



Continuous numerical methods for ODEs with defect control[☆]

W.H. Enright

Department of Computer Science, University of Toronto, Toronto, Ontario, Canada M5S 3G4

Received 7 July 1999; received in revised form 3 December 1999

Abstract

Over the last decade several general-purpose numerical methods for ordinary differential equations (ODEs) have been developed which generate a continuous piecewise polynomial approximation that is defined for all values of the independent variable in the range of interest. For such methods it is possible to introduce measures of the quality of the approximate solution based on how well the piecewise polynomial satisfies the ODE. This leads naturally to the notion of “defect-control”. Numerical methods that adopt error estimation and stepsize selection strategies in order to control the magnitude of the associated defect can be very effective and such methods are now being widely used. In this paper we will review the advantages of this class of numerical methods and present examples of how they can be effectively applied. We will focus on numerical methods for initial value problems (IVPs) and boundary value problems (BVPs) where most of the developments have been introduced but we will also discuss the implications and related developments for other classes of ODEs such as delay differential equations (DDEs) and differential algebraic equations (DAEs). © 2000 Elsevier Science B.V. All rights reserved.

1. Introduction

Consider the ordinary differential equation

$$y' = f(t, y), \quad (1)$$

with exact solution denoted by $y(t)$. Traditional discrete numerical methods partition the interval of interest, $[t_0, t_F]$, by introducing a mesh $t_0 < t_1 < \dots < t_N = t_F$ and generate a discrete approximation $y_i \approx y(t_i)$ for each associated meshpoint. The number of meshpoints, N , and the distribution of these meshpoints is usually determined adaptively by the method in an attempt to deliver acceptable accuracy at a minimum cost. These methods generally accomplish this objective by keeping N as small as possible subject to a constraint that an indirect measure of $\max_{i=1, \dots, N} \|y(t_i) - y_i\|$ be kept

E-mail address: enright@cs.toronto.edu (W.H. Enright).

[☆] This research was supported by the Natural Science and Engineering Research Council of Canada and by Communication and Information Technology Ontario.

small (relative to an accuracy parameter TOL). (We use $\|\cdot\|$ to represent the max–norm for vectors and the induced matrix norm for matrices.) Different methods implement very different strategies in an attempt to achieve this indirect error control and this can make it particularly challenging to interpret the accuracy of a numerical solution.

In recent years, the notion of a continuous extension (of an underlying discrete method) has received considerable attention. With this approach, one associates with each discrete approximation, $\{t_i, y_i\}_{i=1}^N$, a piecewise polynomial, $u(t)$, defined for all t in the interval $[t_0, t_F]$ and satisfying $u(t_i) = y_i$ for $i = 1, 2, \dots, N$. (For a detailed discussion of the advantages and costs of generating such extensions see [7,13,17,23,26]). For example, consider applying a method based on a standard s -stage, p th-order Runge–Kutta formula to (1). The corresponding discrete approximations will then satisfy

$$y_i = y_{i-1} + h_i \sum_{j=1}^s \omega_j k_j, \tag{2}$$

where

$$k_j = f \left(t_{i-1} + h_i c_j, y_{i-1} + h_i \sum_{r=1}^s a_{jr} k_r \right),$$

for $j = 1, 2, \dots, s$; $i = 1, 2, \dots, N$ and $h_i = t_i - t_{i-1}$. A continuous extension is derived by introducing the additional stages, $k_{s+1}, k_{s+2}, \dots, k_{\bar{s}}$ and polynomials $b_j(\tau)$ of degree $\leq p$ for $j = 1, 2, \dots, \bar{s}$ such that the polynomial $u_i(t)$ defined by

$$u_i(t) = y_{i-1} + h_i \sum_{j=1}^{\bar{s}} b_j \left(\frac{t - t_{i-1}}{h_i} \right) k_j \tag{3}$$

satisfies $u_i(t) = y(t) + O(h^p)$ for $t \in (t_{i-1}, t_i)$ and $h = \max_{i=1}^N h_i$. Formulas of this type are called continuous Runge–Kutta (CRK) formulas. The polynomials $\{u_i(t)\}_{i=1}^N$ then define a piecewise polynomial, $u(t)$, which will be continuous on $[t_0, t_F]$ and interpolate the underlying discrete approximation if $b_j(1) = w_j$ for $j = 1, 2, \dots, s$ and $b_{s+1}(1) = b_{s+2}(1) \cdots = b_{\bar{s}}(1) = 0$.

In deriving CRK formulas of order p , several issues arise which can have a significant impact on the implementation and hence on the performance of the resulting continuous method. If the extra stages are restricted to be “explicit”, that is if

$$k_{s+j} = f \left(t_{i-1} + c_{s+j} h_i, y_{i-1} + h_i \sum_{r=1}^{s+j-1} a_{s+j,r} k_r \right)$$

for $j = 1, 2, \dots, (\bar{s} - s)$, then implementation is straightforward and the cost of obtaining the continuous approximation is only the additional $\bar{s} - s$ evaluations of the differential equation on each step. It is, therefore, generally preferable to derive and implement explicit CRK methods although there are classes of problems where one can only use implicit CRK formulas. Differential algebraic problems, DAEs, are one such class and we will consider them in more detail in Section 5.

Another issue that can be important is the magnitude of the defect or residual that is associated with the continuous approximation. This quantity can be interpreted as a measure of the quality of the numerical solution. For a piecewise polynomial approximation $u(t)$ associated with the ODE (1) one defines the defect, $\delta(t)$, for $t \in (t_0, t_F)$ by

$$\delta(t) = u'(t) - f(t, u(t)).$$

Note that, in deriving CRK formulas, our assumption that $u(t)$ be order p implies, for sufficiently smooth problems,

$$u'(t) = y'(t) + O(h^{p-1}). \tag{4}$$

Furthermore, since

$$\begin{aligned} \delta(t) &= u'(t) - f(t, u(t)) - y'(t) + f(t, y(t)) \\ &= (u'(t) - y'(t)) + (f(t, y(t)) - f(t, u(t))), \end{aligned} \tag{5}$$

then for differential equations that are Lipschitz continuous, $\|\delta(t)\|$ will be at worst $O(h^{p-1})$.

If we let $z_i(t)$ be the solution of the local initial value problem

$$z'_i = f(t, z_i), \quad z_i(t_{i-1}) = y_{i-1},$$

for $i = 1, 2, \dots, N$; then the local error, $le_i(t)$, associated with the continuous approximation on step i is defined for $t \in (t_{i-1}, t_i)$ to be

$$le_i(t) = z_i(t) - u_i(t).$$

It is well known that the discrete local error of a p th-order Runge–Kutta formula (2) must be $O(h_i^{p+1})$. That is

$$\begin{aligned} le_i(t_i) &= z_i(t_i) - u_i(t_i) \\ &= z_i(t_i) - y_i \\ &= O(h_i^{p+1}). \end{aligned} \tag{6}$$

If we derive order p CRK formulas which satisfy the additional constraint that the associated local error be $O(h_i^{p+1})$ for $t \in (t_{i-1}, t_i)$ then we will have

$$u'_i(t) = z'_i(t) + O(h_i^p) \tag{7}$$

and therefore, for $t \in (t_{i-1}, t_i)$,

$$\begin{aligned} \delta(t) &= (u'_i(t) - z'_i(t)) + (f(t, z_i(t)) - f(t, u_i(t))) \\ &= O(h_i^p). \end{aligned} \tag{8}$$

Furthermore, if we derive order p CRK formulas with the stronger additional constraint that, for $t \in (t_{i-1}, t_i)$, the local error be $O(h_i^{p+1})$ and satisfy

$$le_i(t) = \psi_i(\tau) D(t_{i-1}) h_i^{p+1} + O(h_i^{p+2}), \tag{9}$$

where $\tau = (t - t_{i-1})/h_i$, $D(t)$ is a function depending only on the problem, and $\psi_i(\tau)$ is a polynomial in τ whose coefficients are independent of the problem and the stepsize, it can be shown from (8) (see [3] for details) that

$$\delta(t) = \psi'_i(\tau) D(t_{i-1}) h_i^p + O(h_i^{p+1}). \tag{10}$$

In this paper we are considering continuous methods which are designed to directly monitor and control the maximum magnitude of an estimate of the defect of the piecewise polynomial $u(t)$ that is delivered as the approximate solution. We will focus on methods based on an order p CRK formula but most of the discussion and analysis will apply to other classes of continuous methods such as those based on multistep formulas. Note that an order p CRK will always satisfy (5) but, without

additional constraints, such formulas may be more difficult to implement in an effective and reliable way (when the objective is to control the magnitude of $\|\delta(t)\|$) since:

- The magnitude of the associated defect will generally only be $O(h_i^{p-1})$,
- Although the defect can be sampled at any point in the interval of interest it may not be easy to justify a rigorous, inexpensive estimate of its maximum magnitude over each step.

The first of these difficulties can be overcome by considering only order p CRK formulas with local error that is $O(h_i^{p+1})$ and both difficulties can be overcome by considering only those order p CRK formulas that satisfy (9). In this latter case as $h_i \rightarrow 0$, for any $t \in (t_{i-1}, t_i)$, $\delta(t)$ will satisfy (10) and therefore for $\bar{t} \in (t_{i-1}, t_i)$ (corresponding to a local maximum of $|\psi'_i((t - t_{i-1})/h_i)|$), $\|\delta(\bar{t})\|$ will be an asymptotically correct estimate of the maximum magnitude of the defect on the i th step. Note that since $\psi'_i(\tau)$ is a polynomial which is independent of the problem and the stepsize, the location of its local maximum magnitude, $\bar{\tau}$ (for $\tau \in (0, 1)$), is known and the corresponding value for \bar{t} is $\bar{t} = t_{i-1} + \bar{\tau}h_i$.

In Section 2, we will consider continuous methods for IVPs based on defect-control and in subsequent sections we will consider such methods for BVPs and DDEs. In each of these areas there are some general purpose software packages available. Finally, we will consider the development of methods for DAEs based on defect-control. We will discuss some prototype and experimental codes that implement this approach.

2. Initial value methods

The development of software based on defect-control for the numerical solution of IVPs in ODEs has a history that goes back several decades and is closely related to the notion of backward error analysis. Consider the standard IVP

$$y' = f(t, y), \quad y(t_0) = y_0, \quad t \in [t_0, t_F]. \tag{11}$$

Hull [18] and Stetter [24] investigated the reliability (or “effectiveness”) of various error control strategies for discrete methods applied to (11) by establishing conditions which would guarantee the existence of a piecewise approximating function, $\hat{u}(t) \in C^0[t_0, t_F]$, which interpolates the discrete approximate solution y_i and satisfies a slightly perturbed IVP

$$\hat{u}' = f(t, \hat{u}) + \hat{\Delta}(t), \quad \hat{u}(t_0) = y_0, \tag{12}$$

where $\hat{\Delta}(t) \in C^0[t_0, t_F]$ and satisfies

$$\|\hat{\Delta}(t)\| \leq \hat{k} \text{ TOL}, \tag{13}$$

for some modest value of \hat{k} (independent of both the problem and the method), and TOL is the specified error tolerance. In these investigations, $\hat{u}(t)$ and $\hat{\Delta}(t)$ were generally not computable but one could use standard results from mathematics, such as the Grobner–Aleksiev variation of constants formula (see [22] for details), to obtain appropriate global error bounds. For example, if (11), (12) and (13) are satisfied then it is straightforward to show

$$y(t_i) - y_i = y(t_i) - \hat{u}(t_i) = \int_{t_0}^{t_i} K(t, s) \hat{\Delta}(s) ds, \tag{14}$$

Table 1
Cost per step of relaxed and strict defect control for some CRK formulas

Formula	p	s	\bar{s}	\tilde{s}
CRK4	4	4	6	7
CRK5	5	6	9	11
CVSS6B	6	7	11	14
CVSS7	7	9	15	20
CVSS8	8	13	21	28

where $K(t, s)$ is a variational matrix that depends only on the problem, and this implies

$$\|y(t_i) - y_i\| \leq (t_i - t_0) \hat{k} K_{\max} \text{TOL}, \tag{15}$$

where K_{\max} is a bound on $\|K(s, t)\|$. Note that this is an example of backward error analysis where the computed solution is guaranteed to exactly satisfy a slightly perturbed IVP and one can interpret K_{\max} as a type of condition number for the problem (which quantifies how sensitive the solution can be to small changes in the problem specification).

Subsequently, several investigators, who were primarily interested in dense output (or off-mesh approximations), analysed and developed computable piecewise polynomial interpolants, $u(t)$, which could be efficiently implemented with new or existing discrete methods (see, for example, [7,17,23,25,27]). It was soon recognized that, once $u(t)$ was generated, the corresponding defect could be sampled and new reliable error and stepsize-control strategies could be developed with the objective of directly satisfying relationships analogous to (12) and (13).

As we have noted earlier, when one considers CRK formulas, the requirement that $\|\delta(t)\|$ be of optimal order and easy to bound by an inexpensive error estimate (for $t \in (t_{i-1}, t_i)$) will generally impose additional constraints on what would comprise a suitable local interpolating polynomial, $u_i(t)$. In a series of investigations [1–3,14–16], several order p CRK formulas have been developed and compared. These investigations have also considered the relative advantages of several alternative defect-control strategies for $3 \leq p \leq 8$. We will now consider two of the most promising strategies.

The first strategy assumes that the local interpolant $u_i(t)$ defined by (3) satisfies (7) and that the estimate of the maximum magnitude of the defect is obtained by the heuristic

$$\begin{aligned} \text{est}_i &= \|\delta(t_{i-1} + \hat{\tau}h_i)\| \\ &= \|u'_i(t_{i-1} + \hat{\tau}h_i) - f(t_{i-1} + \hat{\tau}h_i, u_i(t_{i-1} + \hat{\tau}h_i))\|, \end{aligned} \tag{16}$$

where $\hat{\tau}$ is chosen in a careful way (see [1] for a detailed discussion of how $\hat{\tau}$ can be chosen). We will refer to this strategy as the “relaxed” defect-control strategy. This heuristic for controlling the maximum magnitude of the defect works well on most problems but the strategy is not asymptotically justified (as $h \rightarrow 0$) and it can severely underestimate the size of the defect for some problems. Table 1 reports for $4 \leq p \leq 8$ the value of s and \bar{s} for some order p CRK which have been found to be particularly effective. The higher-order formulas are members of families of formulas derived in [26] (the particular coefficients are identified in [3]). Note that the cost per step of a method using this error-control strategy for these particular CRKs is \bar{s} derivative evaluations. This follows since

although one must perform an additional derivative evaluation to sample the defect (at $t_{i-1} + \hat{\tau}h_i$), each of the formulas is designed to ensure that $u(t) \in C^1[t_0, t_F]$ by requiring that $u'_i(t_{i-1} + h_i) = f(t_{i-1} + h_i, y_i)$. This implies that, on all but the first step, k_1 will be available (as an internal stage of the previous step).

A more rigorous error-control strategy (which we will refer to as the “strict” defect-control strategy) for the same set of underlying discrete formulas can be developed by requiring that the corresponding continuous extension satisfy (9) as well as (7). With this additional constraint we have available (as discussed earlier) an asymptotically correct estimate of the maximum magnitude of the defect (on the i th step) given by $est_i = \|\delta(t_{i-1} + \bar{\tau}h_i)\|$ where $\bar{\tau}$ is fixed and independent of the problem or stepsize. One way to generate such a CRK (although it may not be optimal) is to begin with an \bar{s} -stage order p CRK, $u_i(t)$, satisfying (3) and (7) and replace the “extra” stages, k_{s+j} , $j = 1, \dots, (\bar{s} - s)$ (whose corresponding $c_{s+j} \neq 1$) with the more accurate values

$$\tilde{k}_{s+j} = \begin{cases} f(t_{i-1} + c_{s+j}h_i, u_i(t_{i-1} + c_{s+j}h_i)) & \text{if } c_{s+j} \neq 1, \\ k_{s+j} & \text{if } c_{s+j} = 1. \end{cases} \tag{17}$$

A new, more suitable interpolant, $\tilde{u}_i(t)$ can then be defined by

$$\tilde{u}_i(t) = y_{i-1} + h_i \sum_{j=1}^s b_j \left(\frac{t - t_{i-1}}{h_i} \right) k_j + \sum_{j=1}^{\bar{s}-s} b_{j+s} \left(\frac{t - t_{i-1}}{h_i} \right) \tilde{k}_{s+j}. \tag{18}$$

It can easily be shown that $\tilde{u}_i(t)$ will satisfy (7), (9) and, when $c_j = 1$ for one value of j in the range $1 \leq j \leq (\bar{s} - s)$ (as is the case for each of the CRK formulas identified in Table 1), it will be an \tilde{s} -stage order p CRK with $\tilde{s} = s + 2(\bar{s} - s) - 1$. Table 1 also reports the value of \tilde{s} for this strict defect-control strategy using the CRK corresponding to (18).

Clearly a trade-off between efficiency and reliability needs to be addressed when choosing which defect-control strategy should be used to solve a particular problem. Fortunately, it is convenient and straightforward to implement a numerical method which can apply either strategy (using either (3) with $\hat{\tau}$ or (18) with $\bar{\tau}$ to define the respective interpolants and defect estimates) and thus the choice can be an option which can be selected by the user. From Table 1 we can observe that the cost per step of using the strict defect-control strategy can be between 20 and 33% more than that for the relaxed strategy but better control of the size of the defect should be realised.

To illustrate and quantify this trade-off we have run both versions of a numerical method based on the CRK formula CVSS6B, on the 25 standard nonstiff problems of the DETEST package [11] at nine error tolerances and assessed how well the defect was controlled. The performance on a typical problem, problem B4, is summarised in Table 2.

These results are typical of that observed on all problems. Both methods are robust in that they are able to deliver, over a wide range of tolerances, a close and consistent relationship between the size of the defect and the specified error tolerance.

Both versions of CVSS6B required over 13 500 steps to solve all of the problems at all tolerances. With the strict defect-control strategy the maximum magnitude of the defect rarely (on fewer than 1.8% of the steps) exceeded TOL and never exceeded 7 TOL. With the relaxed defect-control strategy the maximum magnitude of the defect exceeded TOL on 20% of the steps but it rarely exceeded 5 TOL (on fewer than 1.2% of the steps) and it never exceeded 18 TOL.

An alternative rigorous defect control strategy based on the use of a different norm has been proposed and justified in [19]. On each step one introduces a weighted L_2 -norm (which can be

Table 2
Performance of CVSS6B on problem B4 of DETEST

Strategy	TOL ^a	Time ^b	FCN ^c	Steps ^d	GL Err ^e	Max Def ^f	% Succ ^g
Relaxed defect control	10 ⁻²	0.010	166	15	4.7	1.9	73
	10 ⁻³	0.015	254	22	18.0	1.2	91
	10 ⁻⁴	0.023	375	32	28.4	1.3	84
	10 ⁻⁵	0.031	507	46	34.5	1.7	74
	10 ⁻⁶	0.046	749	68	37.6	1.5	62
	10 ⁻⁷	0.067	1101	100	39.0	1.5	49
	10 ⁻⁸	0.099	1618	147	40.5	1.5	45
	10 ⁻⁹	0.145	2377	216	41.1	1.5	45
Strict defect control	10 ⁻²	0.014	253	17	1.2	0.88	100
	10 ⁻³	0.021	379	24	9.5	0.94	100
	10 ⁻⁴	0.027	491	35	15.6	0.98	100
	10 ⁻⁵	0.041	743	53	16.6	0.76	100
	10 ⁻⁶	0.061	1093	78	16.6	0.66	100
	10 ⁻⁷	0.090	1625	116	16.5	0.61	100
	10 ⁻⁸	0.131	2381	170	16.5	0.58	100
	10 ⁻⁹	0.193	3501	250	16.5	0.56	100

^aTOL, specified error tolerance.

^bTime, computer time required to solve the problem measured in seconds on a SUN Sparc4.

^cFCN, number of derivative evaluations required to solve the problem.

^dSteps, number of time steps required to solve the problem.

^eGL Err, maximum observed global error measured in units of TOL and determined by measuring the global error at 100 equally spaced points per step.

^fMax Def, maximum magnitude of the defect measured in units of TOL and determined by sampling the defect at 100 equally spaced points per step.

^g% Succ, percentage of steps where the magnitude of the defect is less than TOL.

interpreted as an average magnitude of the defect),

$$\|\delta_i(t)\|_2 = 1/h_i \left(\int_{t_{i-1}}^{t_i} \|\delta_i(s)\|^2 ds \right)^{1/2}, \tag{19}$$

A method can then be developed with an error control strategy that attempts to ensure that

$$\|\delta_i(t)\|_2 \leq \text{TOL}. \tag{20}$$

As is pointed out in [19], $\delta_i(t)$ is usually known at the meshpoints and therefore, for sufficiently smooth problems, one can derive a low cost, asymptotically correct estimate of $\|\delta_i(t)\|_2^2$ using a suitably chosen Lobatto quadrature formula (which would require only a few additional evaluations of the defect). This approach has been implemented and shown to be very effective for a large class of problems.

3. Boundary value methods

Numerical methods for BVPs of the form

$$y' = f(t, y), \quad t \in [t_0, t_F], \quad g(y(t_0), y(t_F)) = 0 \tag{21}$$

generally produce a discrete approximation on a mesh $t_0 < t_1 < \dots < t_N = t_F$ by solving a large coupled nonlinear system of equations. If the underlying formula that determines the discrete solution is a Runge–Kutta or collocation formula then it is straightforward to introduce a continuous extension $u(t)$ and the associated defect $\delta(t)$ (as we have done for IV methods). From Table 1 we see that the cost per step to compute $u(t)$ and estimate the size of the corresponding defect can be as great as applying the underlying discrete formula. For BV methods the cost per step to determine $u(t)$ and $\delta(t)$, after the discrete solution has been computed, remains the same while the cost of solving for the discrete solution is generally much greater. A consequence of this is that once a converged discrete solution is determined by a BV method (based on the use of a CRK or collocation formula), a continuous extension with an optimum $O(h^p)$ defect can be computed at very little incremental cost (see, for example, [10,12]).

When these formulas are used to determine the discrete solution, defect-based error control and mesh-refinement strategies can be particularly attractive. This approach has been followed in the development of the methods MIRKDC [9] and bvp4c [19] which have been found to be effective for solving a wide class of problems.

In the numerical solution of BVPs, one often encounters difficulties with convergence of the iteration scheme that is used to solve the nonlinear system associated with the discrete mesh. This can be the result of a poor choice of mesh and/or a poor initial guess for the discrete solution. In either case, if the method has available a piecewise polynomial approximation $\bar{u}(t)$ with an associated defect $\bar{\delta}(t)$ (as would be the case, for example, if $\bar{u}(t)$ were associated with a mesh and previously computed discrete solution that was judged not to be sufficiently accurate), then these deficiencies can often be overcome by using the size of the defect to help guide the mesh refinement and using $\bar{u}(t)$ to generate the required initial guess. With this approach one can also use the estimates of the maximum magnitude of the defect to help ensure that the approximate solution that is ultimately delivered by the method, $u(t)$, satisfies

$$\|\delta(t)\| = \|u'(t) - f(t, u(t))\| \leq \text{TOL}, \tag{22}$$

and

$$g(u(t_0), u(t_F)) = 0.$$

Note that, with this approach, one could consider using inexpensive interpolants for the mesh refinement strategies, and the more expensive rigorous interpolants for assessing the accuracy of the numerical solution.

When such strategies are adopted by a method, one only has to compute the interpolant and defect estimate on the final iteration after the underlying discrete approximation has converged. Intermediate calculations associated with preliminary coarse meshes or initial iterations (before convergence) either would not require the determination of any interpolant or would only require the less expensive relaxed interpolant.

Numerical experience reported in [9,19] shows that BV methods that implement such defect-based strategies can outperform methods based on more traditional strategies especially when a strongly

nonuniform mesh is appropriate. Even on problems where asymptotic analysis is not necessarily relevant, carefully designed defect-based strategies can quickly lead to a suitable mesh and rapid convergence on that mesh.

4. Delay differential equation methods

A class of numerical methods based on CRK formulas with defect control has been analysed [6] for systems of retarded and neutral DDEs of the form

$$y' = f(t, y, y(t - \sigma_1(t, y(t))), y'(t - \sigma_2(t, y(t))), \quad t \in [t_0, t_F]. \quad (23)$$

$$y(t) = \phi(t), \quad t \leq t_0, \quad (24)$$

where $\sigma_1(t, y)$ and $\sigma_2(t, y)$ are positive scalar functions. One particular sixth-order formula from this class (the formula CVSS6B discussed in Section 2) has been implemented in a software package, DDVERK [4], and shown to be effective for these DDEs [5].

For this class of problems a discrete method must be able to approximate the solution at off-mesh points in order to evaluate the differential equation at an arbitrary point, $t \in [t_0, t_F]$. Therefore, the requirement that the numerical solution be a piecewise polynomial, $u(t)$, does not impose any extra cost and one can associate an approximation $u(t)$ (with corresponding defect $\delta(t)$) with any numerical method.

To be effective for this class of problems a numerical method must be able to detect and efficiently handle the discontinuities that inevitably arise and are propagated as the integration proceeds. Automatic techniques for detecting and accurately crossing points of discontinuity for standard IVPs based on monitoring changes in the associated defect have been proposed and justified in [8]. This technique has been adapted and refined for DDEs in the solver DDVERK (see [4] for a discussion of the details) where it has proved to be very effective for a wide class of problems. It is certainly competitive with the alternative strategy which explicitly checks for discontinuities at all possible locations where propagation is possible. This is particularly true for systems of equations with multiple delays where the number of potential points of discontinuity can be quite large relative to the number of significant or relevant points of discontinuity. The defect based strategy for coping with discontinuities essentially adjusts the stepsize selection (as well as the error-control mechanism) only on those steps where the presence of a point of discontinuity has severely reduced the stepsize.

5. Differential algebraic equation methods

In recent years, there has been considerable interest and progress made in the development of numerical methods for special classes of DAEs. Nevertheless, very few methods can be applied directly to a system of DAEs in the most general form

$$F(t, y, y') = 0, \quad y(t_0) = y_0, \quad t \in [t_0, t_F], \quad (25)$$

with $(\partial F / \partial y')$ known to be singular. Note that if this matrix is nonsingular for all $t \in [t_0, t_F]$ the problem is said to have index 0. In this case one can solve the nonlinear system associated with

(25) to determine $y'(t)$ for any prescribed value of t and $y(t)$. Any initial value method can be applied and special DAE methods are not necessary.

In general the “index” of a problem of the form (25) refers to the minimum number of times one has to differentiate the equation, $F(t, y, y')=0$, in order to derive an equivalent initial value problem where $y'(t)$ can be determined uniquely in terms of $t, y(t), F$ and various partial derivatives of F . The higher the index of a problem, the more sensitive the solution can be to perturbations in the data and the more difficult it becomes to develop reliable numerical methods. Currently, there are several reliable general-purpose numerical methods for index 1 problems and other reliable methods designed for special classes of index 2 and index 3 problems.

The DAEs that arise in application areas, such as the modelling of constrained mechanical systems or the design of electrical circuits often are of index 2 or index 3 but they possess special structure and numerical methods which exploit this structure have been developed and have received wide acceptance. For example, the algebraic constraints can often be explicitly identified and the system decoupled, $y(t) = [y_1(t), y_2(t)]^T$, and written in the semi-explicit form

$$y_1'(t) = f(t, y_1(t), y_2(t)), \tag{26}$$

$$0 = g(t, y_1(t), y_2(t)). \tag{27}$$

When one considers the development of defect-based error control for DAE methods two key questions must first be answered:

- (1) How does one define a sufficiently accurate continuous extension, $u_i(t)$, of the discrete approximation (for $t \in [t_{i-1}, t_i]$)?
- (2) What measure of the size of the defect is appropriate to control? That is, can one introduce a measure $\mu_i(\delta)$ such that for $t \in [t_{i-1}, t_i]$ the condition that $\mu_i(\delta(t)) \leq \text{TOL}$ will ensure that the global error will be proportional to TOL and $\mu_i(\delta)$ will be inexpensive to estimate on each step?

These questions were considered in [21] where defect-based error-control strategies suitable for important classes of index 2 and index 3, semi-explicit problems were introduced and justified. The approach that was introduced can be applied with any discrete, order p , implicit Runge–Kutta formula to generate, on each step, interpolating polynomials $u_i(t)$ and $v_i(t)$ that approximate $y_1(t)$ and $y_2(t)$, respectively. If one defines the defect of the resulting vector of piecewise polynomials associated with $u(t)$ and $v(t)$ we have (from (26) and (27))

$$\delta_1(t) = u'(t) - f(t, u(t), v(t)), \tag{28}$$

$$\delta_2(t) = g(t, u(t), v(t)). \tag{29}$$

The global errors $\|y_1(t) - u(t)\|$ and $\|y_2(t) - v(t)\|$ were analysed and shown to be bounded by a suitable multiple of TOL provided $\delta_1(t), \delta_2(t)$, and $\delta_2'(t)$ were all suitably bounded in norm. Corresponding measures $\mu_i(\delta)$ were proposed and associated estimates introduced which could be the basis for an effective defect-based numerical method for semi-explicit DAEs.

Another approach has been considered in [20] where no assumptions are made on the structure of the DAE. In order to determine the piecewise polynomial, $u(t)$ which approximates the solution to (25) one begins with an implicit continuous extension of a discrete, order p , implicit Runge–Kutta formula. One then introduces an associated overdetermined system of nonlinear equations on each time step by requiring that the corresponding approximating polynomial, $\tilde{u}_i(t)$ satisfy (in a least

squares sense) the defining equations of the underlying continuous extension as well as additional “collocation” equations (which are equivalent to asking that (25) be satisfied at a prescribed set of sample points). The defect, $\tilde{\delta}(t)$, associated with the resulting piecewise polynomial, $\tilde{u}(t)$, is defined by

$$\tilde{\delta}(t) = F(t, \tilde{u}(t), \tilde{u}'(t)). \quad (30)$$

Conditions on the choice of underlying implicit CRK formulas and on the number and choice of collocation points are identified which result in $\|\tilde{\delta}(t)\|$ being $O(h^p)$ for sufficiently differentiable index 1 and index 2 problems. Estimates of $\|\tilde{\delta}(t)\|$ are justified and an experimental code introduced to illustrate the validity of this approach. A general-purpose numerical method based on this approach is under development.

6. Summary and conclusions

As is clear from our discussion so far there are now several general purpose numerical methods for important classes of ODEs that produce piecewise polynomial approximate solutions and attempt to directly control the magnitude of the associated defect. These methods, although more costly than the classical discrete methods, can be efficiently implemented and they produce solutions whose accuracy can be more readily interpreted and compared.

We have also shown that when implementing numerical methods using defect control one must address a trade-off between reliability and efficiency. This trade-off arises from a choice between the use of an inexpensive heuristic or a more expensive (but asymptotically correct) estimate of the maximum magnitude of the defect. This choice can be left to the user but the implications must be understood when interpreting the numerical results.

There are two difficulties that have not been discussed which limit the applicability of this class of methods and which should be addressed in future investigations. If the underlying problem is not sufficiently smooth, then one is restricted to the use of lower-order methods and defect control can be less competitive with the more classical approach at low orders. Also, at limiting precision, where the effect of round-off error may dominate the local error, the currently employed defect estimates are unreliable. More research is required to develop effective strategies for detecting and coping with this situation.

References

- [1] W.H. Enright, A new error control for initial value solvers, *Appl. Math. Comput.* 31 (1989) 288–301.
- [2] W.H. Enright, Analysis of error control strategies for continuous Runge–Kutta methods, *SIAM J. Numer. Anal.* 26 (3) (1989) 588–599.
- [3] W.H. Enright, The relative efficiency of alternative defect control schemes for high order Runge–Kutta formulas, *SIAM J. Numer. Anal.* 30 (5) (1993) 1419–1445.
- [4] W.H. Enright, H. Hayashi, A delay differential equation solver based on a continuous Runge–Kutta method with defect control, *Numer. Algorithms* 16 (1997) 349–364.
- [5] W.H. Enright, H. Hayashi, The evaluation of numerical software for delay differential equations, in: R. Boisvert (Ed.), *The quality of Numerical Software: Assessment and Enhancement*, Chapman & Hall, London, 1997, pp. 179–197.

- [6] W.H. Enright, H. Hayashi, Convergence analysis of the solution of retarded and neutral delay differential equations by continuous methods, *SIAM J. Numer. Anal.* 35 (2) (1998) 572–585.
- [7] W.H. Enright, K.R. Jackson, S.P. Nørsett, P.G. Thomsen, Interpolants for Runge–Kutta formulas, *ACM Trans. Math. Software* 12 (1986) 193–218.
- [8] W.H. Enright, K.R. Jackson, S.P. Nørsett, P.G. Thomsen, Effective solution of discontinuous IVPs using a Runge–Kutta formula pair with interpolants, *Appl. Math. Comput.* 27 (1988) 313–335.
- [9] W.H. Enright, P.H. Muir, A Runge–Kutta type boundary value ODE solver with defect control, *SIAM J. Sci. Comput.* 17 (1996) 479–497.
- [10] W.H. Enright, P.H. Muir, Superconvergent interpolants for the collocation solution of BVODEs, *SIAM J. Sci. Comput.* 21 (2000) 227–254.
- [11] W.H. Enright, J.D. Pryce, Two FORTRAN packages for assessing initial value methods, *ACM Trans. Math. Software* 13 (1) (1987) 1–27.
- [12] W.H. Enright, R. Sivasothinathan, Superconvergent interpolants for collocation methods applied to mixed order BVODEs, *ACM Trans. Math. Software* (2000), to appear.
- [13] I. Gladwell, L.F. Shampine, L.S. Baca, R.W. Brankin, Practical aspects of interpolation in Runge–Kutta codes, *SIAM J. Sci. Statist. Comput.* 8 (2) (1987) 322–341.
- [14] D.J. Higham, Defect estimation in Adams PECE codes, *SIAM J. Sci. Comput.* 10 (1989) 964–976.
- [15] D.J. Higham, Robust defect control with Runge–Kutta schemes, *SIAM J. Numer. Anal.* 26 (1989) 1175–1183.
- [16] D.J. Higham, Runge–Kutta defect control using Hermite–Birkhoff interpolation, *SIAM J. Sci. Comput.* 12 (1991) 991–999.
- [17] M.K. Horn, Fourth- and fifth-order scaled Runge–Kutta algorithms for treating dense output, *SIAM J. Numer. Anal.* 20 (3) (1983) 558–568.
- [18] T.E. Hull, The effectiveness of numerical methods for ordinary differential equations, *SIAM Stud. Numer. Anal.* 2 (1968) 114–121.
- [19] J. Kierzenka, L.F. Shampine, A BVP solver based on residual control and the MATLAB PSE SMV Math., Report 99–001.
- [20] C. MacDonald, A new approach for DAEs, Ph.D. Thesis, Department of Computer Science, University of Toronto, 1999 (also appeared as DCS Technical Report No. 317/19).
- [21] H. Nguyen, Interpolation and error control schemes for algebraic differential equations using continuous implicit Runge–Kutta methods, Ph.D. Thesis, Department of Computer Science, University of Toronto, 1995 (also appeared as DCS Technical Report No. 298/95).
- [22] S.P. Nørsett, G. Wanner, Perturbed collocation and Runge–Kutta methods, *Numer. Math.* 38 (1981) 193–208.
- [23] L.F. Shampine, Interpolation for Runge–Kutta methods, *SIAM J. Numer. Anal.* 22 (5) (1985) 1014–1027.
- [24] H.J. Stetter, Considerations concerning a theory for ODE-solvers, in: R. Bulirsch, R.D. Grigorieff, J. Schroder (Eds.), *Lecture notes in Mathematics*, Vol. 631, Numerical Treatment of Differential Equations, Springer, Berlin, 1978, pp. 188–200.
- [25] H.J. Stetter, Interpolation and error estimation in Adams PC-codes, *SIAM J. Numer. Anal.* 16 (2) (1979) 311–323.
- [26] J.H. Verner, Differentiable interpolants for high-order Runge–Kutta methods, *SIAM J. Numer. Anal.* 30 (5) (1993) 1446–1466.
- [27] M. Zennaro, Natural continuous extensions of Runge–Kutta methods, *Math. Comput.* 46 (1986) 119–133.