the third reaction rate should read $0.123_{10^5}$ instead of $0.120_{10^5}$). We solved the corrected system in the integration interval [5,60]. Table 3.4 lists results for the following four cases:

$$
\begin{aligned}
&\text{I:} \quad \sigma = 1\,, && \mathbf{d} = (20)\,, && \mathbf{p} = (1, \ldots, 20)^\mathsf{T}\,; \\
&\text{II:} \quad \sigma = 3\,, && \mathbf{d} = (7, 5, 8,)^\mathsf{T}\,, && \mathbf{p} = (1, \ldots, 20)^\mathsf{T}\,; \\
&\text{III:} \quad \sigma = 8\,, && \mathbf{d} = (3, 3, 3, 3, 2, 2, 3, 1)^\mathsf{T}\,, && \\
&&& \mathbf{p} = (16, 17, 18, 5, 6, 8, 9, 10, 11, 12, 13, 14, 15, 7, 19, 20, 3, 1, 4, 2)^\mathsf{T}\,; \\
&\text{IV:} \quad \sigma = 20\,, && \mathbf{d} = (1, \ldots, 1)^\mathsf{T}\,, && \mathbf{p} = (1, \ldots, 20)^\mathsf{T}\,.
\end{aligned}
$$

The Diag $J$ and Trian $J$ modes produced almost the same accuracies. For $h\tilde{h}^{-1} = 2$ and $h\tilde{h}^{-1} = 4$, the cases II and III lead to a satisfactory speed-up factor $S \approx 1.45 + 2.2 \text{ m}^{-1}$ and $S \approx 1.85 + 3.1 \text{ m}^{-1}$, respectively. The extremely cheap, but over-optimistic case IV leads to a rather poor convergence behaviour.

Table 3.4. ATMOS20 problem of Verwer [10]

| Version | $N$ | Case | $m = 1$ | $m = 2$ | $m = 3$ | $m = 4$ | ... $m = 10$ | $m = \infty$ |
|---|---|---|---|---|---|---|---|---|
| Full $J$ | 5 | I | 3.4 | 4.9 | 7.0 | 6.8 | ... 8.7 | 11.0 |
| Trian $J$ | 5 | II | 3.4 | 5.1 | 5.2 | 6.1 | ... 8.2 | |
| | | III | 2.5 | 3.3 | 4.0 | 4.7 | ... 7.7 | |
| | | IV | 2.2 | 2.8 | 3.4 | 4.1 | ... 4.1 | |
| Full $J$ | 10 | I | 3.7 | 5.5 | 7.6 | 8.3 | ... 11.5 | 12.3 |
| Trian $J$ | 10 | II | 3.8 | 5.7 | 6.0 | 7.2 | ... 10.3 | |
| | | III | 2.8 | 3.8 | 4.6 | 5.4 | ... 9.5 | |
| | | IV | 2.4 | 3.2 | 4.1 | 4.7 | ... 4.5 | |
| Full $J$ | 20 | I | 4.0 | 6.2 | 8.2 | 10.0 | ... 12.1 | 12.1 |
| Trian $J$ | 20 | II | 4.1 | 6.4 | 6.7 | 7.7 | ... 11.9 | |
| | | III | 3.0 | 4.3 | 5.3 | 6.1 | ... 11.0 | |
| | | IV | 2.7 | 3.7 | 4.7 | 5.1 | ... 4.9 | |

## References

1 BUTCHER, J. C.: On the implementation of implicit Runge-Kutta methods. BIT 16 (1976), 237–240.
2 DAVISON, E. J.: An algorithm for the computer simulation of very large dynamical systems. Automatica 9 (1973), 665–675.
3 ENRIGHT, W. H.: Improving the efficiency of matrix operations in the numerical solution of stiff ordinary differential equations. ACM Trans. Math. Software 4 (1978), 127–136.
4 GOLUB, G. H.; VAN LOAN, C. F.: Matrix computations. Johns Hopkins University Press, Baltimore, MD 1989, 2nd ed.
5 HAIRER, E.; WANNER, G.: Solving ordinary differential equations. II: Stiff and differential-algebraic problems. Springer Series in Comput. Math., Vol. 14, Springer-Verlag, Berlin 1991.
6 HOUWEN, P. J. VAN DER; SOMMEIJER, B. P.: Iterated Runge-Kutta methods on parallel computers. SIAM J. Sci. Statist. Comput. 12 (1991), 1000–1028.
7 HOUWEN, P. J. VAN DER; SWART J. J. B. DE: Triangularly implicit iteration methods for ODE-IVP solvers. CWI Report NM-R9510, submitted for publication 1995.
8 SCHÄFER, E.: A new approach to explain the 'High Irradiance Responses' of photomorphogenesis on the basis of phytochrome. J. of Math. Biology 2 (1975), 41–56.
9 STREHMEL, K.; WEINER, R.: Numerik gewöhnlicher Differentialgleichungen. Teubner Studienbücher: Mathematik, Teubner, Stuttgart 1995.
10 VERWER, J. G.: Gauss-Seidel iteration for stiff ODEs from chemical kinetics. SIAM J. Sci. Comput. 15 (1994), 1243–1250.

*Address*: Prof. Dr. P. J. VAN DER HOUWEN, Dr. B. P. SOMMEIJER, CWI, P.O. Box 94079, NL-1090 GB Amsterdam, The Netherlands