# Implicit-Explicit Integral Deferred Correction Methods for Stiff Problems

Sebastiano Boscarino [1], Jing-Mei Qiu[2], Giovanni Russo [3]

**Abstract.**

The main goal of this paper is to investigate the order reduction phenomenon that appears in the integral deferred correction (InDC) methods based on implicit-explicit (IMEX) Runge-Kutta (R-K) schemes when applied to a class of stiff problems characterized by a small positive parameter $\varepsilon$, called singular perturbation problems (SPPs). In particular, an error analysis is presented for these implicit-explicit InDC (InDC-IMEX) methods when applied to SPPs. In our error estimate, we expand the global error in powers of $\varepsilon$ and show that its coefficients are global errors of the corresponding method applied to a sequence of differential algebraic systems. A study of these errors in the expansion yields error bounds and it reveals the phenomenon of order reduction. In our analysis we assume uniform quadrature nodes excluding the left-most point in the InDC method and the globally stiffly accurate property for the IMEX R-K scheme. Numerical results for the Van der Pol equation and PDE applications are presented to illustrate our theoretical findings.

**Keywords:** Stiff Problems, Implicit-Explicit, Runge-Kutta methods, Integral deferred correction methods, Differential algebraic systems.

**AMS Subject Classification Indices:** Stiff equations 65L04, Multistep, Runge-Kutta and extrapolation methods 65L06, Extrapolation to the limit, deferred corrections 65B05, Methods for differential-algebraic equations 65L80.

---

[1]Department of Mathematics and Computer Science, University of Catania, Catania, 95127, E-mail: boscarino@dmi.unict.it
[2]Department of Mathematics, University of Houston, Houston, 77004. E-mail: jingqiu@math.uh.edu. Research supported by Air Force Office of Scientific Computing grant FA9550-16-1-0179, NSF grant DMS-1522777 and University of Houston.
[3]Department of Mathematics and Computer Science, University of Catania, Catania, 95125, E-mail: russo@dmi.unict.it

# 1 Introduction

Several physical phenomena of great relevance for applications are described by autonomous stiff systems of ordinary differential equations (ODEs) of the form

$$U' = F(U) + \frac{1}{\varepsilon}G(U), \quad U(t_0) = U_0, \tag{1.1}$$

where $F, G : \mathbb{R}^n \to \mathbb{R}^n$ are sufficiently smooth functions with different stiff properties and $\varepsilon$ is the stiffness parameters. Usually, system (1.1) with a large numbers of equations may arise from spatial discretization of a system of partial differential equations, such as convection-diffusion problems, diffusion-reaction ones and hyperbolic systems with relaxation, ([1], [7], [28], [9], [25], [32]), where the method of lines approach is usually used.

In order to be able to treat problems of the form (1.1), 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 $G(U)$ contains the stiffness $\varepsilon$. Then it is desirable to develop numerical methods which are explicit in $F$ and implicit in $G$. For the numerical integration of (1.1), additive Runge-Kutta (R-K) methods are proposed, see for instance [9, 28, 7, 32], and are chosen with the aim of efficiently integrating the system (1.1). Additive R-K methods combining implicit and explicit methods are also known in the literature as implicit-explicit (IMEX) R-K methods [1, 9, 28].

We observe that system (1.1) can be written as a system of $2n$ equations in the form

$$\begin{aligned} y' &= f(y, z), \quad y(t_0) = y_0, \\ \varepsilon z' &= g(y, z), \quad z(t_0) = z_0, \end{aligned} \tag{1.2}$$

once we set $U = y + z$, $F(U) = f(y, z)$ and $G(U) = g(y, z)$. When splitting the system, one has to assign initial conditions for both $y$ and $z$, while only initial conditions on the sum are specified in the problem. Any choice of the initial value for $y$ and $z$, such $y_0 + z_0 = U_0$ will give the same dynamics for the sum $y + z$. Note however that such a splitting is adopted here for the purpose of analyzing the scheme. In practice, once the scheme is constructed, it is applied directly to systems in the form (1.1), and no arbitrariness in the initial condition remains.

System (1.2) is known in the literature as *singular perturbetion problems* (SPPs). Classical books on SPPs are [31, 27] and we also refer the reader to the book [19]. System (1.2) has several origins in applied mathematics and allows us to understand many phenomena observed for very stiff problems, for example as mentioned in [19]: from fluid dynamics and results in linear boundary value problems containing a small parameter $\varepsilon$ (the viscosity coefficient), in the study of stiff nonlinear oscillators (such as the Van der Pol system with large value of the parameter) or in the study of chemical kinetics with slow and fast reactions. In fact, 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. In [19] and in the original paper [17], the authors, by studying these differential algebraic systems, offered a general framework in order to derive rigorous convergence estimates of most of the implicit R-K methods presented in the literature when applied to SPPs (1.2) and showed that these methods suffer from the phenomenon of order reduction in the stiff regime.

Due to the equivalence of the two systems (1.1) and (1.2), the error estimate for (1.2) offers a theoretical foundation to understand the error behavior of (1.1). In this paper, we perform our error estimate base on the system (1.2), while the application problems we considered are often in the form of (1.1).

The main goal in this paper is to investigate the order reduction phenomenon that appears in the InDC framework when combined with IMEX R-K methods (InDC-IMEX) when applied to SPPs. The novelty here is to provide rigorous and careful convergence analysis for the global error of InDC-IMEX method. As far as we know, InDC methods constructed using IMEX R-K methods, have not been seriously studied in the literature when applied to SPPs. Furthermore, the error estimates of InDC-IMEX methods given in this paper for SPPs, explain theoretically, for the first time, the order reduction phenomenon that usually occurs for mildly stiff or stiff problems.

We recall that the InDC method [14, 13], along the development of deferred correction (DC) [2, 30] and spectral deferred correction (SDC) [16, 26, 24, 22, 23, 20] methods, is an automatic technique of building up very high order numerical integrators based on lower order ones for ODEs. The procedure consists of one prediction and iterations of correction steps. The high order accuracy is accomplished by using a lower order numerical method to solve a series of error equations in each correction step. Compared with the classical DC method, the recently developed SDC and InDC methods are based on Picard integral equation and a deferred

correction procedure is applied to an integral formulation of the error equation in DC methods. It has been shown that SDC and InDC outperforms DC in many problems with better stability and accuracy properties [16, 14]. The main difference between SDC and InDC is the distribution of quadrature nodes: the SDC method uses Gaussian/Lobatto/Radau points for better stability and accuracy properties, while InDC method uses uniform quadrature points to guarantee high order accuracy increase when high order R-K methods are applied in correction steps [14]. The InDC method can be considered as a R-K type method by assembling the corresponding Butcher table, see [13] and Section 3.3 later on in this paper. In fact, R-K or additive R-K methods of third order or less have been very well developed and optimized for their efficiency and effectiveness [9]. The new aspect, that the InDC method offers, is a systematical way in constructing very high order method (e.g. up to 12th order [12]) without the need to working out order conditions for both explicit and implicit methods. Moreover, the InDC method can be used to improve the operator splitting error [11], which is not possible by using a R-K method.

The development of this theoretical study of IMEX InDC method is aided by the knowledge of some results obtained in the following papers: [3] and [6]. In [3], the author studied the global error behavior of IMEX R-K methods existing in the literature, presenting convergence proofs for different types of IMEX R-K methods, and gave error bounds for such methods when applied to the SPP system (1.2). In particular, this study revealed that IMEX R-K methods suffer from the phenomenon of order reduction in the stiff regime, i,e. $\varepsilon \to 0$. In a similar fashion, in [6] the authors studied the global error behaviour of InDC method constructed by using implicit R-K methods when applied to the SPP system (1.2). This study gave error bounds for such methods. In particular, it revealed the phenomenon of order reduction in the stiff regime, when we increase the order of accuracy by series of correction equations.

This paper is a natural continuation of the research started in [6]. Then by combining the results presented in [3] and [6], we study the error bounds of InDC-IMEX methods for SPPs (1.2) in order to seek an understanding on the order reduction phenomenon. Notice that the main idea is to expand the error in powers of $\varepsilon$ whose coefficients are called error terms, and show convergence results for these error terms, as done in [3]. In our analysis we consider $h \gg \varepsilon$, with $h$ being the time step size. We consider some suitable assumption for our analysis: *the globally stiff accuracy* (GSA) for the IMEX R-K methods (we will define it in the next), and the use of uniform quadrature nodes excluding the leftmost endpoint in the InDC framework. The first assumption guarantees that if a high order globally stiffly accurate IMEX R-K method is used to construct the InDC IMEX R-K method, the assembled matrix is invertible and by this result, the order of convergence for the leading order term in the $\varepsilon$-expansion of global error increases with high order in the correction steps. Furthermore, we will show by a counter example that, if the assumption of GSA for the IMEX R-K method is not satisfied, the corresponding InDC IMEX R-K method does not have the order increase as expected for the leading order term in the $\varepsilon$-expansion of the global error. Note that the assumption of GSA provides a sufficient conditions to guarantee the invertibility of assembled matrix. In general we do not know whether GSA is also a necessary condition. To find this out requires a more detailed investigation.

Finally, note that the uniform distribution of nodes is important to increase accuracy when a high order R-K method is applied in the correction steps for classical problems and the use of quadrature nodes excluding the left-most endpoint leads to an important stability condition for stiff problems, i.e. the method becomes L-stable if A-stable; we refer readers to [23] and [14] for details.

The paper is organized in the following way. In Section 2, we briefly present existing local and global error estimates of IMEX R-K methods for SPPs [3]. In Section 3, we introduce InDC-IMEX methods applied to SPPs (1.2). A reformulation of InDC methods constructed with IMEX R-K methods is obtained, with assembled matrices in double Bouchet tableau as in a classical IMEX method. In Section 4, main theoretical results are stated in the form of a theorem, which is proved via the classical error estimate of IMEX methods presented in [3]. Numerical evidences, supporting these theoretical results, are summarized and presented in Section 5. Conclusions are given in Section 6.

## 2 SPPs and IMEX R-K methods

In this section, we review classical concepts and results of the R-K and IMEX R-K methods when applied to SPPs (1.2).In system (1.2) we assume that $0 < \varepsilon \ll 1$ and $f$ and $g$ are sufficiently differentiable vector-valued functions. The functions $f$, $g$ and the initial values $y(0)$, $z(0)$ may depend smoothly on $\varepsilon$. For simplicity of notation, we suppress such dependence. We require that system (1.2) satisfies

$$\mu(g_z(y, z)) \leq -1, \tag{2.1}$$

3

in an $\varepsilon$-independent neighbourhood of the solution, where $\mu$ denotes the logarithmic norm with respect to some inner product, see for details [19].

When $\varepsilon = 0$, the corresponding *reduced* system is the differential algebraic equation (DAE)

$$\begin{aligned} y' &= f(y, z), \\ 0 &= g(y, z), \end{aligned} \qquad (2.2)$$

whose initial values are *consistent* if $0 = g(y(0), z(0))$. We assume that the Jacobian $g_z(y, z)$ is invertible in a neighborhood of the solution of (2.2). This assumption guarantees the solvability of the second equation in (2.2) and that the equation $g(y, z) = 0$ possesses a locally unique solution $z = \mathcal{G}(y)$ (implicit function theorem). Then insert it into the first equation of (2.2), this gives

$$y' = f(y, \mathcal{G}(y)). \qquad (2.3)$$

The same assumption guarantees that system (2.2) is a differential-algebraic one of index 1, (see [19] for more details).

From a classical result in SPPs theory, condition (2.1) guarantees the existence of an $\varepsilon$-expansion, as the sum of a smooth function of the independent variable $t$ and an exponentially decaying function of the stretched variable $\tau = t/\varepsilon$ (initial layer). The exponentially decaying function is not present if the initial values of system (1.2) (which depend on $\varepsilon$) are on the smooth solution, see Chap. VI.3 in [19] for more details.

Thus in our analysis we seek mainly $\varepsilon$-asymptotic expansion of the exact smooth solutions of the problem (1.2) of the form

$$\begin{aligned} y(t) &= y_0(t) + y_1(t)\varepsilon + y_2(t)\varepsilon^2 + \cdots, \\ z(t) &= z_0(t) + z_1(t)\varepsilon + z_2(t)\varepsilon^2 + \cdots. \end{aligned} \qquad (2.4)$$

Inserting (2.4) into (1.2) and collecting equal power of $\varepsilon$ we have, [19]:

$$g(y_0, z_0) = 0, \qquad (2.5)$$

$$y_0' = f(y_0, z_0), \qquad (2.6)$$

$$\varepsilon^1 : \quad \begin{cases} y_1' = f_y(y_0, z_0)y_1 + f_z(y_0, z_0)z_1, \\ z_0' = g_y(y_0, z_0)y_1 + g_z(y_0, z_0)z_1, \end{cases} \qquad (2.7)$$

$$\cdots$$

$$\varepsilon^\nu : \quad \begin{cases} y_\nu' &= f_y(y_0, z_0)y_\nu + f_z(y_0, z_0)z_\nu + \phi_\nu(y_0, z_0, \cdots, y_{\nu-1}, z_{\nu-1}), \\ z_{\nu-1}' &= g_y(y_0, z_0)y_\nu + g_z(y_0, z_0)z_\nu + \psi_\nu(y_0, z_0, \cdots, y_{\nu-1}, z_{\nu-1}). \end{cases} \qquad (2.8)$$

Eq. (2.5) tells us that $z_0$ is algebraically related to $y_0$. Differentiating (2.5) we have

$$g_y(y_0, z_0)y_0' + g_z(y_0, z_0)z_0' = 0.$$

Compatibility with (2.6) implies

$$g_y(y_0, z_0)f(y_0, z_0) + g_z(y_0, z_0)(g_y(y_0, z_0)y_1 + g_z(y_0, z_0)z_1) = 0.$$

Therefore also $z_1$ is linked to $y_0$, $z_0$, and $y_1$, by an algebraic relation. Then we have the following definition:

**Definition 2.1.** An initial condition $(y(0), z(0))$ for the system (1.2) is *well-prepared* to order $\nu$ in $\varepsilon$ if the expansion

$$\begin{aligned} y(0) &= y_0(0) + y_1(0)\varepsilon + y_2(0)\varepsilon^2 + \cdots, \\ z(0) &= z_0(0) + z_1(0)\varepsilon + z_2(0)\varepsilon^2 + \cdots, \end{aligned} \qquad (2.9)$$

satisfies the algebraic conditions (2.5), (2.6), (2.7) and (2.8) up to order $\nu$ in $\varepsilon$.

Note that arbitrary initial values introduce an initial layer in the solution. One possible way to overcome this difficulty is simply to use a small time step $\Delta t = \mathcal{O}(\varepsilon)$ during the initial transient. After a short time, the initial layer is damped out, and the time step can be chosen larger than $\varepsilon$. All this is usually performed by a suitable time step control [19]. Here we are interested in the behaviour of the scheme for time much larger

4

that $\varepsilon$, after the effect of the initial layer has damped out, and this is why we assume that our initial condition is well-prepared. In practice, any initial condition chosen as the solution of the system at a given time large enough compared to $\varepsilon$ will be well prepared.

A sequence of differential-algebraic systems arises in the study of SPPs, i.e. (2.5), (2.6), (2.7) and (2.8), under the assumption of a smooth solutions (2.4) of system (1.2) (a general and detailed investigation about the $\varepsilon$-expansion is given in [19]).

As we will show next, by inserting the ansatz (2.4) into system (1.2) and collecting equal power of $\varepsilon$, one obtain a set of conditions where the coefficients in the expansion (2.4) are the solutions of differential algebraic systems.

We remind that in order to prove the convergence of R-K methods for problem of the form (1.2), in [19] and in the original paper [17], the authors showed that most of the implicit R-K methods presented in the literature suffer from the phenomenon of the order reduction when applied to (1.2) in the stiff regime ($\varepsilon \to 0$). In [3], similar results of convergence are obtained for IMEX R-K methods.

Now we consider an $s$-stage IMEX R-K method with a double *tableau* in the usual Butcher notation,

$$\begin{array}{c|c} \tilde{c} & \tilde{A} \\ \hline & \tilde{b}^T \end{array} \qquad \begin{array}{c|c} c & A \\ \hline & b^T \end{array}. \tag{2.10}$$

where $\tilde{A} = (\tilde{a}_{ij})$ is an $s \times s$ matrix for an explicit scheme, with $\tilde{a}_{ij} = 0$ for $j \geq i$ and $A = (a_{ij})$ is an $s \times s$ matrix for an implicit scheme. For the implicit part of the methods, we use a diagonally implicit scheme for the function $g$, i.e. $a_{ij} = 0$, for $j > i$, in order to guarantee simplicity and efficiency in solving the algebraic equations corresponding to the implicit part of the discretization. The vectors $\tilde{c} = (\tilde{c}_1, ..., \tilde{c}_s)^T$, $\tilde{b} = (\tilde{b}_1, ..., \tilde{b}_s)^T$, and $c = (c_1, ..., c_s)^T$, $b = (b_1, ..., b_s)^T$ complete the characterization of the scheme. The coefficients $\tilde{c}$ and $c$ are given by the usual relation

$$\tilde{c}_i = \sum_{j=1}^{i-1} \tilde{a}_{ij}, \quad c_i = \sum_{j=1}^{i} a_{ij}. \tag{2.11}$$

When such IMEX scheme is applied to the SPP (1.2), we have

$$\begin{pmatrix} y_{n+1} \\ \varepsilon z_{n+1} \end{pmatrix} = \begin{pmatrix} y_n \\ \varepsilon z_n \end{pmatrix} + h \sum_{i=1}^{s} \begin{pmatrix} \tilde{b}_i k_i \\ b_i \ell_i \end{pmatrix} \tag{2.12}$$

where

$$\begin{pmatrix} k_i \\ \ell_i \end{pmatrix} = \begin{pmatrix} f(Y_i, Z_i) \\ g(Y_i, Z_i) \end{pmatrix} \tag{2.13}$$

and the internal stages are given by

$$\begin{pmatrix} Y_i \\ \varepsilon Z_i \end{pmatrix} = \begin{pmatrix} y_n \\ \varepsilon z_n \end{pmatrix} + h \begin{pmatrix} \sum_{j=1}^{i-1} \tilde{a}_{ij} k_j \\ \sum_{j=1}^{i} a_{ij} \ell_j \end{pmatrix}. \tag{2.14}$$

The following definition will be also useful to characterize properties of an IMEX R-K method in the sequel.

**Definition 2.2.** We say that an IMEX R-K scheme is *globally stiffly accurate* (GSA) if $b^T = e_s^T A$, and $\tilde{b}^T = e_s^T \tilde{A}$, with $e_s = (0, ..., 0, 1)^T$, and $c_s = \tilde{c}_s = 1$, i.e. the numerical solution is identical to the last internal stage value of the scheme.

This definition means that the IMEX R-K scheme is a *stiffly accurate*, (i.e. $a_{si} = b_i$, for $i = 1, ..., s$) in the implicit part [19] and a FSAL (*First Same As Last*) R-K method [18] in the explicit part, (i.e. $\tilde{a}_{si} = \tilde{b}_i$, for $i = 1, ..., s - 1$). It is known in the literature that FSAL R-K methods are a special class of $s$-stage explicit R-K schemes. Such schemes have the advantage that they require only $s - 1$ function evaluations per time step, because the last stage of step $n$ coincides with the first step of the step $n + 1$ (see [18] for details). Note that IMEX R-K schemes with this property were already introduced in [8, 5]. Next we shall show the importance of this property for the convergence of the IMEX R-K method.

In our analysis, we use $R(\infty) = \lim_{z \to -\infty} R(z)$, with $R(z)$ *the stability function* of the implicit scheme, given by $R(z) = 1 + z b^T (I - zA)^{-1} \mathbf{1}$, with $b^T = (b_1, ..., b_s)$ and $\mathbf{1} = (1, ..., 1)^T$. From the expression of $R(z)$

we have $R(\infty) = 1 - b^T A^{-1} \mathbf{1}$. An implicit R-K method is said to be $L$-stable if it is $A$-stable and $R(\infty) = 0$. The $L$-stability property is important when one treats with stiff problem.(Chap. IV.3 in [19]). Furthermore, in the special case when the implicit method is stiffly accurate and the matrix $A$ is invertible, one has always $R(\infty) = 0$, and this makes an $A$-stable method $L$-stable.

Now we assume that the implicit R-K matrix $(a_{ij})$ is invertible, so that we get

$$h\ell_{ni} = \varepsilon \sum_{j=1}^{s} \omega_{ij}(Z_{nj} - z_n),$$

where $\omega_{ij}$ are the elements of the inverse of $(a_{ij})$. Note that this assumption is important to provide the error estimates in [3]. Then inserting this into $z_{n+1}$, and setting $\varepsilon = 0$ we obtain

$$Y_i = y_n + h\sum_{j=1}^{i-1} \tilde{a}_{ij} f(Y_j, Z_j) \tag{2.15}$$

$$0 = g(Y_i, Z_i) \tag{2.16}$$

$$y_{n+1} = y_n + h\sum_{i=1}^{s} \tilde{b}_i f(Y_i, Z_i) \tag{2.17}$$

$$z_{n+1} = R(\infty)z_n + \sum_{i,j=1}^{s} b_i \omega_{ij} Z_j. \tag{2.18}$$

In our analysis it is useful to characterize the different type of IMEX R-K methods existing in the literature we will consider in the sequel accordingly to the structure of the implicit part of the method. Following [3] we have

**Definition 2.3.** We call an IMEX R-K method of type A, if the matrix $A \in R^{s \times s}$ is invertible and $\tilde{c} \neq c$.

**Definition 2.4.** We call an IMEX R-K method of type CK, if the matrix $A \in R^{s \times s}$ can be written as

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

with the submatrix $\hat{A} \in R^{(s-1)\times(s-1)}$ invertible and and $\tilde{c} = c$. In particular when $a = 0$, the IMEX R-K method is of type ARS.

In [3], the author presented the error analysis of different types of IMEX R-K schemes when applied to SPP (1.2), some of which is summarized in Section 4 below. Now briefly we review a couple of main results from [3], which we will use to prove the main theorem in this paper.

In order to obtain the error estimate for IMEX R-K methods, we start from the $\varepsilon$-expansion of the exact and numerical solutions of the problem (1.2), (see [3] and [19] for details). Then, by assuming that the initial values are well-prepared, Definition 2.1, an $\varepsilon$-asymptotic expansion of smooth solutions of the system (1.2) is given by the ansatz (2.4):

$$y(t) = \sum_{\nu \geq 0} y_\nu(t)\varepsilon^\nu, \quad z(t) = \sum_{\nu \geq 0} z_\nu(t)\varepsilon^\nu, \tag{2.19}$$

and, similarly, for the numerical solutions by:

$$y_n = \sum_{\nu \geq 0} y_{n,\nu}\varepsilon^\nu, \quad z_n = \sum_{\nu \geq 0} z_{n,\nu}\varepsilon^\nu, \tag{2.20}$$

approximating exact solutions at $t_n$. Then, the errors of the $y$ and $z$-component are formally considered as

$$y(t_n) - y_n = \sum_{\nu \geq 0} \varepsilon^\nu(y_\nu(t_n) - y_{n,\nu}), \quad z(t_n) - z_n = \sum_{\nu \geq 0} \varepsilon^\nu(z_\nu(t_n) - z_{n,\nu}). \tag{2.21}$$

The values $y_\nu(t)$, $z_\nu(t)$ in (2.19) are the coefficients of the $\varepsilon$-expansion of the smooth solution for (1.2) and values $y_{n,0}, z_{n,0}, y_{n,1}, z_{n,1}, \cdots$, in (2.20) represent the numerical solutions of the IMEX R-K method applied to differential algebraic equations (DAEs) of arbitrary order $\nu$. Furthermore, the first differences $y_0(t_n) - y_{n,0}$ and

$z_0(t_n) - z_{n,0}$ in the expansion (2.21) are the global errors of IMEX R-K method applied to the reduced system (2.2), i.e. system of index $\nu = 1$. The other differences for $\nu > 0$ in (2.21) are related to the numerical solutions of the IMEX R-K method when applied to the DAEs of higher index. In order to study the error (2.21), in [3] the author investigated the differences: $y_\nu(t_n) - y_\nu^n$, $z_\nu(t_n) - z_\nu^n$, for $\nu \geq 0$. For details, see [3].

Finally for the next analysis a couple of remarks are in order for GSA IMEX R-K of a given type.

**Remark 2.5.** a) By the Implicit Function Theorem applied to (2.16), we have $Z_i = \mathcal{G}(Y_i)$ for $i = 1, ..., s$ and the internal stages $Z_i$ depend on the internal stages of the explicit part of $Y_i$. Furthermore, in general, the numerical solutions $y_{n+1}$, $z_{n+1}$ do not lie on the manifold $g(y, z) = 0$. However, if the scheme is GSA, then $y_{n+1} = Y_{ns}$, $z_{n+1} = Z_{ns}$ and therefore equation

$$g(y_{n+1}, z_{n+1}) = 0, \tag{2.22}$$

is satisfied anyway.

b) In this case the application of (2.15)-(2.16)-(2.17) and (2.22) to system (2.2) is equivalent to the $s$-stage explicit R-K method (2.15)-(2.17) applied to (2.3) when $\varepsilon = 0$. Then by (2.22), we get $z_{n+1} = \mathcal{G}(y^{n+1})$ and we obtain for the $y$ and $z$ component the following estimates:

$$y_n - y(t_n) = \mathcal{O}(h^p), \quad z_n - z(t_n) = \mathcal{O}(h^p),$$

with $p$ is the classical order of the explicit R-K method. Then the $z$-component possesses the same asymptotic error estimate as the $y$-component in the case $\varepsilon = 0$. Note that if the scheme is not GSA, the equation (2.22) is not satisfied and a loss of accuracy is observed for the variable $z$ (for details, see [3]).

c) The globally stiffly accurate property for an IMEX R-K scheme implies also that $\tilde{b}_i \neq b_i$, for $i = 1, \cdots, s$. Note that in [3] the assumption of $\tilde{b}_i = b_i$ for all $i$ for the IMEX R-K methods represents the only remedy for preserving the order for the differential component $y$ when we consider high index DAEs, i.e. $\nu > 0$, otherwise the order drops to first one (for details, see Theorem 5.2, Theorem 6.1 and Theorem 6.2 in [3]).

# 3   InDC-IMEX R-K formulations applied to SPPs

In this section, we consider the InDC methods constructed using IMEX schemes applied to SPPs (1.2). First of all, we introduce the specific strategy of InDC methods. Later, we show as a InDC IMEX R-K method can be re-written as a standard IMEX R-K one. Clearly this formulation provides a systematic way to construct arbitrary-order IMEX R-K schemes without solving complicated order conditions. The formulation of InDC IMEX R-K methods as a standard IMEX R-K gives us the possibility to estimate the error of the schemes by using some classical results of the error analysis developed in the paper [3] for standard IMEX R-K.

## 3.1   InDC framework

The time interval $[0, T]$ is uniformly discretized into intervals $[t_n, t_{n+1}]$, $n = 0, 1, ..., N - 1$ such that

$$0 = t_0 < t_1 < t_2 < ... < t_n < ... < t_N = T,$$

with the step size $H$. Each interval $[t_n, t_{n+1}]$ is discretized again into $M$ uniform subintervals with quadrature nodes denoted by

$$t_n \doteq \tau_0 < \tau_1 < \cdots < \tau_M \doteq t_{n+1}, \tag{3.1}$$

with $h = \frac{H}{M}$ being the size of a substep. In this paper, the interval $[t_n, t_{n+1}]$ will be referred to as a time step while a subinterval $[\tau_m, \tau_{m+1}]$ will be referred to as a substep. We assume the InDC quadrature nodes are uniform, which is a crucial assumption for high order improvement in accuracy, when consider applying general high order R-K methods in prediction and correction steps for classical ODE system, (see discussions in [14]). Similar conclusions hold when the high order IMEX R-K methods are applied in prediction and correction steps. Since $H = Mh$, we will use $\mathcal{O}(h^p)$ and $\mathcal{O}(H^p)$ interchangeably throughout the paper. The use of quadrature nodes excluding the left-most endpoint (i.e. $\tau_0$) leads to an important stability condition for stiff problems, i.e. the method is L-stable (for details see in [23]). With the above considerations, in this paper, we consider the InDC methods with uniform quadrature nodes excluding the left-most endpoint.

Let's assume we have obtained numerical solutions $\hat{y}_m^{(0)}$ and $\hat{z}_m^{(0)}$ approximating the exact solution at $\tau_m$ by using a low order numerical method for (1.2). Here superscript (0) is used to denote the prediction step in

the InDC method. We build continuous polynomial interpolants $\hat{y}^{(0)}(t)$ and $\hat{z}^{(0)}(t)$ interpolating these discrete values ($\hat{y}_m^{(0)}$ and $\hat{z}_m^{(0)}$ for $m = 1, \cdots M$). Now we define the error functions

$$e^{(0)}(t) = y(t) - \hat{y}^{(0)}(t), \quad d^{(0)}(t) = z(t) - \hat{z}^{(0)}(t). \tag{3.2}$$

Note that $e^{(0)}(t)$ and $d^{(0)}(t)$ are not polynomials in general. We specify the residual function with respect to $y$ and $z$ via the following set of differential equations

$$\begin{aligned}
\delta^{(0)}(t) &= f(\hat{y}^{(0)}(t), \hat{z}^{(0)}(t)) - (\hat{y}^{(0)})'(t), \\
\rho^{(0)}(t) &= g(\hat{y}^{(0)}(t), \hat{z}^{(0)}(t)) - (\varepsilon\hat{z}^{(0)})'(t).
\end{aligned} \tag{3.3}$$

Thus, by subtracting (3.3) from (1.2), and recalling that we restrict our analysis to autonomous systems, the error equations about the error functions (3.2) become

$$\begin{aligned}
(e^{(0)})'(t) - \delta^{(0)}(t) &= \Delta f^{(0)}, \\
\varepsilon(d^{(0)})'(t) - \rho^{(0)}(t) &= \Delta g^{(0)},
\end{aligned} \tag{3.4}$$

where

$$\begin{aligned}
\Delta f^{(0)} &\doteq f(e^{(0)}(t) + \hat{y}^0(t), d^{(0)}(t) + \hat{z}^{(0)}(t)) - f(\hat{y}^{(0)}(t), \hat{z}^{(0)}(t)), \\
\Delta g^{(0)} &\doteq g(e^{(0)}(t) + \hat{y}^0(t), d^{(0)}(t) + \hat{z}^{(0)}(t)) - g(\hat{y}^{(0)}(t), \hat{z}^{(0)}(t)).
\end{aligned} \tag{3.5}$$

The initial conditions for these error equations are zero, i.e. $e^{(0)}(0) = 0$ and $d^{(0)}(0) = 0$.

Suppose that we have obtained approximate solutions $\hat{e}_m^{(0)}$ and $\hat{d}_m^{(0)}$ at $\tau_m$ by using a low order IMEX method for error equations (3.4), the numerical solution can then be improved as

$$\hat{y}_m^{(1)} = \hat{y}_m^{(0)} + \hat{e}_m^{(0)}, \quad \hat{z}_m^{(1)} = \hat{z}_m^{(0)} + \hat{d}_m^{(0)}, \quad \forall m = 0, \cdots M.$$

Such correction procedures can be repeated in each local time step $[t_n, t_{n+1}]$. How the scheme is actually implemented will be illustrated in details in the next section. In summary, the strategy of InDC methods is to use a simple numerical method to compute numerical solutions $\hat{y}^{(0)}(t)$ and $\hat{z}^{(0)}(t)$ as prediction, and then to solve a series of correction equations in the integral form based on equations (3.4). Each correction improves the accuracy of numerical solutions from the previous iteration.

In our description of InDC, we let $\hat{y}_m^{(k)}$, $\hat{z}_m^{(k)}$, $\hat{e}_m^{(k)}$, $\hat{d}_m^{(k)}$ denote the numerical approximations (with hat) to the exact solutions and error functions. We use subscript $m$ to denote the location $\tau_m$ of the quadrature points and use superscript $(k)$ to denote the prediction ($k = 0$) and correction loops ($k = 1, \cdots$). We let $\underline{\cdot}$ denote the vector on InDC quadrature nodes ($\tau_1, \cdots \tau_M$), for example, $\underline{y} = (y_1, \cdots, y_M)$.

## 3.2 InDC-IMEX1 method

In this subsection, we consider InDC method constructed using implicit-explicit Euler method, applied to (1.2).

We use uniformly distributed quadrature nodes $\tau_1, ..., \tau_M$ from equation (3.1) excluding the left-most endpoint.

1. (Prediction step) Use the first implicit-explicit Euler method to compute

$$\underline{\hat{y}}^{(0)} = (\hat{y}_1^{(0)}, ..., \hat{y}_m^{(0)}, ..., \hat{y}_M^{(0)}), \quad \underline{\hat{z}}^{(0)} = (\hat{z}_1^{(0)}, ..., \hat{z}_m^{(0)}, ..., \hat{z}_M^{(0)})$$

   as the approximations of the exact solution $\underline{y} = (y_1, \cdots, y_m, \cdots, y_M)$ and $\underline{z} = (z_1, \cdots, z_m, \cdots, z_M)$ at quadrature nodes $\tau_1, ..., \tau_M$. This gives

$$\begin{aligned}
\hat{y}_{m+1}^{(0)} &= \hat{y}_m^{(0)} + hf(\hat{y}_m^{(0)}, \hat{z}_m^{(0)}), \\
\varepsilon\hat{z}_{m+1}^{(0)} &= \varepsilon\hat{z}_m^{(0)} + hg(\hat{y}_{m+1}^{(0)}, \hat{z}_{m+1}^{(0)}),
\end{aligned} \tag{3.6}$$

   for $m = 0, 1, ...M - 1$.

2. (Correction loop). For $k = 1, ..., K$ ($K$ is the number of the correction step), let $\hat{y}^{(k-1)}$ and $\hat{z}^{(k-1)}$ denote the numerical solutions at the $(k-1)^{th}$ correction.

8

(a) Denote the error function at the $(k-1)^{th}$ correction $e^{(k-1)}(t) = y(t) - \hat{y}^{(k-1)}(t)$, where $y(t)$ is the exact solution and $\hat{y}^{(k-1)}(t)$ is a $(M-1)^{th}$ order polynomial interpolating $\bar{\hat{y}}^{(k-1)}$ at quadrature nodes $\tau_1, ..., \tau_M$. Similarly denote $d^{(k-1)}(t) = z(t) - \hat{z}^{(k-1)}(t)$. The initial conditions for these error functions are zero, i.e. $\hat{e}_0^{(k-1)} = 0$ and $\hat{d}_0^{(k-1)} = 0, \forall k$. Let $\delta^{(k-1)}(t)$ and $\rho^{(k-1)}(t)$, $\Delta f^{(k-1)}$ and $\Delta g^{(k-1)}$ be defined by equations (3.3) and (3.5) respectively, but with the upper script $(0)$ replaced with $(k-1)$. We compute the numerical error vector $\bar{\hat{e}}^{(k-1)} = (\hat{e}_1^{(k-1)}, ..., \hat{e}_M^{(k-1)})$ with $\hat{e}_m^{(k-1)}$ approximating $e^{(k-1)}(\tau_m)$ by applying the first order IMEX R-K method with Butcher table specified in (3.14) to the integral form of (3.4),

$$
\begin{aligned}
\hat{e}_{m+1}^{(k-1)} &= \hat{e}_m^{(k-1)} + h\Delta f_m^{(k-1)} + \int_{\tau_m}^{\tau_{m+1}} \delta^{(k-1)}(s)ds, \\
\varepsilon \hat{d}_{m+1}^{(k-1)} &= \varepsilon \hat{d}_m^{(k-1)} + h\Delta g_{m+1}^{(k-1)} + \int_{\tau_m}^{\tau_{m+1}} \rho^{(k-1)}(s)ds,
\end{aligned}
\tag{3.7}
$$

where

$$
\begin{aligned}
\int_{\tau_m}^{\tau_{m+1}} \delta^{(k-1)}(s)ds &= \int_{\tau_m}^{\tau_{m+1}} f(\hat{y}^{(k-1)}(s), \hat{z}^{(k-1)}(s))ds - \hat{y}_{m+1}^{(k-1)} + \hat{y}_m^{(k-1)}, \\
\int_{\tau_m}^{\tau_{m+1}} \rho^{(k-1)}(s)ds &= \int_{\tau_m}^{\tau_{m+1}} g(\hat{y}^{(k-1)}(s), \hat{z}^{(k-1)}(s))ds - \varepsilon\hat{z}_{m+1}^{(k-1)} + \varepsilon\hat{z}_m^{(k-1)}.
\end{aligned}
\tag{3.8}
$$

Integral terms $\int_{\tau_m}^{\tau_{m+1}}$ on the right hand side of equations (3.8) are approximated by a numerical quadrature, in particular by an interpolary quadrature formula. Especially, let $S$ be the integration matrix, its $(m, l)$ element is

$$
S^{m,l} = \frac{1}{h} \int_{\tau_m}^{\tau_{m+1}} \alpha_l(s)ds, \quad \text{for} \quad m = 0, \cdots, M-1, \quad l = 1, \cdots M,
$$

where $\alpha_l(s)$ is the Lagrangian basis function based on the node $\tau_l$, $l = 1, \cdots M$. Analytical expression of $S^{m,l}$ can be given, which does not depends on $h$,

$$
S^{m,l} = \int_m^{m+1} \prod_{j \neq l} \frac{\theta - j}{l - j} d\theta, \quad \text{for} \quad m = 0, \cdots, M-1, \quad l = 1, \cdots M,
$$

Let

$$
S^m(\underline{f}) = \sum_{j=1}^M S^{m,j} f(y_j, z_j),
\tag{3.9}
$$

then

$$
hS^m(\underline{f}) - \int_{\tau_m}^{\tau_{m+1}} f(y(s), z(s))ds = \mathcal{O}(h^{M+1}),
$$

for any smooth function $f$. In other words, the quadrature formula given by $hS^m(\underline{f})$ approximates the exact integration with $(M+1)^{th}$ order of accuracy *locally*.

(b) Update the approximate solutions

$$
\underline{\hat{y}}^{(k)} = \underline{\hat{y}}^{(k-1)} + \underline{\hat{e}}^{(k-1)}, \quad \underline{\hat{z}}^{(k)} = \underline{\hat{z}}^{(k-1)} + \underline{\hat{d}}^{(k-1)}.
\tag{3.10}
$$

Note that from equations (3.7), (3.8), (3.10) and using the notation introduced in equation (3.9), we get

$$
\begin{aligned}
\hat{y}_{m+1}^{(k)} &= \hat{y}_m^{(k)} + h\Delta f_m^{(k-1)} + hS^m(\underline{\hat{f}}^{(k-1)}), \\
\varepsilon\hat{z}_{m+1}^{(k)} &= \varepsilon\hat{z}_m^{(k)} + h\Delta g_{m+1}^{(k-1)} + hS^m(\underline{\hat{g}}^{(k-1)}).
\end{aligned}
\tag{3.11}
$$

This completes the description of the InDC-IMEX1 method.

Note that in practice $k$ is chosen in such a way that there is no further improvement in order of convergence when increasing it. The precise relation between $k$ and $M$ and the order of accuracy of the method will be reported later.

**Remark 3.1.** The assumption of GSA guarantees, in the case $\varepsilon = 0$ (as in (3.11)), that the approximate solutions in the prediction step $k = 0$ satisfy the equation $g(\hat{y}_m^{(0)}, \hat{z}_m^{(0)}) = 0$, and this implies $\hat{z}_m^{(0)} = \mathcal{G}(\hat{y}_m^{(0)})$, $\forall m$ (see Remark 2.5).

Consequently in the first correction step for $k = 1$, by (3.5), the second equation in system (3.11) is reduced to

$$\Delta g_m^{(0)} = 0 \quad \forall m = 0, \cdots, M. \tag{3.12}$$

and this guarantees that $\hat{\underline{g}}^{(1)} = 0$, i.e. $g(\hat{y}_m^{(1)}, \hat{z}_m^{(1)}) = 0, \forall m$. By Remark 2.5, we get $\hat{z}_m^{(1)} = \mathcal{G}(\hat{y}_m^{(1)}), \forall m = 0, \cdots, M$. A similar argument can be given to show that $\hat{\underline{g}}^{(k)} = 0$, for $k = 2, \cdots, K$. Hence, again by Remark 2.5, we have $g(\hat{y}_{m,0}^{(k)}, \hat{z}_{m,0}^{(k)}) = 0$, i.e.,

$$\hat{z}_{m,0}^{(k)} = \mathcal{G}(\hat{y}_{m,0}^{(k)}), \ \forall k = 0, \cdots K, \ \forall m = 0, \cdots, M.$$

Therefore, the updating of $\hat{y}_{m+1}^{(k)}$ represents the correction step for a non-stiff ordinary differential equation of the form (2.3), i.e.,

$$\hat{y}_{m+1}^{(k)} = \hat{y}_m^{(k)} + h(f(\hat{y}_m^{(k)}) - \hat{f}(y_m^{(k-1)})) + hS^m(\bar{\hat{f}}^{(k-1)})$$

where $f(\cdot) = f(\cdot, \mathcal{G}(\cdot))$. Then by Remark 2.5 b) and by classical results in the InDC framework for non-stiff ordinary differentia equations (Theorem 4.1 in [14] for details), we obtain the classical increasing order of accuracy for the $y$ and $z$-component, i.e.,

$$e_n^{(K)} \doteq \hat{y}_n^{(K)} - y(t_n) = \mathcal{O}(H^{\min(k+2,M)}), \quad d_n^{(k)} \doteq \hat{z}_n^{(K)} - z(t_n) = \mathcal{O}(H^{\min(k+2,M)}). \tag{3.13}$$

with $k \leq M$. In a straightforward manner, in the next section we will generalize this result.

In the next section we show that InDC IMEX R-K methods can be rewritten as standard IMEX R-K schemes.

## 3.3 Reformulation of InDC-IMEX methods

In [13], the InDC method constructed with explicit or implicit R-K methods in the prediction and correction steps has been reformulated as a high-order explicit or implicit R-K method whose Butcher tableau is explicitly constructed.

Similarly in this section we show as an InDC method constructed by a standard IMEX R-K method can be re-written as an IMEX R-K one. We first present the reformulation of InDC-IMEX1 methods constructed from first order classical IMEX R-K schemes as an IMEX method with the corresponding Butcher tableau for the explicit and implicit parts.

We start to give a list of a classical first order IMEX R-K schemes given by a double Butcher tableau (2.10) corresponding to the implicit and explicit part of the method for different types, (2.3) and (2.4). Furthermore we emphisize which of these schemes are GSA.

- First order IMEX R-K method of type ARS (particular case of CK), it is GSA, [1]:

$$\begin{array}{c|cc} 0 & 0 & 0 \\ 1 & 1 & 0 \\ \hline & 1 & 0 \end{array} \qquad \begin{array}{c|cc} 0 & 0 & 0 \\ 1 & 0 & 1 \\ \hline & 0 & 1 \end{array}. \tag{3.14}$$

- First order IMEX R-K method of type A, it is GSA:

$$\begin{array}{c|cc} 0 & 0 & 0 \\ 1 & 1 & 0 \\ \hline & 1 & 0 \end{array} \qquad \begin{array}{c|cc} 1 & 1 & 0 \\ 1 & 0 & 1 \\ \hline & 0 & 1 \end{array}. \tag{3.15}$$

- First order IMEX R-K method of type A, it is not GSA, [28]:

$$\begin{array}{c|c} 0 & 0 \\ \hline & 1 \end{array} \qquad \begin{array}{c|c} 1 & 1 \\ \hline & 1 \end{array}. \tag{3.16}$$

Note that scheme (3.14) is the implicit-explicit Euler method, i.e. Eq. (3.6), applied to problem (1.2). Scheme (3.15) is not used in the literature, because it not efficient. It is reported here only as an example of first order type A scheme which is GSA. As we shall see, InDC IMEX methods based on GSA schemes maintain the

accuracy as $\varepsilon \to 0$, while the InDC IMEX R-K method based on a non GSA scheme (3.16) shows a degradation of accuracy.

We first present the reformulation of InDC-IMEX1 method constructed by scheme (3.14), and prove that the corresponding InDC is an IMEX R-K method of type ARS and it is GSA. We call it InDC-IMEX1-GSA-ARS.

First we present the Butcher tableau for the InDC-IMEX1-GSA-ARS method with one loop of correction step as an example. This takes the form

$$
\begin{array}{c|ccc}
0 & 0 & \mathbf{0}^T & \mathbf{0}^T \\
\mathbf{c} & \frac{1}{M} & \tilde{T} & O \\
\mathbf{c} & \mathbf{0} & \tilde{P} & \tilde{T} \\
\hline
 & 0 & \tilde{\mathbf{b}}_1^T & \tilde{\mathbf{b}}_2^T
\end{array}
\qquad
\begin{array}{c|ccc}
0 & 0 & \mathbf{0}^T & \mathbf{0}^T \\
\mathbf{c} & \mathbf{0} & T & O \\
\mathbf{c} & \mathbf{0} & P & T \\
\hline
 & 0 & \mathbf{b}_1^T & \mathbf{b}_2^T
\end{array}
\tag{3.17}
$$

where $\mathbf{c} = \frac{1}{M}[1, \cdots, M]^T$, $\mathbf{0}$ and $\mathbf{1}$ are vectors with 0 and 1 entries respectively, having the same size as $\mathbf{c}$. $O$ is a $M \times M$ matrix of zeros, $T$, $\tilde{T}$, $P$ and $\tilde{P}$ are $M \times M$ matrices, with

$$
\tilde{T} = \frac{1}{M}
\begin{bmatrix}
0 & 0 & 0 & \dots & 0 \\
1 & 0 & 0 & \dots & 0 \\
\vdots & \vdots & \ddots & \vdots & \vdots \\
1 & 1 & 1 & \dots & 0
\end{bmatrix},
\qquad
T = \frac{1}{M}
\begin{bmatrix}
1 & 0 & 0 & \dots & 0 \\
1 & 1 & 0 & \dots & 0 \\
\vdots & \vdots & \ddots & \vdots & \vdots \\
1 & 1 & 1 & \dots & 1
\end{bmatrix},
$$

and $\tilde{P} = S - \tilde{T}$, $P = S - T$, where the matrix $S$ is constructed such that its entries $S_{ij} = \int_{t_0}^{t_i} \alpha_j(s)ds = \sum_{k=0}^{i-1} S^{k,j}$ with $\alpha_j(s)$ the Lagrangian basis functions for the node $\tau_j$. The vectors $\tilde{\mathbf{b}}_1^T$, $\tilde{\mathbf{b}}_2^T$, $\mathbf{b}_1^T$ and $\mathbf{b}_2^T$ are taken so that they are the last rows of the matrices $\tilde{P}$, $\tilde{T}$, $P$ and $T$. Furthermore, the method (3.17) has the same nodes $\mathbf{c}$ in implicit and explicit parts and the weights $\tilde{\mathbf{b}}_i^T \neq \mathbf{b}_i^T$, for $i = 1, 2$. As an example, we show the Butcher table for InDC-IMEX1-GSA-ARS method with $M = 2$ and with one correction loop below. Note that $S = [3/4, -1/4; 1, 0]$ for $M = 2$.

$$
\begin{array}{c|ccccc}
0 & 0 & 0 & 0 & 0 & 0 \\
1/2 & 1/2 & 0 & 0 & 0 & 0 \\
1 & 1/2 & 1/2 & 0 & 0 & 0 \\
1/2 & 0 & 3/4 & -1/4 & 0 & 0 \\
1 & 0 & 1/2 & 0 & 1/2 & 0 \\
\hline
 & 0 & 1/2 & 0 & 1/2 & 0
\end{array}
\quad
\begin{array}{c|ccccc}
0 & 0 & 0 & 0 & 0 & 0 \\
1/2 & 0 & 1/2 & 0 & 0 & 0 \\
1 & 0 & 1/2 & 1/2 & 0 & 0 \\
1/2 & 0 & 1/4 & -1/4 & 1/2 & 0 \\
1 & 0 & 1/2 & -1/2 & 1/2 & 1/2 \\
\hline
 & 0 & 1/2 & -1/2 & 1/2 & 1/2
\end{array}.
\tag{3.18}
$$

The same conclusion holds for the InDC method constructed with the first order IMEX R-K schemes (3.15). We obtain a InDC IMEX R-K scheme denoted as InDC-IMEX1-GSA-A. It can be shown, by similar argument presented for the previous scheme, that the assembled matrix $\mathbb{A}$ of the implicit part in the table Butcher tableau, for the InDC-IMEX1-GSA-A method, is invertible, and, at the end, we obtain a IMEX R-K method of type A which is GSA.

As an example, we show below the double Butcher tableaus for the InDC-IMEX1-GSA-A method with $M = 2$ and with one loop of correction only. The tableaus, obtained by an using an algebraic manipulator from Eq. (3.11), show that the InDC IMEX R-K scheme is of type A and is GSA.

$$
\begin{array}{c|cccccccc}
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
1/2 & 1/2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
1/2 & 1/2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
1 & 1/2 & 0 & 1/2 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
1/2 & 0 & 3/4 & 0 & -1/4 & 0 & 0 & 0 & 0 \\
1/2 & 0 & 3/4 & 0 & -1/4 & 0 & 0 & 0 & 0 \\
1 & 0 & 1/2 & 0 & 0 & 0 & 1/2 & 0 & 0 \\
\hline
1 & 0 & 1/2 & 0 & 0 & 0 & 1/2 & 0 & 0
\end{array}
\quad
\begin{array}{c|cccccccc}
1/2 & 1/2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
1/2 & 0 & 1/2 & 0 & 0 & 0 & 0 & 0 & 0 \\
1 & 0 & 1/2 & 1/2 & 0 & 0 & 0 & 0 & 0 \\
1 & 0 & 1/2 & 0 & 1/2 & 0 & 0 & 0 & 0 \\
1/2 & 0 & 1/4 & 0 & -1/4 & 1/2 & 0 & 0 & 0 \\
1/2 & 0 & 1/4 & 0 & -1/4 & 0 & 1/2 & 0 & 0 \\
1 & 0 & 1/2 & 0 & -1/2 & 0 & 1/2 & 1/2 & 0 \\
1 & 0 & 1/2 & 0 & -1/2 & 0 & 1/2 & 0 & 1/2 \\
\hline
 & 0 & 1/2 & 0 & -1/2 & 0 & 1/2 & 0 & 1/2
\end{array}.
\tag{3.19}
$$

Finally, we consider an InDC IMEX R-K method constructed by using the first order IMEX scheme (3.16) which is not GSA. The InDC method is denoted by InDC-IMEX1-NGSA-A. For such an scheme, we show that the assembled matrix $\mathbb{A}$ of the implicit part in the Butcher table is non-invertible.

Similaly as an example, we consider InDC-IMEX1-NGSA-A constructed by using $M = 2$ and with one loop of correction only.

Then, the corresponding Butcher tables can be assembled as

| $c$ | | | | | | | | | $c$ | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1/2 | 1/2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1/2 | 1/2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1/2 | 1/2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1/2 | 1/2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1/2 | 0 | 1/2 | 0 | 0 | 0 | 0 | 0 |
| 1 | 1/2 | 0 | 1/2 | 0 | 0 | 0 | 0 | 0 | 1 | 1/2 | 0 | 1/2 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1/2 | $-1/2$ | 3/4 | 0 | $-1/4$ | 1/2 | 0 | 0 | 0 |
| 1/2 | $-1/2$ | 3/4 | 0 | $-1/4$ | 1/2 | 0 | 0 | 0 | 1/2 | $-1/2$ | 3/4 | 0 | $-1/4$ | 1/2 | 0 | 0 | 0 |
| 1/2 | $-1/2$ | 3/4 | 0 | $-1/4$ | 1/2 | 0 | 0 | 0 | 1 | $-1/2$ | 1 | $-1/2$ | 0 | 1/2 | 0 | 1/2 | 0 |
| 1 | $-1/2$ | 1 | $-1/2$ | 0 | 1/2 | 0 | 1/2 | 0 | 1 | $-1/2$ | 1 | $-1/2$ | 0 | 1/2 | 0 | 1/2 | 0 |
| | $-1/2$ | 1 | $-1/2$ | 0 | 1/2 | 0 | 1/2 | 0 | | $-1/2$ | 1 | $-1/2$ | 0 | 1/2 | 0 | 1/2 | 0 |

$$(3.20)$$

Notice that the size of the Butcher matrices (3.19) and (3.20) is $2*M*(K+1)$, in comparison with $M*(K+1)$ for the $\hat{\mathbb{A}}$ matrix for InDC-IMEX1-GSA-ARS methods, where the factor 2 is due to the fact that we need an extra row to compute solutions at InDC quadrature points. When the method is not GSA, the last stage of IMEX is not the updated solution at quadrature nodes, hence we need one more row to represent solutions at quadrature nodes, see Row 2, 4, 6 and 8 in the implicit matrix presented in (3.20). Besides these extra rows, the construction of the assembled Butcher tableau is in a similar way as those in eq. (3.18) and (3.19). Notice that in matrix (3.20) for row 2, 4, 6, 8, the diagonal entries are 0 for the implicit part of the Butcher Tableau, i.e. the assembled matrix of the implicit part is not invertible.

Then for the InDC-IMEX1-GSA schemes constructed with first order IMEX schemes GSA (3.14)-(3.15), we have the following proposition regarding the invertibility of the assembled matrix of the implicit part.

**Proposition 3.2.** (Invertibility of implicit assembled matrix) Consider the InDC method constructed with a first order IMEX R-K methods of type A or CK.If the quadrature nodes used in the InDC method exclude the left-most point and the first order IMEX R-K method is GSA, then the InDC-IMEX1-GSA method is an IMEX R-K of type A or CK and it is GSA with the assembled matrix $\mathbb{A}$ or $\hat{\mathbb{A}}$ in the implicit part being invertible.

*Proof.* First the statement is true for the scheme with only one correction loop, as seen from the Butcher tables presented in eq. (3.18). Such results can be generalized to the InDC method with general $K$ correction loops, with a lower triangular matrix $\hat{\mathbb{A}}$ of size $M(K+1)$ and

$$
\hat{\mathbb{A}} = \begin{pmatrix} T & & & & \\ P & T & & & \\ O & P & T & & \\ & & & \dots & \\ O & \cdots O & P & T & \end{pmatrix}, \tag{3.21}
$$

whose diagonal entries are diagonal entries of $T$ and are $\frac{1}{M}$, i.e. nonzero. Hence we have the invertibility of $\hat{\mathbb{A}}$. Notice that the InDC method is also GSA, i.e. the last rows of the assembled matrices $\tilde{\mathbb{A}}$ and $\hat{\mathbb{A}}$ are the same as the $\tilde{\mathbf{b}}^T$ and $\mathbf{b}^T$ vectors.

In a similar fashion, we prove for IMEX schemes of type A, where the assembled matrix $\mathbb{A}$ is similar to (3.21), with size $2M(K+1)$. ∎

In the next section we extend this previous result for InDC IMEX R-K schemes constructed by considering arbitrary order $s$-stage IMEX R-K.

## 3.4 InDC Methods constructed with IMEX R-K methods

We describe in general how we apply $s$-stage IMEX R-K methods in correction loops in an InDC framework. We assume double Butcher tables for the IMEX R-K method as specified in Section 2. For the internal stages, we introduce the integration and interpolation matrices, which are based on the use of the same set of nodes for each interval, more precisely, in interval $[\tau_m, \tau_{m+1}]$ we assume that the nodes are located at $\tau_m + c_i h$, where $c_i$, $i = 1, \cdots s$, do not depend on the interval $m$. Then we have:

$$
hS^{c_i,k} = \int_{\tau_m}^{\tau_m + c_i h} \alpha_k(s)ds, \quad P^{c_i,k} = \alpha_k(\tau_m + c_i h), \tag{3.22}
$$

Similarly:

$$hS^{\tilde{c}_i,k} = \int_{\tau_m}^{\tau_m+\tilde{c}_i h} \alpha_k(s)ds, \quad P^{\tilde{c}_i,k} = \alpha_k(\tau_m + \tilde{c}_i h), \tag{3.23}$$

$\forall m = 0, \cdots, M-1, \quad \forall k = 1, \cdots M, \quad \forall i = 1, \cdots s,$ where $\alpha_j(s)$ is the Lagrangian basis function based on the node $\tau_j$. Let

$$S^{c_i}(\underline{f}) = \sum_{j=1}^{M} S^{c_i,j} f(y_j, z_j), \quad P^{c_i}(\underline{f}) = \sum_{j=1}^{M} P^{c_i,j} f(y_j, z_j),$$

$$S^{\tilde{c}_i}(\underline{f}) = \sum_{j=1}^{M} S^{\tilde{c}_i,j} f(y_j, z_j), \quad P^{\tilde{c}_i}(\underline{f}) = \sum_{j=1}^{M} P^{\tilde{c}_i,j} f(y_j, z_j).$$

Then

$$hS^{c_i}(\underline{f}) - \int_{\tau_m}^{\tau_m+c_i h} f(y(s), z(s))ds = \mathcal{O}(h^{M+1}), \quad P^{c_i}(\underline{f}) - f(y(\tau_m + c_i h), z(t_m + c_i h)) = \mathcal{O}(h^M),$$

and similarly,

$$hS^{\tilde{c}_{mi}}(\underline{f}) - \int_{\tau_m}^{\tau_m+\tilde{c}_i h} f(y(s), z(s))ds = \mathcal{O}(h^{M+1}), \quad P^{\tilde{c}_{mi}}(\underline{f}) - f(y(\tau_m + \tilde{c}_i h), z(\tau_m + \tilde{c}_i h)) = \mathcal{O}(h^M),$$

for any smooth function $f$. In other words, the quadrature formula given by $hS^{c_i}(\underline{f})$ approximates the exact integration with $(M+1)^{th}$ order accuracy locally, while the interpolation formula given by $P^{c_i}(\underline{f})$ approximates the exact solution at R-K internal stages with $M^{th}$ order accuracy.

We compute the numerical errors approximating the error functions $e^{(k-1)}(\tau_m)$, $d^{(k-1)}(\tau_m)$ with an IMEX R-K method to (3.4),

$$\begin{pmatrix} \hat{e}_{m+1}^{(k-1)} \\ \varepsilon \hat{d}_{m+1}^{(k-1)} \end{pmatrix} = \begin{pmatrix} \hat{e}_m^{(k-1)} + h \int_0^1 \delta(\tau_m + \tau h)d\tau \\ \varepsilon \hat{d}_m^{(k-1)} + h \int_0^1 \rho(\tau_m + \tau h)d\tau \end{pmatrix} + h \begin{pmatrix} \sum_{i=1}^{s} \tilde{b}_i \, \Delta\hat{\mathcal{K}}_{mi}^{(k-1)} \\ \sum_{i=1}^{s} b_i \, \Delta\hat{\mathcal{L}}_{mi}^{(k-1)} \end{pmatrix}, \tag{3.24}$$

with

$$\begin{pmatrix} \Delta\hat{\mathcal{K}}_{mi}^{(k-1)} \\ \Delta\hat{\mathcal{L}}_{mi}^{(k-1)} \end{pmatrix} \doteq \begin{pmatrix} f(\hat{Y}_{mi}^{(k)}, \hat{Z}_{mi}^{(k)}) - P^{\tilde{c}_i}(\underline{\hat{f}}^{(k-1)}) \\ g(\hat{Y}_{mi}^{(k)}, \hat{Z}_{mi}^{(k)}) - P^{c_i}(\underline{\hat{g}}^{(k-1)}) \end{pmatrix}. \tag{3.25}$$

Here we set

$$\hat{Y}_{mi}^{(k)} = P^{\tilde{c}_i}(\underline{\hat{y}}^{(k-1)}) + \hat{E}_{mi}^{(k-1)}, \quad \hat{Z}_{mi}^{(k)} = P^{c_i}(\underline{\hat{z}}^{(k-1)}) + \hat{D}_{mi}^{(k-1)}, \tag{3.26}$$

with

$$\begin{pmatrix} \hat{E}_{mi}^{(k-1)} \\ \varepsilon \hat{D}_{mi}^{(k-1)} \end{pmatrix} = \begin{pmatrix} \hat{e}_m^{(k-1)} + h \int_0^{\tilde{c}_i} \delta(\tau_m + \tau h)d\tau \\ \varepsilon \hat{d}_m^{(k-1)} + h \int_0^{c_i} \rho(\tau_m + \tau h)d\tau \end{pmatrix} + h \begin{pmatrix} \sum_{j=1}^{i} \tilde{a}_{ij} \, \Delta\hat{\mathcal{K}}_{mj}^{(k-1)} \\ \sum_{j=1}^{i} a_{ij} \, \Delta\hat{\mathcal{L}}_{mj}^{(k-1)} \end{pmatrix}. \tag{3.27}$$

The integral terms in the system (3.24) can be approximated by quadrature rules, then we can rewrite (3.24) and (3.27) as

$$\begin{pmatrix} \hat{y}_{m+1}^{(k)} - hS_{\underline{\hat{f}}}^{m,(k-1)} \\ \varepsilon \hat{z}_{m+1}^{(k)} - hS_{\underline{\hat{g}}}^{m,(k-1)} \end{pmatrix} = \begin{pmatrix} \hat{y}_m^{(k)} \\ \varepsilon \hat{z}_m^{(k)} \end{pmatrix} + h \begin{pmatrix} \sum_{i=1}^{s} \tilde{b}_i \, \Delta\hat{\mathcal{K}}_{mi}^{(k-1)} \\ \sum_{i=1}^{s} b_i \, \Delta\hat{\mathcal{L}}_{mi}^{(k-1)} \end{pmatrix}, \tag{3.28}$$

13

$$
\begin{pmatrix}
\hat{Y}_{mi}^{(k)} - hS_{\hat{\underline{f}}}^{\tilde{c}_i,(k-1)} \\
\varepsilon\hat{Z}_{mi}^{(k)} - hS_{\hat{\underline{g}}}^{c_i,(k-1)}
\end{pmatrix}
=
\begin{pmatrix}
\hat{y}_m^{(k)} \\
\varepsilon\hat{z}_m^{(k)}
\end{pmatrix}
+ h
\begin{pmatrix}
\sum_{j=1}^{i} \tilde{a}_{ij}\Delta\hat{\mathcal{K}}_{mj}^{(k-1)} \\
\sum_{j=1}^{i} a_{ij}\Delta\hat{\mathcal{L}}_{mj}^{(k-1)}
\end{pmatrix},
\tag{3.29}
$$

with

$$
\begin{pmatrix}
S_{\hat{\underline{f}}}^{m,(k-1)} = S^m(\hat{\underline{f}}^{(k-1)}) \\
\varepsilon S_{\hat{\underline{g}}}^{m,(k-1)} = S^m(\hat{\underline{g}}^{(k-1)})
\end{pmatrix},
\quad
\begin{pmatrix}
S_{\hat{\underline{f}}}^{\tilde{c}_i,(k-1)} = S^{\tilde{c}_i}(\hat{\underline{f}}^{(k-1)}) \\
\varepsilon S_{\hat{\underline{g}}}^{c_i,(k-1)} = S^{c_i}(\hat{\underline{g}}^{(k-1)})
\end{pmatrix}.
\tag{3.30}
$$

The double Butcher tableaus for the InDC methods constructed with high order IMEX R-K can be assembled just as done for the InDC-IMEX1 method. They are of size related to the number of quadrature points $M$ as well as the number of stages $s$ of the IMEX R-K method. Below we present Proposition 3.3 as a general result regarding the invertibility of the implicit part of assembled Butcher tableau. We choose not to present the assembled Butcher table for the InDC method constructed with high order IMEX R-K to save space. For example, to construct a Butcher table of InDC IMEX2 method with four quadrature points and one correction loop, the assembled matrix will be at least of size $17 \times 17$. For the construction of classical InDC methods using high order R-K methods, the reader can consult the reference [13].

**Proposition 3.3.** (Invertibility of the implicit assembled matrices $\mathbb{A}$ $\hat{\mathbb{A}}$) Consider the InDC method constructed with IMEX R-K methods of type A or CK with invertible matrix $A$ or $\hat{A}$ in implicit part of Butcher tableau, respectively. If the quadrature nodes used in the InDC method exclude the left-most point and the IMEX R-K method is GSA, then the InDC method is a IMEX R-K method of type A or CK and it is GSA with the assembled matrix $\mathbb{A}$ or $\hat{\mathbb{A}}$ in the implicit part being invertible.

*Proof.* The proof is similar to that for Proposition 3.2. We first consider the case of type ARS, which is a special case of type CK. We let $\hat{\mathbb{A}}$ with size $\sigma$ be the invertible sub-matrix in the implicit part of the InDC Butcher table and $\hat{A}$ with size $s$ be the invertible matrix in the implicit part of IMEX method used to construct the InDC-IMEX method. Then $\sigma = s * M * (K+1)$, where $M$ is the number of quadrature points and $K$ is the number of correction loops. For example, for the InDC method constructed with IMEX2 method (5.1) with $M = 2$ and $K = 1$, $s = 2$, then $\sigma = 8$. $\hat{\mathbb{A}}$ can be constructed with a similar structure as shown in eq. (3.21), where $T$ is of size $s * M$ and it's block triangular. Its diagonal blocks ($M$ of them) are matricies $\hat{A}$ of the IMEX R-K scheme (of size $s$) scaled by the size of subinterval $1/M$. Recall from classical linear algebra that, the determinant of a block triangular matrix is the product of the determinants of diagonal blocks. Therefore,

$$
det(T) = (\frac{1}{M}det(\hat{A}))^M \neq 0,
$$

due to the invertibility of the matrix $A$ of the IMEX R-K method used to construct the InDC IMEX one. Hence, $det(\hat{\mathbb{A}}) = det(T)^{K+1} \neq 0$, i.e. $\hat{\mathbb{A}}$ is invertible. Finally, $\tilde{\mathbf{b}}^T$ and $\mathbf{b}^T$ vectors come from the last row of Butcher tableaus $\tilde{\mathbb{A}}$ and $\hat{\mathbb{A}}$, therefore the InDC method is GSA too. Similar results can be obtained for the type A and type CK. $\square$

**Remark 3.4.** Note that the assumptions in Proposition 3.2 and 3.3 provide a sufficient condition to guarantee the invertibility of implicit assembled matrix $\mathbb{A}$ or $\hat{\mathbb{A}}$. In particular, it has been pointed out in [3] that the invertibility of matrix $A$ or submatrix $\hat{A}$ of the IMEX R-K method, is an essential hypothesis for the error analysis for IMEX R-K methods of different types, otherwise error estimates no longer holds.

In the following proposition, by Proposition 3.3 and Remark 3.1, we can generalize the estimates (3.13) for the corresponding InDC method constructed with GSA IMEX R-K type A or type CK method, in the case of $\varepsilon = 0$, (reduced problem).

**Proposition 3.5.** Consider that the reduced system (2.2) with $\varepsilon = 0$ satisfies (2.1) and with consistent initial values. Consider the InDC IMEX R-K method constructed by GSA IMEX R-K methods of type CK or A. Then the global error after $K$ correction loops satisfies the following estimates

$$
e_n^{(K)} \doteq \hat{y}_n^{(K)} - y(t_n) = \mathcal{O}(H^{\min(s_K,M)}) \quad d_n^{(k)} \doteq \hat{z}_n^{(K)} - z(t_n) = \mathcal{O}(H^{\min(s_K,M)}),
\tag{3.31}
$$

with $\hat{y}_n^{(K)}$ and $\hat{z}_n^{(K)}$ being the numerical solution of the InDC methods at time $t_n$ after $K$ corrections, $s_K = \sum_{k=0}^{K} p^{(k)}$, and $H = Mh$ is one InDC time step. The estimates hold uniformly for $H \leq H_0$ and $nH \leq Const$.

14

*Proof.* From Proposition 3.3 and Remark 2.5 b), we get the estimates (3.31). □

Note that if the IMEX R-K method is not GSA, the assembled matrix for the InDC-IMEX R-K method is not invertible and by Remark 2.5 and 3.1, the error estimates (3.31) in general are not satisfied and an order reduction for the case reduced problem (2.2) ($\varepsilon = 0$) is observed.

# 4 Error estimates of InDC methods constructed with IMEX R-K

In this section, we present the main theoretical result in the form of a theorem. On the contrary to what has been done in [6], our idea here is to extend the error analysis of IMEX R-K methods applied to SPPs obtained in [3], to InDC methods constructed by IMEX R-K methods. In particular, we use the global error estimate results of IMEX methods for SPPs from [3] to InDC-IMEX methods by the fact that an InDC-IMEX method can be viewed as an IMEX R-K method with assembled Butcher tableaus, as discussed in the previous section, to obtain global error estimates for InDC-IMEX methods.

Then, as in [3], in the main thorem, we estimate the *global* errors of the InDC-IMEX method

$$e_{n,\nu}^{(K)} \doteq \hat{y}_{\nu}^{(K)}(t_n) - y_{\nu}(t_n), \quad d_{n,\nu}^{(K)} \doteq \hat{z}_{\nu}^{(K)}(t^n) - z_{\nu}(t_n), \quad \nu = 0, 1$$

where $\hat{y}_{\nu}^{(K)}(t_n)$ and $\hat{z}_{\nu}^{(K)}(t_n)$ are the $\nu$-th term in the $\varepsilon$-expansion of the numerical solution of the InDC-IMEX method with the $K$ correction steps at some final time $t^n$. Note that the estimates for the case $\nu = 0$ has been provided in the Proposition 3.5.

**Theorem 4.1.** *Consider the stiff system* (1.2), (2.1) *with well-prepared initial values* $y(0)$, $z(0)$ *admitting a smooth solution. Consider the InDC method constructed with $M$ uniformly distributed quadrature nodes excluding the left-most point and a globally stiffly accurate IMEX R-K method of order $p^{(0)}$ of type A or CK. Apply IMEX R-K methods of different classical orders $(p^{(1)}, p^{(2)}, \ldots, p^{(K)})$ in the correction loops, $k = 1, \cdots K$. Assume that each of these IMEX R-K methods in the correction loops are globally stiffly accurate. Then the global error after $K$ correction loops satisfies the following estimates*

$$\begin{array}{rcl} e_n^{(K)} \doteq \hat{y}^{(K)}(t_n) - y(t_n) &=& \mathcal{O}(H^{\min(s_K, M)}) + \mathcal{O}(\varepsilon H), \\ d_n^{(K)} \doteq \hat{z}^{(K)}(t_n) - z(t_n) &=& \mathcal{O}(H^{\min(s_K, M)}) + \mathcal{O}(\varepsilon H), \end{array} \tag{4.1}$$

*where $\hat{y}^{(K)}(t_n)$ and $\hat{z}^{(K)}(t_n)$ are the numerical solutions of the InDC methods at $t^n$, for $\varepsilon \leq cH$ and for any fixed constant $c > 0$, $s_K = \sum_{k=0}^{K} p^{(k)}$, and $H = Mh$ is one InDC time step. The estimates hold uniformly for $H \leq H_0$ and $nH \leq Const$.*

Now we give the following proposition as a consequence of the Lemma 5.1 and theorem 5.2, and 6.2 in [3]. The main Theorem 4.1 follows from this result.

**Proposition 4.2.** Consider a GSA IMEX R-K method of type A or CK, with $\tilde{b}_i \neq b_i$, for $i = 1, \ldots, s$, and let $p$ be the order of the explicit part of the scheme. Apply this method to the general problem (1.2), under the hypothesis (2.1) with initial values consistent and such that the problem (1.2) admits a smooth solution. Then, for any fixed constant $C$, the global error satisfies for $\varepsilon < Ch$

$$\hat{y}_n - y(t_n) = \mathcal{O}(h^p) + \mathcal{O}(\varepsilon h), \quad \hat{z}_n - z(t_n) = \mathcal{O}(h^p) + \mathcal{O}(\varepsilon h), \tag{4.2}$$

These estimates hold uniformly for $h \leq h_0$ and $nh < C$ for any fixed constant $C > 0$.

*Proof.* The proof of this Proposition is obtained by combining Remarks 2.5 a) b) and by the results in Theorem 5.2 and 6.2 in [3]. □

**Proof of Theorem 4.1.** By Proposition 3.3, the InDC IMEX R-K method constructed by a GSA IMEX R-K method of type A or CK can be viewed as an IMEX R-K method of type A or CK and is GSA. Then the hypothesis of Proposition 4.2 are satisfied for the InDC IMEX R-K scheme. The estimates (4.1) are a consequence of Proposition 3.5 and 4.2.

We point out that the simplest InDC scheme is constructed by repeated use of the same IMEX scheme of order $p$, and the optimal choice of $M$ is given by $M = s_K = p(K+1)$. This is our choice in the PDE applications presented in the paper.

# 5 Numerical evidence and PDE applications

In this section, we consider the following InDC IMEX methods for stiff ODEs and PDEs. Below we present a list of IMEX R-K methods that are used in the InDC framework. These include first order GSA IMEX method of type A and ARS, and first order IMEX method of type A but not GSA. For high order IMEX methods, we choose methods of type CK or ARS. We decide not to use high order GSA IMEX R-K methods of type $A$, because they are more expensive compared with methods of type CK and ARS. In fact, the construction of such methods involves many internal stages. For example for a second order GSA IMEX R-K method of type $A$, we require more than three ($s > 3$) internal stages. In fact, in [5] the authors proved that it is not possible to construct second order GSA IMEX R-K methods of type A with three internal stages.

**InDC method embedded with a first order IMEX.**

- We let InDC-IMEX1-GSA-ARS-M-k denote the InDC methods embedded with first order ($p = 1$) GSA IMEX R-K scheme of type ARS (3.14) with $M$ quadrature points and $k$ correction steps.

- We let InDC-IMEX1-NGSA-M-k denote the InDC methods embedded with first order ($p = 1$) non globally stiffly accurate IMEX R-K scheme (3.16) with $M$ quadrature points and $k$ correction steps. The first order non globally stiffly accurate IMEX R-K method of type A (IMEX1-NGSA) has the following double Butcher tableau

- We let InDC-IMEX1-GSA-A-M-k denote the InDC method constructed with first order ($p = 1$) GSA IMEX R-K scheme of type A (3.15) with $M$ quadrature points and $k$ correction steps.

**InDC method embedded with a second order IMEX.**

- We let InDC-IMEX2-ARS-M-k denote the InDC method constructed with the second order globally stiffly accurate IMEX R-K method of type ARS with the following double Butcher tableau

$$
\begin{array}{c|ccc}
0 & 0 & 0 & 0 \\
\gamma & \gamma & 0 & 0 \\
1 & \delta & 1-\delta & 0 \\
\hline
& \delta & 1-\delta & 0
\end{array}
\qquad
\begin{array}{c|ccc}
0 & 0 & 0 & 0 \\
\gamma & 0 & \gamma & 0 \\
1 & 0 & 1-\gamma & \gamma \\
\hline
& 0 & 1-\gamma & \gamma
\end{array},
\tag{5.1}
$$

  where $\gamma = 1 - \frac{\sqrt{2}}{2}$ and $\delta = 1 - 1/(2\gamma)$, with $k$ correction steps and $M$ quadrature points. This method has order $p = 2$.

- We let InDC-IMEX2-CK-M-k denote the InDC method constructed with a second order globally stiffly accurate method of type CK with the following double Butcher tableau and $\gamma = 1 - \sqrt{2}/2$

$$
\begin{array}{c|ccc}
0 & 0 & 0 & 0 \\
2/3 & 2/3 & 0 & 0 \\
1 & 1/4 & 3/4 & 0 \\
\hline
& 1/4 & 3/4 & 0
\end{array}
\qquad
\begin{array}{c|ccc}
0 & 0 & 0 & 0 \\
2/3 & 2/3-\gamma & \gamma & 0 \\
1 & 1/4+\gamma/2 & 3/4-3\gamma/2 & \gamma \\
\hline
& 1/4+\gamma/2 & 3/4-3\gamma/2 & \gamma
\end{array}.
\tag{5.2}
$$

  with $k$ correction steps and $M$ quadrature points. This method has order $p = 2$.

**InDC method embedded with a third order IMEX.**

- We let InDC-IMEX3-ARS-M-k denote the InDC method constructed with a third order globally stiffly accurate method of type ARS with the following double Butcher tableau

$$
\begin{array}{c|ccccc}
0 & 0 & 0 & 0 & 0 & 0 \\
1/2 & 1/2 & 0 & 0 & 0 & 0 \\
2/3 & 11/18 & 1/18 & 0 & 0 & 0 \\
1/2 & 5/6 & -5/6 & 1/2 & 0 & 0 \\
1 & 1/4 & 7/4 & 3/4 & -7/4 & 0 \\
\hline
& 1/4 & 7/4 & 3/4 & -7/4 & 0
\end{array}
\qquad
\begin{array}{c|ccccc}
0 & 0 & 0 & 0 & 0 & 0 \\
1/2 & 0 & 1/2 & 0 & 0 & 0 \\
2/3 & 0 & 1/6 & 1/2 & 0 & 0 \\
1/2 & 0 & -1/2 & 1/2 & 1/2 & 0 \\
1 & 0 & 3/2 & -3/2 & 1/2 & 1/2 \\
\hline
& 0 & 3/2 & -3/2 & 1/2 & 1/2
\end{array}.
\tag{5.3}
$$

  in the prediction and $k$ correction steps and $M$ quadrature points. This method is stiffly accurate with order $p = 3$.

The indicated order of convergence by Theorem 4.1 for the $y$ and $z$ components in the SPPs are summarized in Table 5.1. For the InDC-IMEX1-GSA-ARS-M-k and InDC-IMEX1-GSA-A-M-k methods, the order of convergence will increase with $k$ for the $\varepsilon^0$ error term when $\varepsilon \ll H$ and $k \leq M - 1$, leading to a term of $H^{\min(k+1,M)}$ for the differential and algebraic component.Similar comments apply to the InDC-IMEX2-ARS-M-k, InDC-IMEX2-CK-M-k, as well as the InDC-IMEX3-ARS-M-k with the order of accuracy for the first error term increased by 2 or 3 per correction step respectively. The order of convergence for the $\varepsilon^1$ error term is $\varepsilon H$ for all schemes we consider here. Note that for those InDC-IMEX methods with the same order of accuracy for the index 1 problem (i.e. when $\varepsilon = 0$), the complexity measured by the number of function evaluations is comparable. For example, when $M = 8$, the number of function evaluations for the InDC-IMEX1-GSA-ARS-8-7, InDC-IMEX2-ARS-8-3, InDC-IMEX2-CK-8-3 are the same. Note that all of these methods achieve eighth order accuracy for the index 1 problem when $\varepsilon = 0$. On the other hand, the number of function evaluations for InDC-IMEX1-GSA-A-8-7 is twice as much as InDC-IMEX1-GSA-ARS-8-7, e.g. see the number of stages in Butcher tableaus in (3.18) and (3.20). From the computational cost point of view, the InDC-IMEX1-GSA-ARS method is preferred. For the sake of completeness, we present results for the InDC-IMEX1-GSA-A method for the ODE Van der Pol equation, but not for the PDE examples.

Table 5.1: Global error predicted by Theorem 4.1 with $H \gg \varepsilon$.

| Method | $y-$comp | $z-$comp |
|---|---|---|
| InDC-IMEX1-GSA-ARS-M-k | $H^{\min(k+1,M)} + \varepsilon H$ | $H^{\min(k+1,M)} + \varepsilon H$ |
| InDC-IMEX1-GSA-A-M-k | $H^{\min(k+1,M)} + \varepsilon H$ | $H^{\min(k+1,M)} + \varepsilon H$ |
| InDC-IMEX2-ARS-M-k | $H^{\min(2(k+1),M)} + \varepsilon H$ | $H^{\min(2(k+1),M)} + \varepsilon H$ |
| InDC-IMEX2-CK-M-k | $H^{\min(2(k+1),M)} + \varepsilon H$ | $H^{\min(2(k+1),M)} + \varepsilon H$ |
| InDC-IMEX3-ARS-M-k | $H^{\min(3(k+1),M)} + \varepsilon H$ | $H^{\min(3(k+1),M)} + \varepsilon H$ |

## 5.1 Van der Pol example

For numerical verification, we test a standard nonlinear oscillatory test problem, Van der Pols equation with well-prepared initial data up to $\mathcal{O}(\varepsilon^3)$, [19]:

$$\begin{cases} y' = z \\ \varepsilon z' = (1 - y^2)z - y \end{cases}, \quad \begin{cases} y(0) = 2 \\ z(0) = -\frac{2}{3} + \frac{10}{81}\varepsilon - \frac{292}{2187}\varepsilon^2 \end{cases} \tag{5.4}$$

and $\varepsilon = 10^{-6}$.

Numerical observations in Figures 5.1 and 5.2 are consistent with Theorem 4.1 and Table 5.1. They produce estimates for the $y$ and $z$ component in the form of equation (4.1). Especially, numerical results of InDC method constructed with first order IMEX methods presented in Figure 5.1 and of InDC method constructed with high order (second or third order) InDC-IMEX2-ARS-M-K presented in Figure 5.2, confirm the theoretical prediction (4.1), i.e., if the IMEX method is GSA, the order of convergence for $\varepsilon^0$ term for the $y$ and $z$ component increases with the correction loops (that is $s_K$ with $K$ the number of correction loops). However, such improvement for $\varepsilon^0$ term is not true if the IMEX method is not globally stiffly accurate (see two panels in the top row of Figure 5.1 for InDC-IMEX1-NGSA method). Furthermore, in the estimate (4.1), the InDC-IMEX methods exhibit order reduction both in differential and algebraic components (see Table 5.1 for every type of InDC method). This phenomenon appears, since the $\varepsilon^1$ term of the error behaves like $\mathcal{O}(\varepsilon H)$ in (4.1) for both $y$ and $z$ components (see in Figures 5.1 and 5.2). When $H$ is very small, the $\mathcal{O}(\varepsilon H)$ term is dominant in the estimates of $y$ and $z$-components respectively. Finally, we want to comment that, despite the same order of convergence, the magnitude of errors of InDC-IMEX1-GSA-A-7-3 are much smaller than those of InDC-IMEX1-GSA-ARS-7-3. As we mentioned earlier, the InDC-IMEX1-GSA-ARS method costs less number of function evaluations and is more efficient.

## 5.2 PDEs examples

In this section we consider systems of the form (1.1), which may arise from a method of lines discretization of PDEs. Notice that in PDEs applications, such as advection-diffusion problem, convection-diffusion-reaction
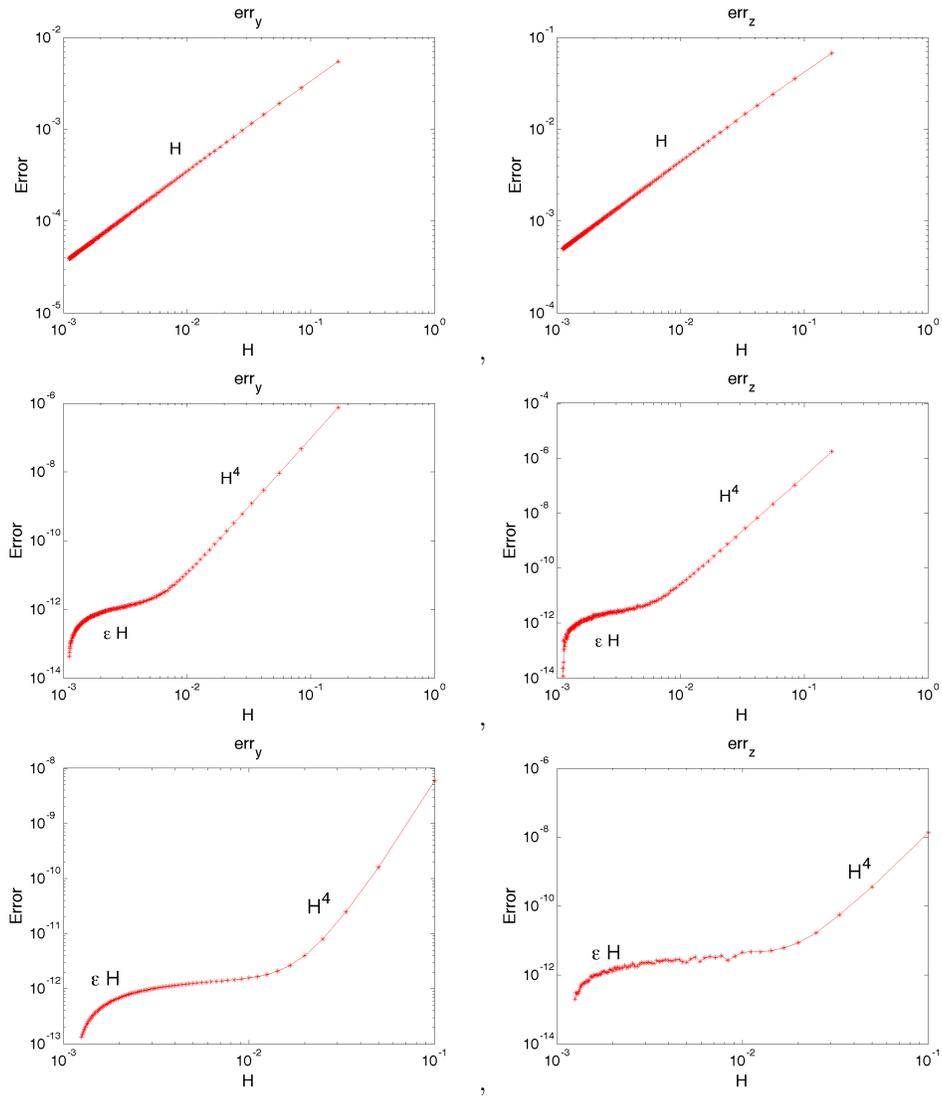
Figure 5.1: Van der Pol equation. Global error ($T = 0.5$) of the InDC-IMEX1-NGSA-5-3 method (top row), InDC-IMEX1-GSA-ARS-7-3 method (middle row), and InDC-IMEX1-GSA-A-7-3 method (bottom row). $\varepsilon = 10^{-6}$.
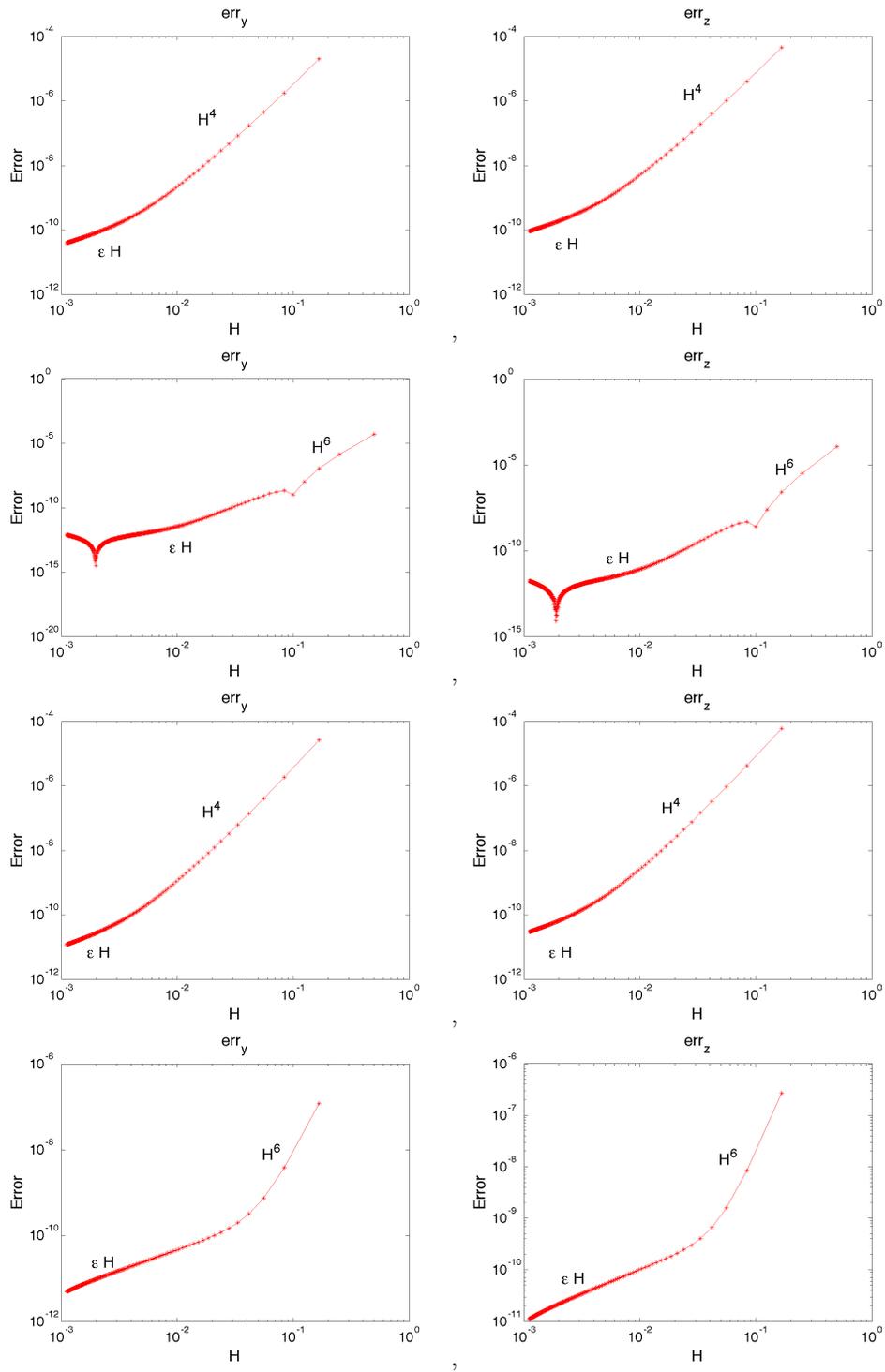
Figure 5.2: Van der Pol equation. Global error ($T = 0.5$) of the InDC-IMEX2-ARS-5-1 method (upper row), of the InDC-IMEX2-ARS-7-2 method (the second row), of the InDC-IMEX2-CK-5-1 method (the third row), and of the InDC-IMEX3-ARS-7-1 method (bottom row). $\varepsilon = 10^{-6}$.

systems and hyperbolic systems with relaxations, the structure of the problems is usually of additive form as (1.1).

**Example 5.1. (Advection-diffusion equation)** As an example of advection-diffusion equation, we consider the 1D viscous Burgers' equation,

$$u_t + \left(\frac{u^2}{2}\right)_x = \frac{1}{R}u_{xx}, \quad R > 0. \tag{5.5}$$

It contains a 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. The solution develops a sharp shock wave front after a certain time of propagation and when $R \to \infty$ discontinuities appear. 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 system (5.5), after a method of lines approach with $U_i(t) = u(x_i, t)$ and $U(t) = (U_1(t), U_2(t), \cdots, U_N(t))^T$, has the form of (1.1). The convection term $-(u^2/2)_x$ is considered as non-stiff and is integrated using the explicit part of the method; it is discretized by high order finite difference discretization with Weighted-Essentially Non Oscillatory (WENO) reconstruction [29]. The diffusion term $u_{xx}$ is considered stiff and is treated implicitly; it is discretized by the standard fourth order finite difference technique, excepted at the nearby boundary points where a 3-rd order formula was implemented.

To check the temporal order of convergence, we consider the periodic boundary condition and a smooth initial condition $u(x, 0) = \sin(\pi x)$ with $x \in [0, 1]$. The equation has been integrated to time $T = 0.02$. Notice that the solution drops dramatically from $sin(\pi x)$ to zero within the first 0.05 second. The exact solution of this problem was obtained by Cole [15] so that numerical comparison can be done. In the Tables 5.2 and 5.3, we show the $L^2$ error and the corresponding order of convergence in time "Order" for a sequence of mesh sizes $\Delta t = \Delta x = 1/N$. We use the InDC-IMEX1 method with one or two corrections, corresponding to a second or third order time integration scheme and InDC-IMEX2-ARS-GSA-4-2 method with one correction, corresponding to a fourth order scheme. The expected temporal order of convergence is numerically observed.

Table 5.2: Example 5.1. Accuracy test with the Burgers' equationfor InDC- IMEX1-M-k GSA, with k = 1, 2 corrections.

| Method | $\Delta t$ | $L_2$ error | Order |
|---|---|---|---|
| InDC-IMEX1-2-1 | 1/96 | $2.4551e - 02$ | $--$ |
| InDC-IMEX1-2-1 | 1/192 | $9.0755e - 03$ | 1.435 |
| InDC-IMEX1-2-1 | 1/384 | $2.7988e - 03$ | 1.697 |
| InDC-IMEX1-2-1 | 1/768 | $7.9568e - 04$ | 1.814 |
| InDC-IMEX1-2-1 | 1/1536 | $2.1206e - 04$ | 1.907 |
| InDC-IMEX1-3-2 | 1/3072 | $5.4748e - 05$ | 1.953 |
| InDC-IMEX1-3-2 | 1/96 | $6.9529e - 04$ | $--$ |
| InDC-IMEX1-3-2 | 1/192 | $2.2174e - 04$ | 1.648 |
| InDC-IMEX1-3-2 | 1/384 | $4.1888e - 05$ | 2.404 |
| InDC-IMEX1-3-2 | 1/768 | $6.5849e - 06$ | 2.669 |
| InDC-IMEX1-3-2 | 1/1536 | $9.1527e - 07$ | 2.846 |
| InDC-IMEX1-3-2 | 1/3072 | $1.2111e - 07$ | 2.920 |

Next we consider hyperbolic systems with stiff relaxation. These systems have the form

$$\mathcal{U}_t + \mathcal{F}(\mathcal{U})_x = \frac{1}{\varepsilon}\mathcal{G}(\mathcal{U}), \tag{5.6}$$

where $\mathcal{U} \in \mathbb{R}^N$, $\mathcal{F}, \mathcal{G} : \mathbb{R}^N \to \mathbb{R}^N$ and the Jacobian matrix $\mathcal{F}'(\mathcal{U})$ has real eigenvalues and admits a basis of eigenvectors for all $\mathcal{U} \in \mathbb{R}^N$ and $\varepsilon > 0$ is the stiffness parameter. The operator $\mathcal{G} : \mathbb{R}^N \to \mathbb{R}^N$ is called a relaxation operator. (5.6) defines a relaxation system. For mathematical results concerning this kind of problems we refer to [25] and [10]. Same as for the previous example, we discretize the system (5.6) by a finite difference scheme with $U_i(t)$ approximates the solution at grid point $x_i$, $i = 1, \cdots N$ and $U(t) = (U_1(t), U_2(t), \cdots, U_N(t))^T$ be solutions at all grid points. With the method of line approach, $U(t)$ satisfies a system in the form of (1.1),

Table 5.3: Example 5.1. Accuracy test with the Burgers' equation examplefor second order IMEX2-ARS GSA method and InDC-IMEX2-ARS-GSA-M-k, with k = 1 correction.

| Method | N | $L_2$ error | Order |
|---|---|---|---|
| InDC-IMEX2-ARS-GSA-4-2 | 1/48 | $7.6079e - 04$ | $--$ |
| InDC-IMEX2-ARS-GSA-4-2 | 1/96 | $2.7095e - 05$ | 4.811 |
| InDC-IMEX2-ARS-GSA-4-2 | 1/192 | $1.1064e - 06$ | 4.614 |
| InDC-IMEX2-ARS-GSA-4-2 | 1/384 | $5.5847e - 08$ | 4.302 |
| InDC-IMEX2-ARS-GSA-4-2 | 1/758 | $3.4640e - 09$ | 4.010 |
| InDC-IMEX2-ARS-GSA-4-2 | 1/1536 | $2.3912e - 10$ | 3.856 |

where the function $F(U)$ in (1.1) as a discretization of the convective term $-\mathcal{F}(\mathcal{U}))_x$ is being treated explicitly, and $G(U)$ as a discretization of the source term $\mathcal{G}(\mathcal{U})$ is being treated implicitly.

A prototype example that we will use to illustrate our theoretical findings is the following $2 \times 2$ nonlinear hyperbolic system with relaxation ([7], [28], [10])

$$
\begin{aligned}
u_t + f_1(u, v)_x &= 0, \\
v_t + f_2(u, v)_x &= \tfrac{1}{\varepsilon} g(u, v).
\end{aligned} \tag{5.7}
$$

System (5.7) is a particular case of (5.6) with $\mathcal{U} = (u, v)^T$ , $u \in \mathbb{R}^n$, $v \in \mathbb{R}^{N-n}$, $n < N$, $\mathcal{F} = (f_1, f_2)^T$ and $\mathcal{G} = (0, g)^T$. Note that system (5.7) possesses a unique local equilibrium, namely, $g(u, v) = 0$ implies $v = q(u)$ and, at the local equilibrium, one has the macroscopic system

$$
u_t + f_1(u, q(u))_x = 0. \tag{5.8}
$$

Equation (5.8) can be derived by sending $\varepsilon$ in (5.7) to zero, the so-called zero relaxation limit, ([7], [28], [10]). If we use a method of lines to (5.7), we get a ODE system of the form (1.1).

**Example 5.2.** Consider a linear system with stiff relaxation source term

$$
\begin{aligned}
\partial_t u + \partial_x v &= 0, \\
\partial_t v + \partial_x u &= -\tfrac{1}{\varepsilon}(v - bu), \quad x \in [0, 2], \quad t \in [0, T],
\end{aligned} \tag{5.9}
$$

with $\varepsilon > 0$ and $b$ constant. Here $\mathcal{F}(\mathcal{U}) = (v, u)^T$ and $\mathcal{G}(\mathcal{U}) = (0, (v - bu))^T$.

To gain an understanding of the system (5.9), we consider a formal expansion of solutions in the form

$$
\begin{aligned}
u &= u_0 + u_1 \varepsilon + \mathcal{O}(\varepsilon^2), \\
v &= v_0 + v_1 \varepsilon + \mathcal{O}(\varepsilon^2),
\end{aligned} \tag{5.10}
$$

and insert it in (5.9). Collect leading terms in approximating the second equation of (5.9) (order $\varepsilon^0$), we get $v_0 = bu_0$, plug which into the first equation, we obtain $\partial_t u_0 + b\partial_x u_0 = 0$. This equation is formally obtained as $\varepsilon \to 0$, and we call it *reduced* equation. Next, we consider the first-order correction to the leading term approximation. By looking for the order $\varepsilon$ correction to the approximation (5.9) we get Then, by keeping first-order terms in the expansion (5.10) and neglecting second and higher order terms, we obtain a dissipative evolutionary equation

$$
\begin{aligned}
v &= bu - (1 - b^2)\partial_x u \\
\partial_t u + b\partial_x u &= \varepsilon \partial_x((1 - b^2)\partial_x u).
\end{aligned} \tag{5.11}
$$

The second equation is a convection-diffusion equation with viscosity coefficient $\nu = \varepsilon(1 - b^2)$ and it is dissipative if $|b| \leq 1$ (subcharacteristc condition of Liu [10] for (5.9)).

Then motivated by this analysis we perform an accurate test for this problem considering well-prepared initial data. Note that the initial conditions for the coefficients of $u_i(0)$ in (5.10) can be chosen arbitrarily, but there is no freedom in the choice of $v_i(0)$. We let $u_0(x, 0) = \sin(2\pi x)$, and $u_1(x, 0) = 0$ with $v_0(x, 0) = bu_0(x, 0)$ and $v_1(x, 0) = (b^2 - 1)\partial_x u(x, 0)$. Such initial data is consistent, i.e., $g(u_0(x, 0), v_0(x, 0)) = (bu_0(x, 0) - v_0(x, 0)) = 0$ for $\varepsilon = 0$.

We consider a periodic smooth solution and we set $b = 0.5$ and $T = 0.2$. The spatial discretization of the domain has a fixed space mesh $\Delta x = 0.02$. Our test problems are computed with coarse temporal grids that
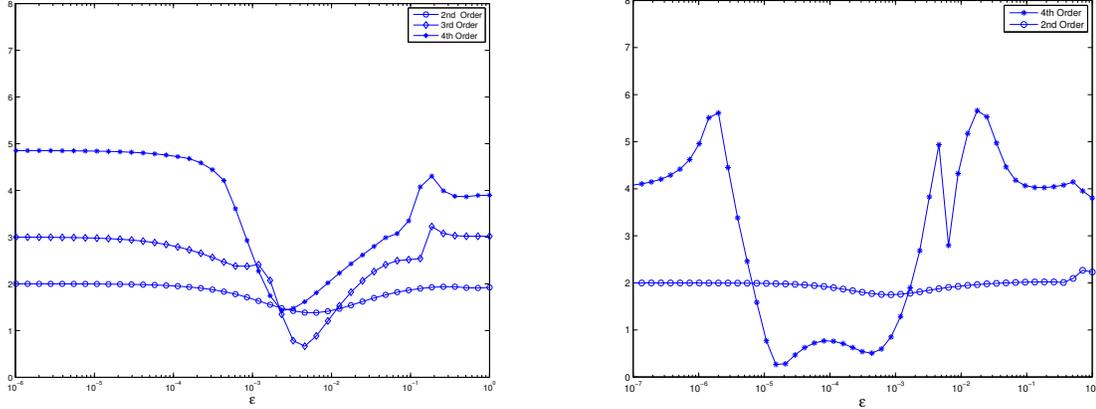
Figure 5.3: Example 5.2. Convergence rate for $u$-component versus $\varepsilon$. Left: for the method InDC-IMEX1-M-k, with $M = 2, 3, 4$ corresponding to $k = 1, 2, 3$ corrections respectively. Right: second order IMEX2-ARS method GSA and InDC-IMEX2-ARS-GSA-M-k, with $M = 4$ and $k = 0, 1$ correction.

do not resolve small scales. High accuracy in space is obtained by finite difference discretization with WENO reconstruction [29]. To compute the error, the exact solution is obtained by the Fourier series

$$U_{ex}(x, t) = \sum_{\ell=-\infty}^{+\infty} U_\ell(t)e^{i\ell x}, \quad V_{ex}(x, t) = \sum_{\ell=-\infty}^{+\infty} V_\ell(t)e^{i\ell x}$$

with $U_\ell(t)$ and $V_\ell(t)$ satisfying the ODE system

$$\begin{aligned}
\dot{U}_\ell &= -i\ell V_\ell, \\
\dot{V}_\ell &= -i\ell U_\ell - \tfrac{1}{\varepsilon}(V_\ell - bU_\ell)
\end{aligned} \tag{5.12}$$

For each $\ell$, system (5.12) can be written as a 2x2 constant coefficient homogeneous system which can be solved exactly. For the initial condition $u(x, 0) = \sin(2\pi x)$ and the corresponding one for $v$, it is sufficient to consider $\ell = 1$ only to obtain the exact solution.

In Figure 5.3, we plot the temporal convergence rate of the $u$-component as a function of $\varepsilon$. Here order of accuracy curves are obtained in the following way: on each curve and for each value of $\varepsilon$, the order of accuracy is computed from measuring $L^1$ errors computed by the difference between the exact solution and the numerical one for two different time step sizes $(\Delta t, \Delta t/2)$. Expected high order convergence rate is observed in the limiting cases of $\varepsilon \to 0$ and $\varepsilon \approx 1$ for the InDC-IMEX1-M-k when $M = 2, 3, 4$ and $k = 1, 2, 3$ and for the InDC-IMEX2-ARS-GSA-M-k when $M = 4$ and $k = 0, 1$ respectively. We observe that the convergence rate is increased by InDC correction iterations for sufficiently stiff parameters; however, for intermediate values of the parameter $\varepsilon$, e.g., $10^{-4} < \varepsilon < 10^{-2}$, we have a deterioration of the accuracy, as we expect. This is an indication that these InDC IMEX methods suffer from the phenomenon of order reduction in the mildly stiff regime ($\Delta t = \mathcal{O}(\varepsilon)$), when the classical order is greater than two [3, 4, 9]. 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. In the next example we will give a brief explanation of this lack of accuracy for mildly stiff regime.

**Example 5.3.** Consider a nonlinear hyperbolic system with relaxation [21, 7]

$$\begin{aligned}
h_t + w_x &= 0, \\
w_t + (h + 0.5h^2)_x &= -\tfrac{1}{\varepsilon}(w - 0.5h^2),
\end{aligned} \tag{5.13}$$

where $\mathcal{U} = (h, w)$, $\mathcal{F}(\mathcal{U}) = (w, (h + 0.5h^2))^T$ and $\mathcal{G}(\mathcal{U}) = (0, -(w - 0.5h^2))^T$. We consider a *well-prepared* initial data given by $h(0, x) = 1 + 0.2\sin(8\pi x)$ and $w = w_0 + \varepsilon w_1$ with $w_0 = f(h(0, x)) = 0.5h^2(0, x)$ and $w_1 = (f'(h(0, x)) - p'(h(0, x))\partial_x h(0, x)$ where $p(h) = (h + 0.5h^2)$. The initial conditions are designed to be well-prepared in the same way as the previous example. The boundary condition is periodic. We evolve the solution to $T_{final} = 0.1$ before the shock forms. We perform our numerical test using a fixed space mesh with $\Delta x = 0.01$ for $x \in [0, 1]$.
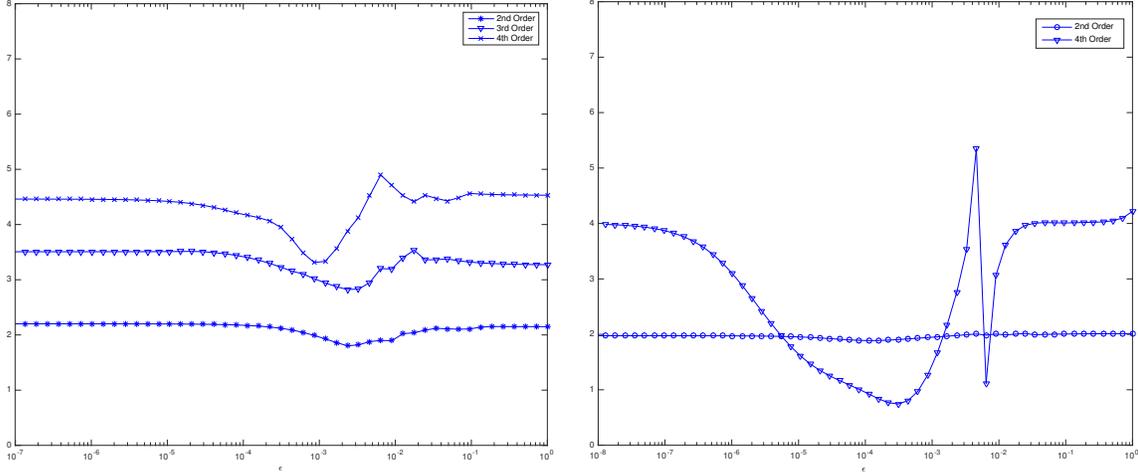
Figure 5.4: Example 5.3. Convergence rate for $h$-component versus $\varepsilon$. Left: for the method InDC-IMEX1-M-k, with $M = 2, 3, 4$ corresponding to $k = 1, 2, 3$ corrections respectively. Right: second order IMEX2-ARS method GSA and InDC-IMEX2-ARS-GSA-M-k, with $M = 4$ and $k = 1$ correction.

In Fig. 5.4, we plot the temporal convergence rate as a function of $\varepsilon$, which is computed in the same fashion as those in Fig. 5.3. Each point of the graph marked by a triangle in Figure 5.4 is given by the expression

$$\text{order}_h(\varepsilon) = \log_2(E_{h/2}(\varepsilon)/E_h(\varepsilon))$$

where the errors $E_{h/2}$ and $E_h$ are computed for the same value of $\varepsilon$. For example, the values along the dashed line in Figure 5.5 are used to compute the point on Figure 5.4 corresponding to $\varepsilon = 10^{-4}$. A similar behavior is observed as that in Fig. 5.3: order of convergence strongly depends on $\varepsilon$. In particular, for small and large values of $\varepsilon$, the order of convergence is increased by InDC correction iterations, while for intermediate values of $\varepsilon$, we observe that the order of convergence goes down to small values, even below 1.

We show now that this observation is not in contradiction with the theoretical prediction. We start observing that if $\varepsilon \gg h$ then a classical analysis can be used to estimate the global error, i.e. the classical global error due to the presence of the small parameter $\varepsilon$ is of the order $\mathcal{O}(h/\varepsilon)^p$, with $p$ the order of the scheme. On the other hand, as a natural consequence of the Proposition 4.2, if $\varepsilon \ll h$, the global error satisfies (4.2). As both estimates about the global error have to be satisfied (see Figure 5.6), we have:

$$\max_{\varepsilon} err_{\varepsilon} \lesssim \max_{\varepsilon} \left( \min \left( \mathcal{O}\left( \frac{h}{\varepsilon} \right)^p, \mathcal{O}(h^p) + \mathcal{O}(\varepsilon h) \right) \right). \tag{5.14}$$

Since the estimate is based on classical order (when $h \ll \varepsilon$) and on the asymptotic expansion in $\varepsilon$ (when $h \gg \varepsilon$), we do not expect it to be sharp in intermediate regime, when $\varepsilon \approx h$. Then a simple calculation shows that the uniform order is $\mathcal{O}(h^{\frac{2p}{p+1}})$, where the worst case takes place where $\mathcal{O}(\varepsilon h) = \mathcal{O}((h/\varepsilon)^p)$, i.e., for $\varepsilon = \varepsilon^*$ with

$$\varepsilon^* = \mathcal{O}(h^{\frac{p-1}{p+1}}). \tag{5.15}$$

Therefore we get $\max_{\varepsilon} err_{\varepsilon} = \mathcal{O}(h^{\frac{2p}{p+1}})$. This argument can be extended to InDC IMEX R-K methods using the estimates (4.1) given in the Theorem 4.1 In this case we get that the uniform order is

$$\max_{\varepsilon} err_{\varepsilon} = \mathcal{O}(h^{\frac{2r}{r+1}}) \tag{5.16}$$

where $r = \min(M, s_k)$ is the classical order of the InDC IMEX R-K.

In order to show that our numerical observations in Fig. 5.4 are not in contradiction with the theoretical prediction eq. (5.16) of the uniform order, we produce two error plots of the fourth order method (InDC-IMEX2-ARS-GSA-M-k, with $M = 4$ and $k = 1$) in Fig. 5.5. On the left panel, for each curve and for each value of $\varepsilon$, we plot the $L_1$ error, which is computed by comparing solutions for different values of time step, i.e. $E_i = \|S_{h/2^i} - S_{h/2^{i-1}}\|_{L^1}$, for $i = 1, 2, 3, 4$. The blue circles are the values corresponding to $\max_{\varepsilon} err_{\varepsilon}$. A lack
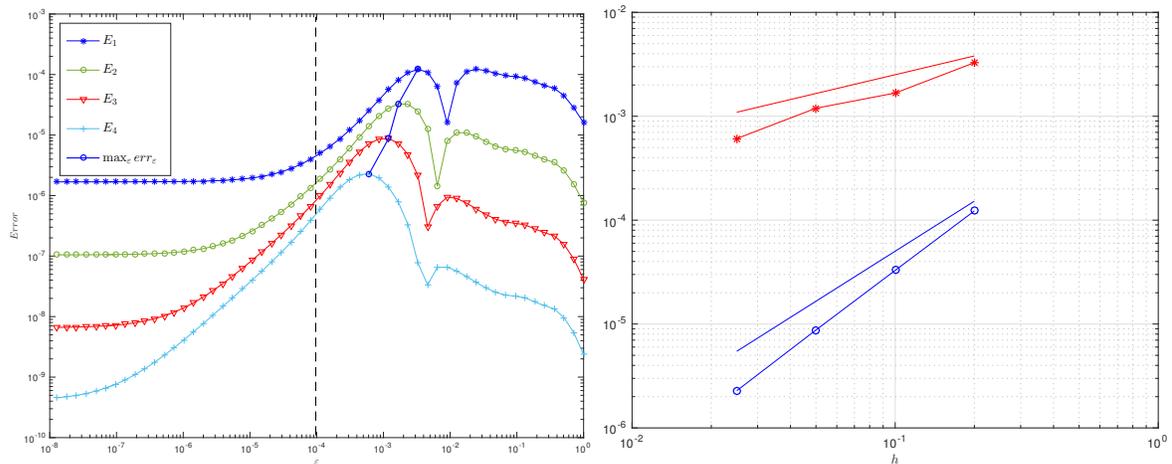
23

Figure 5.5: Example 5.3. Logarithm scale for the error versus the parameter $\varepsilon$ with $10^{-8} \leq \varepsilon \leq 1$ (left panel) and $\varepsilon^* = arg(\max_\varepsilon err_\varepsilon)$ versus $h$ (identified by $-*$) and $err(\varepsilon^*)$ versus $h$ (identified by "o") (right panel). The reference lines are linear curves with reference slopes 0.6 and 1.6 respectively. We used InDC-IMEX2-ARS-GSA-M-k, with $M = 4$ and $k = 1$ correction.

of accuracy is observed in the intermediate regime when $h$ is approximately of the same order of $\varepsilon$ On the right panel of Fig. 5.5, we act differently: for each $h$ we plot $\varepsilon^* = arg(\max_\varepsilon err_\varepsilon(h))$ (identified by stars) and $err_{\varepsilon^*}$ (identify by circles) as a function of $\Delta t$. We compare it with the theoretical estimate (5.16). Using least square best fit on the four blue points, the computed uniform order appears to be 1.92 while theoretical estimate (5.16) predicts a uniform order 1.6 (blue continuous line in Fig. 5.5) for $r = 4$. Likewise, a best fit for four red stars $\varepsilon^* \propto h^\alpha$, gives $\alpha = 0.78$, while the theoretical prediction given in Eq. (5.15) gives $\alpha = (r-1)/(r+1) = 0.6$ (red continuous line in Fig. 5.5) with $r = 4$. A similar behavior is obtained when adopting a third order scheme $r = 3$ using a InDC-IMEX1- M-k, with $M = 3$ and $k = 2$. The computed uniform order is approximately 1.9, while theoretical prediction gives 1.5. There is a small discrepancy between the computed values and the theoretical prediction. Such small discrepancy suggests that the estimate for the uniform order is not sharp. The reason is not fully understood and deserves deeper analysis.

# 6 Conclusions

This paper studies the order of convergence of InDC-IMEX methods when applied to SSPs, using uniform distribution of quadrature points excluding the left-most point. Since InDC methods have a similar structure to R-K methods [13], we construct the InDC-IMEX R-K methods as IMEX R-K methods with enlarged assembled Butcher tables and apply the convergence results in [3] directly to the InDC-IMEX R-K methods. Theoretical results on global error estimates in the form of $\varepsilon$-expansion are presented. The InDC-IMEX schemes are applied to the classical Van der Pol equations and PDE systems in order to illustrate our theoretical findings. In particular the order reduction phenomenon is observed as expected for intermedia $\varepsilon$, while high order convergent rates are observed in the asymptotic limit when $\varepsilon \to 0$ and when $\varepsilon = \mathcal{O}(1)$ (similarly as in [7]).

We also pointed out that the globally stiffly accurate property of the IMEX R-K scheme is an important assumption: the fact that an IMEX RK is GSA guarantees that the assembled matrix of the corresponding InDC IMEX R-K scheme based on it is invertible, and this in turn implies that the high order accuracy in the limit as $\varepsilon \to 0$ is maintained. Note that the assumption of GSA provides a sufficient condition to guarantee the invertibility of the assembled matrix. Although we do not know whether GSA property is necessary, we showed an example of InDC scheme based on a non GSA IMEX R-K, which lacks such asymptotic accuracy property.

Furthermore, we showed that even if InDC IMEX R-K can in principle be re-written as IMEX-RK with many stages, the actual implementation of such schemes as InDC IMEX R-K illustrated in the paper provides a systematic way to construct high order IMEX RK schemes, which would be far too complicated to write down as standard IMEX-RK schemes.
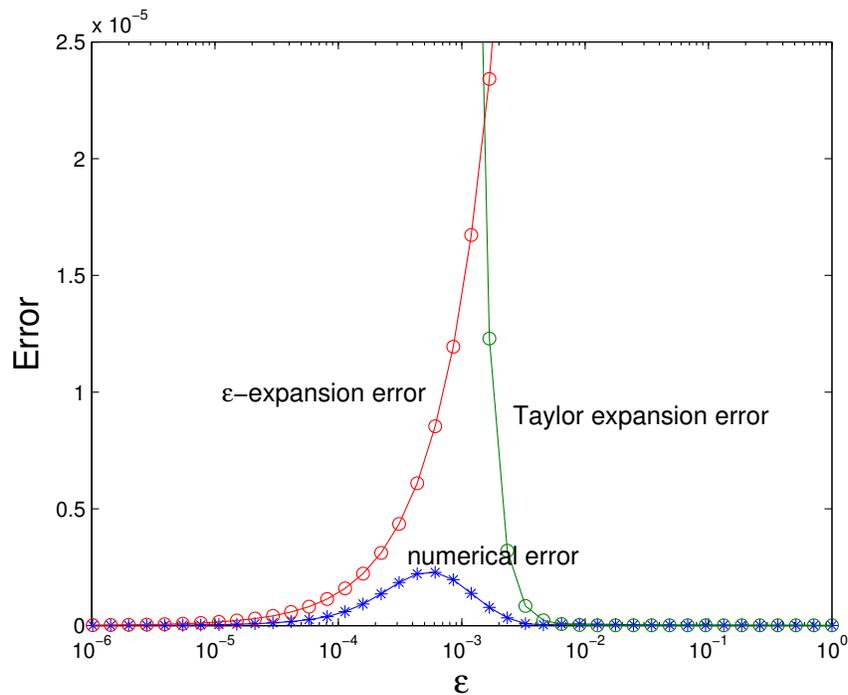
Figure 5.6: Numerical global error (∗) and theoretical global errors (○) versus $\varepsilon$. In this Figure "Taylor expansion error" represents the global error given by the classical error analysis when $\varepsilon \gg h$, and the "$\varepsilon$-expansion error" the global one given by the asymptotic analysis when $\varepsilon \ll h$. The "numerical error" has been computed by a fourth order scheme, i.e. InDC-IMEX2-ARS-GSA-M-k, with M = 4 and k = 1 correction. This plot is consistent with the estimate (5.14) with $p = r = \min(M, s_k) = 4$.

# Acknowledgments

# References

[1] U. M. ASCHER, S. J. RUUTH, AND R. J. SPITERI, *Implicit-explicit Runge-Kutta methods for time-dependent partial differential equations*, Applied Numerical Mathematics, 25 (1997), pp. 151–167.

[2] K. BÖHMER AND H. STETTER, *Defect correction methods. Theory and applications*, (1984).

[3] S. BOSCARINO, *Error analysis of IMEX Runge-Kutta methods derived from differential-algebraic systems*, SIAM Journal on Numerical Analysis, 45 (2008), pp. 1600–1621.

[4] ——, *On an accurate third order implicit-explicit Runge-Kutta method for stiff problems*, Applied Numerical Mathematics, 59 (2009), pp. 1515–1528.

[5] S. BOSCARINO, L. PARESCHI, AND G. RUSSO, *Implicit-explicit Runge–Kutta schemes for hyperbolic systems and kinetic equations in the diffusion limit*, SIAM Journal on Scientific Computing, 35 (2013), pp. A22–A51.

[6] S. BOSCARINO AND J.-M. QIU, *Error estimates of integral deferred correction methods for stiff problems*, ESAIM: Mathematical Modelling and Numerical Analysis, 50 (2016), pp. 1137–1166.

[7] S. BOSCARINO AND G. RUSSO, *On a class of uniformly accurate IMEX Runge-Kutta schemes and applications to hyperbolic systems with relaxation*, SIAM Journal on Scientific Computing, 31 (2010), p. 1926.

[8] S. BOSCARINO AND G. RUSSO, *Flux-explicit IMEX Runge–Kutta schemes for hyperbolic to parabolic relaxation problems*, SIAM Journal on Numerical Analysis, 51 (2013), pp. 163–190.

[9] M. CARPENTER AND C. KENNEDY, *Additive Runge-Kutta schemes for convection-diffusion-reaction equations*, Applied Numerical Mathematics, 44, pp. 139–181.

[10] G. Q. CHEN, C. D. LEVERMORE, AND T.-P. LIU, *Hyperbolic conservation laws with stiff relaxation terms and entropy*, Communications on Pure and Applied Mathematics, 47 (1994), pp. 787–830.

[11] A. CHRISTLIEB, W. GUO, M. MORTON, AND J.-M. QIU, *A high order time splitting method based on integral deferred correction for semi-lagrangian vlasov simulations*, Journal of Computational Physics, 267 (2014), pp. 7–27.

[12] A. CHRISTLIEB, M. MORTON, B. ONG, AND J.-M. QIU, *Semi-implicit integral deferred correction constructed with additive runge–kutta methods*, Commun. Math. Sci, 9 (2011), pp. 879–902.

[13] A. CHRISTLIEB, B. ONG, AND J. QIU, *Comments on high order integrators embedded within integral deferred correction methods*, Communications in Applied Mathematics and Computational Science, 4 (2009), pp. 27–56.

[14] ——, *Integral deferred correction methods constructed with high order Runge-Kutta integrators*, Mathematics of Computation, 79 (2009), p. 761.

[15] J. D. COLE, *On a quasi-linear parabolic equation occurring in aerodynamics*, Quarterly of applied mathematics, 9 (1951), pp. 225–236.

[16] A. DUTT, L. GREENGARD, AND V. ROKHLIN, *Spectral deferred correction methods for ordinary differential equations*, BIT Numerical Mathematics, 40 (2000), pp. 241–266.

[17] E. HAIRER, C. LUBICH, AND M. ROCHE, *Error of Runge-Kutta methods for stiff problems studied via differential algebraic equations*, BIT Numerical Mathematics, 28 (1988), pp. 678–700.

[18] E. HAIRER, S. NØRSETT, AND G. WANNER, *Solving ordinary differential equations: Nonstiff problems*, vol. 1, Springer Verlag, 1993.

[19] E. HAIRER AND G. WANNER, *Solving ordinary differential equations II: stiff and differential algebraic problems*, vol. 2, Springer Verlag, 1993.

[20] J. HUANG, J. JIA, AND M. MINION, *Accelerating the convergence of spectral deferred correction methods*, Journal of Computational Physics, 214 (2006), pp. 633–656.

[21] S. JIN, *Runge-kutta methods for hyperbolic conservation laws with stiff relaxation terms*, Journal of Computational Physics, 122 (1995), pp. 51–67.

[22] A. LAYTON, *On the choice of correctors for semi-implicit picard deferred correction methods*, Applied Numerical Mathematics, 58 (2008), pp. 845–858.

[23] A. LAYTON AND M. MINION, *Implications of the choice of quadrature nodes for picard integral deferred corrections methods for ordinary differential equations*, BIT Numerical Mathematics, 45 (2005), pp. 341–373.

[24] ———, *Implications of the choice of predictors for semi-implicit picard integral deferred corrections methods*, Communications in Applied Mathematics and Computational Science, 1 (2007), pp. 1–34.

[25] T.-P. LIU, *Hyperbolic conservation laws with relaxation*, Communications in Mathematical Physics, 108 (1987), pp. 153–175.

[26] M. MINION, *Semi-implicit spectral deferred correction methods for ordinary differential equations*, Communications in Mathematical Sciences, 1 (2003), pp. 471–500.

[27] R. O'MALLEY JR, *Introduction to singular perturbations. volume 14. applied mathematics and mechanics.*, tech. rep., DTIC Document, 1974.

[28] L. PARESCHI AND G. RUSSO, *Implicit-explicit Runge-Kutta schemes and applications to hyperbolic systems with relaxation*, Journal of Scientific computing, 25 (2005), pp. 129–155.

[29] C.-W. SHU, *Essentially non-oscillatory and weighted essentially non-oscillatory schemes for hyperbolic conservation laws*, Advanced Numerical Approximation of Nonlinear Hyperbolic Equations, (1998), pp. 325–432.

[30] R. D. SKEEL, *A theoretical framework for proving accuracy results for deferred corrections*, SIAM J. Numer. Anal., 19 (1982), pp. 171–196.

[31] A. TIKHONOV, B. VASL'EVA, AND A. SVESHNIKOV, *Differential Equations*, Springer Verlag, 1985.

[32] X. ZHONG, *Additive semi-implicit runge–kutta methods for computing high-speed nonequilibrium reactive flows*, Journal of Computational Physics, 128 (1996), pp. 19–31.