



ELSEVIER

Applied Numerical Mathematics 28 (1998) 127–141



APPLIED
NUMERICAL
MATHEMATICS

Parallel implementation of block boundary value methods on nonlinear problems: theoretical results [☆]

Luigi Brugnano ^{a,*}, Donato Trigiante ^{b,1}

^a *Dipartimento di Matematica “U. Dini”, Viale Morgagni 67/A, I-50134 Firenze, Italy*

^b *Dipartimento di Energetica “S. Stecco”, Via C. Lombroso 6/17, I-50134 Firenze, Italy*

Received 4 September 1997; received in revised form 16 September 1997; accepted 16 September 1997

Dedicated to Professor Mario R. Occorsio in occasion of his 65th birthday

Abstract

Recently a new parallel ODE solver implementing a “parallelism across the steps” has been proposed (Amodio and Brugnano, 1997; Brugnano and Trigiante, 1998). In the mentioned references, the attention was devoted to some essential features of the parallel method, which are already present in the case where it is used to approximate linear continuous problems. In this paper, the previous analysis is completed by discussing questions which typically arise when approximating nonlinear continuous problems. In particular, we shall study, for the parallel solver, the problem of the mesh selection and the convergence of the nonlinear iteration. © 1998 Elsevier Science B.V. and IMACS. All rights reserved.

Keywords: Parallel methods for ODEs; Nonlinear Gauss–Seidel iteration; Mesh selection; Modified Newton method; Boundary Value Methods

1. Introduction

The study of new methods for the numerical solution of initial value problems (IVPs) for ODEs,

$$y' = f(t, y), \quad t \in (t_0, T], \quad y(t_0) = \eta \in \mathbb{R}^m, \quad (1)$$

suitable for parallel computers has been the object of much research in the last thirty years (see [11] for a complete overview). In particular, methods trying to approximate the solution of problem (1) simultaneously at different grid points, are classified as having a *parallelism across the steps*. A new method implementing this kind of parallelism has been recently proposed in [1,2,9]. The parallel solver is based on block Boundary Value Methods (B₂VMs), recently introduced by the authors for

^{*} Work supported by CNR (contract no. 96.00243.CT01) and MURST.

^{*} Corresponding author. E-mail: na.brugnano@na-net.ornl.gov.

¹ E-mail: na.dtrigiante@na-net.ornl.gov.

the approximation of Hamiltonian problems [8–10]. The same methods have been also considered for constructing a very efficient sequential code [12], which compares well with the most reliable existing ones (for example, RADAU5).

The analysis carried out in [1,2] concerns the application of the methods to linear continuous problems. In such case, in fact, one may not consider, in a simplified analysis, matters which arise when dealing with nonlinear problems. Among them, the most important are the mesh selection and the choice of the starting approximation to the solution for the iterative process involved. The aim of the present paper is the study of such problems. As a result, a complete scheme for the parallel solver will be derived. The actual implementation of the proposed method on a parallel computer, along with the analysis of its parallel performances, will be considered in a companion paper [3].

The structure of the paper is the following: in Section 2 we recall the main facts about B_2 VMs and their parallel implementation; Section 3 is devoted to the mesh selection and the derivation of the starting approximation for the nonlinear iteration, whose convergence is studied in Section 4. Finally, Section 5 contains some numerical tests carried out by using a Matlab prototype of the parallel solver.

2. Block Boundary Value Methods (B_2 VMs)

For a complete description of BVMs we refer to the book [9]. Here we shall confine ourselves to a very short introduction.

2.1. Boundary Value Methods (BVMs)

The shortest way of introducing Boundary Value Methods (BVMs) is to consider the approximation of problem (1) by using a k -step LMF,

$$\sum_{i=0}^k \alpha_i y_{n+i} = h \sum_{i=0}^k \beta_i f_{n+i}, \quad (2)$$

where, as usual, $h = (T - t_0)/N$ is the stepsize, y_n is the discrete approximation to $y(t_n)$ and $f_n = f(t_n, y_n)$, $t_n = t_0 + nh$. To get a discrete problem from Eq. (2), k independent conditions must be imposed. They are chosen by fixing the first $k_1 \leq k$ initial points, $y_0, y_1, \dots, y_{k_1-1}$, of the discrete solution, and $k_2 = k - k_1$ final ones, y_{N-k_2+1}, \dots, y_N . In other words, we are approximating the continuous IVP (1) by means of a discrete boundary value problem. The latter defines a BVM with (k_1, k_2) -boundary conditions, which may be stable and have arbitrarily high order [6]. As matter of fact, the Dahlquist barriers are overcome, when $k_1 < k$.

In particular, we shall consider the so called *Generalized Adams Methods* (GAMs) [4,7,9], which are methods having the following form:

$$y_n - y_{n-1} = h \sum_{i=-v}^{k-v} \beta_{i+v} f_{n+i}, \quad n = v, \dots, N - k + v, \quad (3)$$

where

$$v = \begin{cases} (k+1)/2, & \text{for odd } k, \\ k/2, & \text{for even } k. \end{cases}$$

For all $k \geq 1$, the coefficients $\{\beta_i\}$ are uniquely determined by imposing that the formula has order $k + 1$. The resulting method must be used with $(\nu, k - \nu)$ -boundary conditions. We observe that for $k = 1$ we get the usual trapezoidal rule.

In principle, the boundary conditions for (3) require the following values:

$$y_0, y_1, \dots, y_{\nu-1}, y_{N-k+\nu+1}, \dots, y_N,$$

of the discrete solution. Of such values, only y_0 is available from the initial condition of the continuous problem. However, we can treat the remaining values as unknowns, provided that an equal number of equations is introduced. This is done by means of a suitable set of *additional initial* and *final methods*. For GAMs they are conveniently chosen in the following form:

$$y_j - y_{j-1} = h \sum_{i=0}^k \beta_i^{(j)} f_i, \quad j = 1, \dots, \nu - 1, \tag{4}$$

$$y_j - y_{j-1} = h \sum_{i=0}^k \beta_{k-i}^{(j)} f_{N-i}, \quad j = N - k + \nu + 1, \dots, N. \tag{5}$$

The coefficients of the additional methods (4), (5) are uniquely determined by imposing that each formula has the same order $k + 1$ of the *main method* (3).

2.2. Block version of the methods and parallel implementation

From the above arguments one has that a BVM can be considered as a composite method made up by a main method coupled with an appropriate set of additional methods. An instance is given by (3)–(5). In this formulation, the method only requires one value, i.e., y_0 . This feature has led to define a block version of BVMs [8], which has been successfully used for approximating Hamiltonian problems.

Thus, the block version consists in discretizing the interval $[t_0, T]$ by using two different meshes: a coarser one and a finer one. Let the coarser mesh contain the $p + 1$ points

$$\tau_i = \tau_{i-1} + \widehat{h}_i, \quad i = 1, \dots, p, \quad \tau_0 \equiv t_0, \quad \tau_p \equiv T.$$

Then, on each subinterval $[\tau_{i-1}, \tau_i]$, $i = 1, \dots, p$, we apply the same (composite) BVM as described above, by using the *finer stepsize* $h_i = \widehat{h}_i/s$.

As a consequence, the points in the finer mesh belonging to the subinterval $(\tau_{i-1}, \tau_i]$, are given by

$$t_{ji} = \tau_{i-1} + jh_i, \quad j = 1, \dots, s, \quad i = 1, \dots, p. \tag{6}$$

They are called *internal steps* and the rightmost lower index of t_{ji} identifies the i th subinterval. In the following it will be convenient to use the notation t_{0i} to denote τ_{i-1} . In this case, it must be observed that

$$t_{01} \equiv \tau_0, \quad t_{0i} \equiv \tau_{i-1} \equiv t_{s,i-1}, \quad i = 1, \dots, p. \tag{7}$$

The global discrete problem is then given by

$$(A \otimes I_m) \mathbf{y} - (H \cdot B \otimes I_m) \mathbf{f} = \boldsymbol{\eta}, \tag{8}$$

where the following notation has been used: I_m is the identity matrix of size m (the size of the continuous problem), $\boldsymbol{\eta} = \mathbf{e}_1 \otimes \boldsymbol{\eta}$, \mathbf{e}_1 is the first unit vector in \mathbb{R}^{ps+1} , and

$$\mathbf{y} = (y_0^T \quad y_1^T \quad \dots \quad y_p^T)^T, \quad \mathbf{f} = (f_0^T \quad f_1^T \quad \dots \quad f_p^T)^T,$$

the number of steps k of the main method and/or the blocksize s are large enough. This is welcome since the higher k , the higher the order of the method.

When problem (1) is nonlinear, a corresponding linear system having the same structure as in the linear case can still be obtained by applying, for example, the Newton method for solving (8). In this case, however, the following two problems arise:

- (1) determine an efficient way for choosing the appropriate mesh (i.e., to define the finer stepsizes h_1, \dots, h_p);
- (2) get a good starting point, say $\mathbf{y}^{(0)}$, for the nonlinear iteration.

As matter of fact, both requirements are crucial for the convergence of the nonlinear iteration.

3. Diagonally Linearized Gauss–Seidel Method (DLGSM)

In order to solve the above mentioned problems, we propose to use a very cheap sequential method. This allows to get the mesh and the starting approximation $\mathbf{y}^{(0)}$. After that, a more accurate discrete solution is obtained by applying the simplified Newton iteration,

$$M\mathbf{z}^{(i)} = \mathbf{g}^{(i)}, \quad \mathbf{y}^{(i+1)} = \mathbf{y}^{(i)} - \mathbf{z}^{(i)}, \quad i = 0, 1, \dots, \tag{11}$$

where (see (8))

$$\begin{aligned} \mathbf{g}^{(i)} &= \mathbf{g}(\mathbf{y}^{(i)}) \equiv (A \otimes I_m)\mathbf{y}^{(i)} - (H \cdot B \otimes I_m)\mathbf{f}^{(i)} - \boldsymbol{\eta}, \\ M &= (A \otimes I_m) - (H \cdot B \otimes I_m)\mathbf{J}^{(0)}, \end{aligned}$$

and $\mathbf{J}^{(0)}$ is the block diagonal matrix with the Jacobians of the function f at $\mathbf{y}^{(0)}$.

In Section 4 we study the convergence of the *outer iteration* (11). For the moment, we analyze the starting procedure for determining the mesh points and the vector $\mathbf{y}^{(0)}$. This is achieved by using the trapezoidal rule, which has stability properties similar to those of higher order GAMs, implemented in block form. Namely (see (6)), by considering s consecutive steps with the same stepsize h_i . Since this is a sequential procedure, we can study it over any subinterval $[\tau_{i-1}, \tau_i]$ of the coarse mesh. For this reason, in order to simplify the notation, we shall skip the index for the subinterval itself.

Let τ be a generic point where the approximate solution η_τ is known, and let h be the stepsize. The discrete problem defined by the trapezoidal rule at the points $t_i = \tau + ih, i = 0, \dots, s$, is then given by

$$(A_s \otimes I_m)\mathbf{y}_\tau - h(B_s \otimes I_m)\mathbf{f}_\tau = e_1 \otimes \eta_\tau, \tag{12}$$

where e_1 is the first unit vector in \mathbb{R}^{s+1} , and

$$\mathbf{y}_\tau = \begin{pmatrix} y_0 \\ \vdots \\ y_s \end{pmatrix}, \quad \mathbf{f}_\tau = \begin{pmatrix} f_0 \\ \vdots \\ f_s \end{pmatrix}, \quad B_s = \frac{1}{2} \begin{pmatrix} 0 & & & & \\ 1 & 1 & & & \\ & \ddots & \ddots & & \\ & & & \ddots & \ddots \\ & & & & 1 & 1 \end{pmatrix}_{(s+1) \times (s+1)}$$

Eq. (12), which we assume to have a unique solution, will be solved by using a Gauss–Seidel method with diagonal linearization, DLGSM hereafter (Diagonally Linearized Gauss–Seidel Method). By introducing the matrices

$$\tilde{I}_s = I_{s+1} - e_1 e_1^T, \quad L_s = 2B_s - \tilde{I}_s,$$

and denoting by J_0 the Jacobian of f at the initial point, it is defined by the following splitting:

$$\left(A_s \otimes I_m - \frac{h}{2} \tilde{I}_s \otimes J_0 \right) \mathbf{y}_\tau - \frac{h}{2} (L_s \otimes I_m) \mathbf{f}_\tau = \mathbf{e}_1 \otimes \eta_\tau + \frac{h}{2} ((\tilde{I}_s \otimes I_m) \mathbf{f}_\tau - (\tilde{I}_s \otimes J_0) \mathbf{y}_\tau). \quad (13)$$

We then obtain the iterative scheme

$$\begin{aligned} & \left(A_s \otimes I_m - \frac{h}{2} \tilde{I}_s \otimes J_0 \right) \mathbf{y}_\tau^{(j)} - \frac{h}{2} (L_s \otimes I_m) \mathbf{f}_\tau^{(j)} \\ & = \mathbf{e}_1 \otimes \eta_\tau + \frac{h}{2} ((\tilde{I}_s \otimes I_m) \mathbf{f}_\tau^{(j-1)} - (\tilde{I}_s \otimes J_0) \mathbf{y}_\tau^{(j-1)}), \quad j = 1, \dots, g, \end{aligned} \quad (14)$$

where the vectors $\mathbf{y}_\tau^{(j)}$ and $\mathbf{f}_\tau^{(j)}$ contain the approximations to \mathbf{y}_τ and \mathbf{f}_τ (see (12)) at the j th iterate, and g is the maximum number of iterations allowed (see later).

3.1. Convergence of DLGSM

Let us now study the convergence of the iteration (14). By defining

$$\mathbf{e}_\tau^{(j)} = \mathbf{y}_\tau - \mathbf{y}_\tau^{(j)}, \quad \delta \mathbf{f}_\tau^{(j)} = \mathbf{f}_\tau - \mathbf{f}_\tau^{(j)},$$

from (13) and (14) we obtain

$$\begin{aligned} & \left(A_s \otimes I_m - \frac{h}{2} \tilde{I}_s \otimes J_0 \right) \mathbf{e}_\tau^{(j)} - \frac{h}{2} (L_s \otimes I_m) \delta \mathbf{f}_\tau^{(j)} = \frac{h}{2} ((\tilde{I}_s \otimes I_m) \delta \mathbf{f}_\tau^{(j-1)} - (\tilde{I}_s \otimes J_0) \mathbf{e}_\tau^{(j-1)}), \\ & j = 1, \dots, g. \end{aligned} \quad (15)$$

By the way, we observe that if the function $f(t, y)$ is independent of y , namely (1) is a pure quadrature problem, then $\delta \mathbf{f}_\tau^{(j)} = \mathbf{0}$ and $J_0 = 0$. Consequently, one has $\mathbf{e}_\tau^{(1)} = \mathbf{0}$, independently of the stepsize h . The same happens for linear autonomous problems.

Assuming now that the function f is suitably smooth, it follows that

$$\delta \mathbf{f}_\tau^{(j)} = \begin{pmatrix} J_0 & & \\ & \ddots & \\ & & J_s \end{pmatrix} \mathbf{e}_\tau^{(j)} + O(\|\mathbf{e}_\tau^{(j)}\|^2) = (I_{s+1} \otimes J_0 + h D_s \otimes J_0' + O(h^2)) \mathbf{e}_\tau^{(j)} + O(\|\mathbf{e}_\tau^{(j)}\|^2),$$

where J_i is the Jacobian of f at $(t_i, y(t_i))$, J_0' is the derivative of the Jacobian at $(t_0, y(t_0))$, and $D_s = \text{diag}(0, 1, \dots, s)$. One then concludes that, if the procedure converges, Eq. (15) can be approximated as

$$(A_s \otimes I_m - h B_s \otimes J_0) \mathbf{e}_\tau^{(j)} = \frac{h}{2} ((h D_s \otimes J_0') \mathbf{e}_\tau^{(j-1)} + O(\|\mathbf{e}_\tau^{(j-1)}\|^2)), \quad j = 1, \dots, g. \quad (16)$$

The coefficient matrix $(A_s \otimes I_m - h B_s \otimes J_0)$ is nonsingular and has a bounded inverse, when J_0 has no eigenvalues with positive real part. The same applies when J_0 has eigenvalues with positive real part, provided that h is small enough. For sake of simplicity, we shall then assume that the above mentioned matrix is nonsingular and with bounded inverse. By considering norms, we finally obtain

$$\|\mathbf{e}_\tau^{(j)}\| \leq c_1 h^2 \|\mathbf{e}_\tau^{(j-1)}\| + c_2 h \|\mathbf{e}_\tau^{(j-1)}\|^2, \quad j \geq 1, \quad (17)$$

for suitable constants c_1 and c_2 independent of h .

Eq. (17) allows us to derive the following conclusions concerning the convergence of DLGSM:

1. When $\|\mathbf{e}_\tau^{(j-1)}\|$ is large with respect to h , then the first term on the right hand side is negligible.

Quadratic convergence is then to be expected in the initial iterations.

2. Conversely, as soon as the first term dominates, the iteration converges linearly.

By taking into account that the local error of the trapezoidal rule is $O(h^3)$, the above arguments suggest that the number g of iterations for the procedure (14) must be of moderate size. In fact, if we consider the starting vector

$$\mathbf{y}_\tau^{(0)} = (1, 1, \dots, 1)^T \otimes \eta_\tau, \quad (18)$$

the assumptions on the function f imply that $\|\mathbf{e}_\tau^{(0)}\| \approx O(h)$. Consequently, from (17) we obtain $\|\mathbf{e}_\tau^{(j)}\| \approx O(h^{2j+1})$, $j = 0, 1, \dots, g$. Therefore, by using $g = 2$ we can expect that the error associated with the convergence of iteration (14) is smaller than the local error. However, the stepsize control described later, based on the estimates

$$\|\mathbf{e}_\tau^{(j)}\| \approx \|\mathbf{y}_\tau^{(j+1)} - \mathbf{y}_\tau^{(j)}\|, \quad (19)$$

needs $g = 3$ iterations for the DLGSM.

Concerning the convergence of DLGSM, it can be established by using the comparison principle [13]. The following sufficient condition is then derived.

Theorem 1. *For all sufficiently small $h > 0$, the DLGSM iteration is convergent.*

In order to meet the requirements of Theorem 1, we start DLGSM by using a suitably small initial stepsize. This stepsize will be then changed, in the subsequent subintervals, depending on the convergence of the method. This is discussed in the following section.

3.2. Convergence of DLGSM and stepsize selection

In this section the problem of stepsize selection is analyzed in more detail. Namely, we study the determination of the new stepsize h_{new} , depending on the current stepsize h and on the convergence of the DLGSM iteration. The following arguments assume (see (18)) that, for a suitable $\beta > 0$, $\|\mathbf{e}_\tau^{(0)}\| \approx \beta h$. In our analysis, we shall distinguish the two cases: (i) $\beta \approx 1$, and (ii) $\beta \ll 1$. Moreover, by considering the comparison equation associated with (17), namely

$$x_j = c_1 h^2 x_{j-1} + c_2 h x_{j-1}^2, \quad j \geq 1, \quad x_0 = \|\mathbf{e}_\tau^{(0)}\|, \quad (20)$$

we shall suppose that (see (19)), $x_j \approx \|\mathbf{e}_\tau^{(j)}\| \approx \|\mathbf{y}_\tau^{(j+1)} - \mathbf{y}_\tau^{(j)}\|$.

In case (i), by performing $g = 3$ iterations of DLGSM, we obtain estimates for x_0 , x_1 and x_2 . Consequently, from (20) we get

$$x_3 \approx \beta c_1^2 (c_1 + \beta c_2) h^7 \approx \frac{x_2^2}{x_1}.$$

Having fixed a tolerance ε for x_3 , one then obtains that the stepsize

$$h_{\text{new}} = h \left(\varepsilon \frac{x_1}{x_2} \right)^{1/7} \quad (21)$$

can be used for the DLGSM iteration in the subsequent subinterval (or to repeat the current iteration, if $x_3 > \varepsilon$).

In case (ii), again from (20), we obtain that, for $j \geq 1$, $x_j \approx c_1 h^2 x_{j-1}$. Consequently, one has

$$x_3 \approx c_1^3 h^6 x_0 \approx \frac{x_1 x_2}{x_0}. \quad (22)$$

This allows to predict the new stepsize as

$$h_{\text{new}} = h \left(\varepsilon \frac{x_0}{x_1 x_2} \right)^{1/6}. \quad (23)$$

We observe that, if x_0 is small, one could have $x_j \leq \varepsilon$ for $j \leq 2$. In the case $j = 2$, the estimate (22) allows to predict the new stepsize by means of (23). Obviously, we must perform at least two iterations, in order to predict the new stepsize. This may be not possible when both c_1 and c_2 are close to zero, as in the case, for example, of pure quadrature problems (i.e., when $c_1 = c_2 = 0$). This matter will be discussed later. For the moment, we assume that at least two iterations are performed. In such a case, we obtain

$$x_3 \approx \frac{x_1^3}{x_0^2},$$

and, consequently, the new stepsize is predicted as

$$h_{\text{new}} = h \left(\varepsilon \frac{x_0^2}{x_1^3} \right)^{1/6}. \quad (24)$$

In the previous relations, by considering $\varepsilon = \text{tol} \cdot x_0$, one obtains a control of the relative growth of the error. Namely, the stepsize is predicted by imposing

$$\frac{x_3}{x_0} = \text{tol}. \quad (25)$$

3.3. Accuracy control

The most commonly used strategy for the stepsize selection is based on the control of the truncation error. In the case of the trapezoidal rule, it amounts to determine the stepsize in order to have

$$\|T_\tau\| \leq \widetilde{\text{tol}}, \quad (26)$$

where T_τ is the vector whose entries are the truncation errors at t_0, \dots, t_s , and $\widetilde{\text{tol}}$ is the given tolerance. By using the estimate

$$\|T_\tau\| \approx \frac{s}{12} \max_i \|y^{(3)}(t_i)\| h^3,$$

we then obtain the new stepsize

$$\widetilde{h}_{\text{new}} = \left(\frac{12 \cdot \widetilde{\text{tol}}}{s \cdot \max_i \|y^{(3)}(t_i)\|} \right)^{1/3}. \quad (27)$$

Approximations for the values $y^{(3)}(t_i)$ are obtainable from componentwise second divided differences of the f_i , $i = 0, \dots, s$, computed during the DLGSM iteration.

The previous estimate (27) should then be compared with the stepsize predicted by the convergence of DLGSM (see (25)). Since the convergence must be always guaranteed, it should be

$$h_{\text{new}} \leq \widetilde{h}_{\text{new}}, \quad (28)$$

for a suitable tolerance $\widetilde{\text{tol}}$. If this were always the case, the cost for computing $\widetilde{h}_{\text{new}}$ could be avoided. Unfortunately, there are some cases where the above condition is not satisfied. For example, when the function f is independent of y , or it is autonomous and linear, one obtains $x_1 = 0$, independently of the stepsize used, whereas there is a restriction on $\widetilde{h}_{\text{new}}$. It is then evident that, in such cases, the strategy (27) is more restrictive.

In the remaining part of this section, we shall study the cases where (28) holds true. Let us first suppose, for simplicity, that $f = f(y)$ and nonlinear. Then, from the previous analysis, (25) implies that

$$c_1^3 h^6 \approx \frac{x_3}{x_0} = \text{tol},$$

where, for h sufficiently small (see (16) and (17)),

$$c_1 \approx \|(A_s \otimes I_m - h B_s \otimes J_0)^{-1}\| \frac{1}{2} \|D_s \otimes J'_0\| \approx \frac{s^2 m}{2} \|J'_0\|.$$

It follows that

$$\frac{s^6 m^3}{8} h^6 \left\| \frac{\partial J_0}{\partial y} f_0 \right\|^3 \approx \text{tol},$$

that is,

$$h^3 \approx \frac{\sqrt{8}}{s^3 m^{3/2}} \left\| \frac{\partial J_0}{\partial y} f_0 \right\|^{-3/2} \sqrt{\text{tol}}. \tag{29}$$

On the other hand, (26) would approximately require

$$\frac{s}{12} h^3 \left\| \left(\frac{\partial J_0}{\partial y} f_0 \right) f_0 + J_0^2 f_0 \right\| \leq \widetilde{\text{tol}},$$

which, from (29), gives

$$\frac{\sqrt{2}}{6 \cdot s^2 m^{3/2}} \frac{\|((\partial J_0 / \partial y) f_0) f_0 + J_0^2 f_0\|}{\|(\partial J_0 / \partial y) f_0\|^{3/2}} \sqrt{\text{tol}} \leq \widetilde{\text{tol}}.$$

It is then evident that the previous inequality is satisfied with $\widetilde{\text{tol}} \approx \sqrt{\text{tol}}$, provided that

$$\frac{\sqrt{2}}{6 \cdot s^2 m^{3/2}} \frac{\|((\partial J_0 / \partial y) f_0) f_0 + J_0^2 f_0\|}{\|(\partial J_0 / \partial y) f_0\|^{3/2}} \tag{30}$$

is not too large.

It can be shown that the above result continues to hold when f is not autonomous, but with a “moderate” dependence on t . We then conclude that inequality (25) is sufficient for selecting the stepsize when problem (1) is autonomous (or “almost” autonomous), provided that (30) is not very large. In all other cases, the new stepsize will be that given by (27).

The above arguments can be summarized as follows:

(1) when both

$$\frac{df}{dt}(t, y) \approx J(t, y) f(t, y), \quad \frac{\|((\partial J_0 / \partial y) f_0) f_0 + J_0^2 f_0\|}{\|(\partial J_0 / \partial y) f_0\|^{3/2}} \leq \nu_1, \tag{31}$$

hold true, where ν_1 is suitably chosen, we use the requirement (25) for the mesh selection;

(2) conversely, we use (26).

In practice, (31) is assumed to hold true when both

$$\frac{\|f_1 - f_0\|}{h} \leq 1.1 \|J_0 f_0\|, \quad \frac{\|f_0''\|}{(\|f_0'' - J_0^2 f_0\| / \|f_0\|)^{3/2}} \leq \nu_1, \quad (32)$$

are satisfied. Here f_0'' is the second derivative of f at (t_0, y_0) . Conversely, (31) is assumed to be false. As a matter of fact, for the numerical tests in Section 5 the two above requirements are fulfilled at almost all grid points. Consequently, the mesh is essentially always selected through the convergence of DLGSM.

We observe that the checks (32) only require to know information at the initial point. This implies that, in the actual implementation of the above procedure on a parallel computer [3], the corresponding computational load can be done in parallel by a different processor. On the contrary, the computation of \tilde{h}_{new} needs to be done sequentially.

4. Convergence of the nonlinear iteration

Next theorem is a useful tool in studying the convergence of the iteration (11), used for solving

$$\mathbf{g}(\mathbf{y}) \equiv (A \otimes I_m) \mathbf{y} - (H \cdot B \otimes I_m) \mathbf{f} - \boldsymbol{\eta} = \mathbf{0}.$$

Theorem 2. Assume that

- (a) \mathbf{g} is differentiable in $\mathcal{D}_\alpha = \{\mathbf{x}: \|\mathbf{x} - \bar{\mathbf{y}}\| \leq \alpha\}$;
- (b) $\mathbf{y}^{(0)} \in \mathcal{D}_\alpha$;
- (c) $\|M^{-1}(\mathbf{g}'(\mathbf{y}) - \mathbf{g}'(\mathbf{x}))\mathbf{v}\| \leq \gamma \|\mathbf{y} - \mathbf{x}\| \|\mathbf{v}\|$, for all $\mathbf{x}, \mathbf{y} \in \mathcal{D}_\alpha$ and \mathbf{v} such that $\|\mathbf{y} - \mathbf{x} - \mathbf{v}\| \leq \alpha$;
- (d) $\theta \equiv \frac{5}{2}\alpha\gamma < 1$.

Then the simplified Newton iteration (11) converges to the solution $\bar{\mathbf{y}}$ at least linearly.

Proof. We recall that the matrix M in (11) is nothing but $\mathbf{g}'(\mathbf{y}^{(0)})$. The thesis is then proved by showing that

$$\|\mathbf{y}^{(i)} - \bar{\mathbf{y}}\| \leq \theta^i \|\mathbf{y}^{(0)} - \bar{\mathbf{y}}\|.$$

For $i = 1$, in fact, we have

$$\begin{aligned} \|\mathbf{y}^{(1)} - \bar{\mathbf{y}}\| &= \|\mathbf{y}^{(0)} - \bar{\mathbf{y}} - M^{-1}(\mathbf{g}(\mathbf{y}^{(0)}) - \mathbf{g}(\bar{\mathbf{y}}))\| \\ &= \|M^{-1}(\mathbf{g}'(\mathbf{y}^{(0)})(\mathbf{y}^{(0)} - \bar{\mathbf{y}}) - (\mathbf{g}(\mathbf{y}^{(0)}) - \mathbf{g}(\bar{\mathbf{y}})))\| \\ &\leq \int_0^1 \|M^{-1}(\mathbf{g}'(\mathbf{y}^{(0)}) - \mathbf{g}'(\bar{\mathbf{y}} + t(\mathbf{y}^{(0)} - \bar{\mathbf{y}})))(\mathbf{y}^{(0)} - \bar{\mathbf{y}})\| dt \\ &\leq \int_0^1 \gamma(1-t) \|\mathbf{y}^{(0)} - \bar{\mathbf{y}}\|^2 dt \leq \frac{1}{2}\alpha\gamma \|\mathbf{y}^{(0)} - \bar{\mathbf{y}}\| < \theta \|\mathbf{y}^{(0)} - \bar{\mathbf{y}}\|. \end{aligned}$$

By induction, we now suppose true the result for i . Then, for $i + 1$, we get

$$\begin{aligned} \|\mathbf{y}^{(i+1)} - \bar{\mathbf{y}}\| &= \|\mathbf{y}^{(i)} - \bar{\mathbf{y}} - M^{-1}(\mathbf{g}(\mathbf{y}^{(i)}) - \mathbf{g}(\bar{\mathbf{y}}))\| \\ &= \|M^{-1}(\mathbf{g}'(\mathbf{y}^{(0)})(\mathbf{y}^{(i)} - \bar{\mathbf{y}}) - (\mathbf{g}(\mathbf{y}^{(i)}) - \mathbf{g}(\bar{\mathbf{y}})))\| \end{aligned}$$

$$\begin{aligned}
&\leq \|M^{-1}(\mathbf{g}'(\mathbf{y}^{(0)}) - \mathbf{g}'(\mathbf{y}^{(i)}))(\mathbf{y}^{(i)} - \bar{\mathbf{y}})\| \\
&\quad + \|M^{-1}(\mathbf{g}'(\mathbf{y}^{(i)})(\mathbf{y}^{(i)} - \bar{\mathbf{y}}) - (\mathbf{g}(\mathbf{y}^{(i)}) - \mathbf{g}(\bar{\mathbf{y}})))\| \\
&\leq \gamma \|\mathbf{y}^{(0)} - \mathbf{y}^{(i)}\| \|\mathbf{y}^{(i)} - \bar{\mathbf{y}}\| \\
&\quad + \int_0^1 \|M^{-1}(\mathbf{g}'(\mathbf{y}^{(i)}) - \mathbf{g}'(\bar{\mathbf{y}} + t(\mathbf{y}^{(i)} - \bar{\mathbf{y}})))(\mathbf{y}^{(i)} - \bar{\mathbf{y}})\| dt \\
&\leq \gamma (\|\mathbf{y}^{(0)} - \bar{\mathbf{y}}\| + \|\mathbf{y}^{(i)} - \bar{\mathbf{y}}\|) \|\mathbf{y}^{(i)} - \bar{\mathbf{y}}\| + \frac{1}{2} \gamma \|\mathbf{y}^{(i)} - \bar{\mathbf{y}}\|^2 \\
&< \frac{5}{2} \gamma \|\mathbf{y}^{(0)} - \bar{\mathbf{y}}\| \|\mathbf{y}^{(i)} - \bar{\mathbf{y}}\| \leq \theta \|\mathbf{y}^{(i)} - \bar{\mathbf{y}}\|. \quad \square
\end{aligned}$$

Remark 1. It is worth noting that, since all GAMs have similar stability properties, the conditioning of the matrix M is almost independent of the particular method considered in this class. Consequently, it is reasonable to expect that the above parameter θ does not vary sensibly for such methods. This fact is actually observed in practice. In fact, the simplified Newton iteration (11) requires approximately the same number of steps to have the same stopping criterion satisfied, independently of the method chosen.

The result of Theorem 2 is of practical interest, since we can get cheap estimates for the two parameters α and γ during the execution of DLGSM (see next section). If necessary, this allows to split the interval of integration into suitable “windows”, where the convergence of the simplified Newton iteration (11) is assured in a given number of steps.

In other words, having fixed an upper bound θ_{\max} for the parameter θ , we use DLGSM until the estimate for θ is smaller than θ_{\max} . As soon as this bound is exceeded, we stop DLGSM, thus determining the first “window”, where the simplified Newton iteration (11) (i.e., the parallel section of the method) is carried out until convergence. After that, we start DLGSM again, thus repeating the above procedure until the whole interval of integration is covered.

4.1. Estimates for the parameters α and γ

In this section, we shall get estimates for both parameters α and γ defined in Theorem 2. They are obtained dynamically during the execution of DLGSM. In particular, we show that to each subinterval $[\tau_{i-1}, \tau_i]$ of the coarse mesh we can associate a couple of parameters (α_i, γ_i) which take into account of all the integration up to τ_i .

Suppose, for simplicity, that only one window is needed. We start considering the first parameter, which, in practice, is the measure of the global error up the current subinterval. In fact, by considering the vector $\bar{\mathbf{y}}$ whose (block) entries are the values of the continuous solution at the grid points, one has

$$\alpha = \|\bar{\mathbf{y}} - \mathbf{y}^{(0)}\| \approx \|\hat{\mathbf{y}} - \mathbf{y}^{(0)}\|.$$

This because $\bar{\mathbf{y}}$ is the solution of the discrete problem given by a higher order method than the trapezoidal rule, which provides us with the starting vector $\mathbf{y}^{(0)}$.

Let δ_i be the vector whose entries are the errors over the i th subinterval. It is then not difficult to realize that, at first order approximation, it satisfies the equation (see (8) and (12))

$$(A_{i:s} \otimes I_m - (H_i B_{i:s} \otimes I_m) \mathbf{J}_i^{(0)}) \begin{pmatrix} 0 \\ \delta_1 \\ \vdots \\ \delta_i \end{pmatrix} = \begin{pmatrix} 0 \\ T_1 \\ \vdots \\ T_i \end{pmatrix}, \tag{33}$$

where $\mathbf{J}_i^{(0)}$ is the block diagonal matrix with the Jacobians of f at $(t_0, y_0^{(0)}), \dots, (t_{si}, y_{si}^{(0)})$, and the vectors $\{T_r\}$ contain the truncation errors over the corresponding subintervals of the coarse mesh. Consequently, by considering that (33) is a lower block bidiagonal linear system, one obtains the following recurrence (hereafter, see (7)), let J_{jr} be the Jacobian of f at $(t_{jr}, y_{jr}^{(0)})$, $J_{0r} \equiv J_{s,r-1}$, $r > 1$, and $J_{01} \equiv J_0$,

$$\begin{pmatrix} \delta_{s,r-1} \\ \delta_r \end{pmatrix} = \left(A_s \otimes I_m - h_r B_s \otimes I_m \begin{pmatrix} J_{0r} & & \\ & \ddots & \\ & & J_{sr} \end{pmatrix} \right)^{-1} \begin{pmatrix} \delta_{s,r-1} \\ T_r \end{pmatrix}, \quad r = 1, 2, \dots, i, \tag{34}$$

where $\delta_{s,r-1}$ is the last block entry of δ_{r-1} , and $\delta_{s0} = 0$. The vectors $\{\delta_r\}$ are then obtained, provided that estimates for the truncation errors $\{T_r\}$ are available. This can be done, as previously said, by means of componentwise second divided difference of the function f . Alternatively, a *deferred correction approach* as described in [5,9] may be used.

In practice, in order to reduce the computational effort, the above recurrence (34) is approximated by

$$\begin{pmatrix} \delta_{s,r-1} \\ \delta_r \end{pmatrix} = (A_s \otimes I_m - h_r B_s \otimes J_{0r})^{-1} \begin{pmatrix} \delta_{s,r-1} \\ T_r \end{pmatrix}, \quad r = 1, 2, \dots, i.$$

The coefficient matrix turns out to be block Toeplitz lower bidiagonal, whose diagonal block, $(I_m - \frac{1}{2}h_r J_{0r})$, has already been factored during the DLGSM iteration.

By using hereafter the infinity norm, we then conclude that to the i th subinterval there corresponds a parameter

$$\alpha_i = \max\{\alpha_{i-1}, \|\delta_i\|\}, \quad i \geq 1, \quad \alpha_0 = 0.$$

A similar approach can be used for the parameter γ . In fact, one verifies that

$$\mathbf{g}'(\mathbf{y}) - \mathbf{g}'(\mathbf{x}) = (H \cdot B \otimes I_m)(\mathbf{J}(\mathbf{x}) - \mathbf{J}(\mathbf{y})) \approx (H \otimes I_m)(\mathbf{J}(\mathbf{x}) - \mathbf{J}(\mathbf{y})), \tag{35}$$

where $\mathbf{J}(\mathbf{x})$ and $\mathbf{J}(\mathbf{y})$ are the block diagonal matrices, whose diagonal blocks are the Jacobians of the function f evaluated at the entries of \mathbf{x} and \mathbf{y} , respectively. The last approximation follows from the fact that the matrix B (see (9), (10)) has unit row sums.

We shall then obtain the estimate for γ by fixing a particular direction $\mathbf{x} - \mathbf{y}$, and a vector \mathbf{v} of unit norm. In practice, in the i th block we assume that the Jacobian matrices all coincide with J_{0i} , which corresponds to the initial condition, y_{0i} (i.e., $y_{s,i-1}$), for the current block. Then, after we have carried out the DLGSM iterations, in order to start the procedure in the subsequent block, we need to compute $J_{0,i+1} \equiv J_{si}$, namely the Jacobian at $y_{0,i+1} \equiv y_{si}$. As a consequence, by setting $E_s = (1, \dots, 1)^T \in \mathbb{R}^s$, we choose

$$\mathbf{z} = \begin{pmatrix} 0 \\ z_1 \\ \vdots \\ z_p \end{pmatrix} \equiv (\mathbf{J}(\mathbf{x}) - \mathbf{J}(\mathbf{y}))\mathbf{v} = \begin{pmatrix} 0 \\ E_s \otimes (J_{01} - J_{s1})v_1 \\ \vdots \\ E_s \otimes (J_{0p} - J_{sp})v_p \end{pmatrix},$$

where the $\{v_i\}$ are the block entries of \mathbf{v} (let $y_{0i} \equiv y_0$), defined as

$$v_0 = 0, \quad v_i = \frac{y_{0i} - y_{si}}{\|y_{0i} - y_{si}\|}, \quad i = 1, \dots, p.$$

Consequently, we determine the approximation for γ such that

$$\gamma \geq \frac{\|\mathbf{w}\|}{\|\mathbf{u}\|},$$

where (see (35))

$$\mathbf{w} \equiv M^{-1}(H \otimes I_m)\mathbf{z} \approx M^{-1}(\mathbf{g}'(\mathbf{y}) - \mathbf{g}'(\mathbf{x}))\mathbf{v},$$

and

$$\mathbf{u} = (0, u_1, \dots, u_p)^T, \quad u_i = E_s \otimes (y_{0i} - y_{si}), \quad i = 1, \dots, p.$$

By using as an approximation for M the matrix computed for the DLGSM iteration, we then obtain the following recurrence for the entries of the vector \mathbf{w} :

$$\begin{pmatrix} w_{s,r-1} \\ w_r \end{pmatrix} = (A_s \otimes I_m - h_r B_s \otimes J_{0r})^{-1} \begin{pmatrix} w_{s,r-1} \\ h_r z_r \end{pmatrix}, \quad r = 1, 2, \dots, i.$$

Consequently, we choose the following approximation to the parameter γ up to the i th subinterval:

$$\gamma_i = \begin{cases} \max \left\{ \gamma_{i-1}, \frac{\|w_i\|}{\|y_{s,i-1} - y_{si}\|} \right\}, & \text{if } \|y_{s,i-1} - y_{si}\| > 0, \\ \gamma_{i-1}, & \text{otherwise,} \end{cases} \quad i \geq 1,$$

having set $\gamma_0 = 0$.

We then conclude that, during the DLGSM iteration over the i th subinterval in the coarse mesh, we can obtain cheap estimates for α_i and γ_i and, therefore, for the parameter θ_i which determines the speed of convergence of the outer iteration (11). Consequently, when such estimate is larger than a fixed θ_{\max} , we end the sequential procedure based on the DLGSM. Namely, we have determined a *window* where we perform the simplified Newton process (11). As previously said, the whole procedure is then repeated until the entire interval of integration $[t_0, T]$ is covered.

5. Numerical examples

In this section we report some numerical tests obtained by using a Matlab prototype implementing the procedure previously described. In all cases, we use the 9th order GAM (GAM9), after the DLGSM iteration. The tolerance used for the DLGSM iteration is $\text{tol} = 10^{-6}$, with tolerance $\text{tol} = 10^{-3}$ in case of switch to the local error control. Moreover, the value θ_{\max} is chosen in order to have convergence for the outer iteration (11) in at most four iteration. The stopping criterion for the Newton iteration is

$$\max_j \frac{|y_j^{(i+1)} - y_j^{(i)}|}{1 + |y_j^{(i+1)}|} \leq 10^{-9},$$

where $y_j^{(i)}$ is the j th entry of $\mathbf{y}^{(i)}$. For comparison, we also report the results on the same problems by using the stiff solver ode23s, taken from the Matlab ODE Suite [15]. In order to have a comparable

Table 1
Numerical results

Number of	Robertson		van der Pol		Hires	
	GAM9	ode23s	GAM9	ode23s	GAM9	ode23s
windows	3	–	9	–	1	–
points	1010	7883	1851	86541	488	18210
fev	2530	23654	6280	259628	1560	54629
Jev	1010	7883	1851	86541	488	18210
flops	3.4×10^6	2.2×10^6	2.7×10^6	14.3×10^6	14.1×10^6	22.4×10^6

accuracy for both solvers, we have used ode23s with the following parameters: $\text{atol} = \text{rtol} = 10^{-11}$, and providing analytically the Jacobian. Such parameters have been obtained by comparing the solutions computed for two of the following problems at selected grid points, where the values of the exact solutions are known. For both solvers we report the number of mesh points, function and Jacobian evaluations, and the number of Matlab flops. For the prototype of the parallel solver we also report the number of windows required to cover the integration interval. Finally, we want to mention that the flops count for GAM9 includes some extra work required for getting an approximation to the global error.

The first problem we consider are the Robertson equations

$$\begin{aligned} y_1' &= -0.04y_1 + 10^4 y_2 y_3, \\ y_2' &= 0.04y_1 - 10^4 y_2 y_3 - 3 \times 10^7 y_2^2, \quad t \in [0, 10^{15}], \\ y_3' &= 3 \times 10^7 y_2^2, \\ y_1(0) &= 1, \quad y_2(0) = y_3(0) = 0. \end{aligned}$$

The second problem is given by the van der Pol equations

$$\begin{aligned} y_1' &= y_2, \\ y_2' &= -y_1 + \mu y_2(1 - y_1^2), \quad t \in [0, \mu], \\ y_1(0) &= 2, \quad y_2(0) = 0, \end{aligned}$$

where $\mu = 10^6$.

Finally, the third problem is the Hires problem, a stiff ODE of dimension 8, taken from [14]. In Table 1 we report the obtained results, as described above. It is worth mentioning that for the van der Pol problem, whose solution has a huge spike near $t = 806853$, more than one window is needed to cover the interval $[0, 10^6]$ (actually 9 of them), with the first window ending at about $t = 806009$. This confirms the reliability of the proposed procedure for detecting the intervals where the Newton iteration does converge.

References

- [1] P. Amodio and L. Brugnano, Parallel implementation of block boundary value methods for ODEs, *J. Comput. Appl. Math.* 78 (1997) 197–211.

- [2] P. Amodio and L. Brugnano, Parallel ODE solvers based on block BVMs, *Adv. Comput. Math.* 7 (1–2) (1997) 5–26.
- [3] P. Amodio and L. Brugnano, ParalleloGAM: a parallel code for ODEs, *Appl. Numer. Math.* 28 (1998) 95–106.
- [4] P. Amodio and F. Mazzia, A boundary value approach to the numerical solution of initial value problems by multistep methods, *J. Differential Equations Appl.* 1 (1995) 353–367.
- [5] L. Brugnano, Boundary value methods for the numerical approximation of ordinary differential equations, in: *Lecture Notes in Computer Science* 1196 (1997) 78–89.
- [6] L. Brugnano and D. Trigiante, Convergence and stability of boundary value methods for ordinary differential equations, *J. Comput. Appl. Math.* 66 (1996) 97–109.
- [7] L. Brugnano and D. Trigiante, Boundary value methods: the third way between linear multistep and Runge–Kutta methods, *Comput. Math. Appl.* (to appear).
- [8] L. Brugnano and D. Trigiante, Block boundary value methods for linear Hamiltonian systems, *Appl. Math. Comput.* 81 (1997) 49–68.
- [9] L. Brugnano and D. Trigiante, *Solving Differential Problems by Multistep Initial and Boundary Value Methods* (Gordon & Breach, Amsterdam, 1998).
- [10] L. Brugnano and D. Trigiante, On the potentiality of sequential and parallel codes based on extended trapezoidal rules (ETRs), *Appl. Numer. Math.* 25 (1997) 169–184.
- [11] K. Burrage, *Parallel and Sequential Methods for Ordinary Differential Equations* (Clarendon Press, Oxford, 1995).
- [12] F. Iavernaro and F. Mazzia, Solving ordinary differential equations by generalized Adams methods: properties and implementation techniques, submitted.
- [13] V. Lakshmikantham and D. Trigiante, *Theory of Difference Equations: Numerical Methods and Applications* (Academic Press, San Diego, CA, 1988).
- [14] W.M. Lioen, J.J.B. de Swart and W.A. van der Veen, Test set for IVP solvers, CWI Report NM-R9615 (August 1996).
- [15] L.F. Shampine and M.W. Reichelt, The Matlab ODE Suite, Report (1995).