# On an accurate third order implicit-explicit Runge–Kutta method for stiff problems

Sebastiano Boscarino

*Department of Mathematics and Computer Science, University of Catania, viale A. Doria 6, 95125, Italy*

## Abstract

Most of the popular implicit-explicit (IMEX) Runge–Kutta (R-K) methods existing in the literature suffer from the phenomenon of order reduction in the stiff regime when applied to stiff problems containing a non-stiff term and a stiff term. Specifically, order reduction is observed when the problem becomes increasingly stiff. In this paper, our motivation is to derive a third-order IMEX R-K method for stiff problems that has a better temporal order of convergence than other well-known IMEX R-K methods. A comparison with other third-order methods shows substantial potential of this new method.

© 2008 IMACS. Published by Elsevier B.V. All rights reserved.

## 1. Introduction

Several physical phenomena of great importance for applications are described by stiff systems of differential equations in the form

$$U' = F(U) + \frac{1}{\varepsilon}G(U) \tag{1}$$

where $U = U(t) \in \mathbb{R}^n$, $F, G : \mathbb{R}^n \to \mathbb{R}^n$, and $\varepsilon > 0$ is the stiffness parameter. Systems of the form (1) with a large number of equations often arise from the discretization of partial differential equations, such as convection–diffusion problems and hyperbolic systems with relaxation [7,20,23,22,1,19] when a method of lines is usually used.

In order to be able to treat problems of the form (1), it could be interesting to separate the non-stiff and the stiff terms. In most cases $F(U)$ is non-linear and non-stiff and $\frac{1}{\varepsilon}G(U)$ contains the stiffness. Then it is desirable to develop numerical methods which are explicit in $F$ and implicit in $G$.

Concerning this, a general approach to the solution of the problem (1) is based on implicit-explicit (IMEX) multistep methods [18,9,3,2] or IMEX R-K methods [7,20,23,22,1]. In this paper we consider IMEX R-K methods. An IMEX R-K method can be also considered both partitioned [17] and additive R-K method [7,1]. This method consists

of applying an implicit discretization for $G$ and a explicit one for $F$. Since systems of the form (1) usually arise from PDEs, simplicity and efficiency in solving the algebraic equations corresponding to the implicit part of the discretization at each step is of fundamental importance. To this aim, it is natural to consider diagonally implicit R-K (DIRK) methods for function $G$. In this work, the method is derived using Singly Diagonally Implicit R-K methods (SDIRK).

Now, we observe that system (1) can be written as a system of $2n$ equations in the form

$$
y' = f(y, z),
$$
$$
\varepsilon z' = g(y, z),
\tag{2}
$$

once we set $U = y + z$, $F(U) = f(y, z)$ and $G(U) = g(y, z)$. On the other hand, system (2) is a particular case of system (1) when $F(U) = (f(y, z), 0)$ and $G(U) = (0, g(y, z))$. System (2), is a singular perturbation problem, (SPP) [25,21]. System (2) allows us to understand many phenomena observed for very stiff problems. When the parameter $\varepsilon$ is small, the corresponding differential equation is stiff, and when $\varepsilon$ tends to zero, the differential equation becomes differential algebraic. A sequence of differential algebraic systems arises in the study of SPPs (see [16,12,4]).

In [4], Boscarino studied the global error behavior of the most popular IMEX R-K methods existing in the literature presenting convergence proofs for different types of IMEX R-K methods which gave sharp error bounds for such methods when applied to system (2). In particular, this study revealed that these methods suffer from the phenomenon of order reduction in the stiff regime ($\Delta t \gg \varepsilon$), when the classical order is greater than two [7,20]. Furthermore, from the practical point of view, the understanding of this phenomenon is essential in situations where one is interested in the construction of higher order methods.

Another interesting aspect is that the order reduction phenomenon still may be considerable about the numerical integration of the semidiscrete equations arising, for instance, after the spatial discretization of convection–diffusion equations. It is well known that the main difficulty in dealing with advection–diffusion problems is that the system becomes stiffer as the spacial mesh is refined.

Then, the goal of this paper is the construction of an IMEX R-K method in such a way that does not suffer from the order reduction phenomenon producing an error which preserves the order of the method and provides a better temporal order convergence than other well-known IMEX R-K methods found in literature.

The development of this method is aided by the knowledge of additional order conditions derived to impose accuracy in time at the various order in the stiffness parameter $\varepsilon$ [5]. Then, we require extra order conditions in addition to the classical ones presented in literature for the IMEX R-K methods [7,20,23,22,1]. We remark that in the classical literature IMEX R-K methods do not satisfy these extra conditions.

The paper is organized as follows. Section 2 will provide a definition and a classification of IMEX R-K methods. After, in Section 3, a brief review of additional order conditions is presented. Section 4 is devoted to the construction of the new IMEX R-K method. In Section 5 we present several test problems, two simple singular perturbation problems and a convection–diffusion equation. Testing of the new IMEX R-K method on these problems is discussed in Section 6 and comparisons are made with other IMEX R-K methods. Appendix is also included.

## 2. IMEX Runge–Kutta methods: definition and classification

An IMEX R-K scheme applied to (1) has the following form

$$
U_{n+1} = U_n + h \sum_{i=1}^{s} \tilde{b}_i F(t_n + \tilde{c}_i h, U^i) + h \sum_{i=1}^{s} b_i \frac{1}{\varepsilon} G(t_n + c_i h, U^i)
\tag{3}
$$

with internal stages given by

$$
U^i = U_n + h \sum_{j=1}^{i-1} \tilde{a}_{ij} F(t_n + \tilde{c}_i h, U^i) + \sum_{j=1}^{i} a_{ij} \frac{1}{\varepsilon} G(t_n + c_i h, U^i).
\tag{4}
$$

The matrices $(\tilde{a}_{ij})$, with $\tilde{a}_{ij} = 0$ for $j \geqslant i$, and $(a_{ij})$ are $s \times s$ matrices such that the resulting method is explicit in $F$, and implicit in $G$. A diagonally implicit method ($a_{ij} = 0$, for $j > i$) for $G$ gives a sufficient condition to guarantee that

$F$ is always evaluated explicitly. The methods are characterized by the vectors $\tilde{c} = (\tilde{c}_1, \ldots, \tilde{c}_s)^T$, $\tilde{b} = (\tilde{b}_1, \ldots, \tilde{b}_s)^T$, $c = (c_1, \ldots, c_s)^T$, $b = (b_1, \ldots, b_s)^T$. They can be represented by a double *tableau* in the usual Butcher notation,

$$\begin{array}{c|c} \tilde{c} & \tilde{A} \\ \hline & \tilde{b}^T \end{array} \qquad \begin{array}{c|c} c & A \\ \hline & b^T \end{array}.$$

The coefficients $\tilde{c}$ and $c$, used for the treatment of non-autonomous systems, are given by the relation

$$\tilde{c}_i = \sum_{j=1}^{i-1} \tilde{a}_{ij}, \qquad c_i = \sum_{j=1}^{i} a_{ij}. \tag{5}$$

The great number of IMEX R-K methods presented in literature leads us to classify them in three different types characterized by the structure of the matrix $A = (a_{ij})_{i,j=1}^{s}$ of implicit method [4].

**Definition 2.1.** We call an IMEX R-K method of type A (see [22]), if the matrix $A \in \mathbb{R}^{s \times s}$ is invertible.

**Definition 2.2.** We call an IMEX R-K method of type CK (see [7]), if matrix $A \in \mathbb{R}^{s \times s}$ can be written as

$$A = \begin{pmatrix} 0 & 0 \\ a & \hat{A} \end{pmatrix}$$

with the submatrix $\hat{A} \in \mathbb{R}^{(s-1) \times (s-1)}$ invertible.

**Remark.** IMEX R-K methods, called of type ARS (see [1]), are a special case of the type CK with the vector $a = 0$.

## 3. Algebraic order conditions

IMEX Runge–Kutta methods can be viewed as a particular class of partitioned Runge–Kutta methods. Therefore, their order conditions can be derived from the general theory of partitioned methods. The derivation of order conditions relies on the theory of rooted trees and already extensively presented in [16,14]. In particular a generalization of the theory of the order conditions for partitioned Runge–Kutta method was derived in [11].
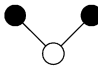
In this section we introduce additional order conditions for IMEX R-K methods when applied to differential algebraic equations (DAEs). We obtained these additional order conditions in a similar way as done in [14,16] by imposing that the IMEX R-K methods were of a given order in the hierarchy of the DAEs [16,13]. In particular, these methods were applied directly to index 1 and 2 differential algebraic systems. We remark that the corresponding reduced system of (2), obtained by considering the limit case $\varepsilon = 0$, is an index 1 differential algebraic system if $g_z(y, z)$ is invertible in a neighborhood of the solution (see [16,11]). For a detail analysis about these additional order conditions the reader is referred to [5]. Instead, for a detail account of the classical order conditions for IMEX R-K methods see, for example [7,20].
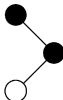
Exactly as done in [14], we derived the order conditions with the help of the so-called bicolour rooted trees [14, 16,6]. Using bicolour rooted trees the Taylor expansion of the solution of DAEs can be easily written in terms of the so-called elementary differentials. For a graphical representations of these formulas, in order to distinguish the derivatives with respect to $y$ and $z$, we need two kind of vertices: *meagre* ($\circ$) and *fat* ($\bullet$). For a full description of the elementary differentials and the bicolour rooted trees we refer to [14–16]. Below we collect several order conditions for index 1 DAEs for the $z$-component and the respective elementary differentials. We note that in this representation we identify the expression $(-g_z)^{-1}g$ with a fat vertex and $f$ with a meagre one.

**Index 1 order conditions.** We have $z(t_0 + h) - z_1 = \mathcal{O}(\Delta t^2)$ if

$$(g_z)^{-1} g_y f, \qquad \sum_{i,j} b_i w_{ij} \tilde{c}_j = 1.$$

We have $z(t_0 + h) - z_1 = \mathcal{O}(\Delta t^3)$ if

$$(g_z)^{-1}g_{yy}(f,f), \qquad \sum_{i,j} b_i w_{ij} \tilde{c}_j^2 = 1,$$



$$(g_z)^{-1}g_y f_y f, \qquad \sum_{i,j,k} b_i w_{ij} \tilde{a}_{jk} \tilde{c}_k = 1/2.$$



Here the coefficients $w_{ij}$ are the elements of the inverse of the matrix $A$ introduced in Section 2.

Furthermore, motivated by convergence results in [4] for different types of IMEX R-K methods when applied to semi-explicit index 2 systems, additional order conditions on the coefficients of this new method are deduced in [5] thus to substantially improve the estimates given in [4]. We remark that the derivation of these order conditions is in complete analogy to the derivation of the order conditions for semi-explicit index 2 systems in [16]. We briefly recall the basic ideas. We start to consider the following system

$$u' = \mathcal{F}(u,v,z), \tag{6a}$$
$$v' = \mathcal{K}(u,v,z), \tag{6b}$$
$$0 = \varphi(u,v). \tag{6c}$$

We assume that $\mathcal{F}$, $\mathcal{K}$ and $\varphi$ are sufficiently differentiable. The system (6) is a differential algebraic system of index 2 if $\varphi_v \mathcal{K}_z$ is invertible in a neighborhood of the solution. For a graphical representation the derivatives of $u$ are characterized by trees with *meagre root*, where the trees will be denoted by $t$ or $t_i$ (the root by $\tau = \bullet$) and the derivatives of $v$ are characterized by trees with *square root*, where the trees are indicated by $w$ or $w_i$ (the root by $\eta = \square$). Instead, the derivatives of $z$ are characterized by trees with a *fat root* and they will be denoted by $r$ or $r_i$. We identify each occurring $\mathcal{F}$ with a *meagre* vertex and each of its derivatives with an upwards leaving branch. The expression $(-\varphi_v \mathcal{K}_z)^{-1}\varphi$ with a *fat* vertex and the derivatives of $\varphi$ therein again with upwards leaving branches. Instead, we identify each occurring $\mathcal{K}$ with a *square* vertex and each of its derivatives with an upwards leaving branch. Let $\mathrm{DAT2} = \mathrm{DAT2}_u \cup \mathrm{DAT2}_v \cup \mathrm{DAT2}_z$ denote a generalization of the set of differential algebraic index 2 tree defined recursively in a similar way as done in [16].
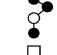
In Table 1 we collect the order conditions explicitly computed for some trees of DAT2. We denote by $\rho(t)$, $\rho(w)$ and $\rho(r)$ the order of the trees. We have not included the trees which have only meagre and square vertices because these order conditions are exactly the classical order condition for IMEX Runge–Kutta methods (see [7,20,22]). The trees of the classical order conditions form a subset of the trees considered here. We also observe that some of the order conditions are identical to those for index 1, i.e., $\sum_{ij} b_i \omega_{ij} c_j = 1$.

From now on, we consider the following conditions (see also [16])

$$B(p): \quad \sum_{i=1}^{s} b_i c_i^{k-1} = \frac{1}{k}, \quad k = 1, \dots, p;$$

$$C(\sigma): \quad \sum_{j=1}^{s} a_{ij} c_j^{k-1} = \frac{c_i^k}{k}, \quad i = 1, \dots, s, \ k = 1, \dots, \sigma;$$

$$D(\zeta): \quad \sum_{i=1}^{s} b_i c_i^{k-1} a_{ij} = \frac{b_j}{k}\left(1 - c_j^k\right), \quad j = 1, \dots, s, \ k = 1, \dots, \zeta,$$

$$\tilde{C}(\sigma): \quad \sum_{j=1}^{s} \tilde{a}_{ij} c_j^{k-1} = \frac{c_i^k}{k}, \quad i = 1, \dots, s, \ k = 1, \dots, \sigma;$$

$$\tilde{D}(\zeta): \quad \sum_{i=1}^{s} b_i c_i^{k-1} \tilde{a}_{ij} = \frac{b_j}{k}\left(1 - c_j^k\right), \quad j = 1, \dots, s, \ k = 1, \dots, \zeta, \tag{7}$$

called simplifying assumptions. Here we have imposed that $\tilde{b}_i = b_i$ and $\tilde{c}_i = c_i$, for $i = 1, \dots, s$. Several consequences we can be derived considering the previous assumptions. Assumption $\tilde{C}(\sigma)$ usually can only be satisfied for the trivial case $\sigma = 1$. On other hand, applying assumptions $D(\xi)$ and $\tilde{D}(\xi)$ produce two inconsistent equations at $j = s$ where

Table 1
Trees and order conditions.

| $\rho(t)$ | graph | order condition |
|---|---|---|
| 2 | $t_1$ | $\sum \tilde{b}_i \omega_{ij} \tilde{c}_j^2 = 1$ |
| 2 | $t_2$ | $\sum \tilde{b}_i \omega_{ij} c_j^2 = 1$ |
| 2 | $t_3$ | $\sum \tilde{b}_i \omega_{ij} \tilde{a}_{jk} \tilde{c}_k = \frac{1}{2}$ |
| 2 | $t_4$ | $\sum \tilde{b}_i \omega_{ij} \tilde{a}_{jk} c_k = \frac{1}{2}$ |
| 2 | $t_5$ | $\sum \tilde{b}_i \omega_{ij} \tilde{c}_j c_j = 1$ |
| $\rho(w)$ | graph | order condition |
| 2 | $w_1$ | $\sum b_i \omega_{ij} \tilde{c}_j^2 = 1$ |
| 2 | $w_2$ | $\sum b_i \omega_{ij} c_j^2 = 1$ |
| 2 | $w_3$ | $\sum b_i \omega_{ij} \tilde{a}_{jk} \tilde{c}_k = \frac{1}{2}$ |
| 2 | $w_4$ | $\sum b_i \omega_{ij} \tilde{a}_{jk} c_k = \frac{1}{2}$ |
| 2 | $w_5$ | $\sum b_i \omega_{ij} \tilde{c}_j c_j = 1$ |
| $\rho(r)$ | graph | order condition |
| 1 | $r_1$ | $\sum b_i \omega_{ij} \omega_{jk} \tilde{c}_k^2 = 2$ |
| 1 | $r_2$ | $\sum b_i \omega_{ij} \omega_{jk} c_k^2 = 2$ |
| 1 | $r_3$ | $\sum b_i \omega_{ij} \omega_{jk} \tilde{a}_{kl} \tilde{c}_l = 1$ |
| 1 | $r_4$ | $\sum b_i \omega_{ij} \omega_{jk} \tilde{a}_{kl} c_l = 1$ |
| 1 | $r_5$ | $\sum b_i \omega_{ij} \omega_{jk} \tilde{c}_k c_k = 1$ |

no values of $c_s$ can satisfy them simultaneously. Hence, if the method is stiffly accurate we get $c_s = 1$ and consequently we can only apply this simplifying assumption to the explicit part of the method. We remark that the condition $C(\sigma)$ is equivalent to $\sum_{j=1}^{s} \omega_{ij} c_j^k = k c_i^{k-1}$ for $k = 1, \ldots, \sigma$ and $D(\zeta)$ is equivalent to $\sum_{j=1}^{s} b_i c_i^k \omega_{ij} = \sum_{i=1}^{s} b_i \omega_{ij} - k b_j c_j^{k-1}$ for $k = 1, \ldots, \zeta$.

These simplifying assumptions will be helpful because they simplify considerably the construction of the new IMEX R-K method.

## 4. Construction of a third order IMEX R-K method

Now we provide a procedure to derive a new third-order IMEX R-K method that works much better than existing IMEX R-K methods in the stiff regime where usually an order reduction phenomenon appears. We start to consider a known IMEX R-K method and modify its coefficients in such a way that the new order conditions, introduced before, are satisfied.

Among the different type of IMEX R-K methods we choose to construct methods of type CK [7,4]. These methods are characterized to have the matrix $A$ according to Definition 2.2, several restrictions on the coefficients as $b_i = \tilde{b}_i$, $c_i = \tilde{c}_i$ and the stiffly accurate condition $a_{si} = b_i$ for $i = 1, \ldots, s$. These restrictions reduce the degree of freedom available to satisfy all classical order conditions and, in particular, the condition $a_{si} = b_i$ facilitates $L$-stability of the implicit part of the method. Furthermore, the implicit part differs from the classical SDIRK one (see for example [16])

by having the element $a_{11} = 0$ (*Explicit*, Singly Diagonally Implicit R-K method, ESDIRK, see [7]) and by getting an explicit first stage so to increase the stage order at two, instead the effective stage order, i.e., one.

Then, based on this knowledge, we require that the new method satisfies the following assumptions:

0. $\tilde{b}_i = b_i$, $\tilde{c}_i = c_i$, for $i = 1, \ldots, s$;
1. $\sum_{j=1}^{i} \hat{a}_{ij} c_j^{k-1} = c_i^k / k$ for $i = 2, \ldots, s-1$ and $k = 1, 2$;
2. $\sum_{j=1}^{i-1} \tilde{a}_{ij} c_j^{k-1} = c_i^k / k$ for $i = 3, \ldots, s$ and $k = 1, 2$;
3. $\sum_{i=1}^{s} b_i = 1$, $\sum_{i=1}^{s} b_i c_i = 1/2$, $\sum_{i=1}^{s} b_i c_i^2 = 1/3$;
4. $\tilde{b}_2 = b_2 = 0$;
5. $\sum_j b_i \hat{\omega}_{ij} \tilde{c}_j = 1$, $\sum_{j,k} b_i \hat{\omega}_{ij} \hat{\omega}_{jk} \tilde{c}_k^2 = 2$ for $i = 2, \ldots, s$;

1. and 2. are the simplifying assumptions $C(2)$, and $\widetilde{C}(2)$ for the implicit and explicit part with $c_i = \tilde{c}_i$ for $i = 1, \ldots, s$. The technical condition 4. is necessary because the assumption 2. cannot be satisfied for $i = 2$. The coefficients $\hat{\omega}_{ij}$ are the elements of the inverse of the matrix $\hat{A}$.

It is worth commenting that considering the assumption 0., the number of classical order conditions, with coupling conditions associated, are reduced only to formulas in 3. (see [7,20,22]). Instead, formulas 5. are the new order conditions. We note that the classical IMEX R-K methods existing in the literature do not satisfy conditions 5.

Formulas 5. appear in a very simple and natural way. For instance, in Table 1 if 0. is satisfied we obtain only $w_1$, $w_3$ and $r_1$ and $r_3$. On the other hand from 2., the order conditions $r_3$ and $w_3$ are the same of $r_1$ and $w_1$. This leaves only two index 2 order conditions that, from

$$\sum_{j=2}^{i} \hat{\omega}_{ij} c_j^k = k c_i^{k-1}, \quad i = 2, \ldots, s, \; k = 1, 2; \tag{8}$$

with $k = 2$, it follows that $w_1$ is equal to $\sum_{i=2}^{s} b_i c_i = 1/2$. Then, we get only the condition $r_1$. The other formula in 5. is obtained analogously by considering the index 1 order conditions and the assumption 2. and (8). Further, it may also be seen that order conditions 5. can be automatically satisfied because if the method is stiffly accurate, from $\sum_{i,j} b_i \hat{\omega}_{ij} \tilde{c}_j = 1$ we get $\tilde{c}_s = c_s = 1$ and from $\sum_{i,j} b_i \hat{\omega}_{ij} \hat{\omega}_{jk} \tilde{c}_k^2 = 2$, considering also (8), we get again $c_s = 1$.

Now a careful study of the type CK method provides two other assumptions that have to be satisfied. We first give a $L$-stability condition for the implicit part of the method.

We remark that, in this paper, the stability properties of the full method are considered as a consequence of the stability properties of the implicit and the explicit methods, separately. More precisely, as in [1], we required that the implicit method is $L$-stable and the explicit one is selected to have the largest stability region.

Of course, this is not enough to ensure good stability properties for the overall IMEX R-K method. Notice that in [6], an appropriate extension of the concept of $L$-stability for linearly implicit R-K methods is given when applied to advection–reaction–diffusion equations and a stability analysis for IMEX schemes applied to stiff systems has been performed by Higueras et al., see [10], using contractivity and monotonicity properties and concept of algebraic stability for additive R-K methods.

In the following we analyze only the stability requirements for the implicit method. We emphasize that the only stiffly accurate condition for ESDIRK part of the method is not enough to guarantee that $R(\infty) = 0$ and an additional condition is required. In this case we have that, if

$$-e_s^T \hat{A}^{-1} a = \sum_{j \geqslant 2} \hat{\omega}_{sj} a_{j1} = 0, \tag{9}$$

then, $R(\infty) = 0$.

Here $R(z)$ is the absolute stability function and $e_s^T = (0, \ldots, 1)^T$. This result is obtained in a very simple way. In fact, to obtain that, we apply one step of the ESDIRK part of the method with matrix given in Definition 2.2 to equation $y' = \lambda y$ and with initial value $y_0 = 1$. The reader is referred to [5], where a proof is given.

However, this is not sufficient for the $L$-stability because the ESDIRK part should also be $A$-stable. Thus, we consider the stability function $R(z)$ for a ESDIRK method that has the form

$$R(z) = \frac{P(z)}{(1 - \gamma z)^{s-1}} \tag{10}$$

with

$$P(z) = (-1)^{s-1} \sum_{j=0}^{s-2} L_{s-1}^{s-1-j} \frac{1}{\gamma} (\gamma z)^j \tag{11}$$

and error constant $C = (-1)^{s-1} L_{s-1}(1/\gamma) \gamma^{s-1}$ where $L_s$ is the $s$-degree Laguerre polynomial and $L_s^{(k)}$ denotes its $k$-th derivative, [16]. Now, concerning the ESDIRK part of the IMEX R-K method (see [16], Section IV.6), the function (10) is analytic in $\mathbb{C}^-$, provided that $\gamma > 0$, and $A$-stability is equivalent to the fact that the polynomial

$$E(y) = |Q(iy)|^2 - |P(0, iy)|^2 \geqslant 0, \quad \text{for all } y, \tag{12}$$

i.e., stability on the imaginary axis ($I$-stability). For completeness, we give the following explicit formula for $E(y)$ where we have $s - 1 = 4$ with $p \geqslant s - 2 = 3$ and $s$ the number of the internal stages,

$$\begin{aligned}
E(y) = {}& y^4 \left( 1/12 - 4\gamma/3 + 6\gamma^2 - 8\gamma^3 + 2\gamma^4 \right) \\
& + y^6 \left( -1/36 + 2\gamma/3 - 6\gamma^2 + 76\gamma^3/3 - 52\gamma^4 + 48\gamma^5 - 12\gamma^6 \right) + y^8 \gamma^8.
\end{aligned}$$

Then (see [16] Table 6.4), the region of $\gamma$ for $A$-stability (and hence $L$-stability) is given by the following interval

$$0.22364780\ldots \leqslant \gamma \leqslant 0.57281606\ldots . \tag{13}$$

Finally, we require that

$$\sum_{j=2}^{s} b_i \hat{\omega}_{ij} \hat{\omega}_{j2} = 0, \quad i = 2, \ldots, s, \tag{14}$$

is included among the assumptions. This condition follows from the proof of Theorem 6.2 in [4,5] about the type CK. We provide this assumption in order to have at least the second order of accuracy in time for the algebraic variable. In fact, without (14) the order of the algebraic variable drops to first order. Hence, an extra care must be taken to properly construct this method.

We remark that if the method is stiffly accurate the assumption (14) is equivalent to $\hat{\omega}_{s2} = 0$ and (9) to $\sum_{j\geqslant 3} \hat{\omega}_{sj} a_{j1} = 0$.

Unfortunately, there are not third-order four-stages IMEX R-K that satisfy all the previous assumptions. This is formally stated in the following.

**Proposition 1.** *Consider IMEX Runge–Kutta method of type CK, stiffly accurate (i.e., $a_{si} = b_i$ for $i = 1, \ldots, s$) in the implicit part, with $b_i = \tilde{b}_i$, $c_i = \tilde{c}_i$ for all $i$ and $a_{11} = 0$. Then there exist no third-order four stage method satisfying assumptions* 1.–4. *with condition* (9) *and assumption* (14).

The proof is very trivial (see [5]), in fact, one can prove easily by algebraic computations several contradictions. Then, a third-order five-stages IMEX R-K method of type CK is designed, i.e., $s = 5$.

In the construction of a third-order type CK IMEX R-K method, we first evaluate the coefficients of a third-order five-stages ESDIRK method, which is stiffly accurate, $a_{ii} = \gamma$, for all $i$, with $\gamma > 0$, and $a_{11} = 0$. Moreover, assumption 0. is imposed. Then, we require to solve the system of equations from 1. to 4. with the condition (9) and the assumption (14). We remark that order conditions in 5. are satisfied as we proved above.

Then, we solve the system

$$1 - b_3 - b_4 - \gamma = b_1, \tag{15a}$$

$$b_3 c_3 + b_4 c_4 + \gamma = 1/2,$$

$$b_3 c_3^2 + b_4 c_4^2 + \gamma = 1/3, \tag{15b}$$

$$\gamma c_2 = c_2^2/2,$$

$$a_{32}c_2 + \gamma c_3 = c_3^2/2,$$

$$a_{42}c_2 + a_{43}c_3 + \gamma c_4 = c_4^2/2, \tag{15c}$$

$$\frac{b_3 a_{32}}{\gamma^3} + \frac{b_4 a_{42}}{\gamma^3} - \frac{b_4 a_{43} a_{32}}{\gamma^4} = 0, \tag{15d}$$

$$\frac{b_4 a_{43} a_{31}}{\gamma^3} - \frac{b_3 a_{31}}{\gamma^2} - \frac{b_4 a_{41}}{\gamma^2} + \frac{b_1}{\gamma} = 0, \tag{15e}$$

with $\sum_{j=1}^{i} a_{ij} = c_i$ for $i = 2, 3, 4$. Now, from the first and second formula in (15c) we get $c_2 = 2\gamma$ and $a_{32} = (c_3(c_3 - 2\gamma))/(2c_2)$. In particular, we get $a_{21} = \gamma$. Using the third formula in (15c), formulas in (15b) and (15d), we compute $b_3$, $b_4$, $a_{42}$ and $a_{43}$ as function of $c_3$, $c_4$ and by trivial computation we obtain $a_{i1}$ for $i = 3, 4$ and $b_1$.

Finally, by using the condition (15e) and substituting the quantities computed before, by algebraic manipulations, we obtain:

$$c_3 = \frac{2(6\gamma^2 - 6\gamma + 1)}{3(2\gamma^2 - 4\gamma + 1)} \tag{16}$$

with $c_4$ a free parameter. Notice that the free parameter $c_4$ has been computed to minimize the fourth-order error terms [15].

Now, concerning the constructing of explicit part of the method, we consider the assumption 2. with 4., i.e., $\sum_{j=1}^{s} \tilde{a}_{ij} c_j^{k-1} = c_i^k/k$ for $i = 3, 4, 5$ with $k = 1, 2$.

Also, we impose that the condition

$$\sum_{i,j,k} b_i \tilde{a}_{ij} \tilde{a}_{jk} c_k = 1/24, \tag{17}$$

is satisfied so that the explicit R-K part of the method has largest possible stability region, i.e., the same stability region as a four-order explicit R-K method.

It is noteworthy that, considering the assumption 2., (17) implies $\sum_{i,j} b_i \tilde{a}_{ij} c_j^2 = 1/12$ if and only if $\sum_i b_i \tilde{a}_{i2} = 0$. Thus, if we choose $\tilde{a}_{42} = 0$ we can immediately compute $\tilde{a}_{52}$ from $b_3 \tilde{a}_{32} + \gamma \tilde{a}_{52} = 0$ and, by using the assumption 2. we get $\tilde{a}_{43} = c_4^2/2c_3$ and $\tilde{a}_{32} = c_3^2/2c_2$ for $i = 3, 4$ and $k = 2$. Finally, we get $\tilde{a}_{21} = 2\gamma$, $\tilde{a}_{31} = c_3 - \tilde{a}_{32}$, $\tilde{a}_{41} = c_4 - \tilde{a}_{43}$, $\tilde{a}_{51} = 1 - \tilde{a}_{52} - \tilde{a}_{53} - \tilde{a}_{54}$ with $\tilde{a}_{53}$ and $\tilde{a}_{54}$ computed using again assumption 2. and $\sum_{i,j} b_i \tilde{a}_{ij} c_j^2 = 1/12$.

## 5. Test problems

In order to test the accuracy of the method described above, we first consider two stiff problems, the Van der Pol's equation [16], and Pareschi Russo's problem [20]. Moreover, as third test problem, we consider an advection–diffusion problem.

The Van der Pol's equation is one of the simplest non-linear equations (describing non-linear oscillations) in the stiff literature. It has the following form

$$\dot{y}(t) = z(t),$$

$$\dot{z}(t) = \varepsilon^{-1}\big((1 - y(t)^2)z(t) - y(t)\big) \tag{18}$$

with $0 \leqslant \varepsilon \ll 1$, as $\varepsilon$ goes to zero this equation becomes increasingly stiff. The IMEX R-K methods considered here treat the first equation explicitly and the second implicitly. This test problem is chosen to compare the accuracy of several types of IMEX R-K methods and the new method when the stiffness parameters $\varepsilon$ is sufficiently small. Numerical results for different types of IMEX R-K methods confirm order reduction especially for the algebraic $z$-component. We will conduct convergence tests with initial conditions, [16]

$$y(0) = 2, \quad z(0) = -\frac{2}{3} + \frac{10}{81}\varepsilon - \frac{292}{2187}\varepsilon^2 - \frac{1814}{19683}\varepsilon^3 + \mathcal{O}(\varepsilon^4),$$

such that the solution is smooth and we choose $\varepsilon$ in a wide range of values from $\varepsilon = 10^{-6}$ to one. In the following figures, we have plotted the relative global error at $t_{\text{end}} = 0.55139$ as a function of the step size $\Delta t$, which was taken constant over the considered interval $[0, t_{\text{end}}]$. We use logarithmic scales in both directions. The relative global error

Fig. 1. Global error versus the stepsize in the Van der Pol equation calculated with $\varepsilon = 10^{-6}$.



Fig. 2. Global error of the third-order BHR(5, 5, 3) method versus the stepsize.

behaves like $C \cdot \Delta t^r$ with $r$ the slope of the straight line and $C$ is a constant. We have indicated this behavior in Figs. 1 and 2.

Pareschi Russo's problem [20], is a simple prototype of stiff system of the form (1), which contains both non-stiff and stiff term. It has the following form

$$\dot{y}(t) = -z(t),$$
$$\dot{z}(t) = y(t) + \varepsilon^{-1}\big(\sin\big(y(t)\big) - z(t)\big). \tag{19}$$

The term multiplied by $\varepsilon^{-1}$ is integrated with the implicit part of the method while other terms with explicit one. The purpose of this test problem is to investigate the numerical convergence rate of each IMEX R-K methods presented in this paper for a while range of values of the stiff parameter $\varepsilon$. This test problem is accomplished with initial values

$$y(0) = \pi/2, \qquad z(0) = \sin y(0) + \varepsilon\big(y(0) + \sin y(0) \cos y(0)\big) + \mathcal{O}\big(\varepsilon^2\big).$$

The numerically observed temporal order of convergence is calculated by

$$p = \frac{\log(E_{\Delta t_1}/E_{\Delta t_2})}{\log(\Delta t_1/\Delta t_2)}$$

Table 2
Convergence rate of $z$-component in $L_\infty$-norm.

| Schemes | $\varepsilon = 1$ | $\varepsilon = 10^{-1}$ | $\varepsilon = 10^{-2}$ | $\varepsilon = 10^{-3}$ | $\varepsilon = 10^{-4}$ | $\varepsilon = 10^{-5}$ | $\varepsilon = 10^{-6}$ |
|---|---|---|---|---|---|---|---|
| ARS | 3.00 | 2.84 | 3.23 | 2.31 | 2.12 | 2.10 | 2.10 |
| MARS | 3.04 | 2.62 | 1.97 | 1.04 | 1.89 | 3.19 | 3.09 |
| ARK3 | 3.05 | 2.95 | 2.45 | 2.12 | 2.02 | 2.01 | 2.01 |
| MARK3 | 3.05 | 2.43 | 2.07 | 1.37 | 1.74 | 2.62 | 2.91 |
| BHR | 2.98 | 2.93 | 2.78 | 3.15 | 3.53 | 3.38 | 3.37 |

with $E_{\Delta t_1}$ and $E_{\Delta t_2}$ the global errors evaluated with steps $\Delta t_1$ and $\Delta t_2 = \Delta t_1/2$. The value $\Delta t_1 = 0.05$ has been used. The system has been integrated for $t \in [0, 5]$. Table 2 summarizes the convergence rates, as a function of $\varepsilon$, of the different types of third order IMEX R-K methods using different values of $\varepsilon$ ranging from $10^{-6}$ to 1.

Finally, as a convection–diffusion problem, we consider the following non-linear partial differential equation

$$u_t + \left(\frac{u^2}{2}\right)_x = \frac{1}{R}u_{xx}, \quad R > 0, \tag{20}$$

called Burgers' equation. It is a mathematical model of free turbulence. It contains non-linear advection term $(u^2/2)_x$ and a dissipation term $u_{xx}/R$, where $R$ is called the Reynolds number. We remark that when the Reynolds number is large the convection term dominates the diffusion term and the solution develops a sharp shock wave front after a certain time of propagation and when $R \to \infty$ discontinuous solutions appear.

For the discretization in space we considered a uniform grid and for the advection term high accuracy in space obtained by finite difference discretization with Weighted-Essentially Non Oscillatory (WENO) reconstruction, [24]. Instead, for the diffusion term we used the standard 4-th order finite difference technique, excepted at the nearby boundary points where a 3-rd order formula was implemented.

The convergence test of the methods is verified using exact and numerical smooth solution. We solve numerically equation (20) under the following boundary conditions $u(0, t) = u(1, t) = 0$ and initial condition $u(x, 0) = \sin(\pi x)$ with $0 \leqslant x \leqslant 1$. The exact solution of this problem was obtained by Cole [8] so that numerical comparison can be done.

In this test problem we choose $R = 0.1$ so that the dissipation term dominates the advection one and a smooth solution is produced. The equation has been integrated from $t = 0$ to $T = 0.02$. Notice that the solution drops dramatically from $\sin(\pi x)$ to zero within the first 0.05 second. In this test problem convection term is treated as non-stiff and is integrated using the explicit part of the method, while diffusion term is treated as stiff and integrated using ESDIRK one.

## 6. Discussion

In order to identify the method derived in this paper we name it BHR($s, \rho, p$), where the triplet $(s, \rho, p)$ characterizes the number $s$ of the stages of the implicit part of the method, $\rho$ the number of the stages of the explicit one and $p$ the order of the method.

In our numerical experiment a good choice for $\gamma$ is, of course, to consider values in interval (13). We took several values of the $\gamma$ ranging from 0.25 to 0.57281606248208 and we have obtained results that are not sensitive to this choice. We only report here the results corresponding to the values $\gamma_1 = 0.435866521508482$ and $\gamma_2 = 0.57281606248208$ which we have considered to be more significant. In order to justify the value of $\gamma_2$ we have to consider again the stability function (10). In fact, if we read the polynomial $P(z)$ in (10) as $P(z) = p_0 + p_1 z + p_2 z^2 + p_3 z^3 + p_4 z^4 + p_5 z^5$ where $p_k = \sum_{j=0}^k \frac{q_j}{(k-j)!}$, for $k = 0, \ldots, 5$, with $q_j = (-\gamma)^j \binom{s-1}{j}$ and $s = 5$ (see [16] Section IV.6), we know that for methods satisfying the stiffly accurate assumption the coefficient of the highest degree of polynomial $P(z)$ is zero. Then in order to obtain that the stability function vanish for the $L$-stability as $z$ tends toward infinity, solving $p_4$ for $\gamma$ results in a equation, $24\gamma^4 - 96\gamma^3 + 72\gamma^2 + 16\gamma + 1 = 0$, having four roots. Only one of these gives $I$-stability resulting in $\gamma_2 = 0.57281606248208$, and leads to a $L$-stable method. Instead, we used the value $\gamma_1$ because it is the same one used in ARS(3,4,3), [1], ARK3(2)4L[2]SA, [7], and in MARS(3,4,3), MARK3(2)4L[2]SA, [4]. Moreover, it belongs to the interval (13). In Appendix we report the coefficients of BHR method for $\gamma_1$ and $\gamma_2$.

Numerical tests are conducted to determine accuracy performance of the new method BHR(5,5,3). Comparisons between BHR(5,5,3) and many existing IMEX R-K methods, as ARS(3,4,3) of Ascher et al. [1], ARK3(2)4L[2]SA of Carpenter and Kennedy [7], and MARS(3,4,3), MARK3(2)4L[2]SA (see [4]), are presented on two different types of ODEs characterized by a stiffness parameter $\varepsilon$ and on a convection–diffusion problem, i.e., Burgers' equation. All methods are formally third order.

We remark that Van der Pol equation and Pareschi–Russo's problem are examples of SPPs and, when $\varepsilon$ decreases these problems become index 1 differential algebraic equations. During this transition, differential variables becomes algebraic ones, an order reduction phenomenon appears and the observed temporal order of convergence fall down the classical order of the method. It is noteworthy that these results are consistent with the theoretical estimates obtained in [4].

To evaluate these two SPPs we have considered the BHR(5,5,3) method 1. in Appendix A. Similarly, if we consider the method 2., we obtain analogous results as in Figs. 1, 2 and Table 2.

About Van der Pol equation, ARS(3,4,3) and ARK3(2)4L[2]SA methods display similar order reduction. In fact, results presented in Fig. 1 show that ARS(3,4,3) and ARK3(2)4L[2]SA methods display an $\mathcal{O}(\Delta t^2)$ error in the second component (algebraic variable) of the solution when $\varepsilon$ is increasingly stiff ($\varepsilon = 10^{-6}$). Instead, for the first component of the solution a third order we observe for the error convergence. We mention that the error of the $y$-component, obtained by ARS(3,4,3) and MARS(3,4,3), is smaller than that obtained with the type CK method as ARK3(2)4L[2]SA, MARK3(2)4L[2]SA and BHR(5,5,3). However, our purpose about these examples is to show how the new method produces a better error estimate for the algebraic component in the stiff regime with respect to the other methods. In fact, it is interesting to note that MARS(3,4,3) and MARK3(2)4L[2]SA methods improve the error estimate in the second component for the solution in the Van der Pol equation displaying an estimate given by $\mathcal{O}(\Delta t^3) + \mathcal{O}(\Delta t)$. A good explanation how to achieve these error estimates is given in [4].

Analogously, an essential gain of accuracy for the $z$-component of the solution is obtained in the Pareshi Russo's problem (see Table 2), where a third order in the limit case ($\varepsilon \to 0$) is provided, in fact, both MARS(3,4,3) and MARK3(2)4L[2]SA show a third-order of accuracy for small and large values of $\varepsilon$ but we have a deterioration of the accuracy for intermediate values. Instead, ARS(3,4,3) and ARK3(2)4L[2]SA show a lost of the accuracy when $\varepsilon$ decreases, i.e., the convergence rate drops to second order for sufficiently stiff parameter ($\varepsilon < 10^{-3}$). This means that a significant benefit for the second (algebraic) component of the solution is observed. But this fact is no enough to ensure an uniform third order in time of the method in a wide range of $\varepsilon$.

Now, taking into account of the assumptions required in the construction of the new method, we obtain the desired results. In fact, BHR(5,5,3) method exhibits the better error estimate with respect to the other methods in the Van der Pol equation. More precisely, Figs. 1 and 2 show better results for this method, a third-order convergence is observed about the algebraic variable for sufficiently stiff parameter ($\varepsilon = 10^{-6}$) and when $\varepsilon = 10^{-k}$, $k = 0, 1, 2, 3, 4, 5, 6$ (see Fig. 2). Results for the $y$-component are similar to those given before. Then, when the parameter $\varepsilon$ decreases the observed order of convergence of the method does not fall below the classical order and the absence of significant order reduction for the BHR(5,5,3) method confirms the gain in term of accuracy in time of this method in the stiff regime.

Analogously, on the Pareschi and Russo's problem, BHR(5,5,3) method exhibits a significant observed third-order of convergence at each value of $\varepsilon$, ranging from $10^{-6}$ to 1. Table 2 shows the corresponding convergence rate for $z$-component in $L_\infty$-norm. Different norms give essentially the same results. Concerning the $y$-component of the solution we do not report the convergence rates for each value of the parameter $\varepsilon$ because a third-order convergence rate for each considered IMEX R-K method is observed.

Now, we use the third test problem to check the temporal order of convergence and the efficiency of the BHR(5,5,3) method. In this example $\Delta x$ decreases with the stepsize $\Delta t$ and, in the Tables we show the quantity $p_2 = \log_2(\|E_{\Delta t_1}\|_2 / \|E_{\Delta t_2}\|_2)$ with $N\Delta t = T$ and $\Delta t_2 = \Delta t_1/2$, i.e., the observed temporal order of convergence, measured using $N = 48$ doubled up to $N = 768$. We have compared the BHR(5,5,3) method with ARS(3,4,3) and ARK3(2)4L[2]SA methods which are formally third-order ones. Concerning the BHR(5,5,3), we remark that the evaluations have been made using method 2. given in Appendix. The results are summarized in Table 3. The second, fourth and sixth columns give the values of the errors and the third, fifth and seventh columns display the order of the accuracy for each method. The observed order for the ARS(3,4,3) and ARK3(2)4L[2]SA methods is about 2.3, respectively, whereas BHR(5,5,3) method converges to 3-rd order. Then, we confirm again a gain in terms of accuracy in time of BHR(5,5,3) method for this test problem.

Table 3
Convergence rate in $L_2$-norm for the Burgers' equation.

| N | ARS(3,4,3) | | ARK3(2)4L[2]SA | | BHR(5,5,3) | |
|---|---|---|---|---|---|---|
| | $L_2$-error | $p_2$ | $L_2$-error | $p_2$ | $L_2$-error | $p_2$ |
| 48 | $1.6 \times 10^{-1}$ | | $1.6 \times 10^{-1}$ | | $1.1 \times 10^{-1}$ | |
| 96 | $3.2 \times 10^{-2}$ | 2.34 | $3.2 \times 10^{-2}$ | 2.36 | $1.6 \times 10^{-2}$ | 2.79 |
| 192 | $6.7 \times 10^{-3}$ | 2.25 | $6.7 \times 10^{-3}$ | 2.58 | $2.3 \times 10^{-3}$ | 2.81 |
| 384 | $1.3 \times 10^{-3}$ | 2.30 | $1.3 \times 10^{-3}$ | 2.30 | $2.9 \times 10^{-4}$ | 3.01 |
| 768 | $2.6 \times 10^{-4}$ | 2.34 | $2.6 \times 10^{-4}$ | 2.34 | $3.2 \times 10^{-5}$ | 3.16 |



Fig. 3. Efficiency diagram for Burgers' equation.

Finally, Fig. 3 gives in the ordinate the global error at the endpoint of integration against in the abscissa the computational cost measured by the computing time (CPU Time) in logarithmic scale. We have applied the methods with fixed step sizes $\Delta t = T/N$ and with $N = 48$ doubled up to $N = 768$. We can observe that among ARS(3,4,3) and ARK3(2)4L[2]SA method, there is not big difference. It is interesting to note that in our test problem the BHR(5,5,3) method shows smaller global errors and needs almost the same CPU time than the other IMEX R-K methods. Then, in term of efficiency, we emphasize that, if the BHR(5,5,3) method has more stages ($s = 5$) than the other comparator methods ($s = 4$), this experiment indicates that this new IMEX R-K method can be considered competitive to each other third-order IMEX R-K method.

## Acknowledgements

# Appendix

1. BHR(5,5,3) IMEX R-K method with $c_2 = 2\gamma$, $\gamma = 0.435866521508482 \approx 424782/974569$.

| 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|
| $c_2$ | $2\gamma$ | 0 | 0 | 0 | 0 |
| $c_3$ | $\tilde{a}_{31}$ | $\tilde{a}_{32}$ | 0 | 0 | 0 |
| $c_4$ | $\tilde{a}_{41}$ | 0 | $\tilde{a}_{43}$ | 0 | 0 |
| 1 | $\tilde{a}_{51}$ | $\tilde{a}_{52}$ | $\tilde{a}_{53}$ | $\tilde{a}_{54}$ | 0 |
| | $b_1$ | 0 | $b_2$ | $b_3$ | $\gamma$ |

| 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|
| $c_2$ | $\gamma$ | $\gamma$ | 0 | 0 | 0 |
| $c_3$ | $a_{31}$ | $a_{32}$ | $\gamma$ | 0 | 0 |
| $c_4$ | $a_{41}$ | $a_{42}$ | $a_{43}$ | $\gamma$ | 0 |
| 1 | $b_1$ | 0 | $b_2$ | $b_3$ | $\gamma$ |
| | $b_1$ | 0 | $b_3$ | $b_4$ | $\gamma$ |

$\tilde{a}_{31} = \gamma$,
$\tilde{a}_{32} = \tilde{a}_{31}$,
$\tilde{a}_{41} = -4758833752202859860333264/59411272693343784570463$,
$\tilde{a}_{42} = 0$,
$\tilde{a}_{43} = 1866233449822026827708736/59411272693343784570463$,
$\tilde{a}_{51} = 6282884581807316958563588168609139173761030824 7/17611291068441210531978163031168634371575305600 0$,
$\tilde{a}_{52} = -b_3$,
$\tilde{a}_{53} = 2623158872930437393370885639960932 07/29742755473037635325208178690649200 0$,
$\tilde{a}_{54} = -9876182318941765814381247170 87/23877337660202969319526901856000$,
$a_{31} = \tilde{a}_{31}$,
$a_{32} = -31733082319927313/45570537722196088937954647102$,
$a_{41} = -3012378541084922027361996617949193605163013778096 10/4512339405658526997790775304503051259795589734581934 9$,
$a_{42} = -6286558929780715329426 8/10255967344161067230558732701909504 7$,
$a_{43} = 41876979692085529960314626700141490094521427700 0/21245436038525770855595459809987481860321716713 9$,
$b_1 = 48769850233674067860351 1/1181159636928185920260208$,
$b_2 = 0$,
$b_3 = 302987763081184622639300143137943089/1535359944203293318639180129368156500$,
$b_4 = -1052359283351006160729382188 63/2282554452064661756575727198000$.
$c_3 = 902905985686/1035759735069$, $c_4 = 2684624/1147171$.

2. BHR(5,5,3) IMEX R-K method with $\gamma = 0.57281606248208 \approx 2051948/3582211$, $c_2 = 2\gamma$.

$\tilde{a}_{31} = 473447115440655855452482357894373/1226306256343706154920072735579148$,
$\tilde{a}_{32} = 1292987660341318823230699787220 19/12263062563437061549200727355 79148$,
$\tilde{a}_{41} = 3749810521082814372451684 8/17264258354639800617 3766007$,
$\tilde{a}_{42} = 0$,
$\tilde{a}_{43} = 762833597425614801408044 16/17264258354639800617 3766007$,
$\tilde{a}_{51} = (-34099758602120646123035398556226393330782744869519)/5886704102363745137792385361113084313351870216475136$,
$\tilde{a}_{52} = (-2374163524338269788569417957340 73)/5546817025768783428914471634 99456$,
$\tilde{a}_{53} = 42981597105462287836382124116507832282 75/2165398513352098924587211488610407046208$,
$\tilde{a}_{54} = 610186561585576085357192228974 9/27286397302587824980364037456 8448$,
$a_{31} = 2592522581696725239027084257804694 69319755/4392887760843243968922388674191715336228$,
$a_{32} = (-172074174703261986564706189586177)/1226306256343706154920072735579148$,
$a_{41} = 110320206157455340528586372919574026878513173939555969375 4/98794577359372770706415224145904934590282646779257673058 37$,
$a_{42} = (-1037545205670589695665425562960873240 94)/45905036388824673483312148227531995 4529$,
$a_{43} = 38632070830699796545968721903772406086027010719471 28/19258690251287609765240683320611425745736762681950551$,
$b_1 = (-2032971420760927701493589)/38017147656515384190997416$,
$b_2 = 0$,
$b_3 = 21976027766516769832652611096438970734 47/94506712327913958354993394737909718 4164$,
$b_4 = (-1281472151942603980706668262353 39)/694684827106875033885629526 26424$.
$c_3 = 12015769930846/24446477850549$, $c_4 = 3532944/5360597$.

## References

[1] U. Ascher, S. Ruuth, R.J. Spiteri, Implicit-explicit Runge–Kutta methods for time dependent partial differential equations, Appl. Numer. Math. 25 (1997) 151–167.

[2] U. Ascher, S. Ruuth, B.T.R. Wetton, Implicit-explicit methods for time dependent PDE's, Appl. Numer. Math. 32 (1995) 797–823.

[3] J.G. Blom, W. Hundsdorfer, J.G. Verwer, An implicit-explicit approach for atmospheric transport-chemistry problems, in: Workshop on the Method of Lines for Time-Dependent Problems, Lexington, KY, 1995, Appl. Numer. Math. 20 (1–2) (1996) 191–209.

[4] S. Boscarino, Error analysis of IMEX Runge–Kutta methods derided from differential algebraic systems, SIAM J. Numer. Anal. 45 (4) (2007) 1600–1621.

[5] S. Boscarino, On the uniform accuracy of implicit-explicit Runge–Kutta methods, Ph.D. Thesis in Mathematics for the Technology, Department of Mathematics and Computer Science, University of Catania, Italy, 2005.

[6] M.P. Calvo, J. de Frutos, J. Novo, Linearly implicit Runge–Kutta methods for advection–reaction–diffusion equations, Appl. Numer. Math. 37 (2001) 535–549.

[7] M.H. Carpenter, C.A. Kennedy, Additive Runge–Kutta schemes for convection–diffusion–reaction equations, Appl. Numer. Math. 44 (1–2) (2003) 139–181.

[8] J.D. Cole, On a quasi-linear parabolic equation occurring in aerodynamics, Quart. Appl. Math. 9 (1978) 225.

 [9] J. Frank, W. Hundsdorfer, J.G. Verwer, On the stability of implicit-explicit linear multistep methods, Appl. Numer. Math. 25 (2–3) (1997) 193–205. Special issue on time integration.
[10] B. Garcia-Celayeta, I. Higueras, T. Roldan, Contractivity/monotonicity for additive Runge–Kutta methods: Inner product norms, Appl. Numer. Math. 56 (2006) 862–878.
[11] E. Hairer, Order conditions for numerical methods for partitioned ordinary differential equations, Numer. Math. 36, 431–445.
[12] E. Hairer, Ch. Lubich, M. Roche, Error of Runge–Kutta methods for stiff problems via differential algebraic equations, BIT 28 (3) (1988) 678–700.
[13] E. Hairer, Ch. Lubich, M. Roche, The Numerical Solution of Differential Algebraic Systems by Runge–Kutta Methods, Lecture Notes in Mathematics, vol. 1409, Springer-Verlag, 1989.
[14] E. Hairer, Ch. Lubich, G. Wanner, Geometric Numerical Integration: Structure-Preserving Algorithms for Ordinary Differential Equations, Springer Series in Computational Mathematics, vol. 31, Springer-Verlag, Berlin, 2002.
[15] E. Hairer, S.P. Nørsett, G. Wanner, Solving Ordinary Differential Equation I: Non Stiff Problems, Springer Series in Computational Mathematics, vol. 8, Springer-Verlag, 1987, second revised edition, 1993.
[16] E. Hairer, G. Wanner, Solving Ordinary Differential Equation II: Stiff and Differential Algebraic Problems, Springer Series in Computational Mathematics, vol. 14, Springer-Verlag, 1991, second revised edition, 1996.
[17] E. Hairer, G. Wanner, On the Butcher group and general multi-value methods, Computing 13 (1) (1974) 1–15.
[18] W. Hundsdorfer, J. Jaffré, Implicit-explicit time stepping with spatial discontinuous finite elements, Appl. Numer. Math. 45 (2–3) (2003) 231–254.
[19] S.F. Liotta, V. Romano, G. Russo, Central schemes for balance laws of relaxation type, SIAM J. Numer. Anal. 38 (2000) 1337–1356.
[20] L. Pareschi, G. Russo, Implicit-explicit Runge–Kutta schemes for stiff systems of differential equations, in: Recent Trends in Numerical Analysis, in: Adv. Theory Comput. Math., vol. 3, Nova Sci. Publ., Huntington, NY, 2001, pp. 269–288.
[21] R.M. O'Malley, Introduction to Singular Perturbations, Academic Press, New York, 1974.
[22] L. Pareschi, G. Russo, High order asymptotically strong-stability-preserving methods for hyperbolic systems with stiff relaxation, in: Hyperbolic Problems: Theory, Numerics, Applications, Springer, Berlin, 2003, pp. 241–251.
[23] L. Pareschi, G. Russo, Implicit-explicit Runge–Kutta schemes and applications to hyperbolic systems with relaxations, J. Sci. Comput. 25 (1) (October, 2005) 129–155.
[24] C.W. Shu, Essentially non-oscillatory and weighted essentially non-oscillatory schemes for hyperbolic conservation laws, in: Advanced Numerical Approximation of Nonlinear Hyperbolic Equation, Lecture Notes in Mathematics, vol. 1967, Springer, New York, 2000.
[25] A.N. Tikhonov, A.B. Vasl'eva, A.G. Sveshnikov, Differential Equations, Springer-Verlag, 1985. Translated from the Russian by A.B. Sossinskij.