

# Cheap One-Step Global Error Estimation for ODEs

A. Murua\*and J. Makazaga†  
Konputazio Zientziak eta A. A. saila  
Informatika Fakultatea, EHU/UPV  
Donostia/San Sebastián, Spain

## Abstract

We propose a class of one-step integrators for ODEs that provide an estimation of the global error along with the numerical solution. The schemes that we propose can be considered as a generalization of the globally embedded RK methods of Dormand, Gilmore and Prince [2]. We present preliminary numerical experiments testing a particular 5th order scheme that we have constructed based on an optimized standard RK method. The numerical results seem to indicate that our method gives useful information on the behaviour of the global error while being practically as efficient as the underlying standard RK scheme.

## 1 Introduction

Very efficient general purpose software is available for the numerical integration of non-stiff ODEs. They more or less succeed in keeping the local error below a prescribed tolerance. However, they do not give any information about the actual global error, unless additional computational effort is done, typically performing a second integration. We believe that some useful information about the actual global error must be provided to the user of general purpose software, since the global error may be much larger than the local error depending on the properties of the system to be integrated, the length of the integration interval, and the properties of the integrator itself.

Our goal is to obtain schemes that give useful information about the propagation of the global error while being as efficient as existing methods that do not provide any sort of global error estimation. Interesting work in this direction has been done by J. R. Dormand, J. P. Gilmore, P. J. Prince [2].

In this paper, we will focus on non-stiff integrators intended to be implemented in sequential computers. We consider ODE systems in autonomous form

$$\frac{dy}{dt} = f(y), \quad y \in R^D, \quad f : R^D \rightarrow R^D. \quad (1)$$

---

\*e-mail: ander@si.ehu.es

†e-mail: ccpmaodj@si.ehu.es

Recall that the  $h$ -flow of the system is a parametric transformation of phase space  $\phi_h : R^D \rightarrow R^D$  such that  $\phi_h(y(t)) = y(t+h)$  for any solution  $y(t)$ .

We wish to integrate initial value problems over an interval  $[0, T]$

$$y(t_0) = y_0, \quad y(t) = ?, \quad t \in [t_0, T].$$

When numerically solving an ODE system by means of a one-step method, we effectively are replacing the flow  $\phi_h$  of the system by a transformation of phase space  $\psi_h : R^D \rightarrow R^D$  that approximates the flow  $\phi_h$ . Then, for a given time discretization  $t_0 < t_1 < \dots < t_N = T$  with  $h_n = t_n - t_{n-1}$ , one computes the numerical solution

$$y_n = \psi_{h_n}(y_{n-1}), \quad n = 1, \dots, N, \tag{2}$$

as approximations of the exact solution values

$$y(t_n) = \phi_{h_n}(y(t_{n-1})), \quad n = 1, \dots, N.$$

The simplest one-step method is the explicit Euler method, where  $\psi_h$  is defined by  $\psi_h(y) = y + hf(y)$ . For the family of explicit Runge-Kutta (RK) methods,  $\psi_h$  is defined as follows.

$$\psi_h(y) = y + h \sum_{i=1}^s b_i f(Y_i),$$

where for  $i = 1, \dots, s$

$$Y_i = y + h \sum_{j=1}^{i-1} a_{ij} f(Y_j).$$

Recall that the local error is defined as  $\delta(y, h) = \psi_h(y) - \phi_h(y)$ , and the method  $\psi_h$  is of order  $p$  if  $\delta(y, h) = O(h^{p+1})$ . In that case we have that the global errors  $e_n$  exhibit convergence of order  $p$  [1, 3], that is,

$$e_n = y_n - y(t_n) = O(H^p), \quad nH \leq \text{Constant}, \quad H = \max_n h_n.$$

## 2 One-Step global error estimation of One-step methods

The global errors  $e_n = y_n - y(t_n)$  of the one-step method (2) satisfy the recurrence

$$e_n = E_{h_n}(y_{n-1}, e_{n-1}),$$

where the mapping  $E_h : R^{2D} \rightarrow R^D$  is defined by  $E_h(y, e) = \psi_h(y) - \phi_h(y - e)$ . Obviously, if the exact flow is not available, that mapping  $E_h$  will not be available either.

Among the different global error estimation techniques proposed in the literature (see [5] for a general survey of different approaches for global error estimation), the procedures that retain the one-step nature of the method itself can be described as a recurrence defined by a mapping  $\tilde{E}_h$  that somehow approximates the true global error mapping  $E_h$ . Thus, the estimates  $\tilde{e}_n$  are obtained as

$$\tilde{e}_n = \tilde{E}_{h_n}(y_{n-1}, \tilde{e}_{n-1}), \quad n = 1, \dots, N, \quad (3)$$

If these estimates  $\tilde{e}_n$  of the global errors are good in some sense, one can extrapolate to obtain a second (hopefully) better approximation  $\bar{y}_n$  as  $y_n - \tilde{e}_n$ . Thus, the process of obtaining the numerical solutions  $y_n$  together with the global error estimates  $\tilde{e}_n$  can be alternatively interpreted as obtaining two approximations  $y_n$  and  $\bar{y}_n$ , and then computing the estimated global error as their difference. More specifically, define  $\bar{\psi}_h(y, \bar{y}) = \psi_h(y) - \tilde{E}_h(y, y - \bar{y})$ , and compute

$$y_n = \psi_{h_n}(y_{n-1}), \quad \bar{y}_n = \bar{\psi}_{h_n}(y_{n-1}, \bar{y}_{n-1}), \quad \tilde{e}_n = y_n - \bar{y}_n. \quad (4)$$

Of course, any process of the form (4) can be interpreted as the application (2) of an one-step method together with a global error estimation of the form (3), where  $\tilde{E}_h(y, e) = \psi_h(y) - \bar{\psi}_h(y, y - e)$ .

**Example 1** Let us consider Richardson extrapolation based on a given method  $\hat{\psi}_h$ . It consists on computing two numerical approximations  $y_n, \hat{y}_n$  ( $n = 1, \dots, N$ ) with  $\hat{y}_0 = y_0$  as follows

$$y_n = \hat{\psi}_{h_n/2}(\hat{\psi}_{h_n/2}(y_{n-1})), \quad \hat{y}_n = \hat{\psi}_{h_n}(\hat{y}_{n-1}), \quad \tilde{e}_n = \frac{1}{2^p - 1}(\hat{y}_n - y_n)$$

and the extrapolated solution will be obtained as  $\bar{y}_n = y_n - \tilde{e}_n$ . This can be interpreted as a process of the form (4) where

$$\begin{aligned} \psi_h(y) &= \hat{\psi}_{h/2}(\hat{\psi}_{h/2}(y)), \\ \bar{\psi}_h(y, \bar{y}) &= \frac{1}{2^p - 1} \left( 2^p \psi_h(y) - \hat{\psi}_h(\hat{y}) \right), \quad \text{with } \hat{y} = \bar{y} - 2^p(\bar{y} - y). \end{aligned}$$

□

We now address the following question: What is meant by  $\tilde{e}_n$  being a good global error estimate? Typically,  $\tilde{e}_n$  is required to be an asymptotically correct global error estimate, in the sense that,

$$\tilde{e}_n = (I + O(H))e_n, \quad \text{for } nH \leq \text{Constant}, \quad H = \max_n h_n,$$

which is obviously equivalent to the extrapolated approximation  $\bar{y}_n$  being of order  $p + 1$ , that is

$$\bar{y}_n - y(t_n) = O(H^{p+1}), \quad nH \leq \text{Constant}.$$

More generally, one could require that  $\tilde{e}_n$  has (for some  $r \geq 1$ )  $r$  asymptotically correct terms, i.e.,

$$\tilde{e}_n = (I + O(H^r))e_n,$$

or equivalently,  $\bar{y}_n - y(t_n) = O(H^{\bar{p}})$ , with  $\bar{p} = p + r$ .

Richardson extrapolation (Example 1), provides in general an asymptotically correct global error estimate with  $r = 1$ , but if the method is symmetric [3], then it gives two asymptotically correct terms of the global error (i.e.,  $r = 2$ ).

An interesting class of schemes that also fits into the format (4) are the so-called *globally embedded RK methods* due to Dormand, Gilmore, and Prince, where the mappings  $\psi_h$  and  $\bar{\psi}_h$  are defined as follows,

$$\psi_h(y) = y + h \sum_{i=1}^s b_i f(Y_i), \quad (5)$$

where

$$Y_i = y + h \sum_{j=1}^{i-1} a_{ij} f(Y_j), \quad i = 1, \dots, s, \quad (6)$$

and

$$\bar{\psi}_h(y, \bar{y}) = \bar{y} + h \sum_{i=s+1}^{\bar{s}} \bar{b}_i f(Y_i), \quad (7)$$

where

$$Y_i = \bar{y} + h \sum_{j=1}^{i-1} a_{ij} f(Y_j), \quad i = s+1, \dots, \bar{s}, \quad (8)$$

The procedure they apply to construct schemes of that kind with different orders  $p$  and  $\bar{p} = p + r$  can be described as follows: They construct an  $s$ -stage RK scheme of order  $p$  with an appropriate continuous extension, and then apply a global error estimation technique known as 'solving for the correction' [5] using RK schemes specially designed for that purpose.

An alternative way of deriving such schemes with prescribed values of  $p$  and  $\bar{p}$  is to study directly which conditions must be satisfied by the mappings  $\psi_h$  and  $\bar{\psi}_h$  so that they provide two approximations of  $y(t_n)$  of order respectively  $p$  and  $\bar{p}$ , and then translate those conditions in terms of the parameters  $b_i$ ,  $\bar{b}_i$ , and  $a_{ij}$ .

### 3 A general class of schemes

Let us now consider the following generalization of processes of the form (4): Take  $\bar{y}_0 = y_0 = y(t_0)$ , and compute for  $n = 1, 2, \dots$ ,

$$y_n = \psi_{h_n}(y_{n-1}, \bar{y}_{n-1}), \quad \bar{y}_n = \bar{\psi}_{h_n}(y_{n-1}, \bar{y}_{n-1}), \quad \tilde{e}_n = y_n - \bar{y}_n. \quad (9)$$

where, for the moment,  $\psi_h$  and  $\bar{\psi}_h$  are arbitrary mappings  $R^{2D} \rightarrow R^D$ . The values  $y_n, \bar{y}_n$  are intended to approximate the solution  $y(t_n)$  of (1) with initial value  $y(t_0) = y_0$ , and  $\tilde{e}_n$  will be an estimate of the global error  $e_n = y_n - y(t_n)$  of the approximation  $y_n$ .

We would like to find which conditions must those mappings satisfy so that the process (9) provides an approximation  $y_n$  to  $y(t_n)$  of order  $p$  together with valid estimates  $\tilde{e}_n = y_n - \bar{y}_n$  of the global errors  $e_n$ . Is it sufficient that  $y_n$  and  $\bar{y}_n$  be approximations of order respectively  $p$  and  $\bar{p} = p + r$ ?

**Example 2** Let  $\hat{\psi}_h$  be a method of order  $p + r$ , and  $C(y, h)$  uniformly bounded for all  $y$  and  $h$ . If we apply (9) with the mappings  $\psi_h, \bar{\psi}_h$  defined by the equations

$$\psi_h(y, \bar{y}) = \hat{\psi}_h(\bar{y}) + C(y, h)h^p, \quad \bar{\psi}_h(y, \bar{y}) = \hat{\psi}_h(\bar{y}),$$

then, we have that

$$\begin{aligned} \bar{e}_n &= \bar{y}_n - y(t_n) = O(H^{p+r}) \\ e_n &= y_n - y(t_n) = O(H^{p+r}) + C(y, h_n)h_n^p \\ \tilde{e}_n &= y_n - \bar{y}_n = C(y, h_n)h_n^p. \end{aligned}$$

We see that  $\tilde{e}_n$  formally is an asymptotically correct estimate of the global error  $e_n$ . But who would accept that as a general purpose global error estimation technique?  $\square$

In order to avoid the kind of poor global error estimation shown in Example 2, we must impose additional conditions to the mappings  $\psi_h$  and  $\bar{\psi}_h$ .

Let us define the *underlying one-step integrators* of the process (9)  $\psi_h(y) := \psi_h(y, y)$ ,  $\bar{\psi}_h(y) := \bar{\psi}_h(y, y)$ . We require that for some  $q, \bar{q} \geq 0$ ,

$$\psi_h(y, y + e) = \psi_h(y, y) + O(h^{q+1}\|e\| + h\|e\|^2), \quad (10)$$

$$\bar{\psi}_h(y + e, y) = \bar{\psi}_h(y, y) + O(h^{\bar{q}+1}\|e\| + h\|e\|^2). \quad (11)$$

The higher  $q$  and  $\bar{q}$ , the more similar is the application of (9) to the independent application of the two underlying one-step methods  $y_{n+1} = \psi_h(y_n)$  and  $\bar{y}_{n+1} = \bar{\psi}_h(\bar{y}_n)$ . Motivated by that, we will refer to (10)–(11) as the *independency conditions*.

In order that the global error estimation in (9) be useful, we want that  $\tilde{e}_n$  and  $e_n$  propagate in a similar way when  $h$  is sufficiently small.

Let us denote the local error of the underlying one-step integrators by

$$\delta(y, h) = \psi_h(y, y) - \phi_h(y), \quad \bar{\delta}(y, h) = \bar{\psi}_h(y, y) - \phi_h(y).$$

**Lemma 1** *If the underlying one-step methods  $\psi_h(y)$  and  $\bar{\psi}_h(y)$  are resp. of order  $p$  and  $p + r$ , and the independency conditions (10)–(11) hold with  $q, \bar{q} \geq 0$ , then*

$$\begin{aligned} e_n &= R_{n,n-1}e_{n-1} + \delta_n + \pi_n, \\ \bar{e}_n &= R_{n,n-1}\bar{e}_{n-1} + \bar{\delta}_n + \bar{\pi}_n, \\ \tilde{e}_n &= R_{n,n-1}\tilde{e}_{n-1} + (\delta_n - \bar{\delta}_n) + (\pi_n - \bar{\pi}_n), \end{aligned} \quad (12)$$

where  $e_n = y_n - y(t_n)$ ,  $\bar{e}_n = \bar{y}_n - y(t_n)$ ,  $\tilde{e}_n = y_n - \bar{y}_n$ , and

$$\begin{aligned}\delta_n &= \delta(y_{n-1}, h_n), & \bar{\delta}_n &= \bar{\delta}(\bar{y}_{n-1}, h_n), \\ \bar{\pi}_n &= O(h_n^{q+1} \|\tilde{e}_{n-1}\| + h_n(\|\tilde{e}_{n-1}\|^2 + \|e_{n-1}\|^2)), \\ \pi_n &= O(h_n^{\bar{q}+1} \|\tilde{e}_{n-1}\| + h_n(\|\tilde{e}_{n-1}\|^2 + \|\bar{e}_{n-1}\|^2)), \\ R_{nk} &= \frac{\partial \phi_{t_n - t_k}}{\partial y}(y(t_k)).\end{aligned}$$

In addition,

$$e_n = \sum_{k=1}^n R_{nk}(\delta_k + \pi_k), \quad \bar{e}_n = \sum_{k=1}^n R_{nk}(\bar{\delta}_k + \bar{\pi}_k), \quad \tilde{e}_n = \sum_{k=1}^n R_{nk}(\delta_k - \bar{\delta}_k + \pi_k - \bar{\pi}_k). \quad (13)$$

**Remark 1** We imply from (12) that, in order that  $e_n$ ,  $\bar{e}_n$ , and  $\tilde{e}_n$  follow similar propagation patterns,  $\pi_n$  and  $\bar{\pi}_n$  should be sufficiently small. In that sense, the higher  $q$  and  $\bar{q}$ , the better. However, it is not obvious in which extent it is important in practice to have higher or lower values of  $q$  and  $\bar{q}$ .  $\square$

**Remark 2** Of course, (13) does not guarantee that the global errors  $e_n$ ,  $\bar{e}_n$  and the estimated global error  $\tilde{e}_n$  propagate in a similar way, even if  $\pi_n$  and  $\bar{\pi}_n$  are negligible. The whole procedure would fail, for instance, in the following situation: Let us consider a two-dimensional system, where the Jacobians  $R_{nk}$  have two eigenvalues  $\lambda_k^1$  and  $\lambda_k^2$  with eigenvectors  $v_k^1$  and  $v_k^2$  such that  $|\lambda_k^1| \gg |\lambda_k^2|$ . Let  $\delta_k = \delta_k^1 v_k^1 + \delta_k^2 v_k^2$  and  $\bar{\delta}_k = \bar{\delta}_k^1 v_k^1 + \bar{\delta}_k^2 v_k^2$  (for each  $k$ ). If  $\delta_k^1 = \bar{\delta}_k^2$ , then  $\tilde{e}_n$  can be much smaller than  $e_n$  and  $\bar{e}_n$ . However, we may hope that, in general, the different way in which each  $R_{nk}$  affects to  $e_n$ ,  $\bar{e}_n$ , and  $\tilde{e}_n$  (depending on the directions of  $\delta_k$ ,  $\bar{\delta}_k$ , and  $\delta_k - \bar{\delta}_k$ ) will tend to compensate for the different  $k = 1, \dots, n$ .  $\square$

**Proof:** We will first prove the first equality in (12). The second equality in (12) can be proven in a completely analogous way, while the third one is obtained by subtracting term by term the first two equalities in (12). We have that

$$\begin{aligned}e_n &= (\psi_{h_n}(y_{n-1}, \bar{y}_{n-1}) - \psi_{h_n}(y_{n-1}, y_{n-1})) \\ &\quad + (\psi_{h_n}(y_{n-1}, y_{n-1}) - \phi_{h_n}(y_{n-1})) \\ &\quad + (\phi_{h_n}(y_{n-1}) - \phi_{h_n}(y(t_{n-1}))).\end{aligned}$$

From definition of  $\delta_n$  and taking into account the independency condition (10), we arrive at

$$e_n = O(h^{q+1} \|\tilde{e}_{n-1}\| + h \|\tilde{e}_{n-1}\|^2) + \delta_n + \left( \frac{\partial \phi_{h_n}}{\partial y}(y(t_{n-1})) e_{n-1} + O(h \|e_{n-1}\|^2) \right),$$

which gives the first equality in (12). The equalities (13) follow from (12) by noting that  $R_{nk} = R_{n,n-1} R_{n-1,k}$ .  $\square$

Adapting the standard techniques of studying the convergence of one-step methods for ODEs [], the following result is obtained.

**Theorem 1** *Under the hypothesis of Lemma 1, if  $\bar{q} \geq r$ , then the numerical approximations provided by the scheme (9) satisfy*

$$e_n = y_n - y(t_n) = O(H^p), \quad \bar{e}_n = \bar{y}_n - y(t_n) = O(H^{p+r}),$$

for  $nH \leq \text{Constant}$  ( $H = \max_n h_n$ ).

## 4 Generalized globally embedded RK schemes

As we have already observed, the globally embedded RK methods (5)–(8) are of the form (4). We now propose the following generalization, which fits in the more general format (9), where the mappings  $\psi_h, \bar{\psi}_h : R^{2D} \rightarrow R^D$  are defined by

$$\psi_h(y, \bar{y}) = y + h \sum_{i=1}^{\bar{s}} b_i f(Y_i), \quad (14)$$

$$\bar{\psi}_h(y, \bar{y}) = \bar{y} + h \sum_{i=1}^{\bar{s}} \bar{b}_i f(Y_i), \quad (15)$$

where for  $i = 1, \dots, \bar{s}$ ,

$$Y_i = \mu_i y + (1 - \mu_i) \bar{y} + h \sum_{j=1}^{i-1} a_{ij} f(Y_j). \quad (16)$$

**Remark 3** Schemes that fits into the format (4) can be obtained by requiring in (14)–(16) that, for some  $s < \bar{s}$ ,

$$\mu_i = 1, \quad i = 1, \dots, s, \quad b_i = 0, \quad i = s + 1, \dots, \bar{s}.$$

If in addition the following conditions are required,

$$\bar{b}_i = 0, \quad i = 1, \dots, s, \quad \mu_i = 0, \quad i = s + 1, \dots, \bar{s},$$

then the family of globally embedded RK schemes (5)–(8) is obtained.  $\square$

**Remark 4** The underlying one-step methods of (14)–(16) are obviously standard RK methods.  $\square$

In order to apply Lemma 1 and Theorem 1, we need to know the order of the underlying RK methods, and in addition, to obtain the values of  $q$  and  $\bar{q}$  for which the independency conditions (10)–(11) hold. It is well known [1, 3] how to obtain systematically the conditions on the coefficients of a RK method to achieve a prescribed order. But how do the independency conditions translate in terms of the coefficients  $b_i, \bar{b}_i, a_{ij}, \mu_i$  of (14)–(16)? It turns out that the resulting equalities in terms of the parameters of the method are similar to those that arise when writing down the order conditions for standard RK methods, only that the parameters  $\mu_i$  come now into play. There is also a nice correspondence with certain kinds of rooted trees. We omit the general formulation of these conditions for lack of space. The conditions for  $q$  and  $\bar{q}$  being  $\geq 3$  are displayed in Table 1. Hereafter, we use the notation  $c_i = \sum a_{ij}$ .

Table 1: Independency conditions for  $q, \bar{q} \leq 3$

$k$	tree	$q = k$	$\bar{q} = k$
1		$\sum_i b_i(1 - \mu_i) = 0$	$\sum_i \bar{b}_i \mu_i = 0$
2		$\sum_i b_i c_i(1 - \mu_i) = 0$	$\sum_i \bar{b}_i c_i \mu_i = 0$
2		$\sum_{i,j} b_i a_{ij}(1 - \mu_j) = 0$	$\sum_i \bar{b}_i a_{ij} \mu_j = 0$
3		$\sum_i b_i c_i^2(1 - \mu_i) = 0$	$\sum_i \bar{b}_i c_i^2 \mu_i = 0$
3		$\sum_{i,j} b_i c_i a_{ij}(1 - \mu_j) = 0$	$\sum_{i,j} \bar{b}_i c_i a_{ij} \mu_j = 0$
3		$\sum_{i,j} b_i a_{ij} c_j(1 - \mu_i) = 0$	$\sum_{i,j} \bar{b}_i a_{ij} c_j \mu_i = 0$
3		$\sum_{i,j} b_i a_{ij} c_j(1 - \mu_j) = 0$	$\sum_{i,j} \bar{b}_i a_{ij} c_j \mu_j = 0$
3		$\sum_{i,j,k} b_i a_{ij} a_{jk}(1 - \mu_k) = 0$	$\sum_{i,j,k} \bar{b}_i a_{ij} a_{jk} \mu_k = 0$

**Remark 5** One might also be interested in reducing the contribution of the  $O(h\|e\|^2)$  terms in (10)–(10). Additional independency conditions (corresponding to rooted trees with more than one white vertex) should be considered in that case. For instance, if  $\sum \bar{b}_i \mu_i^2 = 0$ , then the  $O(h\|e\|^2)$  term in (11) could be replaced by  $O(h^2\|e\|^2 + h\|e\|^3)$ .  $\square$

## 5 Practical considerations

Let us assume that Theorem 1 holds with  $r \geq 1$  for the scheme (14)–(16). Then, it makes sense to give, instead of  $y_n$ , the higher order approximation  $\bar{y}_n$  as the numerical solution. In that case,  $\tilde{e}_n$  is no longer an asymptotically correct global error estimate, but an *uncertainty estimate* [5]. We then expect that,  $\tilde{e}_n = y_n - \bar{y}_n$  is bigger than  $\bar{e}_n = \bar{y}_n - y(t_n)$  for sequences of sufficiently small  $h_n$ , and the exact error  $\bar{e}_n$  and the estimated  $\tilde{e}_n$  (hopefully) propagate in a similar way. In that sense, there is no point in taking  $\bar{p} = p + r$  much higher than  $p$ , since in that case  $\tilde{e}_n$  will be an excessively conservative estimate of  $\bar{e}_n$  for sequences of sufficiently small  $h_n$ .

If we compare the application of the scheme (9) defined by (14)–(16) with the application  $\hat{y}_n = \bar{\psi}_{h_n}(\hat{y}_{n-1}, \hat{y}_{n-1})$  of the underlying  $\bar{p}$ th order RK scheme, we observe on one hand that, they need practically the same computational effort per step. And on the other hand, the higher  $\bar{q}$ , the more similar is  $\bar{y}_n$  to  $\hat{y}_n$  (for small  $h_n$  and  $\tilde{e}_n$ ). Therefore, we may expect that, for reasonable values of  $\bar{q}$  in the independency condition (11),  $\bar{y}_n$  and  $\hat{y}_n$  have approximately the same accuracy. In particular, it can be proven that under the hypothesis



of Theorem 1,

$$\bar{y}_n - y(t_n) = (I + O(H^{\bar{q}-r}))(\hat{y}_n - y(t_n)), \quad \text{for } nH \leq \text{Constant}, \quad H = \max_n h_n. \quad (17)$$

These considerations induce us to expect that computing  $\bar{y}_n$  together with the uncertainty estimate  $\tilde{e}_n$  of the global error is practically as efficient as computing the RK solution  $\hat{y}_n$  alone.

## 6 Construction of a method of order 5(4)

In order to see whether generalized globally embedded schemes (14)–(16) can really be as efficient as standard RK methods, we have constructed a method of order 5 (more specifically,  $\bar{s} = 8$ ,  $\bar{p} = 5$ ,  $p = 4$ ,  $\bar{q} = 2$ ,  $q = 1$ ) based on a very efficient RK scheme, namely, the 5th order ERK method of Bogacki and Shampine implemented in the code RKSUITE. We will refer to that RK scheme as BSRK5.

We have determined the coefficients  $\bar{b}_i$  ( $1 \leq i \leq 7$ ), and  $a_{ij}$  ( $1 \leq i, j \leq 7$ ) so that the underlying RK method of order 5 is the scheme BSRK5. Similarly to what is standard in the construction of (locally) embedded RK schemes, we take  $\bar{y}_{n-1}$  and respectively  $\bar{y}_n$  as the first and respectively last stages of the scheme (14)–(16). This is achieved taking  $a_{8i} = \bar{b}_i$  ( $1 \leq i \leq 7$ ),  $\mu_1 = 0$  and  $\mu_8 = 0$ . Thus, although the resulting method is formally a 8th stage order, it only requires 7 evaluations per step.

The remaining parameters, that is,  $b_i$  ( $1 \leq i \leq 8$ ) and  $\mu_i$  ( $2 \leq i \leq 7$ ) are chosen in such a way that the underlying RK method  $\psi_h(y) = \psi_h(y, y)$  is of order 4 and the independency conditions (10)–(11) are satisfied for  $q = 1$  and  $\bar{q} = 2$ . Due to the properties of the coefficients of the method BSRK5 (i.e., the simplifying assumptions made when constructing the scheme), these conditions read as follows: For  $y_n$  being of order 4,

$$b_2 = 0, \quad \sum_i b_i = 1, \quad \sum_i b_i c_i = \frac{1}{2}, \quad \sum_i b_i c_i^2 = \frac{1}{3}, \quad \sum_i b_i c_i^3 = \frac{1}{4}, \quad \sum_i b_i a_{i2} = 0,$$

and the independency conditions for  $q = 1$ ,  $\bar{q} = 2$ ,

$$\sum_i \bar{b}_i \mu_i = 0, \quad \sum_i \bar{b}_i c_i \mu_i = 0, \quad \sum_i b_i (1 - \mu_i) = 0.$$

That leaves five free parameters, which we have chosen as follows:

$$\mu_2 = \frac{800}{261}, \quad \mu_3 = \frac{1469}{500}, \quad \mu_4 = \frac{-520}{101}, \quad \mu_5 = \frac{-2379}{401}, \quad b_8 = \frac{-26}{225}.$$

If we compare our method to the standard RK solution  $\hat{y}_n = \bar{\psi}_{h_n}(\hat{y}_{n-1}, \hat{y}_{n-1})$  (i.e. the solution given by the method BSRK5), we have, according to (17), that

$$\bar{y}_n - y(t_n) = (I + O(H))(\hat{y}_n - y(t_n)), \quad nH \leq \text{Constant}, \quad H = \max_n h_n.$$

## 7 Numerical experiments

We next present some preliminary numerical experiments in constant step-size mode (that is  $h_n = h$  for every  $n$ ). We have considered two initial valued problems:

1. We have taken the first example from [4]. The system is of dimension  $D = 1$ , and the initial value problem is defined as

$$y' = \cos(t)y, \quad y(0) = 1, \quad t \in [0, 94.6].$$

Its solution is  $y(t) = e^{\sin(t)}$ , a periodic function with period  $2\pi$ . We refer to that example as 'expsin'.

2. The second example, referred as 'Arenstorff', is a 4th dimensional initial value problem taken from [3, pp.129–130]. It corresponds to a periodic solution of the restricted 3-body problem.

The aim of these experiments is to check in which extent our method of order 5(4), gives use-full global error estimates, and is as accurate as its underlying standard Runge-Kutta method of order 5 (i.e. the BSRK5 method). For each of the problems and a given step-size  $h$  we display two plots, both showing time versus the infinity norm of the error in double logarithmic scale: The first one compares the exact global error of our numerical solution  $\bar{y}_n$  of order 5 (dashed line) with the estimated global error (continuous line), while the second one compares the exact global error of  $\bar{y}_n$  (dashed line) with the exact global error of the solution  $\hat{y}_n$  provided by the underlying standard RK scheme (BSRK5, in continuous line).

We have integrated the 'expsin' example over  $[0, 30\pi]$  (fifteen periods) with constant step-size  $h = 2\pi/7$ . The corresponding two plots are displayed in Figure 1. We believe that the result are quite encouraging. The estimated global error remains bigger than the exact global error as expected for sufficiently small values of  $h$  (first plot), and our new method clearly behaves as a small perturbation of the BSRK5 solution (second plot), which shows that in that case our generalized globally embedded RK scheme is as efficient as its underlying RK method.

We have integrated the 'Arenstorff' problem over one period  $[0, T]$  ( $T = 17.0652165\dots$ ), first with step-size  $h = T/14000$  (Figure 2). The results are as encouraging as in the 'expsin' problem.

We have made an additional experiment for the 'Arenstorff' problem, this time with a bigger step-size  $h = T/3500$  (Figure 3), in order to check what happens when the numerical integration process gives completely wrong results. The first plot shows that the estimated global error satisfactorily reflects the propagation of the true global error. However, the second plot shows that the 5th order numerical approximation  $\bar{y}_n$  gives completely wrong results before the end of the integration interval, while the BSRK5 solution  $\hat{y}_n$  gives considerably better results. This is apparently due to the fact that the 4th order approximation  $y_n$  degrades before than the 5th order BSRK5 solution  $\hat{y}_n$ , so that  $\tilde{e}_n = y_n - \bar{y}_n$  is no longer

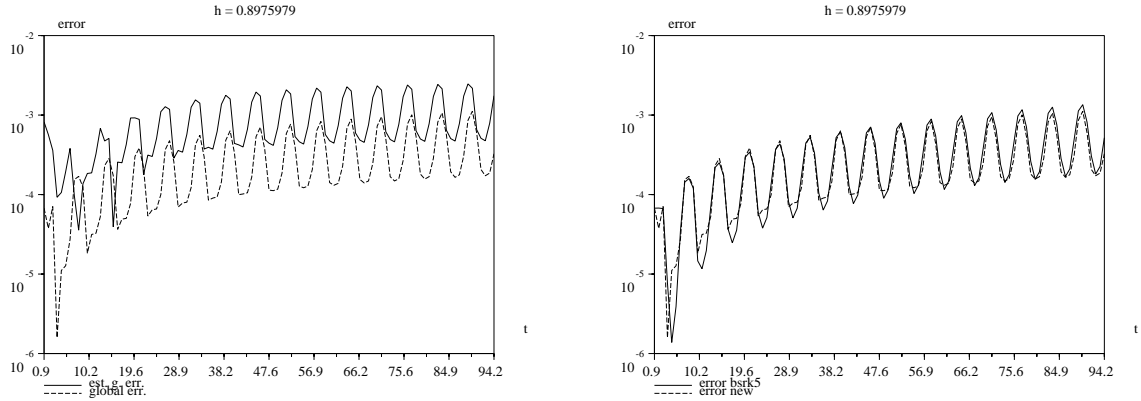


Figure 1: The 'expsin' example for  $h = 2\pi/7$

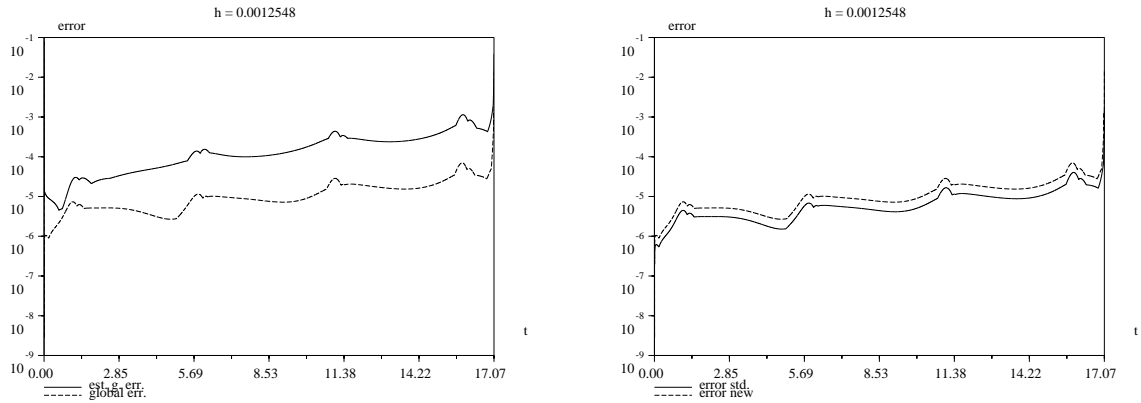


Figure 2: The 'Arenstorf' example for  $h = T/14000$

small, and consequently the  $O(h\|\tilde{e}_{n-1}\|^2)$  term in  $\bar{\pi}_n$  (Lemma 1) dominates the propagation of the error  $\tilde{e}_n$ .

We are currently working on a variable step-size implementation of our scheme. The details about the step-size control strategy, its implementation, and numerical experiments will be reported elsewhere.

**Acknowledgements** This work has been partially supported by grant UPV 140.226-EA 142/98 of the University of the Basque Country.

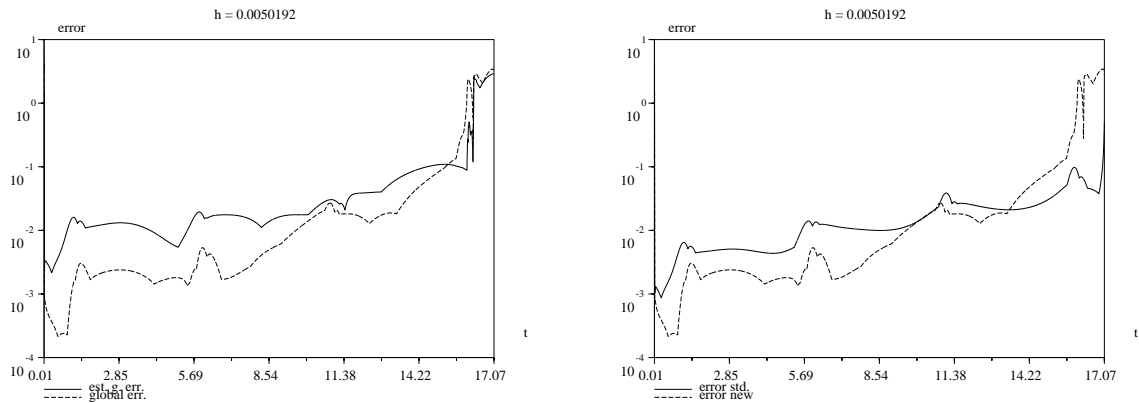


Figure 3: The 'Arenstorff' example for  $h = T/3500$

## References

- [1] J. C. Butcher, *The numerical analysis of ordinary differential equations*, John Wiley & Sons (1987).
- [2] J. R. Dormand, J. P. Gilmore, and P. J. Prince, *Globally Embedded Runge-Kutta Schemes*, *Annals of Numerical Mathematics* 1 (1994), 97–106.
- [3] E. Hairer, S.P. Nørset, G. Wanner: *Solving ordinary differential equations I. Non-stiff problems*, Second Edition, Springer-Verlag (1993).
- [4] L. F. Shampine, *Numerical Solution of Ordinary Differential Equations*, Chapman & Hall (1994).
- [5] R. D. Skeel, *Thirteen Ways to Estimate Global Error*, *Numer. Math.* 48, (1986), 1–20.