

ESTIMATING THE GLOBAL ERROR OF RUNGE-KUTTA  
APPROXIMATIONS FOR ORDINARY DIFFERENTIAL EQUATIONS

K. Dekker and J.G. Verwer

The user of a code for solving the initial value problem for ordinary differential systems is normally left with the difficult task of assessing the accuracy of the numerical result returned by the code. Even when the code reports an estimate of the global error, the question may remain whether this estimate is correct, i.e. whether the user can rely on the estimate. This paper proposes a simple idea of measuring the reliability of the global error estimate with the aim of assisting the user in the validation of the numerical result. The idea is put into practice with the existing code GERK (ACM Algorithm 504) developed by Shampine and Watts. This code uses global Richardson extrapolation for the error estimation, which in many cases can also be applied to Runge-Kutta codes for delay equations.

## 1. INTRODUCTION

This paper deals with the problem of computing reliable estimates for the global error of numerical approximations to the exact solution  $y(x)$  of the initial value problem

$$(1.1) \quad \dot{y} = f(x,y), \quad a \leq x \leq b, \quad y(a) = y_a,$$

where  $f$  is supposed to be a sufficiently smooth, real-valued vector function. We restrict ourselves to non-stiff systems and classical explicit Runge-Kutta approximations (see e.g. [4,7,11]).

Let us first introduce some notations and definitions. The initial

value problem (1.1) is integrated on a grid

$$(1.2) \quad G_N = \{x_n \in [a,b], n = 0(1)N, \text{ with } x_0 = a, x_{n-1} < x_n, x_N = b\}$$

to obtain the approximations  $y_n$ , where  $y_0 = y_a$  and, for  $n = 0, 1, \dots, N-1$ ,

$$(1.3) \quad \begin{aligned} y_{n+1} &= y_n + h_n \sum_{i=1}^m b_i k_i, \quad h_n = x_{n+1} - x_n, \\ k_i &= f(x_n + c_i h_n, y_n + h_n \sum_{j=1}^{i-1} a_{ij} k_j). \end{aligned}$$

The scalar parameters  $a_{ij}$ ,  $b_i$  and  $c_i$  define the Runge-Kutta scheme. The grid  $G_N$  needs not to be uniform and, as is common practice, may be determined during the integration process through some stepsize control mechanism. It will be assumed that for  $N$  sufficiently large the minimal and maximal step-lengths behave like  $O(N^{-1})$ . More specifically, we assume the existence of a piecewise constant function  $\theta: [a,b] \rightarrow [\theta_{\min}, \theta_{\max}]$ ,  $0 < \theta_{\min} \leq \theta_{\max}$ , such that for  $N$  sufficiently large  $h_n = \theta(x_n) H_N = \theta(x_n) \theta_{\max} / N$ ,  $n = 0(1)N-1$ . If this natural assumption is satisfied, we are assured of the existence of an asymptotic expansion in  $H_N$  for the global discretization errors

$$(1.4) \quad \epsilon_n := y_n - y(x_n), \quad n = 1, \dots, N.$$

See STETTER [11] for a detailed analysis. If we let  $f$  be  $M$  times differentiable (in some neighbourhood of  $y(x)$ ), then functions  $e_j$ ,  $j = p, \dots, M$ , exist independent of  $H = H_N$  such that

$$(1.5) \quad \epsilon_n = \sum_{j=p}^M H^j e_j(x_n) + O(H^{M+1}).$$

Here  $p$  denotes the order of accuracy of the Runge-Kutta method. The existence of these asymptotic expansions for  $\epsilon_n$  forms the basis for most of the error estimation techniques.

The usual approach in the literature on global error estimation is to compute a first approximation for  $\epsilon_n$ ,  $\text{est}_n^{(1)}$  say, which satisfies a relation of the form

$$(1.6) \quad \text{est}_n^{(1)} = H^p e_p(x_n) + H^{p+1} v(x_n) + O(H^{p+2}).$$

Here  $v(x)$  is some function different from  $e_{p+1}(x)$ . The user of a code which delivers an estimate like  $est_n^{(1)}$  will normally be interested in the global error. Anyhow, it is reasonable to assume that most users wish to rely on the estimate. Otherwise the extra effort spent is of no use. In this respect global error estimation has to be approached in an essentially different way than local error estimation. The importance of local error estimation lies in stepsize control, while the reliability of the local estimate is of less importance than its additional costs. When reporting global error estimates however, one should make higher demands on reliability than on efficiency for the reason just mentioned. In fact, from the user's point of view, the computation of a highly reliable global error estimate might be considered as important as an efficient computation of the approximation itself.

These considerations lead us to the conclusion that it might be desirable to compute a second and more accurate estimate  $est_n^{(2)}$  satisfying

$$(1.7) \quad est_n^{(2)} = H^p e_p(x_n) + H^{p+1} e_{p+1}(x_n) + O(H^{p+2}),$$

and to compare this result with the first estimate  $est_n^{(1)}$ .

One way to do this is to check whether

$$(1.8) \quad r_{est} \stackrel{c}{=} est_n^{(2)} / est_n^{(1)}, \quad c \text{ means componentwise operation,}$$

is sufficiently close to one. The quantity  $r_{est}$  is a first order approximation to the true error ratio  $r_{true}$ , i.e.  $r_{est} = r_{true} + O(H)$ , where

$$(1.9) \quad r_{true} \stackrel{c}{=} est_n^{(2)} / \epsilon_n.$$

If  $r_{est}$  is close to one and  $est_n^{(2)}$  is of an acceptable magnitude, one has a strong indication that  $est_n^{(2)}$  is an accurate estimate. We believe that the reliability of automatic codes for our initial value problem (1.1) is greatly enhanced if the asymptotic quality of the global error estimation can be verified.

The objective of this paper is to put this idea into practice and to show that it is useful. Our starting point is the existing Runge-Kutta code GERK developed by SHAMPINE & WATTS [10]. This code is based on a Fehlberg (4,5)-pair [3] and computes a first estimate  $est_n^{(1)}$  by means of global Richardson extrapolation. The decision to concentrate on GERK is based on

the fact that this code is very suitable for the task we have set ourselves.

When combined with interpolation formulas, Runge-Kutta methods have also been shown to be applicable to delay equations, such as

$$(1.10) \quad \dot{y} = f(x, y, y(x-\tau)), \quad x \geq 0, \quad y(x) = \phi(x), \quad x \leq 0,$$

where  $\tau = \tau(x, y(x))$  (see e.g. ARNDT [1], OPPELSTRUP [9]). If the order of the interpolation is high enough, and if the problem is sufficiently smooth, the global error satisfies a relation like (1.5) where  $M$  is at least equal to  $p$  [9]. This means that global Richardson extrapolation can be used. Normally, however, it seems difficult to prove that  $M \geq p+k$ ,  $k \geq 1$ .

## 2. GERK AND GLOBAL RICHARDSON EXTRAPOLATION

Global Richardson extrapolation involves parallel integration with the same method on different grids. The use of Richardson extrapolation for estimating the global discretization error of one-step integration methods is well-known (see HENRICI [4], p.81, LETHER [8] and STETTER [11], p. 157). When using non-equidistant grids, which we assume, it is only allowed to change the stepsize at grid points where the various approximations are combined in the extrapolation process.

SHAMPINE & WATTS [10] have implemented global Richardson extrapolation on top of the Runge-Kutta code RKF45 which is based on a Fehlberg (4.5)-pair. They called the resulting code GERK. This code computes two parallel solutions and estimates the global error at the finest grid (grid  $G_{2N}$  of fig. 1). By computing a third parallel solution, on a grid  $G_{3N}$  as shown in fig. 1, the same idea can be used for obtaining two estimates  $est_n^{(1)}$  and  $est_n^{(2)}$  of the global error at the grid  $G_{3N}$ . Having two estimates of the global error available we then can measure the accuracy of these estimates as outlined in the introduction.

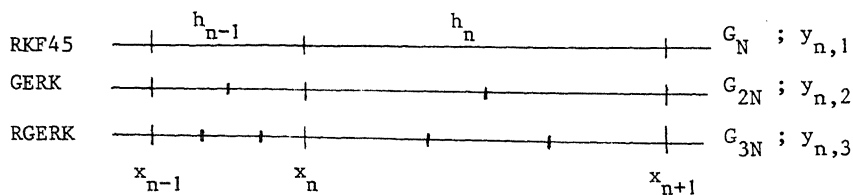


Figure 1.

Let us give some details. Consider three coherent grids as shown in Figure 1. Apply on these grids some Runge-Kutta method of order  $p$  to obtain at the points  $x = x_n$  the approximations  $y_{n,1}$ ,  $y_{n,2}$  and  $y_{n,3}$ . Let  $\epsilon_{n,i} := y_{n,i} - y(x_n)$ . Then

$$(2.1) \quad \epsilon_{n,i} = \sum_{j=p}^M (H/i)^j e_j(x_n) + O(H^{M+1}), \quad i = 1, 2, 3.$$

We now define our estimates by

$$(2.2) \quad \begin{aligned} \text{est}_n^{(1)} &:= (y_{n,2} - y_{n,3}) / (1.5^{p-1}), \\ \text{est}_n^{(2)} &:= (1+\eta)\text{est}_n^{(1)} - \eta(y_{n,1} - y_{n,3}) / (3^{p-1}), \end{aligned}$$

where  $\eta = (1 - (1.5^{p+1} - 1) / (1.5^p - 1)) / ((1.5^{p+1} - 1) / (1.5^p - 1) - (3^{p+1} - 1) / (3^p - 1))$ .

Relations (1.5) - (1.7) are satisfied if  $H$  is replaced by  $H/3$  and  $\epsilon_n$  by  $\epsilon_{n,3}$ . Hence we estimate the error of the most accurate solution  $y_{n,3}$ .

The code GERK computes the solutions  $y_{n,1}, y_{n,2}$  on the grids  $G_N, G_{2N}$  and delivers at the points  $x_n$  the global error estimate  $(y_{n,1} - y_{n,2}) / (2^p - 1)$ , where  $p = 5$ . Thus it also reports the more accurate solution  $y_{n,2}$ . The stepsize selection of GERK is based on a mixed relative-absolute local error control on the coarsest grid  $G_N$  by using the imbedded 4-th order scheme (see [10], section 4 for details). Control on the coarsest grid protects the parallel integration, where no control is performed, against instability. It shall be clear now that it is possible to place our estimation procedure on top of GERK without drastic changes. Only minor modifications are required. These are the implementation of the third parallel integration on  $G_{3N}$  and the implementation of the estimates (2.2). Further, the modified code, which we have named RGERK, should report the numerical solution  $y_{n,3}$ , the estimates  $\text{est}_n^{(2)}$  and  $r_{\text{est}}$ .

We wish to emphasize that we did not modify the stepsize and local error control. This implies that for a given value of the local tolerance input parameters RGERK computes exactly the same solutions  $y_{n,1}$  and  $y_{n,2}$  as GERK does. GERK, in turn, computes the same solution  $y_{n,1}$  as RKF45. It should be noted that in normal situations the global error of  $y_{n,3}$  shall be somewhat smaller than the prescribed local tolerance values. This is because we report the solution computed on the finest grid  $G_{3N}$ , whereas the local error and stepsize control is performed on the coarsest grid  $G_N$ . It

is perhaps clarifying to observe that the grids  $G_N$ ,  $G_{2N}$  and  $G_{3N}$  are determined in the course of the integration, viz. by the stepsize control.

Finally, a few remarks on the cost ratios of RKF45, GERK and RGERK. When we consider the coarsest grid  $G_N$ , RKF45 uses six f-evaluations per step, GERK eighteen, and RGERK thirty-six. However, they report the solution at  $G_N$ ,  $G_{2N}$  and  $G_{3N}$ , respectively. Hence one has to take the accuracy at the three grids into consideration. On the asymptotic basis we thus arrive at the ratios 6 : 9 : 12. In practice the cost ratios, in terms of the numbers of f-evaluations, will slightly differ from the asymptotic ratios. Normally they will be somewhat larger. For further practical information we would like to refer to SHAMPINE & WATTS [10].

### 3. MEASURING THE RELIABILITY OF THE GLOBAL ERROR ESTIMATES

A code like GERK computes a numerical solution of (1.1) and reports at the same time an estimate  $est_n^{(1)}$  of the global discretization error. Experiments reported by Shampine and Watts show that their estimate  $est_n^{(1)}$  will be reliable in many cases. Nevertheless, in real life computation the user of GERK is left with the difficult task of assessing the accuracy of the estimate himself. If it is in doubt, which already may be very difficult to establish, one could apply the code a second time with a more stringent local error tolerance and then by comparison try to get more insight in the accuracy of the reported quantities. The theoretical support of the technique of reintegration is difficult to give, however, when using a stepsize determined by local error control. To assist the user in his validation of the numerical result we prefer to compute the quantity  $r_{est}^c = est_n^{(2)}/est_n^{(1)}$  introduced in equations (1.8), (2.2). From a theoretical point of view, the use of  $r_{est}$  is fully justified. In this section we will consider  $r_{est}$  in some more detail.

For convenience of presentation we now restrict ourselves to a single differential equation. First we introduce the quantity

$$(3.1) \quad r := (1.5^p - 1)(y_{n,1} - y_{n,3}) / (y_{n,2} - y_{n,3})(3^p - 1),$$

and observe that  $r_{est}$  can be written as a function of  $r$ , viz.

$$(3.2) \quad r_{est}(r) = 1 + \eta - \eta r.$$

In fact  $r$  has a similar meaning as  $r_{\text{est}}$  being the quotient of two different estimates of  $\epsilon_n = \epsilon_{n,3}$ . Equation (3.2) shows the range of  $r_{\text{est}}$ . The equality  $r_{\text{est}}(1) = 1$  follows immediately from the observation that both  $r$  and  $r_{\text{est}}$  tend to 1 if  $H \rightarrow 0$ . For  $p = 5$  we have

$$(3.3) \quad r_{\text{est}}(r) = (422 - 121r)/301.$$

Note that  $r_{\text{est}}(0) \simeq 1.4$ , which means that if  $r_{\text{est}}$  is close to 1.4 at least one of the estimates  $\text{est}_n^{(1)}$  or  $\text{est}_n^{(2)}$  is very inaccurate. Generally, too small or too large  $r_{\text{est}}$ -values mean that at least one of these estimates is wrong. One should observe, however, that  $\text{est}_n^{(2)}$  is a more accurate estimate than  $\text{est}_n^{(1)}$  (cf. (1.6), (1.7)). In other words,  $r_{\text{est}}$  normally will be a conservative approximation for  $r_{\text{true}} = \text{est}_n^{(2)}/\epsilon_n$ .

The main question is of course, which range of  $r_{\text{est}}$ -values is still meaningful. We have tried to answer this question in two ways, viz. theoretically and experimentally. The experiments are discussed in the next section. Here we discuss our theoretical answer.

Assume that in equation (2.1) the errors  $\epsilon_{n,i}$  can be represented by infinite series. Let  $e_p(x_n) \neq 0$  and introduce

$$(3.4) \quad \alpha_j := H^{j-p} e_j(x_n)/e_p(x_n), \quad j \geq p.$$

Substitution into (3.1) yields

$$(3.5) \quad r = \frac{1 + (1-3^{-p})^{-1} \sum_{j=p+1}^{\infty} (1-3^{-j}) \alpha_j}{1 + (2^{-p}-3^{-p})^{-1} \sum_{j=p+1}^{\infty} (2^{-j}-3^{-j}) \alpha_j}$$

By imposing bounds for  $\alpha_j$ ,  $j \geq p+1$ , one can obtain bounds for  $r$ ,  $r_{\text{est}}$  and  $r_{\text{true}}$ . The idea is to compare these bounds. We will consider  $\alpha_j$ -values satisfying

$$(3.6) \quad |\alpha_j| \leq c^{j-p}, \quad j \geq p+1 \text{ and } 0 < c < 1.$$

The smaller  $c$ , the more dominance of the error term  $H^p e_p(x_n)$  is supposed by this condition. The following results were obtained.

LEMMA 1. Let  $p = 5$  and denote  $\alpha = (\alpha_6, \alpha_7, \dots)$ ,  $\alpha^- = (-c, -c^2, \dots)$  and

$\alpha^+ = (c, c^2, \dots)$ . Suppose that  $0 < c \leq 2/7$ . For all sequences  $\alpha$  the elements of which satisfy condition (3.6), it then holds that

$$(3.7) \quad r(\alpha^-) = \frac{1-\frac{c}{2}}{1-c} \frac{1-\frac{848}{363}c+\frac{2}{3}c^2}{1-\frac{860}{633}c+\frac{1}{3}c^2} \leq r(\alpha) \leq r(\alpha^+) = \frac{1-\frac{c}{2}}{1-c} \frac{1-\frac{40}{121}c}{1-\frac{65}{211}c}.$$

**PROOF.** Substitute  $p = 5$  into (3.5) and write  $r(\alpha) = N(\alpha)/D(\alpha)$ . Differentiating  $r(\alpha)$  to  $\alpha_k$ ,  $k \geq 6$ , yields

$$\begin{aligned} D^2(\alpha) \frac{\partial r(\alpha)}{\partial \alpha_k} &= \frac{243}{242}(1-3^{-k}) - 32 \frac{243}{211}(2^{-k}-3^{-k}) + \\ & 32 \frac{243}{242} \frac{243}{211} \sum_{j=6} \{(1-3^{-k})(2^{-j}-3^{-j}) - (2^{-k}-3^{-k})(1-3^{-j})\} \alpha_j \geq \\ & \frac{243}{242} - \frac{1}{2} \frac{243}{211} - 32 \frac{243}{242} \frac{243}{211} \sum_{j=6} (2^{-j+2}-3^{-k}) c^{j-5} = \\ & \frac{243}{242} \frac{90}{211} - \frac{243}{242} \frac{243}{211} \left( \frac{c}{2-c} + 2^{5-k} \frac{c}{1-c} \right) \geq \\ & \frac{243}{242 \cdot 211} \left( 90 - 243 \left( \frac{c}{2-c} + \frac{1}{2} \frac{c}{1-c} \right) \right) \end{aligned}$$

for all  $k \geq 6$  and  $0 < c < 1$ . The last expression is positive for all  $c$  between 0 and  $2/7$ . Further,

$$D(\alpha) = 1 + 32 \frac{243}{211} \sum_{j=6} \{(2^{-j}-3^{-j})\alpha_j\} > 1 - \frac{243}{211} \sum_{j=6} 2^{5-j} c^{j-5} > 0$$

if  $0 < c \leq 2/7$ . Hence for all  $k \geq 6$ ,  $\partial r(\alpha)/\partial \alpha_k > 0$  if  $0 < c \leq 2/7$ , which implies that for these values of  $c$ ,  $r(\alpha)$  takes its minimum and maximum at  $\alpha = \alpha^-$  and  $\alpha = \alpha^+$ , respectively.  $\square$

Substitution of (3.4) into  $r_{\text{true}}$  given by (1.9) yields

$$(3.8) \quad r_{\text{true}} = \frac{1+(1+n)(2^{-p}-3^{-p})^{-1} \sum_{j=p+1} (2^{-j}-3^{-j})\alpha_j - n(1-3^{-p})^{-1} \sum_{j=p+1} (1-3^{-j})\alpha_j}{1+3^p \sum_{j=p+1} 3^{-j}\alpha_j}$$

and, for  $p = 5$

$$(3.9) \quad r_{\text{true}} = \frac{1 + 243 \frac{64}{301} \sum_{j=6} (2^{-j}-3^{-j} - \frac{1}{128} + \frac{1}{128} 3^{-j})\alpha_j}{1 + 243 \sum_{j=6} 3^{-j}\alpha_j}.$$



**LEMMA 2.** Let  $p = 5$  and denote  $\alpha = (\alpha_6, \alpha_7, \dots)$ ,  $\alpha^- = (-c, -c^2, \dots)$ ,  $\alpha^* = (-c, c^2, c^3, \dots)$ . Suppose that  $0 < c \leq 602/845$ . For all sequences  $\alpha$  the elements of which satisfy condition (3.6), it then holds that

$$(3.10) \quad r_{\text{true}}(\alpha^*) = 1 - \frac{81}{602} \frac{c^2}{(1-c)(1-\frac{c}{2})(1-\frac{2c}{3}+\frac{2c^2}{9})} \leq r_{\text{true}}(\alpha) \leq \\ \leq 1 + \frac{81}{602} \frac{c^2}{(1-c)(1-\frac{c}{2})(1-\frac{2c}{3})} = r_{\text{true}}(\alpha^-).$$

**PROOF.** We write  $r_{\text{true}}(\alpha) = P(\alpha)/Q(\alpha)$  and note that  $Q(\alpha)$  is positive, because

$$Q(\alpha) = 1 + 243 \sum_{j=6} 3^{-j} \alpha_j \geq 1 - \sum_{j=6} (c/3)^{j-5} > 0,$$

if  $0 < c < 3/2$ . Similarly, we have

$$P(\alpha) = 1 + 243 \frac{64}{301} \sum_{j=6} (2^{-j} 3^{-j} - \frac{1}{128} + \frac{1}{128} 3^{-j}) \alpha_j \geq \\ \geq 1 - \frac{243}{602} \sum_{j=6} c^{j-5} > 0$$

if  $0 < c \leq \frac{602}{845}$ . Differentiating  $r_{\text{true}}(\alpha)$  to  $\alpha_k$ ,  $k \geq 7$ , yields

$$Q^2(\alpha) \frac{\partial r_{\text{true}}(\alpha)}{\partial \alpha_k} = 243 \frac{64}{301} (2^{-k} 3^{-k} - \frac{1}{128} + \frac{1}{128} 3^{-k}) Q(\alpha) - \frac{243}{3^k} P(\alpha) \leq 0,$$

as  $2^{-k} 3^{-k} - \frac{1}{128} + \frac{1}{128} 3^{-k} < 0$  for  $k \geq 7$  and both  $Q(\alpha)$  and  $P(\alpha)$  are positive. Finally,

$$Q^2(\alpha) \frac{\partial r_{\text{true}}(\alpha)}{\partial \alpha_6} = \frac{1}{3} (Q(\alpha) - P(\alpha)).$$

At the maximum we have  $P(\alpha)/Q(\alpha) > 1$ , so that the derivative with respect to  $\alpha_6$  is negative; thus the maximum is obtained for  $\alpha_6 = -c$  and  $\alpha_k = -c^{k-6}$ ,  $k \geq 7$ . At the minimum  $P(\alpha)/Q(\alpha) < 1$ , so the derivative with respect to  $\alpha_6$  is positive and the minimum is obtained for  $\alpha_6 = -c$ ,  $\alpha_k = c^{k-6}$ ,  $k \geq 7$ .  $\square$

Using (3.3) Lemma 1 yields bound for  $r_{\text{est}}$  under condition (3.6), where  $0 < c \leq 2/7$ . Under the same condition Lemma 2 yields bounds for  $r_{\text{true}}$ , but

now for  $0 < c \leq 602/845$ . Table 2 shows these bounds for a number of values for  $c \leq 602/845$ .

c	$r_{\text{est}}$		$r_{\text{true}}$	
	lower	upper	lower	upper
$\frac{1}{100}$	.998	1.002	1.000	1.000
$\frac{1}{30}$	.993	1.007	1.000	1.000
$\frac{1}{10}$	.978	1.023	.998	1.002
$\frac{1}{5}$	.951	1.060	.991	1.009
$\frac{1}{4}$	.936	1.087	.985	1.015
$\frac{2}{7}$	.923	1.110	.978	1.022
$\frac{1}{2}$			.876	1.135
$\frac{7}{10}$			.474	1.634

Table 2. Bounds for  $r_{\text{est}}$  and  $r_{\text{true}}$ .

#### 4. PERFORMANCE OF THE MODIFIED GERK CODE

The purpose of this section is to give practical evidence to our view that the use of a second estimate  $\text{est}_n^{(2)}$  greatly enhances the reliability of the global error estimation procedure. Further we want to give an answer to the question of section 3 how to interpret a reported  $r_{\text{est}}$ -value.

We have subjected the code RGERK to various experiments. In sections 4.1 - 4.3 we present results, in some detail, for three different example problems. In section 4.4 we have collected some statistics on the well-known test set of Hull et al. [5]. All computations have been carried out on a CDC Cyber 750. The arithmetic precision of this computer is about 14 decimal digits (48 bits). We refer to [2] for a more detailed discussion.

##### 4.1. A problem with a peaked solution.

Consider the initial value problem

$$(4.1) \quad \dot{y} = -32xy \ln 2, \quad -1 \leq x \leq 1, \quad y(-1) = 2^{-10},$$

with the peaked solution  $y(x) = 2^{6-16x^2}$ . We have taken this problem from LETHER [8]. For  $x < 0$  the problem is unstable. Hence for  $x < 0$  we will find global errors which increase with  $x$  due to unstable growth. On the other hand, for  $x > 0$  the problem becomes highly stable for increasing  $x$ . Hence for  $x > 0$  the errors should decrease again, as  $x$  increases.

We have solved this problem using pure relative local error control. For the tolerance  $10^{-4}$  we have tabulated  $\epsilon_n$ ,  $r_{est}$ ,  $r_{true}$  and ND = number of f-evaluations, for several values of  $x_n \in [-1,+1]$  (see Table 3, which for the sake of comparison also contains results of GERK). For the remaining grid points similar  $r_{true}$  and  $r_{est}$ -values were found. Hence it can be concluded that the estimation procedure delivers a true copy of the global error behaviour over the complete integration interval.

$x_n$	RGERK				GERK		
	$\epsilon_n$	$r_{true}$	$r_{est}$	ND	$\epsilon_n$	$r_{true}$	ND
-0.884	-7.6(-9)	.99	1.08	155	-5.3(-8)	.78	83
-0.604	-2.0(-6)	.99	1.08	443	-1.4(-5)	.79	227
-0.409	-2.2(-5)	.99	1.08	587	-1.5(-4)	.79	299
.078	-1.5(-4)	.98	1.04	767	-1.1(-3)	.90	389
.307	-5.1(-5)	1.00	1.11	844	-3.6(-4)	.70	430
.617	-2.7(-6)	1.00	1.07	1080	-2.0(-5)	.82	558
.817	-1.4(-7)	1.00	1.04	1296	-1.0(-6)	.90	666
1.000	-4.0(-9)	.99	1.02	1584	-3.0(-8)	.95	810

Table 3. Results for problem (4.1)

#### 4.2. The restricted 3-body problem

Our second example is the restricted 3-body problem

$$(4.2) \quad \begin{aligned} \ddot{u}_1 &= 2\ddot{u}_2 + u_1 - \mu^*(u_1+\mu)/r_1^3 - \mu(u_1-\mu^*)/r_2^3, \\ \ddot{u}_2 &= -2\ddot{u}_1 + u_2 - \mu^*u_2/r_1^3 - \mu u_2/r_2^3, \\ r_1 &= [(u_1+\mu)^2+u_2^2]^{\frac{1}{2}}, r_2 = [(u_1-\mu^*)^2+u_2^2]^{\frac{1}{2}}, \mu = 1/82.45, \mu^* = 1 - \mu, \\ u_1(0) &= 1.2, \dot{u}_1(0) = 0, u_2(0) = 0, \dot{u}_2(0) = -1.04935750983032, \end{aligned}$$

which has also been used by Shampine and Watts. Using absolute local error control we have integrated this difficult problem over the first period  $P = 6.19216933131964$ . Table 4 contains results for the endpoint  $x = P$ . These results belong to that component for which the error of RGERK, in absolute value, is maximal. We observe again good results, except for the tolerance values  $10^{-1} - 10^{-3}$ . For  $10^{-1}$  and  $10^{-2}$   $r_{est}$  fails to indicate that the error estimate  $est_n^{(2)}$  is inaccurate. In both cases, however, one can deduce from the magnitude of  $est_n^{(2)}$  and  $est_n^{(1)}$  that the results are unreliable. It remains necessary to consider the magnitude of  $est_n^{(2)}$  and  $est_n^{(1)}$ . Further we see that for  $10^{-3}$  the estimate  $est_n^{(2)}$  is very good, while  $r_{est}$  has nearly the same value as in the first two cases. We have already predicted this situation in section 3 where we established that  $r_{est}$ -values close to 1.4 may be meaningless.

- log of tolerance	RGERK				GERK		
	$\epsilon_n$	$r_{true}$	$r_{est}$	ND	$\epsilon_n$	$r_{true}$	ND
1	-2.1(+1)	-.44	1.40	355	2.3(+1)	-.03	193
2	-1.3(+1)	-.18	1.39	1494	-1.9(+0)	-.03	810
3	1.6(-2)	.95	1.40	2009	8.7(-2)	-.03	1055
4	2.1(-5)	1.05	1.27	2856	1.3(-4)	.30	1506
5	1.9(-6)	1.04	1.14	4171	1.3(-5)	.64	2191
6	1.4(-7)	1.02	1.06	6257	1.0(-6)	.83	3269
7	7.8(-9)	1.03	1.04	9445	5.9(-8)	.89	4873

Table 4. Restricted 3-body problem (4.2)

#### 4.3. Mildly stiff problems

Though explicit Runge-Kutta codes cannot effectively be used for stiff problems, they may be of value for problems exhibiting a mild stiffness. However, when integrating a mathematically stable problem and numerical stability limits the stepsize, global Richardson extrapolation will always perform badly (see also SHAMPINE & WATTS [10]). More precisely, any estimate which makes use of results from the coarsest grid will be conservative. The reason is that the solutions on the finer grids are not troubled by numerical stability because the local error control, which now prevents the computation from becoming unstable, is performed on the coarsest grid. Often this implies that

due to the stability of the problem and of the computation, the true global error at the finer grids is smooth and small when compared with the global error at the coarsest grid. This causes, fortunately enough, conservative estimates, but also large oscillations in  $r_{est}$ , as well as in  $r_{true}$ .

To see how our estimation procedure performs on a mildly stiff problem, we have shown in Table 5 some results for the simple problem

$$(4.3) \quad \dot{y} = -100\left(y - \frac{x}{x+1}\right) + \frac{1}{(x+1)^2}, \quad x \geq 0, \quad y(0) = 0.$$

The general solution is given by  $y(x) = e^{-100x}y(0) + x/(x+1)$ . Since we take  $y(0) = 0$ , only the rather smooth solution  $x/(x+1)$  has to be computed. Table 6 contains results of the  $n$ -th step, where  $n = 10, 19, 30, 39, 50$ , from the integration under absolute local error control with the tolerance  $10^{-3}$ . Recall that  $\epsilon_{n,1}$  denotes the error at the coarsest grid.

We observe that  $\epsilon_{n,1}$  oscillates and slightly increases with  $n$ , whereas  $\epsilon_n$  smoothly decreases with  $n$ . This results in increasing and oscillating  $r_{true}$ -values. Note that  $r_{est}$  detects this behaviour in a satisfactory way. This is because  $r_{est} = r_{true} * \epsilon_n / est_n^{(1)}$  and  $est_n^{(1)}$  is based on results from the second and third grid. Hence for mildly stiff problems  $est_n^{(1)}$  is to be preferred above  $est_n^{(2)}$ , provided of course that the local error control has to prevent the numerical instability. If the code is applied with a maximal stepsize, chosen in such a way that absolute stability is taken care of, the estimation procedure will perform in a normal way.

$n$	$x_n$	$\epsilon_{n,1}$	$\epsilon_n$	$est_n^{(1)}/\epsilon_n$	$r_{true}$	$r_{est}$
10	.347	-1.4(-4)	2.6(-7)	1.19	2.56	2.15
19	.679	2.6(-4)	1.4(-7)	1.20	-1.40	-1.17
30	1.084	-2.0(-4)	5.9(-8)	1.19	7.45	6.25
39	1.418	3.6(-4)	4.5(-8)	1.19	-11.37	-9.57
50	1.823	-3.4(-4)	3.0(-8)	1.19	20.57	17.33

Table 5. Results of RGERK for problem (4.3)

#### 4.4 Performance of RGERK on the test set of Hull et al.

To gain further insight in the use of computing a second and more accurate estimate  $est_n^{(2)}$ , we have applied RGERK to the five problem classes of [5]. For all 25 problems we used, componentwise, the local error criterion

$$|\text{estimated local error}| \leq \text{tolerance} * |\text{solution}| + 10^{-14}$$

For the tolerances  $10^{-3}$ ,  $10^{-5}$  and  $10^{-7}$ . Hence, in normal cases, a relative control. Results, in percentages shown in Table 6, have been collected for five regions in  $(r_{\text{true}}, r_{\text{est}})$ -space.

Region I  $\cup$  II shows the number of times that  $est_n^{(2)}$  approximates  $\epsilon_n$  rather accurately. Region II shows the number of times that  $r_{\text{est}}$  should be considered as too conservative. Region III  $\cup$  IV  $\cup$  V shows the number of times that the accuracy of  $est_n^{(2)}$  is not so good. Fortunately, in most of the cases this has been detected by  $r_{\text{est}}$ , according to the percentages of region III. Regions IV and V show the most interesting information, viz. the percentages of the number of times that  $r_{\text{est}}$  fails to indicate an inaccurate global error estimation. By way of illustration we show these percentages for two different ranges for  $r_{\text{true}}$ . When considering the failure percentages the reader should realize that we deal with a collection of 25 initial value problems which are divided into 5 different classes, each class having its own degree of difficulty.

Since the results are easily contaminated by round-off when the estimates are very small [2,10], we have distinguished between two intervals for  $est_n^{(2)}$ , viz.  $|est_n^{(2)}| > 10^{-10}$  and  $|est_n^{(2)}| \leq 10^{-10}$ .

#### Regions

I	$1/\sqrt{2} \leq r_{\text{true}} \leq \sqrt{2}, .6 \leq r_{\text{est}} \leq 1.3$
II	$1/\sqrt{2} \leq r_{\text{true}} \leq \sqrt{2}, r_{\text{est}} \leq .6, r_{\text{est}} \geq 1.3$
III	$r_{\text{true}} \leq 1/\sqrt{2}, r_{\text{true}} \geq \sqrt{2}, r_{\text{est}} \leq .6, r_{\text{est}} \geq 1.3$
IV	$1/4 \leq r_{\text{true}} \leq 1/\sqrt{2}, \sqrt{2} \leq r_{\text{true}} \leq 4, .6 \leq r_{\text{est}} \leq 1.3$
V	$r_{\text{true}} \leq 1/4, r_{\text{true}} \geq 4, .6 \leq r_{\text{est}} \leq 1.3$

Regions	$ \text{est}_n^{(2)}  > 10^{-10}$			$ \text{est}_n^{(2)}  \leq 10^{-10}$		
	$10^{-3}$	$10^{-5}$	$10^{-7}$	$10^{-3}$	$10^{-5}$	$10^{-7}$
I	46.4	84.2	96.9	57.3	68.3	75.6
II	22.2	6.5	2.5	3.1	3.8	2.3
III	27.1	8.7	.6	30.5	19.2	10.8
IV	3.7	.6		2.2	2.7	5.2
V	.6			6.9	6.0	6.1
totals	75.0	70.1	52.8	25.0	29.9	47.2

Table 6. Percentages for two intervals for  $\text{est}_n^{(2)}$ .

Further, the percentages were first determined per problem, while counting over all components and all gridpoints, and then averaged over the whole problem collection. These averaged percentages are given in the table.

## 5. CONCLUSIONS

The percentages of Table 6 show that in case of extremely small estimates the reliability is insufficient. This cannot be avoided since the failures are due to roundoff effects. Fortunately, one is usually not interested in a very accurate estimate of extremely small errors, so these results are not as bad as they look.

The reliability is much larger if  $\text{est}_n^{(2)}$  itself is not very small. In fact, for  $10^{-7}$  the score for region V is exactly zero, while for region IV only a few failures out of approximately 20.000 data points were found. The score for region II is still too large, however. This is caused by inaccuracies in  $\text{est}_n^{(1)}$ . For  $10^{-3}$  and  $10^{-5}$  the reliability is less, as expected. In particular for  $10^{-3}$ , a current tolerance value for a 5-th order code, the score for region I ( $\text{est}_n^{(1)}$  and  $\text{est}_n^{(2)}$  both accurate) is too low, whereas the score for II, IV, V is too high. Part of the failures for the larger tolerances is of course due to a failure of the asymptotics. However, numerical instability at the coarsest grid  $G_N$  (cf. section 4.3) also influences the results in a negative way.

Therefore we intend to continue our investigations with an estimation procedure which performs local error-stepsize control on the coarsest grid  $G_N$  and which performs global error estimation only on the finer grids

$G_{2N}, \dots, G_{PN}$ ,  $P \geq 3$ . Herewith we avoid non-smoothness effects which might interfere with the estimation procedure. Such a procedure delivers  $P-2$  estimates  $est_n^{(i)}$ ,  $i = 1, \dots, P-2$ , satisfying  $est_n^{(i)} = \epsilon_{n,P} + O(h^{P+i})$ . The results of the present investigation show that a value  $P > 3$  should be investigated. The cost ratio, in terms of  $f(y)$ -evaluations, is given by  $(P+1)/2$ . Hence, for a given accuracy, the additional computer time for the global estimation will be roughly a factor  $(P+1)/2-1$  of the computer time required when no global error estimation is performed. In this respect it is worthwhile to observe that global Richardson extrapolation is uncommonly attractive for users who have a parallel computer at their disposal. Depending on the number of processors, the additional computer time can then be greatly reduced, even to zero (see e.g. JOUBERT & MAEDER [6]).

ACKNOWLEDGEMENT. The authors greatly acknowledge the programming assistance of Mrs. M. Louter-Nool.

#### REFERENCES

- [1] ARNDT, H., *Der Einfluss der Interpolation auf den globalen Fehler bei retardierten Differentialgleichungen*. These proceedings.
- [2] DEKKER, K. & J.G. VERWER, *Estimating the global error of Runge-Kutta approximations*, Report NW 130/82, Mathematical Centre, Amsterdam, 1982.
- [3] FEHLBERG, E., *Low-order classical Runge-Kutta formulas with step-size control and their applications to some heat transfer problems*, NASA Tech. Rep. TR R-315, George C. Marshall Space Flight Center, Marshall, Ala.
- [4] HENRICI, P., *Discrete Variable Methods for Ordinary Differential Equations*, Wiley, New York, 1962.
- [5] HULL, T.E., W.H. ENRIGHT, B.M. FELLEN & A.E. SEDGWICK, *Comparing numerical methods for ordinary differential equation*, SIAM J. Numer. Anal. 9 (1972), 603-637.
- [6] JOUBERT, G. & A. MAEDER, *Solution of differential equations with a simple parallel computer*, ISNM series Vol. 68, pp. 137-144.



- [7] LAMBERT, J.D., *Computational methods in ordinary differential equations*, Wiley, New York, 1974.
- [8] LETHER, F.G., *The use of Richardson extrapolation in one-step methods with variable step-size*, Math. Computation 20 (1966), 379-385.
- [9] OPPELSTRUP, J., *The RKFHB4 method for delay-differential equations*, Lecture Notes in Mathematics 631 (1978), 133-147.
- [10] SHAMPINE, L.F. & H.A. WATTS, *Global error estimation for ordinary differential equations*, ACM Transactions on Mathematical Software 2 (1976) 172-186.
- [11] STETTER, H.J., *Analysis of discretization methods for ordinary differential equations*, Springer-Verlag, Berlin-Heidelberg-New York, 1973.

K. Dekker, J.G. Verwer  
Mathematisch Centrum  
Kruislaan 413  
1098 SJ Amsterdam  
The Netherlands