

Robust defect control for RK-methods using efficiently computed optimal-order interpolants

WAYNE H. ENRIGHT AND WAYNE B. HAYES, {enright,wayne}@cs.toronto.edu, December 9, 2004
Dept. of Computer Science, University of Toronto, Toronto, M5S 3G4, CANADA

Abstract

The quest for reliable integration of *initial value problems* (IVPs) for *ordinary differential equations* (ODEs) is a long-standing problem in numerical analysis. At one end of the reliability spectrum are fixed stepsize methods implemented using standard floating point, where the onus lies entirely with the user to ensure the stepsize chosen is adequate for the desired accuracy. At the other end of the reliability spectrum are rigorous interval-based methods, that can provide provably correct bounds on the error of a numerical solution. This rigour comes at a price, however: interval methods are generally 2-3 orders of magnitude more expensive than fixed stepsize floating-point methods. Along the spectrum between these two extremes lie various methods of different expense that estimate local errors and adjust the stepsize accordingly.

In this paper, we continue previous investigations into a class of interpolants for use in Runge-Kutta methods that have a defect function whose shape is asymptotically independent of the problem being integrated, so that the point at which the maximum defect occurs is known *a priori*. In addition, we find interpolants for which the defect function is almost independent of the timestep, allowing the defect to be monitored and controlled in an efficient and robust manner even for modestly large timesteps, at a fraction of the cost of previously published methods. Our interpolants also offer the highest order possible based on the discrete method. We demonstrate the approach on three Runge-Kutta integrators of order 5, 6, and 8, and provide Fortran and preliminary *Matlab* interfaces to all three new integrators. We also consider how sensitive such methods are to roundoff errors. Numerical results for four problems on a range of accuracy requests are presented.

1 Introduction and Motivation

Consider the *initial value problem* (IVP) for an *ordinary differential equation* (ODE)

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

$$\mathbf{y}(t_0) = \mathbf{y}_0, \quad (2)$$

where \mathbf{y} is an n -dimensional vector and \mathbf{f} an n -dimensional vector-valued function. Standard forward error analysis (*eg.*, Dahlquist and Björck (1974) or Kahaner, Moler, and Nash (1989)) tells us that, for a large class of ODEs, it is impossible in fixed-precision arithmetic to produce a numerical solution to an IVP which remains uniformly close to the exact solution for a long time. However, *backward error analysis* (see for example Corless (1994)) shows that numerical solutions to IVPs of ODEs often exactly solve a “nearby” problem. Here, “nearby” can mean either that the numerical solution follows an exact solution of (1) with slightly different initial conditions (see Hayes (2001) for a review), or that the numerical solution exactly solves a problem with the same initial condition but with a slightly perturbed (1). One form of the latter is *defect-based* backwards error analysis, in which it can be shown that a continuous numerical solution $\tilde{\mathbf{y}}(t)$ exactly satisfies a problem of the form

$$\tilde{\mathbf{y}}'(t) = \mathbf{f}(t, \tilde{\mathbf{y}}(t)) + \delta(t), \quad (3)$$

$$\tilde{\mathbf{y}}(t_0) = \mathbf{y}_0, \quad (4)$$

where $\delta(t)$ is an n -dimensional vector-valued function called the *defect*. Of course, the form and size of the defect depends upon the algorithm used to compute the numerical solution $\tilde{\mathbf{y}}(t)$. In this paper we deal exclusively with p th order Runge-Kutta methods. These methods have the advantage that, with just a few (sometimes zero) extra evaluations of \mathbf{f} , (called *function evaluations*), we can compute a polynomial interpolant between the discrete solution points that is a continuous, piecewise differentiable function which approximates the local solution to $O(h^p)$ or $O(h^{p+1})$ where h is the timestep (Enright, Jackson, Nørsett,

and Thomsen 1986). This sequence of piecewise-differentiable polynomial interpolants forms the continuous solution approximation $\tilde{\mathbf{y}}(t)$. This piecewise polynomial can be differentiated giving $\tilde{\mathbf{y}}'(t)$, and then $\tilde{\mathbf{y}}(t)$ is by definition an exact solution to (3) where

$$\delta(t) \equiv \tilde{\mathbf{y}}'(t) - \mathbf{f}(t, \tilde{\mathbf{y}}(t)). \quad (5)$$

In generating the underlying discrete numerical solution (which is then interpolated), one cannot directly control the forward global error, but must be satisfied with attempting to control some measure of the *local error*. That is, if φ_h is the exact time- h solution operator—ie., $\varphi_h(\mathbf{y}_0)$ is the exact solution to (1,2) at time $t_0 + h$ —and ϕ_h represents an approximate numerical solution computed by some algorithm using stepsize h , then the local error is

$$\Delta(h, \tilde{\mathbf{y}}_0) = \phi_h(\tilde{\mathbf{y}}_0) - \varphi_h(\tilde{\mathbf{y}}_0), \quad (6)$$

where $\tilde{\mathbf{y}}_0$ is the solution value at the beginning of the timestep. We desire that $\|\Delta(h, \tilde{\mathbf{y}}_0)\| \leq \varepsilon$, where ε is the *desired maximum local error*. To control this error, one traditionally can use a variable stepsize integration that varies the stepsize in order to maintain the norm of the estimated local error near to the desired local error. In practice, reliably estimating the local error is relatively inexpensive. For example, if one is using a numerical method which is asymptotically order p , then we expect that $\|\Delta(h, \tilde{\mathbf{y}}_0)\| \sim h^{p+1}$. We can then estimate the local error by comparing two integrations across the same small interval using different timesteps or different formulas. However, it is not difficult to show that this simple strategy can be deceived, in that both integrations can sometimes give similar but inaccurate results. When this happens, their difference is small, leading to a spuriously small estimate of the actual local error. Note that this underestimate of the size of the local error will not necessarily result in an approximation with a large local error in one component on one step being accepted, because other conservative heuristics governing the stepsize-choosing strategy will usually prevent such deceptions.

Another way to control local error is to directly control the defect, which is explicitly computable, at least in principle, as shown above (Enright 1989b; Enright 1989a; Enright 1993). In general, however, the defect can be a complicated function of both the problem and the numerical algorithm, making an estimate of the *maximum* defect across a timestep difficult and expensive to compute. One could, for example, sample the defect at many points across a timestep, and take the maximum of these as an estimate of the maximum defect. The likelihood of this being a severe underestimate of the maximum would depend on the number and location of the sample points, and would be very small as the number of sample points increases. However, this is very *ad hoc* and expensive since it requires extra evaluations of \mathbf{f} at each point we wish to monitor the defect. Furthermore, without sufficient care it is possible to construct examples that deceive such an *ad hoc* method into thinking the maximum defect is smaller than it actually is, just as the above local error estimate can be deceived.

The class of methods we address in this paper are the class of Continuous Runge-Kutta methods (CRKs) where the approximate solution is a piecewise polynomial, $z(t)$, which has an associated defect that is bounded by a small multiple of the tolerance for all problems and whose shape depends only on the method, and not the problem being integrated (Higham 1989; Enright and Higham 1991). Such a method can be developed at a relatively modest increase in cost (over that of standard discrete methods) and they have several advantages. First, since the shape of the defect is independent of the problem being integrated, we can pre-compute (at compile time) the point in the timestep at which the maximum defect is expected to occur (at least asymptotically); at run-time, a reliable estimate of the maximum defect can be obtained with a single evaluation at that point. Secondly, since the maximum defect (and thus the local error per unit step (Stetter 1978)) is bounded by a small multiple of the user-specified tolerance, the methods can be used and the accuracy understood without a user knowing anything about the underlying formula, or the notion of local error or stepsize. Convergence as the user-specified tolerance approaches zero can be proved (Stetter 1978), and interfaces and calling sequences developed which are method independent. Since the methods are continuous, they are useful when visualization or estimates of the global error are of use or when alternative methods are to be investigated for solving the same problem.

Table 1: Number of stages per step for various previous works compared to the present work.

maximum order	5	6	8
Enright (1989a)	4–9	8–10	
Higham (1989)	8–9		
Present work	11	14	24
Enright+Higham (1991)	21	32	60

One way to develop such a class of methods is to use direct defect control and in this manuscript we investigate the validity of this approach. Other approaches could also be used, but the point of this paper is that the cost of developing effective methods need not be that great. The work of Enright (1989a) was preliminary and based on the most natural and inexpensive strategies that could be used for direct defect control. Higham (1989) also considered asymptotically correct estimates. On the other hand, Enright and Higham (1991) considered an expensive technique that assumes parallel processing is available. The work we present here is intermediate: although more expensive than the methods of Enright (1989a) and Higham (1989), the methods we present in this paper achieve robustness comparable to the parallel methods in Enright and Higham (1991) at a fraction of the cost that the parallel algorithms would require if performed sequentially. In particular, Table 1 presents a sample of the number of sequential function evaluations required by the different approaches. Furthermore, we note that the “shape” of the defect curves, although technically only justified asymptotically, are most suitable if they are relatively insensitive as the stepsize ranges over the values that arise in typical applications. The CRKs we analyze and implement here have a relatively constant shape even for modestly large stepsizes, a property which has not been previously demonstrated. (See, for example, Figure 10, which depicts the worst case of non-conformity seen in our problems.)

We also note that, to our knowledge, the three methods described in this paper comprise the first usable, user-friendly implementation of direct defect control. The availability of both Fortran and *Matlab* versions of these integrators makes them good candidates for wide distribution and use.

The reliability of both local error and defect estimates depend on the assumption that truncation error dominates round-off. At stringent accuracy requests, this assumption may not be valid. We test this assumption and conclude that special care is necessary to reduce the effects of round-off when using high-order formulas. This is particularly true if the formula requires many floating-point operations per step.

Finally, we provide implementations of our codes in Fortran and Matlab, and compare them to existing ones using DETEST (Pryce and Enright 1987) and Matlab.

2 Methods

2.1 Continuous, explicit Runge-Kutta methods

We are primarily concerned with explicit Runge-Kutta methods. A traditional explicit s -stage, discrete Runge-Kutta formula is described by its *Butcher tableau*,

$$\begin{array}{c|ccc} & 0 & & \\ c_2 & a_{21} & & \\ c_3 & a_{31} & a_{32} & \\ \vdots & \vdots & \vdots & \ddots \\ c_s & a_{s,1} & a_{s,2} & \dots & a_{s,s-1} \\ \hline & w_1 & w_2 & \dots & w_{s-1} & w_s \end{array}, \quad (7)$$

which has associated with it the discrete Runge-Kutta formula,

$$\varphi_{t_i+h}(\tilde{\mathbf{y}}_i) \approx \phi_{t_i+h}(\tilde{\mathbf{y}}_i) \equiv \tilde{\mathbf{y}}_{i+1} = \tilde{\mathbf{y}}_i + h \sum_{j=1}^s w_j \mathbf{f}(t_i + c_j h, \tilde{\mathbf{Y}}_j),$$

where $\tilde{\mathbf{Y}}_j$ is the j -th stage estimate of the solution at time $t_i + c_j h$,

$$\tilde{\mathbf{Y}}_j = \tilde{\mathbf{y}}_i + h \sum_{r=1}^{j-1} a_{jr} \mathbf{f}(t_i + c_r h, \tilde{\mathbf{Y}}_r).$$

For convenience, we let

$$\mathbf{k}_j = \mathbf{f}(t_i + c_j h, \tilde{\mathbf{Y}}_j),$$

giving

$$\tilde{\mathbf{y}}_{i+1} = \tilde{\mathbf{y}}_i + h \sum_{j=1}^s w_j \mathbf{k}_j.$$

To define a continuous Runge-Kutta method associated with this discrete formula, we determine a different interpolant $\mathbf{u}_i(t)$ over each step of the integration¹ (see, for example, Enright, Jackson, Nørsett, and Thomsen (1986)),

$$\mathbf{u}_i(t) = \tilde{\mathbf{y}}_i + h \sum_{j=1}^{\bar{s}} b_j(\tau) \mathbf{k}_j, \quad t \in [t_i, t_{i+1}], \quad (8)$$

where $\tau = (t - t_i)/h$, $\tau \in [0, 1]$, and

$$b_j(\tau) = \sum_{k=1}^{p+1} \beta_{jk} \tau^k. \quad (9)$$

The extra $(\bar{s} - s)$ stages and the polynomial coefficients β_{jk} are not unique and can be computed as described in Verner (1993). The interpolants that interest us are accurate (in that they agree with the local solution) either to order p or $p + 1$, and their derivatives (and hence their defects) are accurate to order $p - 1$ or p , respectively. In the case that the order of the interpolant is only p , $\beta_{j,(p+1)} \equiv 0$ for all j , and it is possible to increase the order to $p + 1$ by adding a few well-chosen stages $\mathbf{k}_{\bar{s}+1}, \dots, \mathbf{k}_{\bar{s}}$. This gives a modified interpolant $\tilde{\mathbf{u}}(t)$ with coefficients

$$\tilde{b}_j(\tau) = \sum_{k=1}^{p+1} \tilde{\beta}_{jk} \tau^k. \quad (10)$$

Finally, by re-computing some of the \mathbf{k} 's using $\tilde{\mathbf{u}}(t)$ one can derive interpolants $\mathbf{v}(t)$ whose defect is both $O(h^p)$ and whose shape asymptotically approaches a constant polynomial, independent of the problem being integrated. This is accomplished by ensuring that, for those stages \mathbf{k}_j that are not accurate to $O(h^{p+1})$ but which do contribute to the definition of $\tilde{\mathbf{u}}(t)$, we replace \mathbf{k}_j by

$$\tilde{\mathbf{k}}_j = \mathbf{f}(t_i + c_j h, \tilde{\mathbf{u}}(t_i + c_j h)).$$

With this replacement, the coefficients β_{jk} (defining $\mathbf{v}_i(t)$) are identical to the coefficients defining $\tilde{\mathbf{u}}_i(t)$.

2.2 Improved interpolants

For a p th order discrete RK formula, we are primarily interested in interpolants that have an associated p th order defect, such as $\tilde{\mathbf{u}}(t)$ prescribed by (10), since for defect error control (our main concern) to be efficient

¹We sometimes delete the subscript i from $\mathbf{u}_i(t)$ when the meaning is clear from context.

a p th order defect is necessary. In this case we have (Enright 1989b; Enright 1989a) that the leading term in the expansion of the defect satisfies

$$\delta(t) = G(\tau)h^p + O(h^{p+1})$$

where

$$G(\tau) = q_1(\tau)F_1 + q_2(\tau)F_2 + \dots + q_K(\tau)F_K. \quad (11)$$

The q_j are polynomials in τ that depend only on the method and the F_j are constants depending on both the method and the problem.

Now for a given problem and any method, as $h \rightarrow 0$, $G(\tau)$ approaches a constant polynomial and therefore we should observe plots of $\delta(\tau)$ vs. τ approaching a unique polynomial as $h \rightarrow 0$. This polynomial will in general be very different for different problems and this is why it is difficult to choose a fixed sample point τ_* representing the maximum defect across the timestep for use in the error estimate. For the polynomials $\mathbf{u}(t)$ and $\tilde{\mathbf{u}}(t)$ we have used a value for τ_* that is not near any of the zeros of the q_1, q_2, \dots, q_K (Enright 1989b).

There is a special (but more expensive) class of interpolants (including the $\mathbf{v}_i(t)$ discussed in the previous section) which have $K = 1$ and therefore the leading term in the expansion of the defect satisfies

$$G(\tau) = q_1(\tau)F_1 \quad (12)$$

where F_1 is a multiple of the discrete local error evaluated at the right endpoint ($\tau = 1$) (Enright 1989a). That is, there is a direct relationship between the discrete local error and the continuous defect. Since q_1 is independent of the problem we should observe as $h \rightarrow 0$ that the shape of G , and subsequently the (norm of the) local maximum defect, should be proportional to the fixed polynomial q_1 , independent of the problem. The value to use for τ_* is then the location of the maximum of $|q_1(\tau)|, \tau \in [0, 1]$. In the sections below we verify that this is the case. We compute the appropriate $q_1(\tau)$ polynomial theoretically and then verify that its shape is virtually indistinguishable from the numerical defect for a “reasonable” range of stepsizes. Example code in *Maple* for computing the theoretical $q_1(\tau)$ for the 5/6 pair is included in the Appendix.

To obtain such an improved interpolant, $\mathbf{v}_i(\tau)$, we need only recompute the \mathbf{k}_j for $j = s + 1, \dots, \bar{s}$ (corresponding to the points $c_{s+1}, \dots, c_{\bar{s}}$ in the definition of $\tilde{\mathbf{u}}(t)$). In each of these cases we replace \mathbf{k}_j by

$$\mathbf{f}(t_i + c_j h, \tilde{\mathbf{u}}_i(t_i + c_j h))$$

where $\tilde{\mathbf{u}}_i(t)$ is the standard interpolant defined by (10). Then the evaluation of \mathbf{v}_i at a prescribed value of τ is identical to that for $\tilde{\mathbf{u}}_i$, except the re-computed \mathbf{k} 's are used.

In summary, the accuracy of $\mathbf{k}_{s+1}, \dots, \mathbf{k}_{\bar{s}}$ are all $O(h^p)$ or $O(h^{p+1})$ in $\mathbf{u}_i(\tau)$, while with the improved interpolant the order of accuracy of the corresponding stages is increased by one. This is sufficient (as is shown in (Enright 1993)) to ensure that these errors will not contribute to the leading term in the expansion of the defect.

2.3 Testing our methods

We created improved interpolants and implemented the resulting modified methods for several integrators including the 4/5 pair that is `ode45` in *Matlab*, as well as methods developed locally based on 5/6 and 7/8 RK formula pairs. Each such method was extensively tested using the DETEST package (Pryce and Enright 1987), which assesses the performance of a method on a suite of 25 problems.

In addition, we closely studied the following one-dimensional problems taken (along with their names) from DETEST. These were chosen because they have closed-form solutions which facilitate easy evaluation of both local and global solutions:

$$\text{Problem A1: } y' = -y. \quad (13)$$

$$\text{Problem A2: } y' = -y^3/2. \quad (14)$$

$$\text{Problem A4: } y' = \frac{y}{4} \left(1 - \frac{y}{20}\right). \quad (15)$$

$$(16)$$

- (Ss): error/defect evaluated only at τ_* [S], or the max over $\tau \in [0..1]$ [s]. Ideally, we would like τ_* to be near the point of maximum defect.
- (Cc): Use compensated summation when summing Butcher tableau? C=yes, c=no.
- (Bb): pre-computed accurate values of $b_j(\tau)$ of (9) and its derivative at $\tau = \tau_*$? B=yes, b=no.
- (Ee): Minimize roundoff error in (9) and its derivative by ordering terms in increasing absolute value at run-time? E=yes, e=no.
- (Nn): Minimize roundoff error in (9) and its derivative by compile-time ordering by increasing absolute value of the coefficients β_{jk} ? E=yes, e=no.
- (Qq): Use QUAD precision to evaluate (9) and its derivative? Q=yes, q=no.

Table 2: The meanings of each letter used in the names of U-curves in many of the figures in this paper. Note that not all variations were used in all tests. For example, compensated summation was not necessary for the method `ode45`, and so we do not show that variation for `ode45`.

In addition, we used problem D3, which is the planar, two-body Kepler problem, which has four dimensions: position x, y , and velocity v_x, v_y . The Kepler problem also has a closed-form solution.

For each problem and method above, we studied the effect of several variations. Since we want the magnitude of the defect at τ_* to be a good approximation of the maximum defect, we compare the two. Since the coefficients in the Butcher Tableau of the discrete Runge-Kutta formula and the coefficients defining the interpolant can both be quite large, we investigated the effect of *compensated summation*, which attempts to decrease the effects of roundoff error by performing all sums using an algorithm that partially extends precision beyond the hardware floating point precision (see, eg., Higham (2002)). We also studied the effect of ordering the terms β_{jk} in (9) and its derivative (used in computing the defect) both by term size at run-time, and by coefficient size at compile-time. Finally, to compare against the best answer we can expect, we compared each of these implementations to one that uses QUAD precision in all internal calculations. These variations are summarized in Table 2.

Note that for a vector-valued \mathbf{f} such as the Kepler problem, the defect is a also vector-valued function that changes from step to step. Strictly speaking, in many applications we may only be interested in estimating and controlling the magnitude of the defect. After some preliminary tests on many problems across many steps and accuracy requests, we found that looking at a single component of the defect on the first step of an integration reflected the typical situation that arose on each step of a system. For the Kepler problem, we also looked only at the first component of vector, after verifying that all components were acting in a similar fashion.

3 Results

3.1 The 4/5 pair

The tableau for the well-known discrete 4/5 Runge-Kutta pair used by *Matlab*'s `ode45` is shown in Table 3. The first 6 (non-boldfaced) rows represent the tableau used to build the values k_1, \dots, k_6 for the standard, discrete solution.

0	0						
1/5	1/5	0					
3/10	3/40	9/40	0				
4/5	44/45	-56/15	32/9	0			
8/9	19372/6561	-25360/2187	64448/6561	-212/729	0		
1	9017/3168	-355/33	46732/5247	49/176	-5103/18656	0	
1	35/384	0	500/1113	125/192	-2187/6784	11/84	0
	35/384	0	500/1113	125/192	-2187/6784	11/84	0

Table 3: The tableau of *Matlab*'s `ode45`. The non-boldfaced items (all but the last row of a_{ij} and the last entry of b_j) give the standard discrete tableau. The last row gives the extra evaluation from which *Matlab* builds an interpolant with associated local error that is $O(h^5)$ and defect that is $O(h^4)$.

The standard $O(h^5)$ interpolating polynomial (with a defect that is $O(h^4)$) used in `ode45` is

$$\begin{pmatrix} b_1(\tau) \\ b_2(\tau) \\ b_3(\tau) \\ b_4(\tau) \\ b_5(\tau) \\ b_6(\tau) \\ b_7(\tau) \end{pmatrix} = \begin{pmatrix} 1 & -183/64 & 37/12 & -145/128 \\ 0 & 0 & 0 & 0 \\ 0 & 1500/371 & -1000/159 & 1000/371 \\ 0 & -125/32 & 125/12 & -375/64 \\ 0 & 9477/3392 & -729/106 & 25515/6784 \\ 0 & -11/7 & 11/3 & -55/28 \\ 0 & 3/2 & -4 & 5/2 \end{pmatrix} \begin{pmatrix} \tau \\ \tau^2 \\ \tau^3 \\ \tau^4 \end{pmatrix}. \quad (17)$$

This defines the standard interpolating polynomial at step i ,

$$\mathbf{u}(t_i + \tau h) = \tilde{\mathbf{y}}_i + h \sum_{j=1}^7 b_j(\tau) \mathbf{k}_j.$$

To construct $\tilde{\mathbf{u}}(t)$, a standard $O(h^6)$ interpolant, we add 2 more stages (Enright 1989b; Enright 1989a),

$$\mathbf{k}_8 = \mathbf{u}(t_i + 0.86h),$$

$$\mathbf{k}_9 = \mathbf{u}(t_i + 0.93h),$$

and then our $O(h^6)$ polynomial has 9 stages, with corresponding coefficients from (10),

$$\begin{pmatrix} \tilde{b}_1(\tau) \\ \tilde{b}_2(\tau) \\ \tilde{b}_3(\tau) \\ \tilde{b}_4(\tau) \\ \tilde{b}_5(\tau) \\ \tilde{b}_6(\tau) \\ \tilde{b}_7(\tau) \\ \tilde{b}_8(\tau) \\ \tilde{b}_9(\tau) \end{pmatrix} = \begin{pmatrix} 1 & -\frac{1708582621}{524156928} & \frac{1232939669}{262078464} & -\frac{1663764925}{524156928} & \frac{208375}{253952} \\ 0 & 0 & 0 & 0 & 0 \\ 0 & \frac{499875}{94976} & -\frac{1618625}{142464} & \frac{871875}{94976} & -\frac{15625}{5936} \\ 0 & \frac{499875}{65536} & -\frac{1618625}{98304} & \frac{871875}{65536} & -\frac{15625}{4096} \\ 0 & -\frac{26237439}{6946816} & \frac{28319463}{3473408} & -\frac{45762975}{6946816} & \frac{820125}{434176} \\ 0 & \frac{43989}{28672} & -\frac{142439}{43008} & \frac{76725}{28672} & -\frac{1375}{1792} \\ 0 & -\frac{2291427}{100352} & \frac{3838251}{50176} & -\frac{8579075}{100352} & \frac{199625}{6272} \\ 0 & -\frac{47953125}{1078784} & \frac{74828125}{539392} & -\frac{155453125}{1078784} & \frac{78125}{1568} \\ 0 & \frac{8734375}{145824} & -\frac{14359375}{72912} & \frac{31234375}{145824} & -\frac{234375}{3038} \end{pmatrix} \begin{pmatrix} \tau \\ \tau^2 \\ \tau^3 \\ \tau^4 \\ \tau^5 \end{pmatrix}. \quad (18)$$

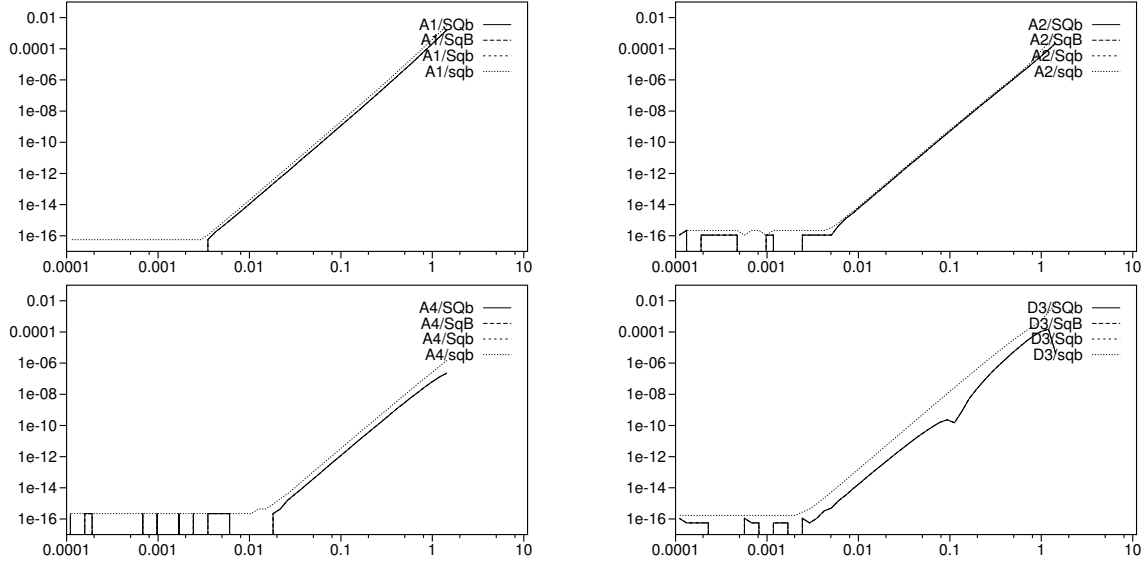


Figure 1: Magnitude of the local error as a function of stepsize for $\mathbf{u}(t)$ as used by the *Matlab* integrator `ode45`. Each figure is for a different problem, A1, A2, A4, and D3, respectively. In all cases the slope of the curve (at least above the machine precision of $1e-16$) is 5.0, which corresponds to the expected local order of the method. We plot the U-curves as computed in several different ways for comparison. The three letters of the name of the curve correspond to the options explained in the text and enumerated in Table 2.

This defines the interpolant $\tilde{\mathbf{u}}(t_i + \tau h)$ at step i , with a defect that is $O(h^5)$,

$$\tilde{\mathbf{u}}(t_i + \tau h) = \tilde{\mathbf{y}}_i + h \sum_{j=1}^9 \tilde{b}_j(\tau) \mathbf{k}_j.$$

Finally, to compute the asymptotically correct interpolant $\mathbf{v}(t)$, we “recompute” \mathbf{k}_8 and \mathbf{k}_9 , giving

$$\hat{\mathbf{k}}_8 = \mathbf{f}(t_i + 0.86h, \tilde{\mathbf{u}}(t_i + 0.86h)),$$

$$\hat{\mathbf{k}}_9 = \mathbf{f}(t_i + 0.93h, \tilde{\mathbf{u}}(t_i + 0.93h)).$$

Let $\hat{\mathbf{k}}_j \equiv \mathbf{k}_j$ for those terms not recomputed, *i.e.*, for $j = 1, \dots, 7$. Then the asymptotically improved interpolant, which also has local error $O(h^6)$, is

$$\mathbf{v}(t_i + \tau h) = \tilde{\mathbf{y}}_i + h \sum_{j=1}^9 \tilde{b}_j(\tau) \hat{\mathbf{k}}_j.$$

By “improved” here we simply mean that the “shape” of the polynomial is asymptotically independent of the problem and depends only on the method, allowing us to pre-compute τ_* , which in the asymptotic limit as $h \rightarrow 0$ is the point at which the maximum defect occurs.

A useful tool to visualize both the order of an integration scheme and the effect of roundoff is the so-called *U-shaped curve*, or U-curve for short, in which we plot the local error or the defect across a single step *vs.* the stepsize. We compute the local error by comparing the approximate solution with the exact solution for each problem. The slope of the U-curve on a log-log plot gives the observed order. In Figure 1 we plot the U-curve of the local error for the standard *Matlab* `ode45` for our 4 problems. The value of τ_* for this method, with interpolating polynomials $\mathbf{u}(t)$ as given in (17), is $\tau_* = 0.23$. (This value is chosen simply to avoid zeros of the polynomials q_1, \dots, q_K (11).) We plot the U-curves as computed in several different ways

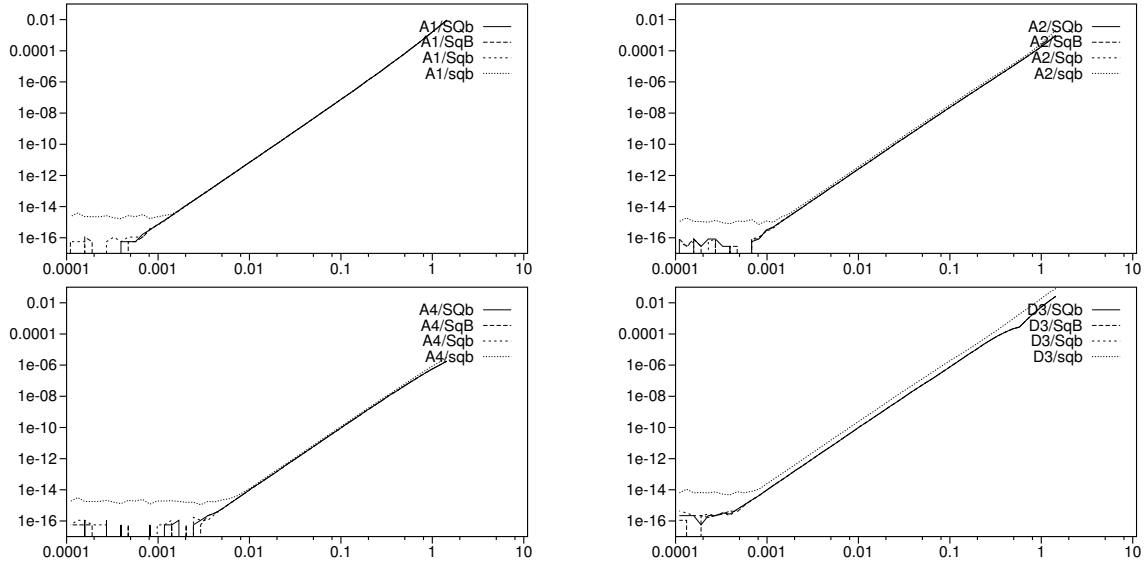


Figure 2: Plots similar to Figure 1, except this time we plot the magnitude of the defect. The slopes of the curves are all about 4.0, because the defect of the *Matlab ode45* polynomial interpolant is $O(h^4)$. Note that for problems A1, A2, and A4, the defect at τ_* (curve **Sqb**) is virtually identical to the maximum defect across $\tau \in [0, 1]$ (curve **sqb**); however, for problem D3, the maximum defect occurs elsewhere and so our estimate at τ_* is a too low by about a factor of 3. Note also that roundoff affects the computation of the maximum defect near limiting precision, so that the levelling of the **sqb** curves near $1e-14$ is due to roundoff. If we pre-compute accurate coefficients for the polynomial at τ_* (curve **SqB**), we eliminate these roundoff effects, allowing computation of the defect at τ_* as accurately as if we had used QUAD precision (curve **SQb**). No compensated summation was used.

for comparison. The first three letters of the name of the curve in the legend correspond to the options described in Table 2. As can be seen, for problems A1 and A2, the local error at τ_* (although we would not be able to measure it in practice) closely approximates the maximum error, but for problems A4 and D3 it is a significant under-estimate. This is not surprising since τ_* is chosen based on attempting to control the defect rather than the local error.

Figure 2 is similar to Figure 1 except we plot the U-curves of the defect rather than the local error. The defect is computed based on evaluation of (5) for $\mathbf{u}_i(t)$. The U-shaped curves for both interpolants $\tilde{\mathbf{u}}(t)$ and $\mathbf{v}(t)$ (not shown) are qualitatively similar, except the order of the error and defect (and hence the slopes of the U-curves) are each higher by one than in Figures 1 and 2. Note that for problem D3, the defect at τ_* is a significant under-estimate of the maximum defect for both $\mathbf{u}(t)$ (shown) and $\tilde{\mathbf{u}}(t)$ (not shown).

Figure 3 plots the shape (normalized so they all have the same magnitude) of the defect of all four problems across $\tau \in [0, 1]$ for a typical stepsize. It clearly demonstrates that the shape of the defect—and in particular its maximum value—is different for the four problems for both $\mathbf{u}(t)$ and $\tilde{\mathbf{u}}(t)$. Thus there is no single choice of τ_* that will reliably estimate the maximum defect for all problems for these interpolants. The figure also demonstrates that the shape of the defect computed from $\mathbf{v}'(t)$ is very similar across all four problems, and that the maximum occurs very close to the expected value of $\tau_* = 0.23$. This allows easy, accurate, robust estimation of the maximum defect with only one evaluation of $\mathbf{v}'(t)$.

3.2 The 5/6 pair

The tableau for the 5/6 pair, along with the derivation of the continuous interpolants $\mathbf{u}(t)$ and $\tilde{\mathbf{u}}(t)$ are described in (Enright 1993). There are 8 stages for the explicit method. One extra function evaluation results in an interpolant of $O(h^6)$ accuracy, and an additional two evaluations provide an interpolant of $O(h^7)$ accuracy. Recomputing these final three stages (giving a total of 6 additional function evaluations

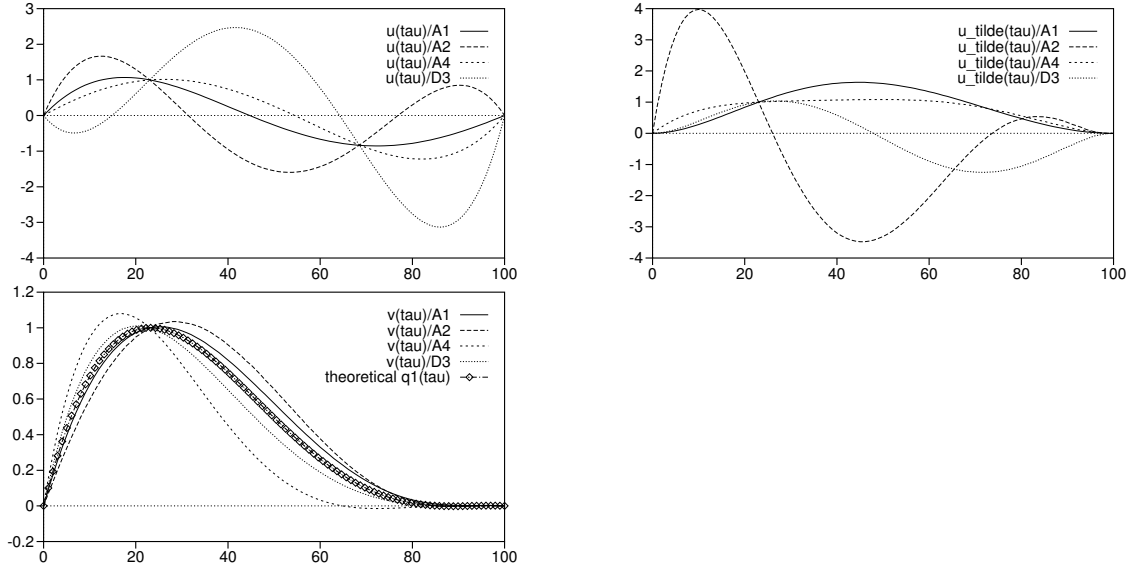


Figure 3: The shape of the 4/5 pair defect curves for all four problems and three interpolants across an entire timestep from $\tau = 0$ to $\tau = 1$. The horizontal axis is labelled by percentage across the timestep, and the curves are normalized so that they are all unity at $\tau_* = 0.23$. The top-left plot is for $\mathbf{u}(t)$ of *Matlab's ode45* which has local error $O(h^5)$, while the right plot is for $\tilde{\mathbf{u}}(t)$ of (18), which has local error $O(h^6)$. Note that the maxima (in magnitude) occur at different places for different problems, and that the max defect is often several times larger in magnitude than that at τ_* ; in some cases (not shown) it was observed to be 40 times larger. The bottom plot is the for the interpolant $\mathbf{v}(t)$ (which has an asymptotically correct $O(h^6)$ defect), showing that it has a very similar shape (for a typical stepsize) for all four problems, so that the defect at τ_* is very close to the maximum defect.

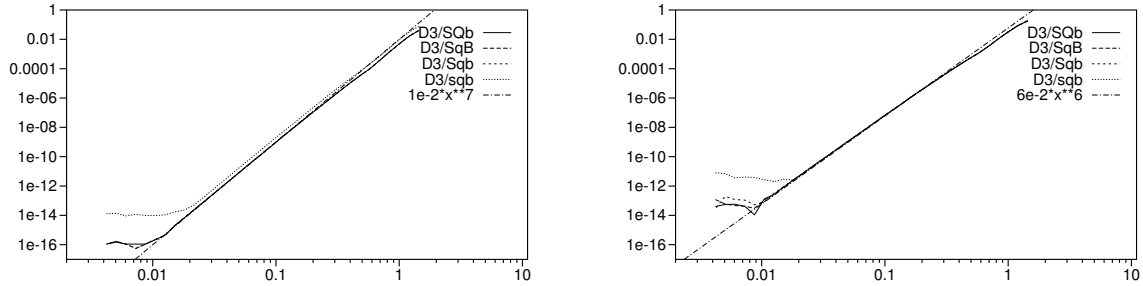


Figure 4: U-shaped-curves for the 5/6 pair interpolant $\tilde{\mathbf{u}}(t)$ operating on the Kepler problem (problem D3); local error is on the left, and defect on the right. The letters describing the curves are from Table 2. Also plotted are straight lines demonstrating that the local error of the improved interpolant is 7th order and the defect 6th order, respectively. As seen for the 4/5 pair, roundoff has a significant effect when interpolating at points other than at τ_* , and pre-computing the accurate polynomial at τ_* virtually eliminates roundoff at τ_* , giving results as good as if QUAD had been used. Similar effects are seen for the other problems.

over the original 8) provides the asymptotically correct interpolant, $\mathbf{v}(t)$. Figure 4 plots the U-curves for the interpolant $\tilde{\mathbf{u}}(t)$ from the 5/6 pair operating on the Kepler problem. Unlike the 4/5 pair, the defect at τ_* (right figure) is only a slight underestimate of the maximum defect across the timestep, at least for these 4 problems. As figures 5, 6, and 7 demonstrate, however, the shape of the defect across a timestep for $\mathbf{u}(t)$ or $\tilde{\mathbf{u}}(t)$ is still a strong function of the problem, and so the sample at τ_* could potentially be worse for problems other than these four. Figures 5, 6, and 7 plot the computed local error and defect (computed by direct evaluation of (5)) across a single timestep for the 5/6 pair operating on all 4 problems. As with the 4/5 pair, we see that the shape of $\mathbf{v}(t)$ is almost problem-independent, allowing a single evaluation of the defect at τ_* to give the defect corresponding to the maximum defect across the timestep. In the appendix, we include *Maple* code for the derivation of $q_1(\tau)$, whose shape corresponds to that of the defect computed from $\mathbf{v}'(t)$, and whose maximum occurs at $\tau_* \approx 0.1509$.

3.3 The 7/8 pair

The derivation of the tableaux for several 7/8 pairs is described in (Verner 1993). We have chosen one member of this family of Runge-Kutta formulas. Figure 8 plots several U-shaped curves for this 7/8 pair operating on problem D3. Local error is on the left, defect is on the right. The top two figures compare the QUAD-computed value of the error and defect at $\tau_* \approx 0.92$ to the maximum error and defect across the timestep, respectively. Two broad observations are clear from these top two plots: first, the error and defect at τ_* is a significant under-estimate of the maximum across the timestep (which explains the gap between the QUAD and all other curves in both top diagrams); and second, there is significant roundoff error associated with evaluation of (9) and its derivative when the timestep is smaller than about 0.2 (which is why the top curves level off well above machine precision). Various methods of compensated summation (represented by the letters C and E in the labels, described in Table 2) have surprisingly little effect in decreasing the roundoff in evaluating (9). The significant roundoff demonstrated in the top left figure (local error) results from evaluation of (9), which has several large ($\sim 5 \times 10^4$) polynomial coefficients for the 7/8 pair; the problem is exacerbated in the right figure (defect), where we evaluate the derivative of that polynomial, which has even larger coefficients (up to 4×10^5).

The bottom two figures in Figure 8 compare the QUAD-computed value of the error and defect at τ_* to the double-precision-computed values at τ_* with pre-computed accurate values for the polynomials b_j (9) at τ_* , including again various methods of compensated summation. These figures demonstrate that (i) pre-computing accurate values of polynomials b'_j at τ_* (lines labelled with a capital B), which is needed to compute the defect, recovers effective QUAD precision (bottom-right figure), and (ii) again none of the types of compensated summation significantly improve on not using compensated summation (both

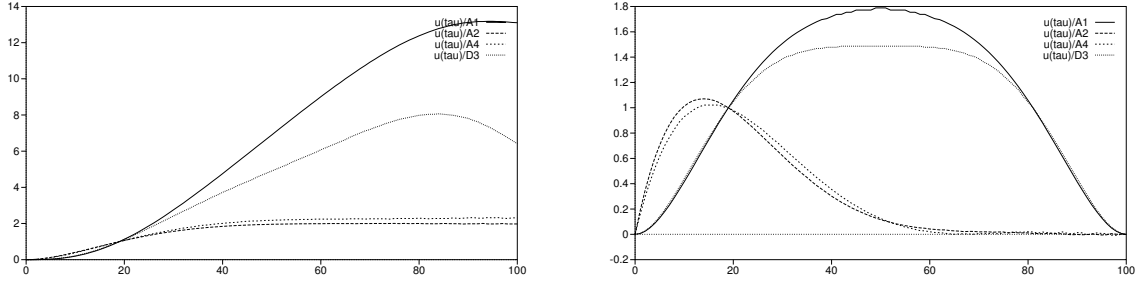


Figure 5: Shape of the local error (left) and defect (right) across a timestep for the $O(h^6)$ interpolant $\mathbf{u}(t)$ for the 5/6 pair, operating on all four problems.

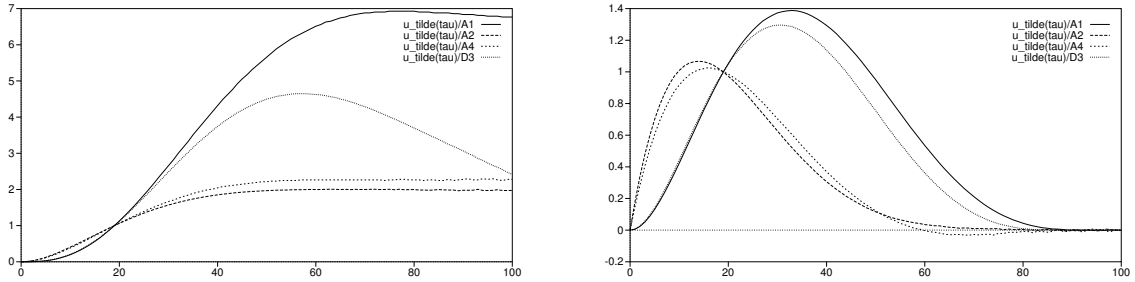


Figure 6: Same as previous figure, but for the $O(h^7)$ interpolant $\tilde{\mathbf{u}}(t)$.

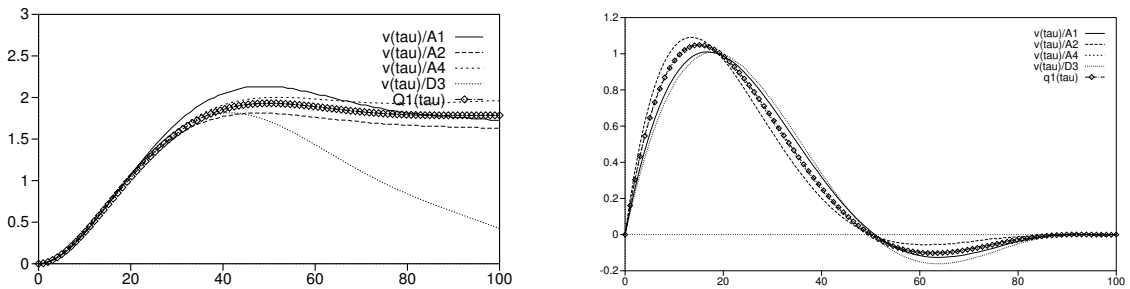


Figure 7: Same as previous figure but for the $O(h^7)$ interpolant $\mathbf{v}(t)$ (whose defect is $O(h^6)$), plus the theoretical $Q_1(\tau)$ and $q_1(\tau)$. Note the shape of the defect is almost problem independent, and corresponds well to the shape of the theoretical defect $q_1(\tau)$.

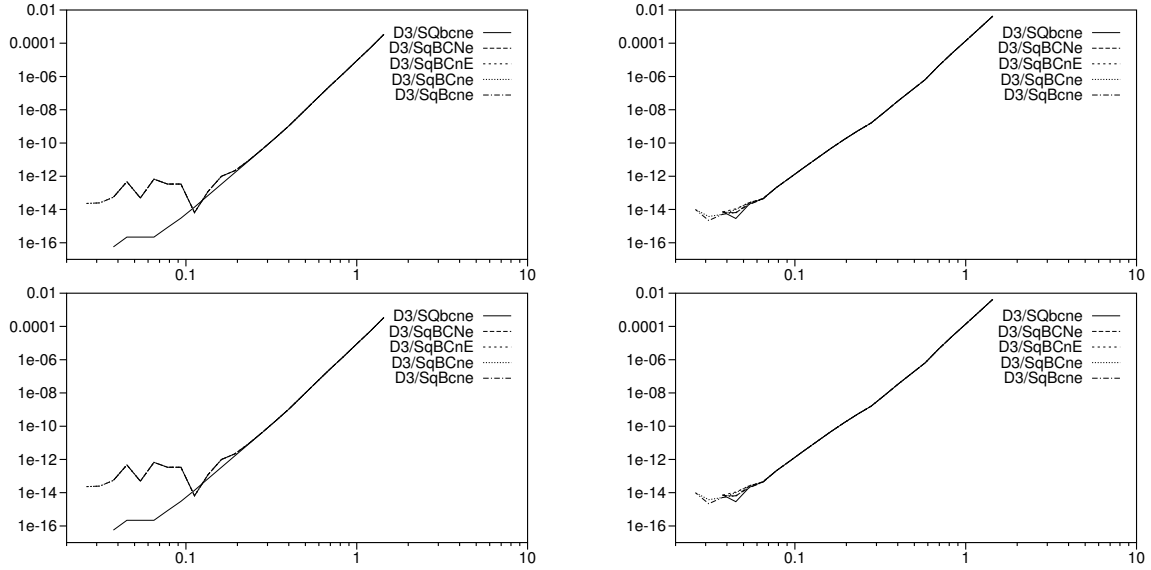


Figure 8: U-shaped-curves for the 7/8 pair interpolant $\mathbf{u}(t)$ operating on the Kepler problem (problem D3). The letters describing the curves are from Table 2. Local error is in the left two figures, defect in the right two figures. The top two figures compare QUAD at τ_* to all cases computing the maximum error and defect across $\tau_* \in [0, 1]$. The bottom two figures compare QUAD at τ_* to all cases computing the error and defect at τ_* with pre-computing accurate polynomials $b_j(\tau_*)$ of (9).

figures). In fact, the local error (bottom-left figure) is identical for all the non-QUAD cases, demonstrating that compensated summation has surprisingly little effect. Similar results are seen for the other problems. DETEST results also confirmed this conclusion, in that no significant change in DETEST results is observable when compensated summation is added.

The shapes of the defects across a timestep for the 7/8 interpolants are plotted in Figure 9. Ironically, the 7/8 pair is so accurate that even when the local error and defect are close to the machine precision for problem A2, the stepsize is too large to reflect the limiting behaviour of the optimal interpolant, although for the other three problems we see that the limiting polynomial is being approached. In the bottom-right figure we resort to QUAD-precision integrations with defect tolerance set to $1\text{e-}25$ in order to force all problems to have a small enough stepsize to demonstrate the limiting behaviour of the defect computed via $\mathbf{v}'(t)$, and in Figure 10 we demonstrate how this limiting polynomial is approached as the timestep decreases for problem A2. This also provides a lesson: if the problem is too “easy” to integrate, then defect measurements at τ_* may be inappropriate even for the “optimal” interpolant, because the stepsize is not small enough for the limiting behaviour to be exhibited.

We note that QUAD precision versions of the interpolant and derivative are available to the user in our package, which is especially useful for the 7/8 pair.

4 Discussion

When using defect control, it is misleading to say we are using a “pair” of RK integrators, because we don’t actually use the lower order integrator of the pair in our error estimates. This allows us to avoid some of the function evaluations related only to the lower order of the pair, but of course we generally add more than we save. Thus these methods tend to be between about 30–70% more expensive than traditional methods, but can be significantly more reliable.

Since satisfying the order conditions when generating a tableau can sometimes be difficult, we briefly studied the effect of errors in the coefficients of the tableau on integrations. In particular, near the beginning

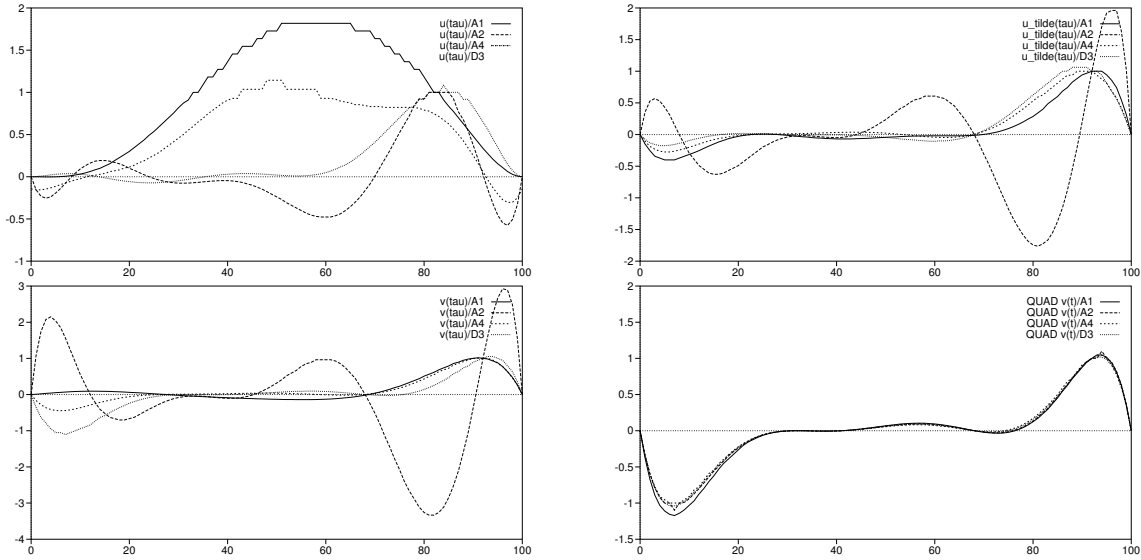


Figure 9: Shape of the defect across a typical timestep for the interpolants of the 7/8 pair operating on all 4 problems. The interpolants are $\mathbf{u}(t)$ (top-left), $\tilde{\mathbf{u}}(t)$ (top-right), and $\mathbf{v}(t)$ (bottom-left) using stepsizes typical of double-precision integration. On the bottom-right, we re-plot the shape of $\mathbf{v}(t)$ using a smaller stepsize giving QUAD-precision accuracy for all 4 problems, to demonstrate the limiting behaviour of $v(t)$.

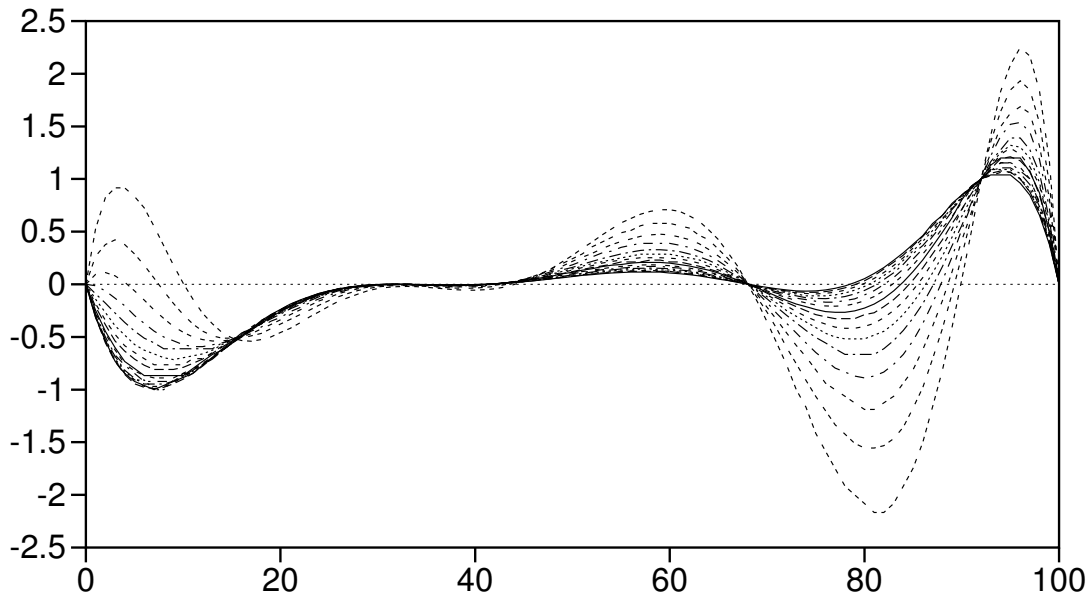


Figure 10: How the shape of A2's defect curve for the 7/8 pair converges on the asymptotic shape as the timestep decreases. Sixteen stepsizes are plotted, ranging from 0.4823 decreasing successively by a factor of 1.2 down to 0.0313.

of this investigation our 7/8 pair had its coefficients computed using only double-precision accuracy, and they were accurate to only 7-8 places. We later re-computed the coefficients in QUAD precision, giving coefficients accurate to 16 digits. We then tested both versions using DETEST (assessing the performance of the defect-based error control) and found very little difference in performance. Only at the most stringent error tolerances was there any difference (in some cases) and these differences were not significant. In particular when a failure occurs in the original version it occurs for the version with more accurate coefficients at the same point or very close by. In some cases there is a bit of improvement in the accuracy obtained or in the number of steps to achieve a given accuracy, but overall the reported summary statistics are not noticeably different. This could be because we are not measuring the right things in our DETEST statistics, but this would require further investigation. It could also be a result of the fact that satisfying the order conditions exactly may not be as important as satisfying a small residual. We know that the original version used had a small residual for the order conditions that we checked (even though the accuracy of the co-efficients was low).

5 Conclusions

In this paper, we have investigated several defect-controlled, explicit, continuous Runge-Kutta methods in which the maximum defect across a timestep can be reliably and efficiently measured, allowing reliable control of the defect. We have created interpolants whose defect, in the limit of small timestep, has a predicted “shape” dependent only on the integration method and not the problem being integrated. Although methods of similar functionality have appeared in the past, ours is unique in that (1) the order of the interpolant is optimal for the order of the discrete method, and (2) it provides effectively equivalent robustness to more expensive such methods, at a fraction of their cost. This gives an integration method in which the error control is both efficient and theoretically justified for all problems.

Appendix A: derivation of $q_1(\tau)$ for the 5/6 pair

The shape of the local error curve for the asymptotically justified interpolant $\mathbf{v}(t)$ for the 5/6 pair is denoted $Q_1(\tau)$. Its derivative, $q_1(\tau) = Q_1'(\tau)$, defines the shape of the defect curve, where $Q_1(\tau)$ is the unique polynomial of degree 6 satisfying

$$Q_1(0) = 1,$$

$$Q_1'(0) = Q_1(1) = Q_1'(1) = Q_1'(0.5) = Q_1'(0.9) = Q_1'(0.95) = 0,$$

where $\tau = 0.5, 0.9,$ and 0.95 are the points at the last stages of the interpolant that are re-computed as described earlier.

Note that if we have

$$Q_1(x) = b_0 + b_1x^1 + \dots + b_6x^6$$

then the linear system that determines the b_i 's (with the constraints defined above at 0 and 1) is simple to write. For the constraints at $1/2, 9/10,$ and $19/20,$ the linear equation contains powers of these points and corresponds to a Vandermond matrix. Using *Maple* to solve for $Q_1(\tau)$ for the 5/6 pair, we have

```
Q1:=x->b0+b1*x+b2*x^2+b3*x^3+b4*x^4+b5*x^5+b6*x^6;
q1:=x->b1+2*b2*x+3*b3*x^2+4*b4*x^3+5*b5*x^4+6*b6*x^5;
solve({q1(0)=0,Q1(0)=1,Q1(1)=0,q1(1)=0,q1(1/2)=0,q1(9/10)=0,q1(19/20)=0},
      {b0,b1,b2,b3,b4,b5,b6});
```

Maple then gives

```
b0 = 1
b1 = 0
b2 = -513.0/17
```

```

b3 = 1766.0/17
b4 = -2478.0/17
b5 = 1608.0/17
b6 = -400.0/17

```

which can be substituted into

$$Q_1(x) = b_0 + b_1x + b_2x^2 + b_3x^3 + b_4x^4 + b_5x^5 + b_6x^6,$$

$$q_1(x) = b_1 + 2b_2x + 3b_3x^2 + 4b_4x^3 + 5b_5x^4 + 6b_6x^5.$$

Now we solve for the zeros of $q_1'(x)$, which tells us where the maxima are:

```

diff(q1(x),x);
# The output of the above diff is: 2*b2+6*b3*x+12*b4*x^2+20*b5*x^3+30*b6*x^4
# Now assign the variables we got above:
b0 := 1;
b1 := 0;
b2 := -513/17;
b3 := 1766/17;
b4 := -2478/17;
b5 := 1608/17;
b6 := -400/17;

Digits:=17;
# And solve for the zeros using the output of the above diff:
fsolve(2*b2+6*b3*x+12*b4*x^2+20*b5*x^3+30*b6*x^4=0,x);

```

Maple's output of the above fsolve is

```
0.15087073414920514, 0.62730617149404269, 0.92210661992260407, 0.9797164744341481
```

Viewing the plot of the right side of Figure 7, we see that it is the peak at ~ 0.15087 which is the maximum of $|q_1(\tau)|$.

Appendix B: Preliminary *Matlab* versions of the integrators

Our implementations of these methods are in Fortran with an interface modelled after DVERK (Hull, Enright, and Jackson 1976). To implement them in *Matlab*, we used *Matlab*'s MEX interface that allows Fortran programs to be directly interfaced with Matlab code.

The mapping between *Matlab*'s ODE options datatype and the options in our current Fortran implementations is not one-to-one; the only options that map cleanly are the suggested initial timestep HSTART and the maximum allowed timestep HMAX. Although our Fortran version can specify a smallest allowable timestep and the maximum number of steps, the *Matlab* options datatype has no facility for specifying these. The tolerance specifications are also different. *Matlab* has both RelTol and AbsTol, and the documentation specifies that the component-by-component error estimate should be the worse of the two:

$$e(i) <= \max(\text{RelTol} * \text{abs}(y(i)), \text{AbsTol}(i)) \quad (19)$$

where AbsTol can be either a scalar or a vector. Our Fortran methods, on the other hand, have only one tolerance, but 5 different ways to interpret it. To account for these differences, ideally, we should add features to our Fortran codes to implement (19) in conformity with Matlab (except we control defect rather than local error). However, we have not yet done this.

In the interim, we have implemented an interface that is as similar to ode23 and ode45 as we can reasonably make it, except for events. We currently have three new integrators called ode45d, ode56d, and ode78d, and the currently implemented features of each are:

- function handles can be used, so either the “old” or “new” interfaces are valid:

```
ode78(@f, tspan, y0)
```

or

```
ode78('f', tspan, y0)
```

- the following options are accepted via `odeset` / `odeget`:
 - `RelTol` becomes DVERK’s `tol`, which is interpreted as DVERK’s default mix of absolute and relative
 - `Refine` is fully implemented
 - `InitialStep` becomes DVERK’s `HSTART`
 - `MaxStep` becomes DVERK’s `HMAX`
 - all other *Matlab* options are ignored.
 - there is no way to use any of DVERK’s other options (*eg.*, no way to communicate with the communication vector)
- no events are implemented
- if there are no output arguments and the user is in the graphical environment, we use `odeplot` to plot the solution, just like `ode23` and `ode45`. The only difference is that we do not return to `Matlab` to plot each step as it is computed, since this would be very inefficient (each task switch between Fortran and Matlab is very expensive). Instead, we wait until DVERK is finished and plot the entire solution vector at once.
- `tspan` is fully implemented. That is, if it has 2 elements, then `Refine` defaults to 1 (rather than `ode45`’s default of 4, since we prefer 1 to mean the natural stepsize of the integrator); or if `tspan` has more than 2 elements, then the solution is returned at exactly these points using, of course, the interpolant to produce any off-mesh solution points.
- we always use DVERK’s natural stepsize, and all off-mesh points are computed using the interpolant (when either `Refine` > 1 or `tspan` has more than 2 elements).

Both the Fortran and *Matlab* implementations are available from the authors on request, including a QUAD precision version of the 7/8 pair.

Acknowledgements We thank Jacek Kierzenka for help with Matlab MEX files, and Philip Sharp for routines to compute coefficients. This research was conducted at the Fields Institute for Research in Mathematical Sciences in Toronto during the *Thematic Year on Numerical Computing*. This research was supported in part by the National Sciences and Engineering Research Council of Canada, and by the Information Technology Research Centre of Ontario.

References

- Corless, R. M. (1994). Error Backward. *Contemporary Mathematics* 172, 31–62.
- Dahlquist, G. and Å. Björck (1974). *Numerical Methods*. Prentice-Hall series in Automatic Computation. Prentice-Hall.
- Enright, W. H. (1989a). Analysis of error control strategies for continuous Runge-Kutta methods. *SIAM Journal on Numerical Analysis* 26(3), 588–599.

- Enright, W. H. (1989b). A new error-control for initial value solvers. *Appl. Math. Comp.* 31, 288–301.
- Enright, W. H. (1993). The relative efficiency of alternative defect control schemes for high-order continuous Runge-Kutta formulas. *SIAM Journal on Numerical Analysis* 30(5), 1419–1445.
- Enright, W. H. and D. J. Higham (1991). Parallel defect control. *BIT* 31, 647–663.
- Enright, W. H., K. R. Jackson, S. P. Nørsett, and P. G. Thomsen (1986). Interpolants for runge-kutta formulas. *ACM Trans. Math. Soft* 12, 193–218.
- Hayes, W. B. (2001). *Rigorous shadowing of numerical solutions of ordinary differential equations by containment*. Ph. D. thesis, Department of Computer Science, University of Toronto. Available on the web as <http://www.cs.toronto.edu/NA/reports.html#hayes-01-phd>.
- Higham, D. J. (1989). Robust defect control with runge-kutta schemes. *SIAM Journal on Numerical Analysis* 26(5), 1175–1183.
- Higham, N. J. (2002). *Accuracy and Stability of Numerical Algorithms, Second edition*. SIAM.
- Hull, T. E., W. H. Enright, and K. R. Jackson (1976, october). User’s guide for DVERK — a subroutine for solving non-stiff ODEs. Technical Report UT-DCS-100, Dept. of Computer Science, University of Toronto. Source code available from the Authors, or from <http://www.netlib.org/> in the ODE directory.
- Kahaner, D., C. Moler, and S. Nash (1989). *Numerical Methods and Software*. Prentice-Hall series in Computational Mathematics. Prentice-Hall.
- Pryce, J. and W. Enright (1987). Two FORTRAN packages for assessing initial value methods. *ACM Trans. Math. Soft.* 13, 1–27.
- Stetter, H. J. (1978). Considerations concerning a theory for ODE-solvers. In *Numerical Treatment of Differential Equations, Proc. Oberwolfach 1976, Lecture notes in Maths 631*, pp. 188–200.
- Verner, J. H. (1993). Differentiable interpolants for high-order Runge-Kutta methods. *SIAM J. Numer. Anal.* 30(5), 1446–1466.