

## ERROR ESTIMATES OF THE INTEGRAL DEFERRED CORRECTION METHOD FOR STIFF PROBLEMS\*

SEBASTIANO BOSCARINO<sup>1</sup> AND JING-MEI QIU<sup>2</sup>

**Abstract.** In this paper, we present error estimates of the integral deferred correction method constructed with stiffly accurate implicit Runge–Kutta methods with a nonsingular matrix  $A$  in its Butcher table representation, when applied to stiff problems characterized by a small positive parameter  $\varepsilon$ . In our error estimates, we expand the global error in powers of  $\varepsilon$  and show that the coefficients are global errors of the integral deferred correction method applied to a sequence of differential algebraic systems. A study of these errors and of the remainder of the expansion yields sharp error bounds for the stiff problem. Numerical results for the van der Pol equation are presented to illustrate our theoretical findings. Finally, we study the linear stability properties of these methods.

**Mathematics Subject Classification.** 65-XX.

Received November 14, 2014. Revised August 27, 2015. Accepted August 28, 2015.

### 1. INTRODUCTION

The deferred correction (DC) method for solving an initial value problem in the form of

$$y'(t) = f(t, y(t)), \quad y(t_0) = y_0 \in \mathbb{R}^N, \quad (1.1)$$

has been investigated intensively [1, 2, 16]. An advantage of the DC method is that one can use a simple numerical method, for instance a first order method, to compute the solution with higher order accuracy. This is accomplished by using a lower order numerical method to solve a series of correction equations during each time step. In each iteration, the order of the method increases. In [6], a new variant of the deferred correction method called the spectral deferred correction (SDC) was proposed. In SDC, a deferred correction procedure is applied to an integral formulation of the error equation in the DC method. It has been shown that the SDC method outperforms DC in many problems with promising numerical results [6]. This is mainly due to the integral formulation of the error equation, as numerical integration is considered to be a more stable and accurate process than numerical differentiation. Moreover, the selection of quadrature nodes plays some role

---

*Keywords and phrases.* Stiff problems, Runge–Kutta methods, integral deferred correction methods, differential algebraic systems.

\* Research supported by Air Force Office of Scientific Research YIP grant FA9550-12-0318, NSF grant DMS-1522777 and DMS-1217008 and University of Houston.

<sup>1</sup> Department of Mathematics and Computer Science, University of Catania, 95125 Catania, Italy. [boscarino@dmf.unict.it](mailto:boscarino@dmf.unict.it)

<sup>2</sup> Department of Mathematics, University of Houston, 77004 Houston, USA. [jingqiu@math.uh.edu](mailto:jingqiu@math.uh.edu).

in the performance of the SDC method [12]. In [6], the quadrature nodes in the proposed SDC method are chosen to be Gauss–Lobatto, Gauss–Radau or Gauss–Legendre points for high order of accuracy. When the quadrature nodes are uniform, the SDC method is called the integral deferred correction (InDC) method. There are various SDC/InDC methods with different implementation strategies, *e.g.* in selecting time integrators in prediction and correction steps [3–5, 10, 11, 13, 14] and in coupling with the Krylov subspace method [10]. Within the InDC framework, it is shown in [4, 5] that if an  $r$ th order integrator is used to solve the error equation, then the accuracy of the scheme increases by  $r$  orders after each correction loop. This analysis has recently been extended in [3] for the InDC method constructed with implicit and semi-implicit integrators. In [4], the InDC method constructed with high order Runge–Kutta (RK) methods has been reformulated as a RK method, whose Butcher tableau has been explicitly constructed.

The main goal of this paper is to study the convergence behavior of the InDC method constructed using implicit RK methods of different orders, when applied to a special class of stiff problems called *singular perturbation problems* (SPPs). A typical SPP has the form

$$\begin{aligned} y'(t) &= f(y(t), z(t)), \\ \varepsilon z'(t) &= g(y(t), z(t)), \end{aligned} \tag{1.2}$$

where  $y$  and  $z$  are vectors in  $\mathbb{R}^N$  with  $N$  being the dimension of the vectors and  $\varepsilon > 0$  is the *stiffness* parameter. We call these vectors the differential component for  $y$  and the algebraic one for  $z$ . Classical books on this subject are [15, 17]. 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

$$\mu(g_z(y, z)) \leq -1, \tag{1.3}$$

in an  $\varepsilon$ -independent neighbourhood of the solution, where  $\mu$  denotes the logarithmic norm with respect to some inner product. From a classical result in SPPs theory, the condition (1.3) guarantees the existence of an  $\varepsilon$ -expansion, whose coefficients are 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 of [8]) for more details. We thus suppose, in our analysis, that the initial values lie on the smooth solution, that  $\varepsilon \ll H$  where  $H$  is the time step size, and that the initial layer is over. In fact, arbitrary initial values introduce an initial layer in the solution. One possible way to overcome this difficulty is simply to ensure that the numerical method resolves the initial layer by taking small step size of  $\mathcal{O}(\varepsilon)$ .

System (1.2) allows us to understand many phenomena observed for very stiff problems. Indeed, in [8] and in the original paper [9], the authors showed that most of the RK methods presented in the literature suffer from the phenomenon of order reduction in the stiff regime. To this aim, we investigate the same phenomenon when it appears in the InDC framework. In the past, such order reduction has been numerically investigated without much theoretical justification [3, 14]. The novelty of this paper is to provide rigorous and careful convergence analysis for the global error of the InDC method and investigate its stability property.

In this paper, we study the global error of the InDC method when it is applied to SPPs in the form of (1.2), in order to seek an understanding on the order reduction phenomenon. First we consider the InDC method constructed with the backward Euler (BE) method, denoted as InDC-BE, and then with implicit RK (IRK) methods, denoted as InDC-IRK.

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. Order reduction phenomenon exists for both differential and algebraic components in the InDC framework. Specifically, under suitable assumptions, the order of convergence for the first term in the  $\varepsilon$ -expansion of global error increases with high order if a high order RK method is applied in the correction steps of the InDC method; whereas the order of convergence for the second term in  $\varepsilon$ -expansion is determined by the stage order of the RK method for the prediction step. We focus our analysis on the InDC

method using uniform quadrature nodes, but excluding the left-most endpoint. The uniform distribution of nodes is important to increase accuracy by the corresponding high order, when a high order RK method is applied in correction steps for classical problems; we refer readers to [5] for details. 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 discuss such stability issues in Section 5. We also remark that important assumptions on the IRK method are that the method is stiffly accurate and has nonsingular matrix  $A$  in its Butcher table representation. We will show that, if these properties are not satisfied, the corresponding InDC method becomes unstable and the numerical solution diverges. A satisfactory explanation of this fact is given in the Appendix.

The paper is organized in the following way. In the rest of this section, we present the basic notations of IRK methods for SPPs in [8] (for more details see [9]). In Section 2, we introduce the InDC-BE method for SPPs (1.2). In Section 3, main theoretical results are stated in the form of two Theorems; numerical evidence supporting these theoretical results are summarized and presented. In Section 4, we prove convergence results for the InDC-BE method. In Section 5, we study the linear stability properties of these InDC methods. Conclusions are given in Section 6. We organize the description of InDC-IRK methods, the  $\varepsilon$ -expansion of the numerical solution, as well as the corresponding error estimates and the estimation of the remainder, into the Appendix for better readability of the paper. Throughout the paper, for classical concepts and convergence results related to RK methods applied to SPPs, we will cite the classical book on the subject [8] (with the Chapter numbering) from time to time.

### 1.1. The IRK method applied to SPPs

In order to get more insight in the convergence estimates of InDC methods, it is useful to consider the convergence results for the RK methods when applied to (1.2). We observe that when the parameter  $\varepsilon$  in system (1.2) is small, the corresponding differential equation is stiff, and when  $\varepsilon$  tends to zero, the differential equations become a differential algebraic system. The corresponding *reduced* system, *i.e.*  $\varepsilon = 0$ , is the differential algebraic equation (DAE)

$$\begin{aligned} y' &= f(y, z), \\ 0 &= g(y, z), \end{aligned} \tag{1.4}$$

whose initial values are *consistent* if  $0 = g(y_0, z_0)$ . We assume that the Jacobian

$$g_z(y, z) \quad \text{is invertible,} \tag{1.5}$$

in a neighbourhood of the solution of (1.4). This assumption guarantees the solvability of (1.4) and that the equation  $g(y, z) = 0$  possesses a locally unique solution  $z = \mathcal{G}(y)$  (Implicit Function Theorem), which inserted into (1.4) gives

$$y' = f(y, \mathcal{G}(y)). \tag{1.6}$$

From now on we assume a Lipschitz condition for  $\mathcal{G}$ . Furthermore, under the assumption (1.5), equation (1.4) is said to be a differential-algebraic equation of index 1. For a definition of the index of differential algebraic problems, we refer to [7, 8].

Now in order to solve system (1.2) we apply an IRK method. This gives

$$\begin{pmatrix} y_{n+1} \\ z_{n+1} \end{pmatrix} = \begin{pmatrix} y_n \\ z_n \end{pmatrix} + h \sum_{i=1}^s b_i \begin{pmatrix} k_{ni} \\ \ell_{ni} \end{pmatrix}, \tag{1.7}$$

where

$$\begin{pmatrix} k_{ni} \\ \varepsilon \ell_{ni} \end{pmatrix} = \begin{pmatrix} f(Y_{ni}, Z_{ni}) \\ g(Y_{ni}, Z_{ni}) \end{pmatrix}, \tag{1.8}$$

and the internal stages are given by

$$\begin{pmatrix} Y_{ni} \\ Z_{ni} \end{pmatrix} = \begin{pmatrix} y_n \\ z_n \end{pmatrix} + h \sum_{j=1}^s a_{ij} \begin{pmatrix} k_{nj} \\ \ell_{nj} \end{pmatrix}. \quad (1.9)$$

Such method is characterized by the coefficient matrix  $A = (a_{ij})$  and vectors  $c = (c_1, \dots, c_s)^T$ ,  $b = (b_1, \dots, b_s)^T$ . They can be represented by a *tableau* in the usual Butcher notation,

$$\left| \begin{array}{c} c \\ \hline A \\ \hline b^T \end{array} \right|. \quad (1.10)$$

The coefficients  $c$  are given by the basic consistency relation  $c_i = \sum_{j=1}^s a_{ij}$ .

We now suppose that the matrix  $A$  is invertible and put  $\varepsilon = 0$ , we obtain by algebraic manipulations from (1.7), (1.8), (1.9) that,

$$\begin{aligned} y_{n+1} &= y_n + h \sum_{i=1}^s b_i f(Y_{ni}, Z_{ni}) \\ z_{n+1} &= R(\infty)z_n + h \sum_{i=1}^s b_i w_{ij} Z_{nj}, \end{aligned} \quad (1.11)$$

where

$$\begin{aligned} Y_{ni} &= y_n + h \sum_{j=1}^s a_{ij} f(Y_{ni}, Z_{ni}) \\ 0 &= g(Y_{ni}, Z_{ni}), \end{aligned} \quad (1.12)$$

with  $R(\infty) = 1 - \sum_{i,j=1}^s b_i w_{ij}$ , where  $R(z)$  is the stability function of the method and  $w_{ij}$  the elements of the inverse of the matrix  $A$ . We note that the numerical solution  $z_{n+1}$  is independent of  $\varepsilon$  and this represents an interesting approach to solve the reduced system (1.4). In general the numerical solutions (1.11) do not lie on the manifold  $g(y, z) = 0$ . Of special importance here is the following definition which will be an important assumption in the next for the analysis.

**Definition 1.1.** An IRK method is called *stiffly accurate* (SA) if  $b^T = e_s^T A$  with  $e_s^T = (0, \dots, 0, 1)$ , *i.e.*, methods for which the numerical solution is identical to the last internal stage.

Now we have a couple of remarks in order here.

**Remark 1.2.** By the non-singularity of the matrix  $A$  and with Definition 1.1, we have  $R(\infty) = 0$  for a SA IRK method. This makes an  $A$ -stable SA IRK method  $L$ -stable. Note that a method is called  $L$ -stable if it is  $A$ -stable and if its stability function  $R(z) \rightarrow 0$  when  $z \rightarrow \infty$ . For details (see Chap. IV.3 in [8]).

**Remark 1.3.** By Definition 1.1, we get for the numerical solutions  $y_{n+1} = Y_{ns}$ , and  $z_{n+1} = Z_{ns}$ , *i.e.* they are identical to the last internal stage of the method. Furthermore, by the second equation in (1.12), we have  $Z_{ni} = \mathcal{G}(Y_{ni})$  and then  $g(y_{n+1}, z_{n+1}) = 0$ , *i.e.* the numerical solutions lie on the manifold and it follows that the numerical solution  $z_{n+1}$  depends on  $y_{n+1}$ , *i.e.*  $z_{n+1} = \mathcal{G}(y_{n+1})$ .

**Remark 1.4.** If the method is stiffly accurate, we say that the numerical solutions of the numerical method (1.11), (1.12), with  $Z_{ni} = \mathcal{G}(Y_{ni})$  and  $z_{n+1} = \mathcal{G}(y_{n+1})$  are identical to the solutions of equation (1.6) with the same Runge–Kutta method, [8].

Now we review the main convergence results of IRK methods for SPPs, for a detailed review we refer the reader to [8, 9]. This result represents the starting point of convergence analysis for InDC methods applied to (1.2).

Under the assumptions of (Thm. 3.8 in Chap. IV.3 in [8]), the global error of an IRK method satisfies the following convergence results

$$y_n - y(t_n) = \mathcal{O}(h^p) + \mathcal{O}(\varepsilon h^{q+1}), \quad z_n - z(t_n) = \mathcal{O}(h^{q+1}).$$

In addition, if the method is stiffly accurate, we have

$$z_n - z(t_n) = \mathcal{O}(h^p) + \mathcal{O}(\varepsilon h^q),$$

where  $p$  is the *classical* order of the method, and  $q$  is the *stage order* of the method, (*i.e.* condition  $C(q)$  of Sect. IV.5 in [9]).

Our idea here is to use the error analysis of IRK methods applied to SPPs obtained in (Chap. VI.3 of [8]), and extend them to the InDC methods. In fact, in order to do that, we perform an asymptotic expansion of smooth solutions of the system (1.2) and similarly for the numerical solutions of an IRK method applied to (1.2). The errors of the  $y$  and  $z$ -component are formally considered as

$$y_n - y(t_n) = \sum_{\nu \geq 0} \varepsilon^\nu (y_{n,\nu} - y_\nu(t_n)), \quad z_n - z(t_n) = \sum_{\nu \geq 0} \varepsilon^\nu (z_{n,\nu} - z_\nu(t_n)), \tag{1.13}$$

where values  $y_\nu(t)$ ,  $z_\nu(t)$  are coefficients of the  $\varepsilon$ -expansion of the smooth solution for (1.2) and  $y_{n,0}$ ,  $z_{n,0}$ ,  $y_{n,1}$ ,  $z_{n,1}, \dots$ , represent the numerical solution of the RK method applied to DAEs of arbitrary order. Furthermore, the first differences  $y_{n,0} - y_0(t_n)$  and  $z_{n,0} - z_0(t_n)$  in the expansion (1.13) are the global errors of the RK method applied to the reduced system (1.4), *i.e.* system of index 1. The other differences for  $\nu > 0$  in (1.13) are related to the numerical solutions of the RK method when applied to the DAEs of higher index. For details, see [8].

## 2. INDC FORMULATIONS APPLIED TO SPPS

In this section, we consider InDC-IRK method for the solution of SPPs written in the form of (1.2). The use of uniform nodes is important for the increase of high order of accuracy, if high order RK methods are used in correction loops. This is related to the concept of “smoothness of the rescaled error vector”, when we apply high order RK methods in correction loops, for more details see [5]. The use of quadrature nodes excluding the left-most endpoint leads to an important stability condition for stiff problems, *i.e.* the method is L-stable if A-stable with  $R(\infty) = 0$ , see [12]. Then, in this paper, we consider the InDC methods with uniform nodes excluding the left-most endpoint.

### 2.1. InDC framework

We consider InDC procedure [6] applied to a SSP,

$$\begin{aligned} y'(t) &= f(y, z), & y(t_0) &= y_0, \\ \varepsilon z'(t) &= g(y, z), & z(t_0) &= z_0. \end{aligned} \tag{2.1}$$

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

$$0 = t_0 < t_1 < t_2 < \dots < t_n < \dots < t_N = T,$$

with the step size  $H$ . Then, each interval  $[t_n, t_{n+1}]$  is discretized again into  $M$  uniform subintervals with quadrature nodes referred to as

$$t_n \doteq \tau_0 < \tau_1 < \dots < \tau_M \doteq t_{n+1}. \tag{2.2}$$

Let  $h = \frac{H}{M}$  be the size of a substep. For simplicity of notation, we assume that  $h$  is constant. 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 remark that the size of time interval  $[t_n, t_{n+1}]$  may vary as the InDC method is a one-step, multi-stage method. We assume the InDC quadrature nodes are uniform, which is a crucial assumption for high order improvement in accuracy, when we apply general high order IRK methods in prediction and correction steps for a classical ODE system (1.1), (see discussions in [5]). We also note that since  $h = \frac{H}{M}$ , we will use  $\mathcal{O}(h^p)$  and  $\mathcal{O}(H^p)$  interchangeably throughout the paper.

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 (2.1) for a single time interval  $[t_n, t_{n+1}]$  with  $m = 1, \dots, M$ . Here superscript (0) is used to denote the prediction step in the InDC method. Let us assume that we build continuous polynomial interpolants  $\hat{y}^{(0)}(t)$  and  $\hat{z}^{(0)}(t)$  interpolating these discrete values. Now we define the error functions

$$e^{(0)}(t) = y(t) - \hat{y}^{(0)}(t), \quad d^{(0)}(t) = z(t) - \hat{z}^{(0)}(t), \quad t \in [t_n, t_{n+1}]. \tag{2.3}$$

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

$$\begin{aligned} \delta^{(0)}(t) &= f(\hat{y}^{(0)}(t), \hat{z}^{(0)}(t)) - (\hat{y}^{(0)})'(t), \\ \rho^{(0)}(t) &= g(\hat{y}^{(0)}(t), \hat{z}^{(0)}(t)) - (\varepsilon \hat{z}^{(0)})'(t). \end{aligned} \tag{2.4}$$

Thus, by subtracting (2.4) from (2.1), the error equations about the error functions (2.3) become

$$\begin{aligned} (e^{(0)})'(t) - \delta^{(0)}(t) &= f(e^{(0)}(t) + \hat{y}^{(0)}(t), d^{(0)}(t) + \hat{z}^{(0)}(t)) - f(\hat{y}^{(0)}(t), \hat{z}^{(0)}(t)), \\ \varepsilon (d^{(0)})'(t) - \rho^{(0)}(t) &= g(e^{(0)}(t) + \hat{y}^{(0)}(t), d^{(0)}(t) + \hat{z}^{(0)}(t)) - g(\hat{y}^{(0)}(t), \hat{z}^{(0)}(t)). \end{aligned} \tag{2.5}$$

A low order numerical method can be used to obtain numerical solutions  $\hat{e}_m^{(0)}$  and  $\hat{d}_m^{(0)}$  at  $\tau_m$  by discretizing the error equations (2.5). Then the numerical solution can be improved as

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

Such correction procedures can be repeated in each local time step  $[t_n, t_{n+1}]$ . 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 (2.5), each correction improves the accuracy of numerical solutions from the previous iteration.

**Remark 2.1** (About notations). In our description of InDC, we let  $y_m, z_m, e_m^{(k)}, d_m^{(k)}$  denote the exact solutions and exact error functions (without hat); and 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  $t = \tau_m$  and use superscript  $(k)$  to denote the prediction ( $k = 0$ ) and correction loops ( $k = 1, \dots$ ). We let  $\bar{\cdot}$  denote the vector on InDC quadrature nodes, for example,  $\bar{y} = (y_1, \dots, y_M)$ .

### 2.2. InDC-BE method

In this subsection, we consider InDC-BE method for the solution of system (2.1). We use uniformly distributed quadrature nodes  $\tau_1, \dots, \tau_M$  given by (2.2) excluding the left-most endpoint.

1. (Prediction step). Use a BE discretization to compute

$$\bar{\hat{y}}^{(0)} = (\hat{y}_1^{(0)}, \dots, \hat{y}_m^{(0)}, \dots, \hat{y}_M^{(0)})$$

as the approximation of the exact solution  $\bar{y} = (y_1, \dots, y_m, \dots, y_M)$  for (2.1) at quadrature nodes  $\tau_1, \dots, \tau_M$ . We make the same for the  $z$ -component. This gives

$$\begin{aligned} \hat{y}_{m+1}^{(0)} &= \hat{y}_m^{(0)} + hf(\hat{y}_{m+1}^{(0)}, \hat{z}_{m+1}^{(0)}), \\ \varepsilon \hat{z}_{m+1}^{(0)} &= \varepsilon \hat{z}_m^{(0)} + hg(\hat{y}_{m+1}^{(0)}, \hat{z}_{m+1}^{(0)}), \end{aligned} \tag{2.6}$$

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

2. (Correction loop). Let  $\hat{y}^{(k-1)}$  and  $\hat{z}^{(k-1)}$  denote the numerical solutions at the  $(k - 1)$ th sequence correction, for  $k = 1, \dots, K$  with  $K$  the number of correction steps.

- (a) Denote the error function at the  $(k - 1)$ th correction by  $e^{(k-1)}(t) = y(t) - \hat{y}^{(k-1)}(t)$ , where  $y(t)$  is the exact solution and  $\hat{y}^{(k-1)}(t)$  is a polynomial of degree  $(M - 1)$  interpolating  $\bar{y}^{(k-1)}$  at quadrature nodes  $\tau_1, \dots, \tau_M$ . Similarly denote  $d^{(k-1)}(t) = z(t) - \hat{z}^{(k-1)}(t)$ . Let  $\delta^{(k-1)}(t)$  and  $\rho^{(k-1)}(t)$  be defined by equation (2.4), but with the upper script (0) replaced with  $(k - 1)$ . We compute the numerical error vector  $\bar{e}^{(k-1)} = (\hat{e}_1^{(k-1)}, \dots, \hat{e}_M^{(k-1)})$  where  $\hat{e}_m^{(k-1)}$  is the approximation of  $e^{(k-1)}(\tau_m)$  by applying a BE method to the integral form of (2.5) with  $\hat{e}_m^{(k-1)}$  approximating  $e^{(k-1)}(\tau_m)$  by applying a BE method to the integral form of (2.5),

$$\begin{aligned} \hat{e}_{m+1}^{(k-1)} &= \hat{e}_m^{(k-1)} + h\Delta f_{m+1}^{(k-1)} + \int_{\tau_m}^{\tau_{m+1}} \delta^{(k-1)}(s)ds, \\ \varepsilon \hat{d}_{m+1}^{(k-1)} &= \varepsilon \hat{d}_m^{(k-1)} + h\Delta g_{m+1}^{(k-1)} + \int_{\tau_m}^{\tau_{m+1}} \rho^{(k-1)}(s)ds, \end{aligned} \tag{2.7}$$

where

$$\begin{aligned} \Delta f_{m+1}^{(k-1)} &= f(\hat{y}_{m+1}^{(k-1)} + \hat{e}_{m+1}^{(k-1)}, \hat{z}_{m+1}^{(k-1)} + \hat{d}_{m+1}^{(k-1)}) - f(\hat{y}_{m+1}^{(k-1)}, \hat{z}_{m+1}^{(k-1)}), \\ \Delta g_{m+1}^{(k-1)} &= g(\hat{y}_{m+1}^{(k-1)} + \hat{e}_{m+1}^{(k-1)}, \hat{z}_{m+1}^{(k-1)} + \hat{d}_{m+1}^{(k-1)}) - g(\hat{y}_{m+1}^{(k-1)}, \hat{z}_{m+1}^{(k-1)}), \end{aligned} \tag{2.8}$$

and

$$\begin{aligned} \int_{\tau_m}^{\tau_{m+1}} \delta^{(k-1)}(s)ds &= \int_{\tau_m}^{\tau_{m+1}} f(\hat{y}^{(k-1)}(s), \hat{z}^{(k-1)}(s))ds - \hat{y}_{m+1}^{(k-1)} + \hat{y}_m^{(k-1)}, \\ \int_{\tau_m}^{\tau_{m+1}} \rho^{(k-1)}(s)ds &= \int_{\tau_m}^{\tau_{m+1}} g(\hat{y}^{(k-1)}(s), \hat{z}^{(k-1)}(s))ds - \varepsilon \hat{z}_{m+1}^{(k-1)} + \varepsilon \hat{z}_m^{(k-1)}. \end{aligned} \tag{2.9}$$

The integral terms  $\int_{\tau_m}^{\tau_{m+1}}$  in equations (2.9) are approximated by a numerical quadrature. Especially, let  $S$  be the integration matrix; its  $(m, k)$  element is

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

where  $\alpha_k(s)$  is the Lagrangian basis function based on the node  $\tau_k$ . Note that  $S^{m,k}$  can be obtained from the computation based on a standard interval  $[0, 1]$ . Let

$$S^m(\bar{f}) = \sum_{j=1}^M S^{m,j} f(y_j, z_j), \tag{2.10}$$

then

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

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

- (b) Update the approximate solutions  $\bar{y}^{(k)} = \bar{y}^{(k-1)} + \bar{e}^{(k-1)}$  and  $\bar{z}^{(k)} = \bar{z}^{(k-1)} + \bar{d}^{(k-1)}$ .

**Remark 2.2.** Using the notation introduced in equation (2.10), we get from equation (2.7) and (2.9),

$$\begin{aligned} \hat{y}_{m+1}^{(k)} &= \hat{y}_m^{(k)} + h\Delta f_{m+1}^{(k-1)} + hS^m(\bar{f}^{(k-1)}), \\ \varepsilon \hat{z}_{m+1}^{(k)} &= \varepsilon \hat{z}_m^{(k)} + h\Delta g_{m+1}^{(k-1)} + hS^m(\bar{g}^{(k-1)}). \end{aligned} \tag{2.11}$$

**Remark 2.3.** Since we consider the nodes excluding the left most quadrature point  $t_0$ , the order of approximation for integration/interpolation will be one order lower than the usual one considered in [5, 6].

**Remark 2.4.** The InDC-BE described above, can be generalized to the InDC-IRK method, for solving SPPs (2.1). To avoid heavy notations from the InDC-IRK method and for a better presentation of the paper, we organize the description of InDC-IRK method and the corresponding error estimates in Appendix.



### 2.3. $\varepsilon$ -asymptotic expansion

In this section, we introduce  $\varepsilon$ -asymptotic expansion of the exact and numerical solution for system (2.1). This  $\varepsilon$ -asymptotic expansion will be useful to study the behavior of the local error for the InDC method.

We are mainly interested in smooth solutions of (2.1) which provide the  $\varepsilon$ -asymptotic expansions for  $t > 0$ , of the form

$$y(t) = \sum_{j=0}^{\infty} y_j(t)\varepsilon^j, \quad z(t) = \sum_{j=0}^{\infty} z_j(t)\varepsilon^j. \tag{2.12}$$

As just pointed out in the introduction, we suppose that the initial values of (2.1) lie on the smooth solution, *i.e.* that an expansion of the form (2.12) holds.

From (2.12) we note that the exact solutions have a power series in  $\varepsilon$  and, considered a truncated series, a remainder after any  $N + 1$  number of terms could be obtained or estimated. In particular, for any  $t \in [0, \bar{t}]$ , the remainder is bounded above by a term  $C_N\varepsilon^{N+1}$  with  $C_N > 0$ , for  $\varepsilon$  small enough, *i.e.*

$$y(t) = \sum_{j=0}^N y_j(t)\varepsilon^j + \mathcal{O}(\varepsilon^{N+1}), \quad z(t) = \sum_{j=0}^N z_j(t)\varepsilon^j + \mathcal{O}(\varepsilon^{N+1}). \tag{2.13}$$

Furthermore we note that a sequence of DAEs arise in the study of (2.1). In fact, the coefficients in the expansion (2.12) are the solutions of DAEs of different indices, for more details (see Chap. VI.3 of [8]). This is obtained by inserting the  $\varepsilon$ -expansion of the exact solution (2.12) into (2.1) and collecting terms of equal powers of  $\varepsilon$ .

$$\varepsilon^0 : \quad \begin{cases} y'_0 = f(y_0, z_0) \\ 0 = g(y_0, z_0) \end{cases}, \tag{2.14}$$

$$\varepsilon^1 : \quad \begin{cases} y'_1 = f_y(y_0, z_0)y_1 + f_z(y_0, z_0)z_1 \doteq \mathbb{F}_1 \\ z'_0 = g_y(y_0, z_0)y_1 + g_z(y_0, z_0)z_1 \doteq \mathbb{G}_1 \end{cases}, \tag{2.15}$$

...

$$\varepsilon^\nu : \quad \begin{cases} y'_\nu = f_y(y_0, z_0)y_\nu + f_z(y_0, z_0)z_\nu + \phi_\nu(y_0, z_0, \dots, y_{\nu-1}, z_{\nu-1}) \doteq \mathbb{F}_\nu \\ z'_{\nu-1} = g_y(y_0, z_0)y_\nu + g_z(y_0, z_0)z_\nu + \psi_\nu(y_0, z_0, \dots, y_{\nu-1}, z_{\nu-1}) \doteq \mathbb{G}_\nu \end{cases}, \tag{2.16}$$

with initial values  $y_\nu(0), z_\nu(0)$  known from (2.12). We observe that system (2.14) under the condition (1.5) is a DAE of index 1. According to [8], if we consider (2.14) and (2.15) together, we have a differential algebraic system of index 2. In general (2.14)–(2.16) is a differential algebraic system of index  $\nu$ .

Now let us look for an  $\varepsilon$ -asymptotic expansion of the numerical solution at the  $k$ th correction step of the InDC-BE method in the form

$$\hat{y}_m^{(k)} = \sum_{\nu=0}^{\infty} \hat{y}_{m,\nu}^{(k)}\varepsilon^\nu, \quad \hat{z}_m^{(k)} = \sum_{\nu=0}^{\infty} \hat{z}_{m,\nu}^{(k)}\varepsilon^\nu. \tag{2.17}$$

The case of  $k = 0$  corresponds to the prediction step of InDC method. Then, inserting the above ansatz (2.17) into the numerical scheme (2.6)–(2.9), and collecting terms of equal powers of  $\varepsilon$ , we have the following:

- for the prediction step ( $k = 0$ )

$$\varepsilon^0 : \quad \begin{cases} \hat{y}_{m+1,0}^{(0)} = \hat{y}_{m,0}^{(0)} + hf(\hat{y}_{m+1,0}^{(0)}, \hat{z}_{m+1,0}^{(0)}), \\ 0 = g(\hat{y}_{m+1,0}^{(0)}, \hat{z}_{m+1,0}^{(0)}), \end{cases} \tag{2.18}$$

$$\varepsilon^1 : \quad \begin{cases} \hat{y}_{m+1,1}^{(0)} = \hat{y}_{m,1}^{(0)} + h\hat{\mathbb{F}}_{m+1,1}^{(0)}, \\ \hat{z}_{m+1,0}^{(0)} = \hat{z}_{m,0}^{(0)} + h\hat{\mathbb{G}}_{m+1,1}^{(0)}, \end{cases} \tag{2.19}$$



where

$$\begin{cases} \hat{\mathbb{F}}_{m+1,1}^{(0)} \doteq f_y(\hat{y}_{m+1,0}^{(0)}, \hat{z}_{m+1,0}^{(0)})\hat{y}_{m+1,1}^{(0)} + f_z(\hat{y}_{m+1,0}^{(0)}, \hat{z}_{m+1,0}^{(0)})\hat{z}_{m+1,1}^{(0)}, \\ \hat{\mathbb{G}}_{m+1,1}^{(0)} \doteq g_y(\hat{y}_{m+1,0}^{(0)}, \hat{z}_{m+1,0}^{(0)})\hat{y}_{m+1,1}^{(0)} + g_z(\hat{y}_{m+1,0}^{(0)}, \hat{z}_{m+1,0}^{(0)})\hat{z}_{m+1,1}^{(0)}, \end{cases} \tag{2.20}$$

- for the correction steps ( $k \geq 1$ ),

$$\varepsilon^0 : \begin{cases} \hat{y}_{m+1,0}^{(k)} = \hat{y}_{m,0}^{(k)} + h\Delta\hat{f}_{m+1,0}^{(k-1)} + hS^m(\bar{f}_0^{(k-1)}), \\ 0 = h\Delta\hat{g}_{m+1,0}^{(k-1)} + hS^m(\bar{g}_0^{(k-1)}), \end{cases} \tag{2.21}$$

$$\varepsilon^1 : \begin{cases} \hat{y}_{m+1,1}^{(k)} = \hat{y}_{m,1}^{(k)} + h\Delta\hat{\mathbb{F}}_{m+1,1}^{(k-1)} + hS^m(\bar{\mathbb{F}}_1^{(k-1)}), \\ \hat{z}_{m+1,0}^{(k)} = \hat{z}_{m,0}^{(k)} + h\Delta\hat{\mathbb{G}}_{m+1,1}^{(k-1)} + hS^m(\bar{\mathbb{G}}_1^{(k-1)}), \end{cases} \tag{2.22}$$

where in (2.21),

$$\begin{cases} \Delta\hat{f}_{m+1,0}^{(k-1)} = f(\hat{y}_{m+1,0}^{(k)}, \hat{z}_{m+1,0}^{(k)}) - f(\hat{y}_{m+1,0}^{(k-1)}, \hat{z}_{m+1,0}^{(k-1)}), \\ \Delta\hat{g}_{m+1,0}^{(k-1)} = g(\hat{y}_{m+1,0}^{(k)}, \hat{z}_{m+1,0}^{(k)}) - g(\hat{y}_{m+1,0}^{(k-1)}, \hat{z}_{m+1,0}^{(k-1)}), \end{cases} \tag{2.23}$$

and in (2.22),

$$\begin{aligned} \Delta\hat{\mathbb{F}}_{m+1,1}^{(k-1)} &= \hat{\mathbb{F}}_{m+1,1}^{(k)} - \hat{\mathbb{F}}_{m+1,1}^{(k-1)} \\ &= \left( f_y(\hat{y}_{m+1,0}^{(k)}, \hat{z}_{m+1,0}^{(k)})\hat{y}_{m+1,1}^{(k)} + f_z(\hat{y}_{m+1,0}^{(k)}, \hat{z}_{m+1,0}^{(k)})\hat{z}_{m+1,1}^{(k)} \right) \\ &\quad - \left( f_y(\hat{y}_{m+1,0}^{(k-1)}, \hat{z}_{m+1,0}^{(k-1)})\hat{y}_{m+1,1}^{(k-1)} + f_z(\hat{y}_{m+1,0}^{(k-1)}, \hat{z}_{m+1,0}^{(k-1)})\hat{z}_{m+1,1}^{(k-1)} \right), \end{aligned} \tag{2.24}$$

where

$$\hat{\mathbb{F}}_{m+1,1}^{(k)} = f_y(\hat{y}_{m+1,0}^{(k)}, \hat{z}_{m+1,0}^{(k)})\hat{y}_{m+1,1}^{(k)} + f_z(\hat{y}_{m+1,0}^{(k)}, \hat{z}_{m+1,0}^{(k)})\hat{z}_{m+1,1}^{(k)}. \tag{2.25}$$

We note that both equations (2.18)–(2.19) for the prediction step ( $k = 0$ ), and equations (2.21)–(2.22) for the correction step ( $k \geq 1$ ), are consistent discretizations of equations (2.14)–(2.15). It is possible to generalize the  $\varepsilon$ -asymptotic expansion to  $\varepsilon^\nu$  ( $\nu \geq 2$ ), but we skip this to avoid heavy notations.

Finally, let  $\varepsilon$ -asymptotic expansion of error functions  $e^{(k)}(t)$ ,  $d^{(k)}(t)$  at the  $k$ th iteration be

$$\begin{pmatrix} e_m^{(k)} \\ d_m^{(k)} \end{pmatrix} = \begin{pmatrix} \sum_{\nu=0}^{\infty} e_{m,\nu}^{(k)}\varepsilon^\nu \\ \sum_{\nu=0}^{\infty} d_{m,\nu}^{(k)}\varepsilon^\nu \end{pmatrix} = \begin{pmatrix} \sum_{\nu=0}^{\infty} (y_{m,\nu} - \hat{y}_{m,\nu}^{(k)})\varepsilon^\nu \\ \sum_{\nu=0}^{\infty} (z_{m,\nu} - \hat{z}_{m,\nu}^{(k)})\varepsilon^\nu \end{pmatrix}. \tag{2.26}$$

Note that in the above ansatz, we consider truncated series of (2.26) with estimate of the remainder as

$$\begin{pmatrix} e_m^{(k)} \\ d_m^{(k)} \end{pmatrix} = \begin{pmatrix} e_{m,0}^{(k)} + e_{m,1}^{(k)}\varepsilon + \dots + e_{m,\nu}^{(k)}\varepsilon^\nu + \mathcal{O}(\varepsilon^{\nu+1}) \\ d_{m,0}^{(k)} + d_{m,1}^{(k)}\varepsilon + \dots + d_{m,\nu}^{(k)}\varepsilon^\nu + \mathcal{O}(\varepsilon^{\nu+1}) \end{pmatrix}, \tag{2.27}$$

where a finite number of terms  $\nu$  are taken. We will see that  $\nu$  is related to the value  $q^{(0)}$ , *i.e.* the stage order of the implicit RK method in the prediction step  $k = 0$ .

In the later part of this paper, our goal is to give rigorous estimates of the coefficients  $e_{m,\nu}^{(k)} = y_{m,\nu} - \hat{y}_{m,\nu}^{(k)}$  and  $d_{m,\nu}^{(k)} = z_{m,\nu} - \hat{z}_{m,\nu}^{(k)}$ , given by (2.12) and (2.17) for  $1 \leq \nu \leq q^{(0)} + 1$  and, finally, estimates of the remainders will be given.

Similarly, we consider the  $\varepsilon$ -asymptotic expansion of numerical approximations of error functions  $\hat{e}^{(k)}(t)$ ,  $\hat{d}^{(k)}(t)$  at the  $k$ th iteration

$$\begin{pmatrix} \hat{e}_m^{(k)} \\ \hat{y}_m^{(k)} \\ \hat{d}_m^{(k)} \end{pmatrix} = \begin{pmatrix} \sum_{\nu=0}^{\infty} \hat{e}_{m,\nu}^{(k)} \varepsilon^\nu \\ \sum_{\nu=0}^{\infty} \hat{y}_{m,\nu}^{(k)} \varepsilon^\nu \\ \sum_{\nu=0}^{\infty} \hat{d}_{m,\nu}^{(k)} \varepsilon^\nu \end{pmatrix} = \begin{pmatrix} \sum_{\nu=0}^{\infty} (\hat{y}_{m,\nu}^{(k+1)} - \hat{y}_{m,\nu}^{(k)}) \varepsilon^\nu \\ \sum_{\nu=0}^{\infty} (\hat{z}_{m,\nu}^{(k+1)} - \hat{z}_{m,\nu}^{(k)}) \varepsilon^\nu \end{pmatrix}. \tag{2.28}$$

Note that combining (2.26) and (2.28), we get with  $k, \nu \geq 0, m = 0, \dots, M$ ,

$$e_{m,\nu}^{(k)} = \hat{e}_{m,\nu}^{(k)} + e_{m,\nu}^{(k+1)}, \quad d_{m,\nu}^{(k)} = \hat{d}_{m,\nu}^{(k)} + d_{m,\nu}^{(k+1)}. \tag{2.29}$$

**Remark 2.5.** Similar  $\varepsilon$ -asymptotic expansions can be given for the numerical solutions of the InDC-IRK method. Again, to avoid heavy notations, we organize them in Appendix.

### 3. MAIN RESULTS AND NUMERICAL EVIDENCE

In this section, we present the main theoretical results in the form of theorems, and provide numerical evidence supporting the main theorems. We will provide a rigorous mathematical proof in the next section.

#### 3.1. Main results

The aim of this section is to present convergence results of the InDC-BE and InDC-IRK method when applied to (2.1).

**Theorem 3.1.** *Consider the stiff system (1.2), (1.3) with initial values  $y(0), z(0)$  admitting a smooth solution. Consider the InDC-BE method constructed with  $M$  uniformly distributed quadrature nodes excluding the left-most point and  $K$  correction steps. Then the global errors after  $K$  correction satisfy,*

$$\begin{aligned} e_n^{(K)} &= \hat{y}_n^{(K)} - y(t_n) = \mathcal{O}(H^{\min\{K+1, M\}}) + \mathcal{O}(\varepsilon H), \\ d_n^{(K)} &= \hat{z}_n^{(K)} - z(t_n) = \mathcal{O}(H^{\min\{K+1, M\}}) + \mathcal{O}(\varepsilon H), \end{aligned} \tag{3.1}$$

for  $\varepsilon \leq cH$  and for any fixed constant  $c > 0$ , where  $H = Mh$  is one InDC time step. The estimates hold uniformly for  $H \leq H_0$  and  $nH \leq \text{Const}$ .

**Theorem 3.2.** *Consider the stiff system (1.2), (1.3) with 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 stiffly accurate IRK method of order  $p^{(0)}$ , stage order  $q^{(0)}$  with  $(q^{(0)} < p^{(0)})$  for the prediction step. Apply IRK methods of different classical orders  $(p^{(1)}, p^{(2)}, \dots, p^{(K)})$  in the correction loops,  $k = 1, \dots, K$ . Assume that each of these IRK methods in the prediction and correction loops are stiffly accurate and the matrices  $A$  are nonsingular. Then the global errors after  $K$  correction loops satisfy the estimates*

$$\begin{aligned} e_n^{(K)} &= \hat{y}_n^{(K)} - y(t_n) = \mathcal{O}(H^{\min\{s_K, M\}}) + \mathcal{O}(\varepsilon H^{q^{(0)}}), \\ d_n^{(K)} &= \hat{z}_n^{(K)} - z(t_n) = \mathcal{O}(H^{\min\{s_K, M\}}) + \mathcal{O}(\varepsilon H^{q^{(0)}}), \end{aligned} \tag{3.2}$$

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  $nH \leq \text{Const}$ .

From the above two theorems, it is observed that the order of convergence for the first terms in (3.1) and (3.2) increases with the correction iteration  $k$ , whereas the order for later terms does not change with the number of corrections  $k$ .

We note that (2.26), but replacing  $m$  with  $n$ , can be adopted to represent the  $\varepsilon$ -asymptotic expansion of the global error functions  $e_n^{(K)}$  and  $d_n^{(K)}$  at the  $K$ th correction step, where  $e_{n,\nu}^{(K)}$  and  $d_{n,\nu}^{(K)}$  for  $\nu = 0, 1, \dots$ , are the global errors of InDC stiffly accurate (SA) IRK method (InDC SA-IRK), applied to the differential algebraic

systems of different indices (2.14)–(2.16). In Section 4, we only prove Theorem 3.1 for estimating  $e_{n,\nu}^{(K)}$  and  $d_{n,\nu}^{(K)}$  with  $\nu = 0, 1$  for the InDC-BE method.

To avoid heavy notations and technical details, we prove Theorem 3.2 for general InDC-IRK methods and estimate the remainder of the expansion (2.27) in Sections A.3 and A.4 in the Appendix.

### 3.2. Numerical evidence

We present some numerical evidence of the estimates given in Theorems 3.1 and 3.2. Below, we consider the following InDC methods constructed with  $M$  quadrature points.

- The InDC-BE method with  $k$  correction steps (InDC-BE-M-k). The BE method has order  $p = 1$  and stage order  $q = 1$ .
- The InDC method constructed with a second order stiffly accurate DIRK method in  $k$  correction steps (InDC-DIRK2-SA-M-k). The second order DIRK method (*DIRK2-SA*) has the Butcher tableau

$$\begin{array}{c|cc} \gamma & \gamma & 0 \\ 1 & 1 - \gamma & \gamma \\ \hline & 1 - \gamma & \gamma \end{array} \tag{3.3}$$

where  $\gamma = 1 - \frac{\sqrt{2}}{2}$ . This method is stiffly accurate with order  $p = 2$  and stage order  $q = 1$ .

- The InDC method constructed with a second order non stiffly accurate midpoint method in  $k$  correction steps (InDC-DIRK2-NSA-M-k). The second order midpoint method (*DIRK2-NSA*) has the Butcher tableau

$$\begin{array}{c|cc} 1/2 & 1/2 & \\ \hline & & 1 \end{array} \tag{3.4}$$

This method is not stiffly accurate, with order  $p = 2$  and stage order  $q = 1$ .

- The InDC method constructed with a second order stiffly accurate Lobatto IIIA method (trapezoidal rule) in  $k$  correction steps (InDC-LobattoIIIA2-M-k). This method has the Butcher tableau

$$\begin{array}{c|cc} 0 & 0 & 0 \\ 1 & 1/2 & 1/2 \\ \hline & 1/2 & 1/2 \end{array} \tag{3.5}$$

It is stiffly accurate with matrix  $A$  singular. It is  $A$ -stable but not  $L$ -stable ( $R(\infty) \neq 0$ ), order  $p = 2$  and stage order  $q = 1$ .

- The InDC method constructed with a third order stiffly accurate Radau IIA method in the prediction step and with the BE method in  $k$  correction steps (InDC-Radau-BE-M-k). The third order Radau IIA method has the Butcher tableau

$$\begin{array}{c|cc} 1/3 & 5/12 & -1/12 \\ 1 & 3/4 & 1/4 \\ \hline & 3/4 & 1/4 \end{array} \tag{3.6}$$

This method is stiffly accurate with order  $p = 3$  and stage order  $q = 2$ .

The indicated order of convergence by Theorems 3.1 and 3.2 for the  $y$  and  $z$  components in the SPPs are summarized in Table 1. Below we discuss the convergence rates specified in Table 1.

- For the InDC-BE-M-k method, the order of convergence will increase with  $k$  for the first error term in equation (2.26) 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 in (3.1). The BE method has stage order  $q = 1$ . The order of convergence for the second error term in equation (2.26) will be determined by the stage order of the prediction  $q^{(0)} = 1$  when  $k$  increases, leading to a term of  $\varepsilon H$  in equation (3.1).
- For the InDC-DIRK2-SA-M-k, the order of convergence will increase with  $k$  by 2 for the first error term in equation (2.26) when  $\varepsilon \ll H$  and  $k \leq M - 1$ , leading to a term of  $H^{\min(2(k+1), M)}$  for the differential and algebraic component in equation (3.1). DIRK2-SA method has stage order  $q = 1$ . The order of convergence

TABLE 1. Global error predicted by Theorem 3.1 and Theorem 3.2 with  $H \gg \varepsilon$ . Note that ‘SA’/‘NSA’ means stiffly accurate/not stiffly accurate.

Method	$y$ -comp	$z$ -comp
InDC-BE-M-k	$H^{\min(k+1,M)} + \varepsilon H$	$H^{\min(k+1,M)} + \varepsilon H$
InDC-DIRK2-SA-M-k	$H^{\min(2(k+1),M)} + \varepsilon H$	$H^{\min(2(k+1),M)} + \varepsilon H$
InDC-DIRK2-NSA-M-k	diverges	diverges
InDC-LobattoIIIA2-SA-M-k	diverges	diverges
InDC-Radau-BE-M-k	$H^{\min(3+k,M)} + \varepsilon H^2$	$H^{\min(3+k,M)} + \varepsilon H^2$

for the second error term in equation (2.26) will be determined by the stage order of the prediction  $q^{(0)} = 1$  when  $k$  increases, leading to a term of  $\varepsilon H$  in equation (3.1).

- An important ingredient, suggested by the analysis, is to require that the methods to be stiffly accurate, *i.e.*  $a_{sj} = b_j$  for  $j = 1, \dots, s$  and that the matrix  $A$  is nonsingular. Such a choice provides a significant benefit for the convergence of the numerical solution, without which the numerical solutions will diverge. For example, if we consider using the second order non stiffly accurate DIRK method in both the prediction and  $k$  correction steps of an InDC framework with  $M$  quadrature points (InDC-DIRK2-NSA-M-k), divergence results are expected (see Fig. 2). Note that in the analysis for InDC-IRK method in the appendix, a satisfactory theoretical explanation of this fact is given. Finally, if we consider using methods with singular matrix  $A$ , as for example the second order Lobatto IIIA method, in both the prediction and  $k$  correction steps of an InDC framework with  $M$  quadrature points (InDC-LobattoIIIA2-M-k, right plot in Fig. 2), again divergence is expected.
- For the InDC-Radau-BE-M-k, the order of convergence will increase with  $k$  by 1 for the first error term in equation (2.26) when  $\varepsilon \ll H$  and  $k \leq M - 1$ , leading to a term of  $H^{\min(3+k,M)}$  for the differential and algebraic component in equation (3.1). Radau IIA method has stage order  $q = 2$ . The order of convergence for the second error term in equation (2.26) will be determined by the stage order of the prediction  $q^{(0)} = 2$  when  $k$  increases, leading to a term of  $\varepsilon H^2$  in equation (3.1).

For numerical verification, we first consider a scalar example [8]

$$\varepsilon z' = -z + \cos(t) \tag{3.7}$$

with the analytical solution

$$z(t) = \frac{\cos(t) + \varepsilon \sin(t)}{1 + \varepsilon^2} + C \exp(-t/\varepsilon),$$

where  $C = z(0) - 1$  is determined by the initial condition. For a consistent initial condition, let  $C = 0$ . This is a good example to investigate the order of convergence for the  $\varepsilon^1$  term in equation (1.13), as the error for  $\varepsilon^0$  is 0. Indeed, for stiff parameter  $\varepsilon = 10^{-6}$  only a region of first order convergence is observed for the BE method, where the global and local error given for the  $z$ -component is  $\mathcal{O}(\varepsilon H)$  (see Cor. 3.10 in [8]). Figure 1 gives the one step error (local error) and global error of BE method, expected  $\mathcal{O}(\varepsilon H)$  is observed. We also test the InDC-DIRK2-NSA-3-1 and InDC-LobattoIIIA2-SA-4-2 method. Numerical results are presented in Figure 2. Divergence results are observed when time step is large compared to  $\varepsilon$  if an InDC-correction is performed.

Now we consider the van der Pol equation [8] with the well-prepared initial data up to  $\mathcal{O}(\varepsilon^3)$

$$\begin{cases} y' = z \\ \varepsilon z' = (1 - y^2)z - y \end{cases}, \quad \begin{cases} y(0) = 2 \\ z(0) = -\frac{2}{3} + \frac{10}{81}\varepsilon - \frac{292}{2187}\varepsilon^2 \end{cases} \tag{3.8}$$

- Numerical results of the InDC-BE-3-2 method are presented in the upper row of Figure 3. The order of convergence for  $\varepsilon^0$  term would increase with the correction loops. The  $\varepsilon^1$  term of error behaves like  $\mathcal{O}(\varepsilon H)$  for both  $y$  and  $z$  components.

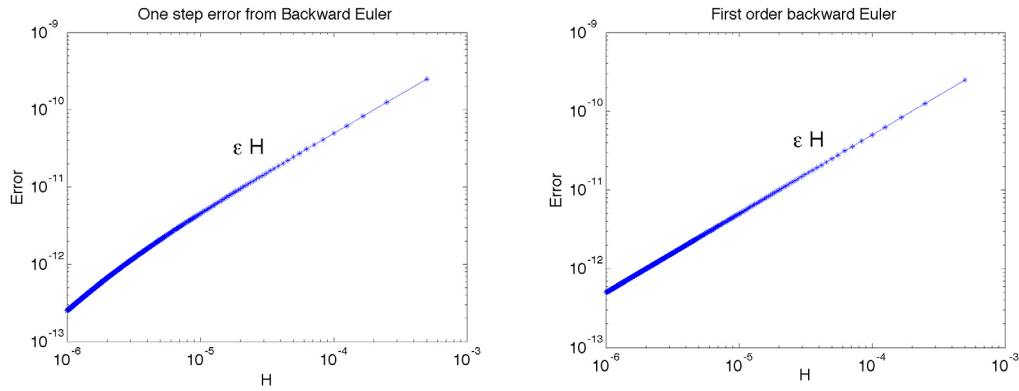


FIGURE 1. Scalar example. Local, *i.e.* one step error (*left plot*) and global error at  $T = 0.5$  (*right plot*) of BE method.  $\mathcal{O}(\epsilon H)$  is observed in both plots with  $\epsilon = 10^{-6}$ .

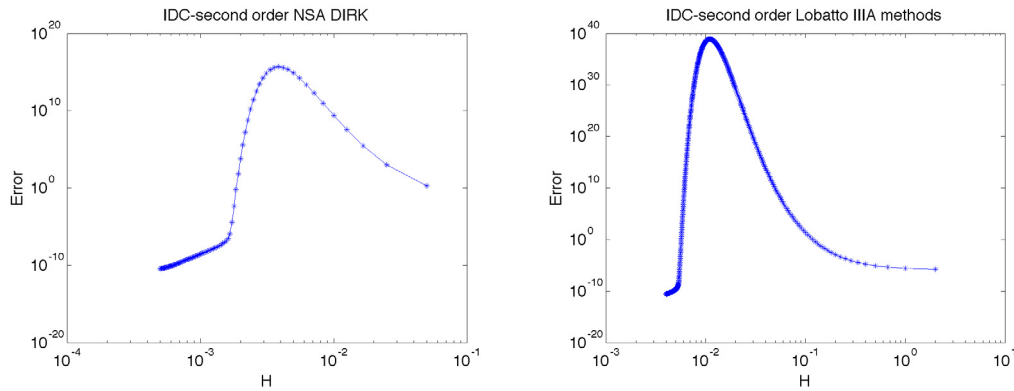


FIGURE 2. Scalar example.  $\epsilon = 10^{-4}$ . *Left*: global error ( $T = 0.1$ ) of the InDC-second order DIRK method that is not SA with three quadrature points and one correction step. *Right*: global error ( $T = 0.1$ ) of the InDC-second order Lobatto IIIA method with matrix  $A$  singular, four quadrature points and one correction step.

- The numerical results of the InDC-DIRK2-SA-4-1 method are presented in the middle row of Figure 3. The order of convergence for  $\epsilon^0$  term would increase with second order with the correction loop. The  $\epsilon^1$  term of error behaves like  $\mathcal{O}(\epsilon H)$  for both  $y$  and  $z$  components.
- The numerical results of the InDC-Radau-BE-6-2 method are presented in the bottom row of Figure 3. The order of convergence for  $\epsilon^0$  term would increase with first order correction loop and is observed to be  $\mathcal{O}(H^5)$ . The  $\epsilon^1$  term of error behaves like  $\mathcal{O}(\epsilon H^2)$  for both  $y$  and  $z$  components.

Numerical observations in Figure 3 are consistent with Theorems 3.1, 3.2 and Table 1. Especially, it is observed that the InDC SA-IRK methods exhibit order reduction both in differential and algebraic components. They produce an estimate for the  $y$  and  $z$  component in the form of equation (3.2). For example, in Figure 3, we observe a behavior like  $e_n^{(k)} = \mathcal{O}(H^3) + \mathcal{O}(\epsilon H)$ . Furthermore, if the step size  $H > \epsilon^{\frac{1}{s_k - q^{(0)}}}$ ,  $\mathcal{O}(H^{s_k})$  is dominant, otherwise the term  $\mathcal{O}(\epsilon H^{q^{(0)}})$  is observed. We observe that in the neighborhood of  $H \approx \epsilon^{\frac{1}{s_k - q^{(0)}}}$ , we have a cancellation of error terms between  $\mathcal{O}(H^{s_k})$  and  $\epsilon \mathcal{O}(H^{q^{(0)}})$ , if error constants are of opposite signs, see for example the plots in middle and bottom rows of Figure 3.

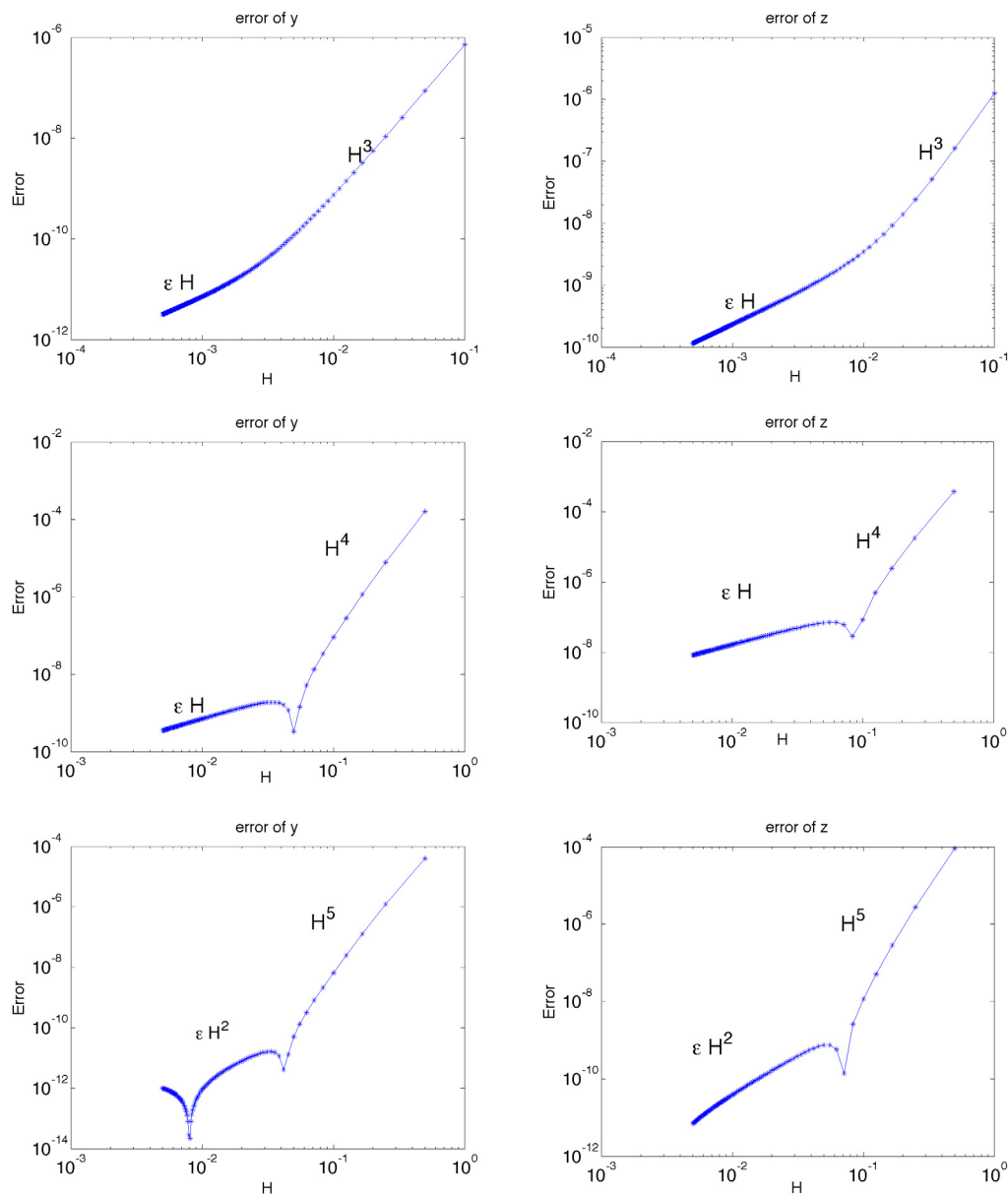


FIGURE 3. Van der Pol equation. Global error ( $T = 0.5$ ) of the InDC-BE-3-2 method (*upper row*); and of the InDC-DIRK2-SA-4-1 method (*middle row*); and of the InDC-Radau-BE-6-2 method (*bottom row*).  $\varepsilon = 10^{-6}$ .

#### 4. PROOFS OF MAIN RESULTS

In this section, we prove Theorem 3.1, which is a special case of Theorem 3.2. By proving Theorem 3.1 through several lemmas, we demonstrate the basic ingredients of the general proof for Theorem 3.2 presented in the Appendix. Our error estimates are based on the  $\varepsilon$ -expansion outlined in Section 2.3.

### 4.1. Error estimates for Theorem 3.1

We perform local error estimate for Theorem 3.1 by two Lemmas. We again note that since  $h = \frac{H}{M}$ , we use  $\mathcal{O}(h^p)$  and  $\mathcal{O}(H^p)$  interchangeably below in our proof. We then prove the global error estimate based on the two Lemmas.

**Lemma 4.1.** [ $\varepsilon^0$  error term] *Let us assume that the reduce system (1.4) with  $\varepsilon = 0$  satisfies (1.5) and that the initial values are consistent. Consider the InDC-BE method constructed with  $M$  uniformly distributed quadrature nodes excluding the left-most point and  $k$  correction steps, i.e. (2.18) for the prediction step and (2.21) for the correction one, with  $k = 1, \dots, K$ . Then the numerical solutions satisfy the following local error estimates at each interior node of InDC  $\tau_m$  with  $m = 0, \dots, M$ ,*

$$e_{m,0}^{(k)} = y_{m,0} - \hat{y}_{m,0}^{(k)} = \mathcal{O}(h^{\min(k+2, M+1)}), \quad d_{m,0}^{(k)} = z_{m,0} - \hat{z}_{m,0}^{(k)} = \mathcal{O}(h^{\min(k+2, M+1)}), \quad (4.1)$$

with

$$g(\hat{y}_{m,0}^{(k)}, \hat{z}_{m,0}^{(k)}) = 0, \quad (4.2)$$

for  $k = 0, \dots, K$ .

*Proof.* For  $k = 0$ , equation (4.2) is a consequence of the consistency of the initial conditions for system (2.14). Then we start to prove the local error estimate (4.1) for the prediction step ( $k = 0$ ). For the exact solution, by equation (2.14) and assumption (1.5), we have for the  $y_0(t)$  component, equation (1.6) and  $g(y_0(t), z_0(t)) = 0$ . By equation (1.5), it follows that  $z_0(t) = \mathcal{G}(y_0(t))$ .

For the numerical solution, we have (2.18). By  $g(\hat{y}_{m,0}^{(0)}, \hat{z}_{m,0}^{(0)}) = 0$ , with  $m = 0, \dots, M$ , we get  $\hat{z}_{m,0}^{(0)} = \mathcal{G}(\hat{y}_{m,0}^{(0)})$  with  $\hat{y}_{m,0}^{(0)}$  being numerical solution of the ordinary differential equation (1.6). Then from classical error estimates for the BE method, we have for the local truncation error  $|\hat{y}_{m,0}^{(0)} - y_{m,0}| \leq C_m h^2$  with  $m = 0, \dots, M$ , for some constant  $C_m$  independent of  $H$ . Therefore,  $|\hat{y}_{m,0}^{(0)} - y_{m,0}| = \mathcal{O}(h^2)$  and by  $\hat{z}_{m,0}^{(0)} = \mathcal{G}(\hat{y}_{m,0}^{(0)})$  and the Lipschitz condition of  $\mathcal{G}$ , it follows that  $|\hat{z}_{m,0}^{(0)} - z_{m,0}| = \mathcal{O}(h^2)$  with  $m = 0, \dots, M$ .

Now we prove the local error estimate (4.1) and equation (4.2) for the correction step  $k = 1$ , assuming a fixed  $M \geq 1$ . By  $g(\hat{y}_{m,0}^{(0)}, \hat{z}_{m,0}^{(0)}) = 0$  in the prediction step, from the second equation in (2.21), we obtain  $g(\hat{y}_{m,0}^{(1)}, \hat{z}_{m,0}^{(1)}) = 0$ , with  $m = 0, \dots, M$ , i.e. equation (4.2) with  $k = 1$ . Then, from the condition (1.5) it follows  $\hat{z}_{m,0}^{(1)} = \mathcal{G}(\hat{y}_{m,0}^{(1)})$ , and this gives from (2.21)

$$\hat{y}_{m+1,0}^{(1)} = \hat{y}_{m,0}^{(1)} + h(\hat{f}(\hat{y}_{m+1,0}^{(1)}) - \hat{f}(\hat{y}_{m+1,0}^{(0)})) + hS^m(\bar{f}_0^{(0)}), \quad (4.3)$$

where  $\hat{f}(\hat{y}_{m+1,0}^{(1)}) = f(\hat{y}_{m+1,0}^{(1)}, \mathcal{G}(\hat{y}_{m+1,0}^{(1)}))$ , and  $S^m(\bar{f}_0^{(0)}) = S^m(\hat{f}(\bar{y}_0^{(0)}, \mathcal{G}(\bar{y}_0^{(0)})))$ . The method (4.3) for updating  $\hat{y}_{m,0}^{(1)}$  represents the first correction step of the InDC-BE method to solve the non-stiff ordinary differential equation (1.6). Therefore, from classical error estimates of InDC-BE method when applied to a non-stiff ordinary differential equation in [5], we have  $|y_{m,0} - \hat{y}_{m,0}^{(1)}| \leq C_m h^3$  for some constant  $C_m$  independent of  $h$  with  $h \leq h_0$ . Therefore  $|y_{m,0} - \hat{y}_{m,0}^{(1)}| = \mathcal{O}(h^3)$  and by  $\hat{z}_{m,0}^{(1)} = \mathcal{G}(\hat{y}_{m,0}^{(1)})$  and Lipschitz condition of  $\mathcal{G}$ , we get  $|z_{m,0} - \hat{z}_{m,0}^{(1)}| = \mathcal{O}(h^3)$ ,  $\forall m = 1, \dots, M$  and  $h \leq h_0$ . The estimate for general  $k > 1$  can be proved in a similar fashion and by mathematical induction with respect to  $k$ .  $\square$

**Lemma 4.2** ( $\varepsilon^1$  error term). *Assume condition (1.3) holds and initial values of the differential algebraic system (2.14)–(2.15) are consistent. Consider the InDC-BE method constructed with  $M$  uniformly distributed quadrature nodes excluding the left-most point, and with (2.18)–(2.19) for the prediction step and (2.21)–(2.22) for the correction step with  $k = 1, \dots, K$  for solving the differential algebraic system (2.14)–(2.15). Then the local error estimates of the InDC-BE method*

$$e_{m,1}^{(k)} = y_{m,1} - \hat{y}_{m,1}^{(k)} = \mathcal{O}(h^2), \quad d_{m,1}^{(k)} = z_{m,1} - \hat{z}_{m,1}^{(k)} = \mathcal{O}(h), \quad (4.4)$$

hold for  $m = 1, \dots, M$  at the interior nodes of InDC, and for  $k = 0, \dots, K$ .



*Proof.* The proof for the case of  $k = 0$  (prediction step) is a consequence of (Lem. 4.4 in Chap. VII. 4 in [8]). We then consider the first correction step with  $k = 1$  and assume a fixed  $M \geq 1$ . We prove (4.4) by mathematical induction w.r.t.  $m$ . Especially, we know  $e_{m,1}^{(1)} = d_{m,1}^{(1)} = 0$ , with  $m = 0$ . We assume (4.4) is valid for  $0, \dots, m$ . We will prove that (4.4) is valid for  $m + 1$ . The integration of (2.15) over  $[\tau_m, \tau_{m+1}]$  gives

$$\varepsilon^1 : \begin{cases} y_{m+1,1} = y_{m,1} + \int_{\tau_m}^{\tau_{m+1}} \mathbb{F}_1(\tau) d\tau, \\ z_{m+1,0} = z_{m,0} + \int_{\tau_m}^{\tau_{m+1}} \mathbb{G}_1(\tau) d\tau, \end{cases} \tag{4.5}$$

with  $\mathbb{F}_1$  and  $\mathbb{G}_1$  defined in (2.15). We consider now

$$e_{m+1,1}^{(1)} = y_{m+1,1} - \hat{y}_{m+1,1}^{(1)}, \quad d_{m+1,1}^{(1)} = z_{m+1,1} - \hat{z}_{m+1,1}^{(1)}, \tag{4.6}$$

*i.e.* the difference between the exact and numerical solution at  $\tau_{m+1}$ . From (2.24), as well as from the estimates (4.1) in Lemma 4.1, we have

$$\Delta \hat{\mathbb{F}}_{m+1,1}^{(k-1)} = f_y \hat{e}_{m+1,1}^{(k-1)} + f_z \hat{d}_{m+1,1}^{(k-1)} + \mathcal{O}(h^{k+1}). \tag{4.7}$$

Here we used the abbreviations  $f_y = f_y(y_{m+1,0}, z_{m+1,0})$  and similarly for  $f_z$ . Equally, we have

$$\Delta \hat{\mathbb{G}}_{m+1,1}^{(k-1)} = g_y \hat{e}_{m+1,1}^{(k-1)} + g_z \hat{d}_{m+1,1}^{(k-1)} + \mathcal{O}(h^{k+1}). \tag{4.8}$$

Then from (4.7) and (4.8) for  $k = 1$  it follows

$$\begin{cases} \Delta \hat{\mathbb{F}}_{m+1,1}^{(0)} = (f_y \hat{e}_{m+1,1}^{(0)} + f_z \hat{d}_{m+1,1}^{(0)}) + \mathcal{O}(h^2), \\ \Delta \hat{\mathbb{G}}_{m+1,1}^{(0)} = (g_y \hat{e}_{m+1,1}^{(0)} + g_z \hat{d}_{m+1,1}^{(0)}) + \mathcal{O}(h^2). \end{cases} \tag{4.9}$$

Now subtracting equation (2.22) from equation (4.5) this gives

$$\varepsilon^1 : \begin{cases} e_{m+1,1}^{(1)} = e_{m,1}^{(1)} - h \Delta \hat{\mathbb{F}}_{m+1,1}^{(0)} - h S^m(\bar{\mathbb{F}}_1^{(0)}) + \int_{\tau_m}^{\tau_{m+1}} \mathbb{F}_1(\tau) d\tau, \\ d_{m+1,0}^{(1)} = d_{m,0}^{(1)} - h \Delta \hat{\mathbb{G}}_{m+1,1}^{(0)} - h S^m(\bar{\mathbb{G}}_1^{(0)}) + \int_{\tau_m}^{\tau_{m+1}} \mathbb{G}_1(\tau) d\tau. \end{cases} \tag{4.10}$$

On the right-hand side of the equations in (4.10) we add and subtract the following quantities:  $h S^m(\bar{\mathbb{F}}_1)$  and  $h S^m(\bar{\mathbb{G}}_1)$ , these are the integrals of  $(M - 1)$ th degree interpolating polynomials on  $(\tau_m, \mathbb{F}_1(\tau_m))_{m=1}^M$  and  $(\tau_m, \mathbb{G}_1(\tau_m))_{m=1}^M$  over the subinterval  $[\tau_m, \tau_{m+1}]$ , hence they are accurate to the order  $\mathcal{O}(h^{M+1})$  locally, *i.e.*  $\int_{\tau_m}^{\tau_{m+1}} \mathbb{F}_1(\tau) d\tau - h S^m(\bar{\mathbb{F}}_1) = \mathcal{O}(h^{M+1})$ . By the local error estimates in Lemma 4.1, as well as equation (4.4) for  $k = 0$ , it follows that  $S^m(\bar{\mathbb{F}}_1) - S^m(\tilde{\mathbb{F}}_1)$  and  $S^m(\bar{\mathbb{G}}_1) - S^m(\tilde{\mathbb{G}}_1)$  are accurate to the order  $\mathcal{O}(h)$ . Thus, from (4.10) we get

$$\begin{cases} e_{m+1,1}^{(1)} = e_{m,1}^{(1)} - h \left( f_y \hat{e}_{m+1,1}^{(0)} + f_z \hat{d}_{m+1,1}^{(0)} \right) + \mathcal{O}(h^2), \\ d_{m+1,0}^{(1)} = d_{m,0}^{(1)} - h \left( g_y \hat{e}_{m+1,1}^{(0)} + g_z \hat{d}_{m+1,1}^{(0)} \right) + \mathcal{O}(h^2). \end{cases} \tag{4.11}$$

Now from (2.29) and (4.1), we have

$$\begin{cases} \hat{e}_{m,1}^{(0)} = \hat{y}_{m,1}^{(1)} - \hat{y}_{m,1}^{(0)} = e_{m,1}^{(0)} - e_{m,1}^{(1)} = -e_{m,1}^{(1)} + \mathcal{O}(h^2), \\ \hat{d}_{m,1}^{(0)} = \hat{z}_{m,1}^{(1)} - \hat{z}_{m,1}^{(0)} = d_{m,1}^{(0)} - d_{m,1}^{(1)} = -d_{m,1}^{(1)} + \mathcal{O}(h), \end{cases} \tag{4.12}$$

and put it into equation (4.11) gives,

$$\begin{cases} e_{m+1,1}^{(1)} = e_{m,1}^{(1)} + h \left( f_y e_{m+1,1}^{(1)} + f_z d_{m+1,1}^{(1)} \right) + \mathcal{O}(h^2), \\ d_{m+1,0}^{(1)} = d_{m,0}^{(1)} + h \left( g_y e_{m+1,1}^{(1)} + g_z d_{m+1,1}^{(1)} \right) + \mathcal{O}(h^2). \end{cases} \tag{4.13}$$

Now using the estimate (4.1) about  $d_{m,0}^{(1)}$ , from the second equation in (4.13) we obtain

$$d_{m+1,1}^{(1)} = -g_z^{-1} g_y e_{m+1,1}^{(1)} + \mathcal{O}(h), \tag{4.14}$$

with the invertibility of  $g_z$ . Inserting this into the first equation in (4.13) gives

$$e_{m+1,1}^{(1)} = (1 - h(f_y - f_z g_z^{-1} g_y))^{-1} e_{m,1}^{(1)} + \mathcal{O}(h^2). \tag{4.15}$$

Finally  $e_{m+1,1}^{(1)} = \mathcal{O}(h^2)$  follows from (4.15), and  $d_{m+1,1}^{(1)} = \mathcal{O}(h)$  from (4.14). The estimate for general  $k > 1$  can be proved in a similar fashion and by mathematical induction with respect to  $k$ .  $\square$

**Remark 4.3.** In [4], the InDC method constructed with explicit RK methods in the prediction and correction steps has been reformulated as a high-order explicit RK method whose Butcher tableau is explicitly constructed. Similarly, the InDC-BE can be viewed as an IRK method with the corresponding Butcher tableau. Below we present the Butcher tableau for the InDC-BE method with one loop of correction step. This takes the form

$$\begin{array}{c|cc} \mathbf{c} & T & Z \\ \mathbf{c} & P & T \\ \hline \mathbf{b}_1^T & \mathbf{b}_2^T & \end{array} \tag{4.16}$$

where  $\mathbf{c} = \frac{1}{M} [1, \dots, M]^T$ ,  $Z$  is a  $M \times M$  matrix of zeros,  $T$  and  $P$  are  $M \times M$  matrices, with

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

$$P = \begin{bmatrix} (\tilde{S}_{11} - \frac{1}{M}) & \tilde{S}_{12} & \dots & \tilde{S}_{1,M-1} & \tilde{S}_{1,M} \\ (\tilde{S}_{21} - \frac{1}{M}) & (\tilde{S}_{22} - \frac{1}{M}) & \dots & \tilde{S}_{2,M-1} & \tilde{S}_{2,M} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ (\tilde{S}_{M,1} - \frac{1}{M}) & (\tilde{S}_{M,2} - \frac{1}{M}) & \dots & (\tilde{S}_{M,M-1} - \frac{1}{M}) & (\tilde{S}_{M,M} - \frac{1}{M}) \end{bmatrix},$$

where the term  $\tilde{S}_{ij} = \int_{t_0}^{t_i} \alpha_j(s) ds$  with  $\alpha_j(s)$  the Lagrangian basis functions for the node  $\tau_j$ , and the vector

$$\mathbf{b}_1^T = \left( \left( \tilde{S}_{M,1} - \frac{1}{M} \right), \left( \tilde{S}_{M,2} - \frac{1}{M} \right), \dots, \left( \tilde{S}_{M,M} - \frac{1}{M} \right) \right), \quad \mathbf{b}_2^T = \frac{1}{M} (1, 1, \dots, 1).$$

Now from remark 4.3 the following proposition follows.

**Proposition 4.4.** *The InDC-BE method with  $K$  correction steps is an implicit stiffly accurate IRK method with an invertible matrix  $A$  in the Butcher tableau (1.10). Especially when  $K = 1$ , we get*

$$A = \begin{pmatrix} T & Z \\ P & T \end{pmatrix}. \tag{4.17}$$

**Remark 4.5.** In the estimates in Lemma 4.2, we show that there is no improvement for  $e_{m,1}^{(k)}$  and  $d_{m,1}^{(k)}$  as  $k$  increases, see equation (4.4). This is consistent with our numerical evidences presented in the previous section. The reason is that *both* the local and global error for the  $z$ -component in the prediction and correction steps is of first order. This sets the bottleneck for the order increase in the second equation of (4.11).

We are now in the position to prove Theorem 3.1 by the local error estimates of the two lemmas above.

*Proof of Theorem 3.1.* Our first step here is to estimate  $e_{n,0}^{(K)}$  and  $d_{n,0}^{(K)}$ . For this, from Lemma 4.1 we have after one step from  $t_0$  to  $t_1$ , the local error estimate

$$y_0(t_1) - \hat{y}_{M,0}^{(K)} = \mathcal{O}(H^{\min(K+2, M+1)}), \quad (4.18)$$

with  $m = M$  and  $\tau_M = t_1$  in equation (4.1). In the estimate of the global error from local error, we obtain

$$e_{n,0}^{(K)} = y_0(nH) - \hat{y}_{n,0}^{(K)} = \mathcal{O}(H^{\min(K+1, M)}).$$

It thus follows from (4.2), and by the Lipschitz condition of  $\mathcal{G}$ , that

$$d_{n,0}^{(K)} = z_0(nH) - \hat{z}_{n,0}^{(K)} = \mathcal{O}(H^{\min(K+1, M)}).$$

Now our next aim is to estimate  $e_{n,1}^{(K)}$  and  $d_{n,1}^{(K)}$ . From Lemma 4.2, we have for the local error estimate

$$y_1(t_1) - \hat{y}_{M,1}^{(k)} = \mathcal{O}(H^2). \quad (4.19)$$

By Lemma 4.2, the proof of the global error estimates for  $y$  and  $z$  is similar to that of (Thm. 4.5 and 4.6 in Chap. VII. 4 of [8]). Thus we obtain

$$e_{n,1}^{(K)} = y_1(nH) - \hat{y}_{n,1}^{(k)} = \mathcal{O}(H), \quad d_{n,1}^{(K)} = z_1(nH) - \hat{z}_{n,1}^{(k)} = \mathcal{O}(H),$$

which proves the statement.  $\square$

**Remark 4.6.** Similar error estimates can be given for the InDC SA-IRK method. We present and prove these error estimates in Appendix.

## 5. STABILITY PROPERTIES

One important aspect of stability of numerical integrators can be visualized by the stability region [8] in a complex plane around the origin. For implicit methods discussed in Section 3, we plot their stability regions in Figure 4. In these plots, the region outside the bounded domains are the stability regions. Note that in these methods, the left-most quadrature point is always excluded in the construction of the InDC method. The following observations can be made.

1. In general, as more InDC corrections are performed, the stability region shrinks.
2. The InDC-BE method appears to be A-stable with  $M$  quadrature nodes and with up to  $M - 1$  correction loops for  $M = 4$  and  $M = 6$ . For the other  $M$ 's, it appears that such conclusion still holds in our tests for  $M \leq 8$ . Note that such methods can be viewed as diagonally implicit RK methods.
3. The InDC-DIRK2-SA method appears to be  $A(\alpha)$ -stable with one and two correction loops, with the angle  $\alpha$  decreasing as more correction loops are taken.
4. Among three different implicit RK methods, the InDC method constructed with the BE method appears to have the largest stability region.

Now we prove the following proposition. We notice that a similar result on the InDC method using the first order BE scheme was established earlier in [12].

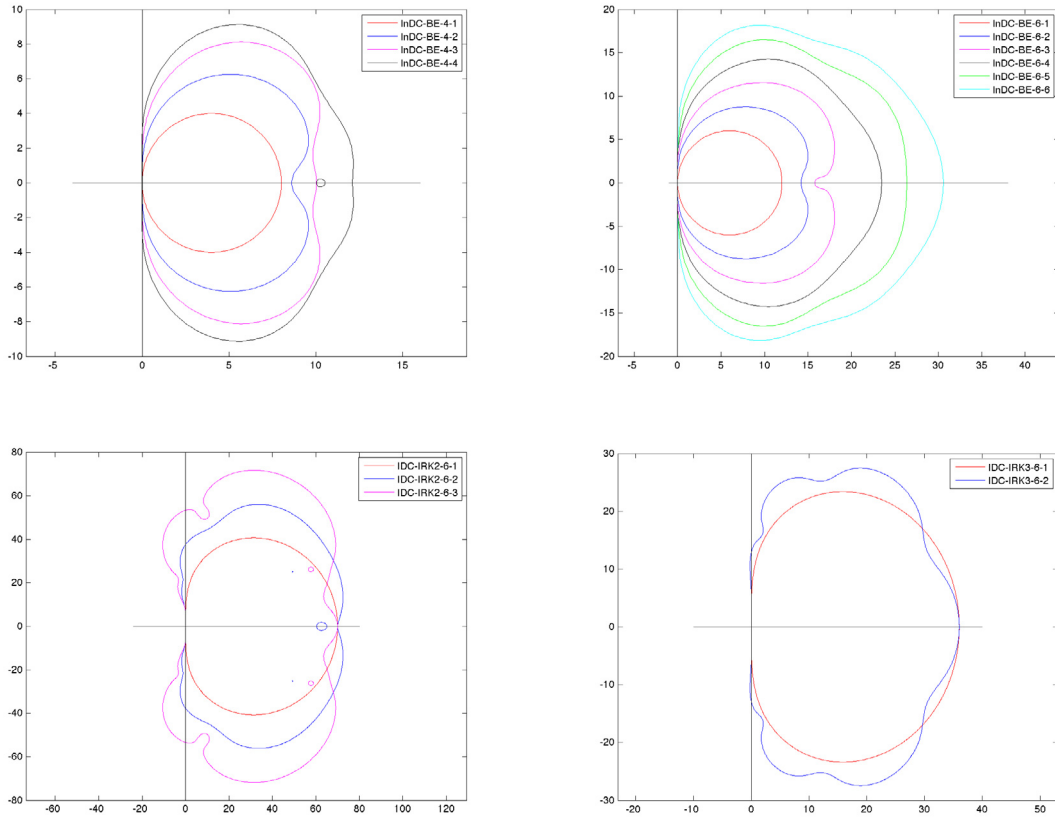


FIGURE 4. Stability regions for InDC-BE method with 4 quadrature points (*upper left*) and 6 quadrature points (*upper right*) with various correction loops as indicated in the legend. Lower left plot shows the stability region for the InDC-IRK2-SA method with 6 quadrature points with zero, one and two correction loops as indicated in the legend. Lower right plot shows the stability region for the InDC-Radau method with 6 quadrature points with zero and one correction loops as indicated in the legend.

**Proposition 5.1.** *Let  $\mathcal{R}(z)$  be the stability function of the InDC method constructed by a stiffly accurate IRK method with nonsingular matrix  $A$  in the corresponding Butcher table. We assume  $M$  uniform quadrature points, but excluding the left-most point, are used. Then  $\lim_{|z| \rightarrow \infty} \mathcal{R}(z) = 0$ . That is, the method is  $L$ -stable, if  $A$ -stable.*

*Proof.* We consider the linear scalar problem  $y' = \lambda y$  with  $z = \lambda h$  and  $y(t = 0) = 1$ . For a stiffly accurate IRK method with nonsingular matrix  $A$  the numerical solution for this linear scalar problem is equal to the last internal stage and the corresponding stability function is

$$R(z) = \mathbf{e}^T (I - Az)^{-1} \mathbf{1} \tag{5.1}$$

with  $\mathbf{e} = (0, \dots, 0, 1)^T$  and  $\mathbf{1} = (1, 1, \dots, 1)^T$  such that for  $z \rightarrow \infty$ , we get  $R(\infty) = 0$ , [8]. Let  $\mathcal{R}_m^{(k)}(z)$  be the amplification factor of the InDC method in the  $k$ th iteration at the  $m$ -th quadrature point. In the prediction step,

$$\mathcal{R}_m^{(0)}(z) = \left( R \left( \frac{z}{M} \right) \right)^m .$$

Hence  $\lim_{|z| \rightarrow \infty} \mathcal{R}_m^{(0)}(z) = 0$ , for  $m = 1, \dots, M$ , as for the IRK method  $\lim_{|z| \rightarrow \infty} R\left(\frac{z}{M}\right) = 0$ . Let  $\mathcal{R}^{(0)} = (\mathcal{R}_1^{(0)}, \dots, \mathcal{R}_M^{(0)})$ , then  $\lim_{|z| \rightarrow \infty} \mathcal{R}^{(0)}(z) = \mathbf{0}$ , where  $\mathbf{0}$  is a zero vector. Note that this is true only if the left-most point is excluded.

In the first correction step, on the first subinterval  $[0, \Delta t/M]$ , the amplification factor of the updated solution at the IRK intermediate stages can be represented as a vector  $\mathbf{r}_1^{(1)}$  with the length of the vector being  $s$ , the stage number of the IRK method. Then, with the help of Butcher table notation in Remark 4.3,

$$\mathbf{r}_1^{(1)}(z) = \mathbf{1} + \frac{z}{M}A(\mathbf{r}_1^{(1)} - P\mathcal{R}^{(0)}) + zS\mathcal{R}^{(0)} = \mathbf{1} + \frac{z}{M}A\mathbf{r}_1^{(1)} + z\left(-\frac{AP}{M} + S\right)\mathcal{R}^{(0)}, \quad \text{with } \mathbf{1} = (1, \dots, 1, 1)'. \tag{5.2}$$

Here we let  $P$  and  $S$  are interpolation and integration matrices of size  $s \times M$ ; they are coefficients that maps the  $M$  function values at quadrature nodes to approximate function values at  $s$  IRK intermediate stages  $(c_i/M, i = 1, \dots, s)$  and over  $s$  integrals  $([0, c_i/M], i = 1, \dots, s)$ . Hence,

$$\mathbf{r}_1^{(1)}(z) = \left(I - \frac{z}{M}A\right)^{-1} \mathbf{1} + \left(I - \frac{z}{M}A\right)^{-1} z\left(-\frac{AP}{M} + S\right)\mathcal{R}^{(0)}.$$

Since the IRK method is stiffly accurate, then

$$\begin{aligned} \mathcal{R}_1^{(1)}(z) &= \mathbf{e}^T \cdot \mathbf{r}_1^{(1)}(z), \quad \text{with } \mathbf{e} = (0, \dots, 0, 1)^T \\ &= \mathbf{e}^T \left(I - \frac{z}{M}A\right)^{-1} \mathbf{1} + \mathbf{e}^T \cdot \left(I - \frac{z}{M}A\right)^{-1} z\left(-\frac{AP}{M} + S\right)\mathcal{R}^{(0)} \\ &\stackrel{(5.1)}{=} R(z_i) + \mathbf{e}^T \cdot \left(I - \frac{z}{M}A\right)^{-1} z\left(-\frac{AP}{M} + S\right)\mathcal{R}^{(0)}. \end{aligned} \tag{5.3}$$

Hence

$$\lim_{|z| \rightarrow \infty} \mathcal{R}_1^{(1)}(z) = \lim_{|z_i| \rightarrow \infty} R(z_i) + \mathbf{e} \lim_{|z| \rightarrow \infty} \left( \left(I - \frac{z}{M}A\right)^{-1} z\right) \left(-\frac{AP}{M} + S\right) \lim_{|z| \rightarrow \infty} \mathcal{R}_0 = 0,$$

since  $\lim_{|z| \rightarrow \infty} \mathcal{R}_0 = 0$  from the prediction step. Similar procedure could be repeated for other subintervals by a mathematical induction argument with respect to the  $m$ , from which we have  $\lim_{|z| \rightarrow \infty} \mathcal{R}_m^{(1)}(z) = 0$ , for  $m = 1, \dots, M$ . Specifically,  $\lim_{|z| \rightarrow \infty} \mathcal{R}_M^{(1)}(z) = 0$  after the first correction loop.

The same conclusion holds for the future correction steps by the mathematical induction argument with respect to the correction loop  $k$ , *i.e.*  $\lim_{|z| \rightarrow \infty} \mathcal{R}_m^{(k)}(z) = 0, m = 1, \dots, M$ . □

**Remark 5.2.** The above result can be generalized to the case when quadrature nodes are not uniformly distributed. On the other hand, the assumption to exclude the left-most point is necessary to guarantee the  $L$ -stability. Specifically, if the left-most point is included, then  $\lim_{|z| \rightarrow \infty} \mathcal{R}_0 \neq 0$ , due to its first component.

**Remark 5.3.** From the stability plots in Figure 4, the InDC-BE methods are  $L$ -stable when  $M \leq 8$  and for the number of iterations  $k \leq M$ .

## 6. CONCLUSIONS

This paper studies the order of convergence of the InDC-BE and InDC-IRK methods when applied to SSPs, using uniform distribution of quadrature points excluding the leftmost point. We applied the technique of asymptotic expansion in powers of  $\varepsilon$  for the smooth exact solution and for the corresponding numerical solution presented in [8, 9]. Two Theorems on global error estimate in the form of  $\varepsilon$ -expansion are presented and proved. Especially, we point out that the InDC methods improve the order of the  $\varepsilon$ -independent error, but there is no order improvement on the higher order terms  $\varepsilon^\nu$  ( $\nu \geq 1$ ). Such asymptotic analysis enables us to understand the phenomenon of order reduction for InDC methods when applied to stiff problems. A solution in order to

solve this problem is not a trivial matter. In fact, as mentioned in Remarks 4.5 and A.8, the bottleneck is the order reduction phenomenon in the prediction step. Further deep studies are required. It is an interesting topic for future investigation but it is beyond our scope in this paper. Numerical results on van der Pol equations confirm these convergence results.

### APPENDIX A.

In the appendix, we extend the error estimates of the InDC-BE method to InDC-IRK method when applied to SPPs. We first describe the InDC-IRK method applied to (1.2), then perform an  $\varepsilon$ -expansion of the numerical solution of this method, and finally we prove Theorem 3.2.

#### A.1. InDC-IRK method

We consider the InDC-IRK method constructed with  $s$ -stage IRK methods, where  $A$  matrices in the Butcher tableau (1.10) are invertible. For the internal stages in the IRK method, we introduce the integration matrix and interpolation matrix as following

$$hS^{c_{mi},k} = \int_{\tau_m}^{\tau_m+c_{mi}h} \alpha_k(s)ds, \quad P^{c_{mi},k} = \alpha_k(\tau_m + c_{mi}h), \tag{A.1}$$

$\forall m = 0, \dots, M - 1, \quad \forall k = 1, \dots, M$  and  $\forall mi = 1, \dots, s$ , where  $mi$  is index used for the  $i$ th-stage of the IRK method over the subinterval  $[\tau_m, \tau_{m+1}]$ . Here  $\alpha_k(s)$  is the Lagrangian basis function based on the node  $\tau_k$ . Let

$$S^{c_{mi}}(\bar{f}) = \sum_{j=1}^M S^{c_{mi},j} f(y_j, z_j), \quad P^{c_{mi}}(\bar{f}) = \sum_{j=1}^M P^{c_{mi},j} f(y_j, z_j),$$

then we have

$$hS^{c_{mi}}(\bar{f}) - \int_{\tau_m}^{\tau_m+c_ih} f(y(s), z(s))ds = \mathcal{O}(h^{M+1}), \tag{A.2}$$

$$P^{c_{mi}}(\bar{f}) - f(y(\tau_m + c_ih), z(\tau_m + c_ih)) = \mathcal{O}(h^M), \tag{A.3}$$

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

To compute the numerical error approximating the error function  $e^{(k-1)}(\tau_m), d^{(k-1)}(\tau_m)$  with a general IRK method to (2.5), we obtain

$$\begin{pmatrix} \hat{e}_{m+1}^{(k-1)} \\ \varepsilon \hat{d}_{m+1}^{(k-1)} \end{pmatrix} = \begin{pmatrix} \hat{e}_m^{(k-1)} + h \int_0^1 \delta(\tau_m + \tau h) d\tau \\ \varepsilon \hat{d}_m^{(k-1)} + h \int_0^1 \rho(\tau_m + \tau h) d\tau \end{pmatrix} + h \sum_{i=1}^s b_i \begin{pmatrix} \Delta \hat{\mathcal{K}}_{mi}^{(k-1)} \\ \Delta \hat{\mathcal{L}}_{mi}^{(k-1)} \end{pmatrix}, \tag{A.4}$$

and

$$\begin{pmatrix} \hat{E}_{mi}^{(k-1)} \\ \varepsilon \hat{D}_{mi}^{(k-1)} \end{pmatrix} = \begin{pmatrix} \hat{e}_m^{(k-1)} + h \int_0^{c_{mi}} \delta(\tau_m + \tau h) d\tau \\ \varepsilon \hat{d}_m^{(k-1)} + h \int_0^{c_{mi}} \rho(\tau_m + \tau h) d\tau \end{pmatrix} + h \sum_{j=1}^s a_{ij} \begin{pmatrix} \Delta \hat{\mathcal{K}}_{mj}^{(k-1)} \\ \Delta \hat{\mathcal{L}}_{mj}^{(k-1)} \end{pmatrix}, \tag{A.5}$$

with

$$\begin{pmatrix} \Delta \hat{\mathcal{K}}_{mi}^{(k-1)} \\ \Delta \hat{\mathcal{L}}_{mi}^{(k-1)} \end{pmatrix} \doteq \begin{pmatrix} f(\hat{Y}_{mi}^{(k)}, \hat{Z}_{mi}^{(k)}) - P^{c_{mi}}(\hat{f}^{(k-1)}) \\ g(\hat{Y}_{mi}^{(k)}, \hat{Z}_{mi}^{(k)}) - P^{c_{mi}}(\hat{g}^{(k-1)}) \end{pmatrix} \tag{A.6}$$

$$= \begin{pmatrix} f(\hat{Y}_{mi}^{(k)}, \hat{Z}_{mi}^{(k)}) - f(P^{c_{mi}}(\hat{y}^{(k-1)}), P^{c_{mi}}(\hat{z}^{(k-1)})) \\ g(\hat{Y}_{mi}^{(k)}, \hat{Z}_{mi}^{(k)}) - g(P^{c_{mi}}(\hat{y}^{(k-1)}), P^{c_{mi}}(\hat{z}^{(k-1)})) \end{pmatrix} + \mathcal{O}(h^M), \tag{A.7}$$

where we put

$$\hat{Y}_{mi}^{(k)} = P^{c_{mi}}(\bar{y}^{(k-1)}) + \hat{E}_{mi}^{(k-1)}, \quad \hat{Z}_{mi}^{(k)} = P^{c_{mi}}(\bar{z}^{(k-1)}) + \hat{D}_{mi}^{(k-1)}, \tag{A.8}$$

and equation (A.7) is due to the high order interpolation accuracy of  $P^{c_{mi}}$ , see equation (A.3). We can rewrite the system (A.4) and (A.5) as

$$\begin{pmatrix} \hat{y}_{m+1}^{(k)} - hS_{\bar{f}}^{m,(k-1)} \\ \varepsilon \hat{z}_{m+1}^{(k)} - hS_{\bar{g}}^{m,(k-1)} \end{pmatrix} = \begin{pmatrix} \hat{y}_m^{(k)} \\ \varepsilon \hat{z}_m^{(k)} \end{pmatrix} + h \sum_{i=1}^s b_i \begin{pmatrix} \Delta \hat{\mathcal{K}}_{mi}^{(k-1)} \\ \Delta \hat{\mathcal{L}}_{mi}^{(k-1)} \end{pmatrix}, \tag{A.9}$$

$$\begin{pmatrix} \hat{Y}_{mi}^{(k)} - hS_{\bar{f}}^{c_{mi},(k-1)} \\ \varepsilon \hat{Z}_{mi}^{(k)} - hS_{\bar{g}}^{c_{mi},(k-1)} \end{pmatrix} = \begin{pmatrix} \hat{y}_m^{(k)} \\ \varepsilon \hat{z}_m^{(k)} \end{pmatrix} + h \sum_{j=1}^s a_{ij} \begin{pmatrix} \Delta \hat{\mathcal{K}}_{mj}^{(k-1)} \\ \Delta \hat{\mathcal{L}}_{mj}^{(k-1)} \end{pmatrix}, \tag{A.10}$$

with

$$\begin{pmatrix} S_{\bar{f}}^{m,(k-1)} = S^m(\bar{f}^{(k-1)}) \\ S_{\bar{g}}^{m,(k-1)} = S^m(\bar{g}^{(k-1)}) \end{pmatrix}, \quad \begin{pmatrix} S_{\bar{f}}^{c_{mi},(k-1)} = S^{c_{mi}}(\bar{f}^{(k-1)}) \\ S_{\bar{g}}^{c_{mi},(k-1)} = S^{c_{mi}}(\bar{g}^{(k-1)}) \end{pmatrix}. \tag{A.11}$$

**Remark A.1.** Under the assumption  $A$  invertible, from the second equation of (A.10) we obtain in vectorial form

$$h\Delta \bar{\mathcal{L}}^{(k-1)} = A^{-1}(\varepsilon \bar{Z}^{(k)} - \varepsilon \hat{z}_m^{(k)} \mathbf{1} - hS^{\bar{c}}(\bar{g}^{(k-1)})),$$

with  $\Delta \bar{\mathcal{L}}^{(k-1)} = (\Delta \hat{\mathcal{L}}_{m1}^{(k-1)}, \dots, \Delta \hat{\mathcal{L}}_{ms}^{(k-1)})^T$ ,  $\mathbf{1} = (1, \dots, 1)^T$  and  $\bar{c} = (c_{m1}, \dots, c_{ms})$ . Inserting this into the second equation of (A.9), we get

$$\varepsilon \hat{z}_{m+1}^{(k)} = \varepsilon R(\infty) \hat{z}_m^{(k)} + \varepsilon b^T A^{-1} \bar{Z}^{(k)} + h(S^m(\bar{g}^{(k-1)}) - b^T A^{-1} S^{\bar{c}}(\bar{g}^{(k-1)})). \tag{A.12}$$

Of special importance now are stiffly accurate RK methods, *i.e.*, methods which satisfy  $b^T A^{-1} = e_s^T$ . This implies  $R(\infty) = 0$  and  $b^T A^{-1} S^{\bar{c}}(\bar{g}^{(k-1)}) = e_s^T S^{\bar{c}}(\bar{g}^{(k-1)}) = S^m(\bar{g}^{(k-1)})$ . Hence by (A.12) we have:  $\hat{z}_{m+1}^{(k)} = \hat{z}_{ms}^{(k)}$ .

### A.2. $\varepsilon$ -asymptotic expansion of InDC-IRK methods

We formally expand the quantities  $\Delta \hat{\mathcal{K}}_{mi}^{(k-1)}, \Delta \hat{\mathcal{L}}_{mi}^{(k-1)}$  from (A.6) and  $\hat{Y}_{mi}^{(k)}, \hat{Z}_{mi}^{(k)}, \hat{y}_{m+1}^{(k)}, \hat{z}_{m+1}^{(k)}$  from (A.8) and (A.9) into powers of  $\varepsilon$  with  $\varepsilon$ -independent coefficients

$$\begin{aligned} \hat{y}_m^{(k)} &= \hat{y}_{m,0}^{(k)} + \varepsilon \hat{y}_{m,1}^{(k)} + \varepsilon^2 \hat{y}_{m,2}^{(k)} + \dots, \\ \hat{Y}_{mi}^{(k)} &= \hat{Y}_{mi,0}^{(k)} + \varepsilon \hat{Y}_{mi,1}^{(k)} + \varepsilon^2 \hat{Y}_{mi,2}^{(k)} + \dots, \\ \Delta \hat{\mathcal{K}}_{mi}^{(k-1)} &= \Delta \hat{\mathcal{K}}_{mi,0}^{(k-1)} + \varepsilon \Delta \hat{\mathcal{K}}_{mi,1}^{(k-1)} + \varepsilon^2 \Delta \hat{\mathcal{K}}_{mi,2}^{(k-1)} + \dots, \\ \hat{z}_m^{(k)} &= \hat{z}_{m,0}^{(k)} + \varepsilon \hat{z}_{m,1}^{(k)} + \varepsilon^2 \hat{z}_{m,2}^{(k)} + \dots, \\ \hat{Z}_{mi}^{(k)} &= \hat{Z}_{mi,0}^{(k)} + \varepsilon \hat{Z}_{mi,1}^{(k)} + \varepsilon^2 \hat{Z}_{mi,2}^{(k)} + \dots, \\ \Delta \hat{\mathcal{L}}_{mi}^{(k-1)} &= \varepsilon^{-1} \Delta \hat{\mathcal{L}}_{mi,-1}^{(k-1)} + \Delta \hat{\mathcal{L}}_{mi,0}^{(k-1)} + \varepsilon \Delta \hat{\mathcal{L}}_{mi,1}^{(k-1)} + \varepsilon^2 \Delta \hat{\mathcal{L}}_{mi,2}^{(k-1)} + \dots \end{aligned} \tag{A.13}$$



Inserting (A.13) into (A.6) we obtain

$$\varepsilon^0 : \quad \Delta \hat{\mathcal{K}}_{mi,0}^{(k-1)} = f(\hat{Y}_{mi,0}^{(k)}, \hat{Z}_{mi,0}^{(k)}) - f(P^{c_{mi}}(\bar{y}_0^{(k-1)}), P^{c_{mi}}(\bar{z}_0^{(k-1)})) + \mathcal{O}(h^M), \tag{A.14}$$

$$\begin{aligned} \varepsilon^1 : \quad \Delta \hat{\mathcal{K}}_{mi,1}^{(k-1)} &= \left( f_y(\hat{Y}_{mi,0}^{(k)}, \hat{Z}_{mi,0}^{(k)}) \hat{Y}_{mi,1}^{(k)} + f_z(\hat{Y}_{mi,0}^{(k)}, \hat{Z}_{mi,0}^{(k)}) \hat{Z}_{mi,1}^{(k)} \right) \\ &\quad - \left( f_y(P^{c_{mi}}(\bar{y}_0^{(k-1)}), P^{c_{mi}}(\bar{z}_0^{(k-1)})) P^{c_{mi}}(\bar{y}_1^{(k-1)}) \right) \\ &\quad + f_z(P^{c_{mi}}(\bar{y}_0^{(k-1)}), P^{c_{mi}}(\bar{z}_0^{(k-1)})) P^{c_{mi}}(\bar{z}_1^{(k-1)}) + \mathcal{O}(h^M). \end{aligned} \tag{A.15}$$

...

ans so on. Similarly, we have

$$\Delta \hat{\mathcal{L}}_{mi,-1}^{(k-1)} = g(\hat{Y}_{mi,0}^{(k)}, \hat{Z}_{mi,0}^{(k)}) - g(P^{c_{mi}}(\bar{y}_0^{(k-1)}), P^{c_{mi}}(\bar{z}_0^{(k-1)})) + \mathcal{O}(h^M), \tag{A.16}$$

$$\begin{aligned} \Delta \hat{\mathcal{L}}_{mi,0}^{(k-1)} &= \left( g_y(\hat{Y}_{mi,0}^{(k)}, \hat{Z}_{mi,0}^{(k)}) \hat{Y}_{mi,1}^{(k)} + g_z(\hat{Y}_{mi,0}^{(k)}, \hat{Z}_{mi,0}^{(k)}) \hat{Z}_{mi,1}^{(k)} \right) \\ &\quad - \left( g_y(P^{c_{mi}}(\bar{y}_0^{(k-1)}), P^{c_{mi}}(\bar{z}_0^{(k-1)})) P^{c_{mi}}(\bar{y}_1^{(k-1)}) \right) \\ &\quad + g_z(P^{c_{mi}}(\bar{y}_0^{(k-1)}), P^{c_{mi}}(\bar{z}_0^{(k-1)})) P^{c_{mi}}(\bar{z}_1^{(k-1)}) + \mathcal{O}(h^M). \end{aligned} \tag{A.17}$$

...

ans so on.

Because of the linearity of relations (A.9) and (A.10), we have to order  $\varepsilon^\nu$  with  $\nu = -1$  in vectorial form

$$hA \Delta \hat{\mathcal{L}}_{m,-1}^{(k-1)} + hS^{\vec{c}}(\bar{g}) = 0, \quad hb^T \Delta \mathcal{L}_{m,-1}^{(k-1)} + hS^m(\bar{g}) = 0, \tag{A.18}$$

and for  $\nu \geq 0$ ,

$$\begin{pmatrix} \hat{y}_{m+1,\nu}^{(k)} - hS_{\mathbb{F}_\nu}^{m,(k-1)} \\ \hat{z}_{m+1,\nu}^{(k)} - hS_{\mathbb{G}_{\nu+1}}^{m,(k-1)} \end{pmatrix} = \begin{pmatrix} \hat{y}_{m,\nu}^{(k)} \\ \hat{z}_{m,\nu}^{(k)} \end{pmatrix} + h \sum_{i=1}^s b_i \begin{pmatrix} \Delta \hat{\mathcal{K}}_{mi,\nu}^{(k-1)} \\ \Delta \hat{\mathcal{L}}_{mi,\nu}^{(k-1)} \end{pmatrix}, \tag{A.19}$$

$$\begin{pmatrix} \hat{Y}_{mi,\nu}^{(k)} - hS_{\mathbb{F}_\nu}^{c_{mi},(k-1)} \\ \hat{Z}_{mi,\nu}^{(k)} - hS_{\mathbb{G}_{\nu+1}}^{c_{mi},(k-1)} \end{pmatrix} = \begin{pmatrix} \hat{y}_{m,\nu}^{(k)} \\ \hat{z}_{m,\nu}^{(k)} \end{pmatrix} + h \sum_{j=1}^i a_{ij} \begin{pmatrix} \Delta \hat{\mathcal{K}}_{mj,\nu}^{(k-1)} \\ \Delta \hat{\mathcal{L}}_{mj,\nu}^{(k-1)} \end{pmatrix}, \tag{A.20}$$

where

$$S_{\mathbb{F}_\nu}^{m,(k-1)} = S^m(\bar{\mathbb{F}}_\nu^{(k-1)}), \quad S_{\mathbb{G}_{\nu+1}}^{m,(k-1)} = S^m(\bar{\mathbb{G}}_{\nu+1}^{(k-1)}). \tag{A.21}$$

Similarly for  $S_{\mathbb{F}_\nu}^{c_{mi},(k-1)}$  and  $S_{\mathbb{G}_{\nu+1}}^{c_{mi},(k-1)}$ .

### A.3. Proof of Theorem 3.2

Before proving Theorem 3.2, we first give some preliminary results as propositions. We remark that the crucial assumption in Theorem 3.2 is that the IRK method is *stiffly accurate*. In the case that this property is not satisfied, the method becomes unstable and the numerical solutions diverge, (see Fig. 2).

In order to justify this, from the invertibility of matrix  $A$  and by the first formula in (A.18) we get

$$\Delta \hat{\mathcal{L}}_{m,-1}^{(k)} = -A^{-1} S^{\vec{c}}(\bar{g}^{(k)}), \tag{A.22}$$

substituting now into the second formula in (A.18) yields

$$-b^T A^{-1} S^{\bar{c}}(\bar{g}^{(k-1)}) + S^m(\bar{g}^{(k-1)}) = 0. \quad (\text{A.23})$$

Then we have the following result as an immediate consequence of the fact that the IRK method is stiffly accurate:

**Proposition A.2.** *Equation (A.23) is automatically satisfied, if the IRK methods in the prediction and correction steps of the InDC method are stiffly accurate.*

*Proof.* An IRK method is stiffly accurate if

$$b^T A^{-1} = e_s^T, \quad (\text{A.24})$$

with  $e_s = (0, \dots, 0, 1)^T$ . From (A.23) we get

$$-e_s^T S^{\bar{c}}(\bar{g}^{(k-1)}) + S^m(\bar{g}^{(k-1)}) = 0. \quad (\text{A.25})$$

Since the last row of the spectral integration matrix is  $s^{m,k} = \int_{\tau_m}^{\tau_m + c_s h} \alpha_k(\tau) d\tau$  by (A.24) we get  $c_s = 1$  and then  $\int_{\tau_m}^{\tau_m + c_s h} \alpha_k(\tau) d\tau = \int_{\tau_m}^{\tau_m + 1} \alpha_k(\tau) d\tau$ . This yields that  $e_s^T S^{\bar{c}}(\bar{g}^{(k-1)}) = S^m(\bar{g}^{(k-1)})$ , and the equation (A.25) is satisfied.  $\square$

Furthermore, similar to the Proposition 4.4, we have the following result for InDC-IRK methods. This Proposition follows from Remark 4.3. In fact, similar reformulation have been performed for the InDC method constructed with explicit RK methods in the prediction and correction steps [4].

**Proposition A.3.** *The InDC method constructed with stiffly accurate IRK methods can be considered again as a stiffly accurate IRK method with a corresponding Butcher Tableau as in (1.10) with the matrix  $A$  invertible.*

Now we are in the position to prove Theorem 3.2 via the following two lemmas.

**Lemma A.4.** ( $\varepsilon^0$  error term) *Consider the reduced system (1.4) satisfying (1.5) with consistent initial values. The numerical solutions of the InDC method after  $k$  correction loops have the following local error estimates at the interior nodes  $\tau_m$ ,  $m = 0, \dots, M$ ,*

$$e_{m,0}^{(k)} = \mathcal{O}(h^{\min(s_k+1, M+1)}), \quad d_{m,0}^{(k)} = \mathcal{O}(h^{\min(s_k+1, M+1)}). \quad (\text{A.26})$$

*Proof.* Since the IRK method in the prediction step is stiffly accurate, by definition (1.1), we have  $b^T A^{-1} = e_s^T$ . This implies that the numerical solution is equal to the last stage of the method, i.e.  $\hat{z}_{m+1,0}^{(0)} = \hat{Z}_{ms,0}^{(0)}$  and  $\hat{y}_{m+1,0}^{(0)} = \hat{Y}_{ms,0}^{(0)}$ . By  $g(\hat{Y}_{mi,0}^{(0)}, \hat{Z}_{mi,0}^{(0)}) = 0$ , we get  $\hat{Z}_{mi,0}^{(0)} = \mathcal{G}(\hat{Y}_{mi,0}^{(0)})$  for all  $mi$  and, in particular,  $\hat{Z}_{ms,0}^{(0)} = \mathcal{G}(\hat{Y}_{ms,0}^{(0)})$ . Then this gives  $\hat{z}_{m+1,0}^{(0)} = \mathcal{G}(\hat{y}_{m+1,0}^{(0)})$ .

Now, by the fact that the IRK method is stiffly accurate and that  $\bar{g}_0^{(0)} = (g(\hat{y}_{1,0}^{(0)}, \hat{z}_{1,0}^{(0)}), \dots, g(\hat{y}_{M,0}^{(0)}, \hat{z}_{M,0}^{(0)})) = \mathbf{0}$  in the prediction step, for the first correction step, i.e.  $k = 1$ , it follows from (A.19) with  $\nu = 0$

$$\begin{aligned} \hat{y}_{m+1,0}^{(1)} &= \hat{y}_{m,0}^{(1)} + h S^m(\bar{f}_0^{(0)}) + h \sum_{i=1}^s b_i \Delta \hat{\mathcal{K}}_{mi,0}^{(0)}, \\ g(\hat{y}_{m+1,0}^{(1)}, \hat{z}_{m+1,0}^{(1)}) &= 0, \end{aligned} \quad (\text{A.27})$$

where, from (A.20), we have for the internal stages

$$\begin{aligned} \hat{Y}_{mi,0}^{(1)} &= \hat{y}_{m,0}^{(1)} + h S^{c_{mi}}(\bar{f}_0^{(0)}) + h \sum_{j=1}^i a_{ij} \Delta \hat{\mathcal{K}}_{mj,0}^{(0)}, \\ g(\hat{Y}_{mi,0}^{(1)}, \hat{Z}_{mi,0}^{(1)}) &= 0. \end{aligned} \quad (\text{A.28})$$

Now, from the invertibility of function  $g_z$ , by (A.27) and (A.28) we get  $\hat{Z}_{mi,0}^{(1)} = \mathcal{G}(\hat{Y}_{mi,0}^{(1)})$  and  $\hat{z}_{m+1,0}^{(1)} = \mathcal{G}(\hat{y}_{m+1,0}^{(1)})$ . Thus the IRK method reads

$$\begin{aligned}\hat{Y}_{mi,0}^{(1)} &= \hat{y}_{m,0}^{(1)} + h \sum_{j=1}^s a_{ij} \Delta \hat{\mathcal{K}}_{mj,0}^{(0)} + h S^{c_{mi}} (\tilde{f}_0^{(0)}), \\ \hat{y}_{m+1,0}^{(1)} &= \hat{y}_{m,0}^{(1)} + h \sum_{i=1}^s b_i \Delta \hat{\mathcal{K}}_{mi,0}^{(0)} + h S^m (\tilde{f}_0^{(0)}),\end{aligned}\tag{A.29}$$

where  $\tilde{f}_0^{(0)} = (f(\hat{y}_{0,0}^{(0)}, \mathcal{G}(\hat{y}_{0,0}^{(0)})), \dots, f(\hat{y}_{M,0}^{(0)}, \mathcal{G}(\hat{y}_{M,0}^{(0)})))$ . The scheme (A.29) of updating  $\hat{y}_{m+1,0}^{(1)}$  can be interpreted as the applying a correction step of the InDC method to the ordinary differential equation (1.6). Therefore applying similar local truncation error estimates as in [4, 5] for InDC frameworks using RK methods when applied to a classical ordinary differential equation, we obtain the local error estimate

$$e_{m,0}^{(1)} = \mathcal{O}(h^{\min(s_2+1, M+1)}),\tag{A.30}$$

for  $m = 0, \dots, M$ , with  $s_2 = p^{(0)} + p^{(1)}$ . By  $\hat{z}_{m,0}^{(1)} = \mathcal{G}(\hat{y}_{m,0}^{(1)})$ , using the Lipschitz condition of  $\mathcal{G}$ , we get

$$d_{m,0}^{(1)} = z_{m,0} - \hat{z}_{m,0}^{(1)} = \mathcal{O}(h^{\min(s_2+1, M+1)}).\tag{A.31}$$

Similarly, at internal stages of the IRK method, by  $\hat{Z}_{mi,0}^{(1)} = \mathcal{G}(\hat{Y}_{mi,0}^{(1)})$ , we have the following local error estimates,

$$E_{mi,0}^{(1)} = y_0(\tau_m + c_i h) - \hat{Y}_{mi,0}^{(1)} = \mathcal{O}(h^{\min(s_1+q^{(1)}+1, M+1)}),\tag{A.32}$$

and

$$D_{mi,0}^{(1)} = z_0(\tau_m + c_i h) - \hat{Z}_{mi,0}^{(1)} = \mathcal{O}(h^{\min(s_1+q^{(1)}+1, M+1)}),\tag{A.33}$$

where  $q^{(1)}$  is the stage order for the IRK method applied to the first correction loop. We note that the proof of the general  $k$  is similar.  $\square$

**Remark A.5.** The local truncation error estimate (A.30) from [5] is quite technically involved; it is related to estimating the smoothness of rescaled error functions. The estimate (A.32) follows a similar fashion. We refer readers to the original paper [5] for details.

**Remark A.6.** With the estimates in the above Lemma, *i.e.* equations (A.30)–(A.33), it follows from equation (A.15)

$$\begin{aligned}\Delta \hat{\mathcal{K}}_{mi,1}^{(k-1)} &= f_y(y_{mi,0}, z_{mi,0}) \hat{E}_{mi,1}^{(k-1)} + f_z(y_{mi,0}, z_{mi,0}) \hat{D}_{mi,1}^{(k-1)} + \mathcal{O}(h^{s_{k-1}+1}), \\ &\doteq \Delta \mathcal{K}_{mi,1}^{(k-1)} + \mathcal{O}(h^{s_{k-1}+1})\end{aligned}\tag{A.34}$$

where  $\hat{E}_{mi,1}^{(k-1)}$  and  $\hat{D}_{mi,1}^{(k-1)}$  are defined by the corresponding  $\varepsilon$ -expansion of equation (A.8),  $s_k = \sum_{r=0}^k p^{(r)}$ , and  $\Delta \mathcal{K}_{mi,1}^{(k-1)} \doteq f_y(y_{mi,0}, z_{mi,0}) \hat{E}_{mi,1}^{(k-1)} + f_z(y_{mi,0}, z_{mi,0}) \hat{D}_{mi,1}^{(k-1)}$ . Here we have used the abbreviations  $y_{mi,0}$  and  $z_{mi,0}$ , *i.e.* the exact solution  $y(t)$  and  $z(t)$  at the position  $t = \tau_m + c_i h$  respectively. We note that, from (A.15), we replaced  $\hat{Y}_{mi,0}^{(k)}$  and  $P^{c_{mi}}(\tilde{y}_0^{(k-1)})$  by adding and subtracting  $y_{mi,0}$  with an error of  $\mathcal{O}(h^{s_{k-1}+q^{(k)}+1})$  and  $\mathcal{O}(h^{s_{k-1}+1})$ , the same for  $\hat{Z}_{mi,0}^{(k)}$  and  $P^{c_{mi}}(\tilde{z}_0^{(k-1)})$ . Similarly, we have from (A.17)

$$\begin{aligned}\Delta \hat{\mathcal{L}}_{mi,0}^{(k-1)} &= g_y(y_{mi,0}, z_{mi,0}) \hat{E}_{mi,1}^{(k-1)} + g_z(y_{mi,0}, z_{mi,0}) \hat{D}_{mi,1}^{(k-1)} + \mathcal{O}(h^{s_{k-1}+1}) \\ &\doteq \Delta \mathcal{L}_{mi,0}^{(k-1)} + \mathcal{O}(h^{s_{k-1}+1}),\end{aligned}\tag{A.35}$$

where  $\Delta \mathcal{L}_{mi,0}^{(k-1)} \doteq g_y(y_{mi,0}, z_{mi,0}) \hat{E}_{mi,1}^{(k-1)} + g_z(y_{mi,0}, z_{mi,0}) \hat{D}_{mi,1}^{(k-1)}$ .

**Lemma A.7.** ( $\varepsilon^\nu$  error term) Consider the same assumptions as in Theorem 3.2 with  $0 < \varepsilon \ll 1$ . Then the numerical solutions of the InDC method after  $k$  correction loops have the following local error estimates at the interior nodes  $\tau_m$  with  $m = 0, \dots, M$

$$e_{m,\nu}^{(k)} = y_{m,\nu} - \hat{y}_{m,\nu}^{(k)} = \mathcal{O}(h^{q^{(0)}+2-\nu}), \quad d_{m,\nu}^{(k)} = z_{m,\nu} - \hat{z}_{m,\nu}^{(k)} = \mathcal{O}(h^{q^{(0)}+1-\nu}), \quad (\text{A.36})$$

with  $1 \leq \nu \leq q^{(0)} + 1$ .

*Proof.* We first prove (A.36) in the case  $\nu = 1$ . In the prediction step ( $k = 0$ ), under the assumption of stiffly accurate IRK method, by the Corollary 3.10 in [8], we get that the error estimates for  $\varepsilon^1$  in (1.13) at the interior nodes of the InDC method with  $m = 0, \dots, M$  satisfy

$$e_{m,1}^{(0)} = y_{m,1} - \hat{y}_{m,1}^{(0)} = \mathcal{O}(h^{q^{(0)}+1}), \quad d_{m,1}^{(0)} = z_{m,1} - \hat{z}_{m,1}^{(0)} = \mathcal{O}(h^{q^{(0)}}). \quad (\text{A.37})$$

We consider  $\varepsilon$ -expansions of  $\hat{y}_m^{(1)}$ ,  $\hat{z}_m^{(1)}$  and  $\hat{E}_{mi}^{(1)}$  and  $\hat{D}_{mi}^{(1)}$  as in (A.13). Inserting them onto equations (A.34), (A.35), from (A.19) and (A.20) for the power  $\varepsilon^1$  with  $k = 1$  and  $\nu = 1$ , we have

$$\begin{pmatrix} \hat{y}_{m+1,1}^{(1)} - hS_{\bar{\mathbb{F}}_1}^{m,(0)} \\ \hat{z}_{m+1,0}^{(1)} - hS_{\bar{\mathbb{G}}_1}^{m,(0)} \end{pmatrix} = \begin{pmatrix} \hat{y}_{m,1}^{(1)} \\ \hat{z}_{m,0}^{(1)} \end{pmatrix} + h \sum_{i=1}^s b_i \begin{pmatrix} \Delta\mathcal{K}_{mi,1}^{(0)} \\ \Delta\mathcal{L}_{mi,0}^{(0)} \end{pmatrix} + \mathcal{O}(h^{p^{(0)}+2}), \quad (\text{A.38})$$

and

$$\begin{pmatrix} \hat{Y}_{mi,1}^{(1)} - hS_{\bar{\mathbb{F}}_1}^{c_{m,i},(0)} \\ \hat{Z}_{mi,0}^{(1)} - hS_{\bar{\mathbb{G}}_1}^{c_{m,i},(0)} \end{pmatrix} = \begin{pmatrix} \hat{y}_{m,1}^{(1)} \\ \hat{z}_{m,0}^{(1)} \end{pmatrix} + h \sum_{j=1}^s a_{ij} \begin{pmatrix} \Delta\mathcal{K}_{mj,1}^{(0)} \\ \Delta\mathcal{L}_{mj,0}^{(0)} \end{pmatrix} + \mathcal{O}(h^{p^{(0)}+2}). \quad (\text{A.39})$$

Now from (4.5) we have for  $\varepsilon^1$ ,

$$y_{m+1,1} = y_{m,1} + \int_{\tau_m}^{\tau_{m+1}} \mathbb{F}_1(t) dt, \quad z_{m+1,0} = z_{m,0} + \int_{\tau_m}^{\tau_{m+1}} \mathbb{G}_1(t) dt. \quad (\text{A.40})$$

We subtract (A.38) from (A.40) and so obtain

$$\begin{pmatrix} e_{m+1,1}^{(1)} + hS_{\bar{\mathbb{F}}_1}^{m,(0)} - \int_{\tau_m}^{\tau_{m+1}} \mathbb{F}_1(t) dt \\ d_{m+1,0}^{(1)} + hS_{\bar{\mathbb{G}}_1}^{m,(0)} - \int_{\tau_m}^{\tau_{m+1}} \mathbb{G}_1(t) dt \end{pmatrix} = \begin{pmatrix} e_{m,1}^{(1)} \\ d_{m,0}^{(1)} \end{pmatrix} - h \sum_{i=1}^s b_i \begin{pmatrix} \Delta\mathcal{K}_{mi,1}^{(0)} \\ \Delta\mathcal{L}_{mi,0}^{(0)} \end{pmatrix} + \mathcal{O}(h^{p^{(0)}+2}). \quad (\text{A.41})$$

From the Corollary 3.10 in [8] and (A.37), we have the following estimates for the local errors

$$\begin{aligned} e_{m,0}^{(0)} = y_{m,0} - \hat{y}_{m,0}^{(0)} &= \mathcal{O}(h^{p^{(0)}+1}), & d_{m,0}^{(0)} = z_{m,0} - \hat{z}_{m,0}^{(0)} &= \mathcal{O}(h^{p^{(0)}+1}), \\ e_{m,1}^{(0)} = y_{m,1} - \hat{y}_{m,1}^{(0)} &= \mathcal{O}(h^{q^{(0)}+1}), & d_{m,1}^{(0)} = z_{m,1} - \hat{z}_{m,1}^{(0)} &= \mathcal{O}(h^{q^{(0)}}). \end{aligned} \quad (\text{A.42})$$

Similarly as done in the proof of Lemma 4.2, on the right hand-side of (A.41) we add and subtract the quantities  $S^m(\bar{\mathbb{F}}_1)$  and  $S^m(\bar{\mathbb{G}}_1)$ , these are the integrals of  $(M - 1)$ th degree interpolating polynomials on  $(\tau_m, \mathbb{F}_1(\tau_m))_{m=1}^M$  and  $(\tau_m, \mathbb{G}_1(\tau_m))_{m=1}^M$  over the subinterval  $[\tau_m, \tau_{m+1}]$ . Hence,  $\int_{\tau_m}^{\tau_{m+1}} \mathbb{F}_1(\tau) d\tau - hS^m(\bar{\mathbb{F}}_1) = \mathcal{O}(h^{M+1})$  and by (A.42), we have  $S^m(\bar{\mathbb{F}}_1) - S_{\bar{\mathbb{F}}_1}^{m,(0)} = \mathcal{O}(h^{q^{(0)}})$  and  $S^m(\bar{\mathbb{G}}_1) - S_{\bar{\mathbb{G}}_1}^{m,(0)} = \mathcal{O}(h^{q^{(0)}})$ . Then we have from (A.41)

$$\begin{aligned} e_{m+1,1}^{(1)} &= e_{m,1}^{(1)} - h \sum_{i=1}^s b_i \Delta\mathcal{K}_{mi,1}^{(0)} + \mathcal{O}(h^{q^{(0)}+1}), \\ d_{m+1,0}^{(1)} &= d_{m,0}^{(1)} - h \sum_{i=1}^s b_i \Delta\mathcal{L}_{mi,0}^{(0)} + \mathcal{O}(h^{q^{(0)}+1}). \end{aligned} \quad (\text{A.43})$$

Now we consider the  $\varepsilon$ -expansion of the error at internal stages  $\tau_m + c_i h$ , and as in equation (4.12) we get

$$E_{mi,1}^{(1)} = P^{c_{mi}}(\bar{e}_1^{(0)}) - \hat{E}_{mi,1}^{(0)}, \quad D_{mi,1}^{(1)} = P^{c_{mi}}(\bar{d}_1^{(0)}) - \hat{D}_{mi,1}^{(0)}, \quad \forall k \geq 0, m, \quad (\text{A.44})$$

where  $\bar{e}_1^{(0)} = (e_{m1,1}^{(0)}, \dots, e_{ms,1}^{(0)})$ ,  $\bar{d}_1^{(0)} = (d_{m1,0}^{(0)}, \dots, d_{ms,0}^{(0)})$ ,  $s$  is the number of internal stages in an IRK method. Especially, by (A.42), it follows from (A.44),

$$\begin{aligned}\hat{E}_{mi,1}^{(0)} &= -E_{mi,1}^{(1)} + \mathcal{O}(h^{q^{(0)}+1}), \\ \hat{D}_{mi,1}^{(0)} &= -D_{mi,1}^{(1)} + \mathcal{O}(h^{q^{(0)}}).\end{aligned}\quad (\text{A.45})$$

Similarly as equations (A.43), from the definition of stage order for the prediction step, we have for the internal stages in vectorial form

$$\begin{aligned}\bar{E}_1^{(1)} &= e_{m,1}^{(1)} \mathbf{1} - hA\Delta\bar{\mathcal{K}}_1^{(0)} + \mathcal{O}(h^{q^{(0)}+1}), \\ \bar{D}_0^{(1)} &= d_{m,0}^{(1)} \mathbf{1} - hA\Delta\bar{\mathcal{L}}_0^{(0)} + \mathcal{O}(h^{q^{(0)}+1}),\end{aligned}\quad (\text{A.46})$$

where  $\bar{E}_1^{(1)} = (E_{m1,1}^{(1)}, \dots, E_{ms,1}^{(1)})$ ,  $\bar{D}_0^{(1)} = (D_{m1,0}^{(1)}, \dots, D_{ms,0}^{(1)})$  and  $\mathbf{1} = (1, 1, \dots, 1)^T$  is a vector of size  $s$ . Now from the second equation in (A.46) and using (A.26) and (A.45), we get

$$A(g_y(y_{mi,0}, z_{mi,0})\hat{E}_{mi,1}^{(0)} + g_z(y_{mi,0}, z_{mi,0})\hat{E}_{mi,1}^{(0)}) = \mathcal{O}(h^{q^{(0)}}), \quad (\text{A.47})$$

Thus, from the invertibility of matrix  $A$  we have

$$\hat{D}_{mi,1}^{(0)} = -(g_z^{-1}g_y)(y_{mi,0}, z_{mi,0})\hat{E}_{mi,1}^{(0)} + \mathcal{O}(h^{q^{(0)}}), \quad (\text{A.48})$$

for  $mi = m1, \dots, ms$ . Plug the above equation (A.48) into equation (A.34) and replace  $\hat{E}_{mi,1}^{(0)}$  by  $E_{mi,1}^{(1)}$  with  $\mathcal{O}(h^{q^{(0)}+1})$  error and  $\hat{D}_{mi,1}^{(0)}$  by  $D_{mi,1}^{(1)}$  with  $\mathcal{O}(h^{q^{(0)}})$  error, by (A.45) we obtain

$$\Delta\mathcal{K}_{mi,1}^{(0)} = (f_y - f_z g_z^{-1} g_y)(y_{mi,0}, z_{mi,0})E_{mi,1}^{(1)} + \mathcal{O}(h^{q^{(0)}}). \quad (\text{A.49})$$

Our next aim now is to prove the local error  $e_{m,1}^{(1)} = \mathcal{O}(h^{q^{(0)}+1})$  by mathematical induction w.r.t.  $m$ . Especially, we would like to show that  $e_{m+1,1}^{(1)} = \mathcal{O}(h^{q^{(0)}+1})$ , if we assume the local error  $e_{l,1}^{(1)} = \mathcal{O}(h^{q^{(0)}+1})$ ,  $\forall l \leq m$ . To show this, we plug equation (A.49) into the first equation (A.46) and obtain  $E_{mi,1}^{(1)} = \mathcal{O}(h^{q^{(0)}+1})$ , for  $mi = m1, \dots, ms$ . From (A.49),  $\Delta\mathcal{K}_{mi,1}^{(0)} = \mathcal{O}(h^{q^{(0)}})$  and plug this estimate into the first equation of (A.43), we obtain the desired estimate of

$$e_{m+1,1}^{(1)} = \mathcal{O}(h^{q^{(0)}+1}). \quad (\text{A.50})$$

Thus, from (A.48) and (A.45), it follows

$$D_{mi,1}^{(1)} = \mathcal{O}(h^{q^{(0)}}). \quad (\text{A.51})$$

Now in order to prove the estimate  $d_{m,1}^{(1)} = \mathcal{O}(h^{q^{(0)}})$ , we start to considering equation (A.12). Since the IRK method is stiffly accurate, from Remark A.1, we have  $\hat{z}_{m+1,1}^{(1)} = \hat{Z}_{ms,1}^{(1)}$ . Hence from (A.51),

$$z_1(\tau_{m+1}) - \hat{z}_{m+1,1}^{(1)} = d_{m+1}^{(1)} = D_{ms,1}^{(1)} \stackrel{(\text{A.51})}{=} \mathcal{O}(h^{q^{(0)}}), \quad m = 0, \dots, M-1. \quad (\text{A.52})$$

The above proof can be generalized for the InDC method with different IRK methods applied to  $k$  correction steps. The local error estimates at the interior nodes of the InDC method  $\tau_m$  with  $m = 0, \dots, M$  are

$$e_{m,1}^{(k)} = \mathcal{O}(h^{q^{(0)}+1}), \quad d_{m,1}^{(k)} = \mathcal{O}(h^{q^{(0)}}).$$

We have thus proved equation (A.36) with  $\nu = 1$ . The general estimates for  $\nu > 1$  in equation (A.36) can be obtained in a similar fashion to the case of  $\nu = 1$ , as in the Theorem 3.4 in (Chap. VI of [8]).  $\square$

*Proof of Theorem 3.2.* The proof is similar to that for Theorem 3.1. In fact, we obtain estimates (3.2) by using the results of Lemmas A.4 and A.7. From Lemma A.4 we have the local error estimate, for one time step  $t_0$  to  $t_1$ ,

$$e_{M,0}^{(K)} = \mathcal{O}(H^{\min(s_K+1, M+1)}).$$

From local to global error, we obtain  $e_{n,0}^{(K)} = \mathcal{O}(H^{\min(s_K, M)})$ . From equation (4.2) and the Lipschitz condition of  $\mathcal{G}$ , we get

$$d_{n,0}^{(K)} = \mathcal{O}(H^{\min(s_K, M)}).$$

Now in order to complete the proof of the theorem, we consider the estimates (A.36) in Lemma A.7 with  $\nu = 1$ . Then we have for the local error estimates after one step (from  $t_0$  to  $t_1$ )

$$e_{M,1}^{(K)} = \mathcal{O}(H^{q^{(0)}+1}), \quad d_{M,1}^{(K)} = \mathcal{O}(H^{q^{(0)}}).$$

Finally, the global estimate (3.2) from the local estimate above is a consequence of Theorems 4.5 and 4.6 in (Chap. VII of [8]). □

**Remark A.8.** We remark that we can not improve the estimate of the global error for the  $y$ -component as done in (Thm. 3.4 in [8]) for high-indices. Indeed the reason for such loss of accuracy is related to the evaluation of the integrals in equation (A.21). These integrals are obtained from the prediction step, and the algebraic variable  $z$  obtained in the prediction step is the cause that reduces the order of the differential variable  $y$  in the correction steps. This can be seen in the evaluation from equation (A.41) to (A.43) due to (A.42). We note that a similar conclusion for the remainder can be drawn.

#### A.4. Estimation of the Remainder

Finally, in order to estimate the remainder for the global error functions  $e_n^{(K)}$  and  $d_n^{(K)}$ , we have the following result.

**Theorem A.9.** *Under the same hypothesis as in Theorem 3.2 for any fixed constant  $C > 0$  and  $\nu \leq q^{(0)} + 1$ , the global error satisfies for  $\varepsilon \leq CH$*

$$e_n^{(K)} = e_{n,0}^{(K)} + \varepsilon e_{n,1}^{(K)} + \dots + \varepsilon^\nu e_{n,\nu}^{(K)} + \mathcal{O}(\varepsilon^{\nu+1}/H), \quad d_n^{(K)} = d_{0,n}^{(K)} + \varepsilon d_{n,1}^{(K)} + \dots + \varepsilon^\nu d_{n,\nu}^{(K)} + \mathcal{O}(\varepsilon^{\nu+1}/H). \quad (\text{A.53})$$

These estimates hold uniformly for  $H \leq H_0$  and  $nH \leq \text{Const}$ .

*Proof.* By the estimates (A.36) it is sufficient to prove the result for  $\nu = q^{(0)} + 1$ .

Through the  $\varepsilon$ -asymptotic expansion (2.26) for the global error functions  $e_n^{(k)}$  and  $d_n^{(k)}$ , by considering estimates (A.36) globally, *i.e.*

$$e_{n,\nu}^{(k)} = y_{n,\nu} - \hat{y}_{n,\nu}^{(k)} = \mathcal{O}(H^{q^{(0)}+1-\nu}), \quad d_{n,\nu}^{(k)} = z_{n,\nu} - \hat{z}_{n,\nu}^{(k)} = \mathcal{O}(H^{q^{(0)}+1-\nu}), \quad (\text{A.54})$$

and  $\nu = q^{(0)} + 1$ , we get:

$$e_n^{(K)} = e_{n,0}^{(K)} + \varepsilon e_{n,1}^{(K)} + \dots + \varepsilon^\nu e_{n,\nu}^{(K)} + \mathcal{O}(\varepsilon^{\nu+1}/H), \quad d_n^{(K)} = d_{0,n}^{(K)} + \varepsilon d_{n,1}^{(K)} + \dots + \varepsilon^\nu d_{n,\nu}^{(K)} + \mathcal{O}(\varepsilon^{\nu+1}/H). \quad (\text{A.55})$$

with  $e_{n,0}^{(K)} = y_{n,0}^{(K)} - \hat{y}_{n,0}^{(K)}$  and  $d_{n,0}^{(K)} = z_{n,0}^{(K)} - \hat{z}_{n,0}^{(K)}$ , ... (see formula (2.26)).

In order to estimate the remainder, we consider the truncated series of the quantities in (A.13):

$$\hat{\underline{y}}_n^{(K)} = \hat{y}_{n,0}^{(K)} + \varepsilon \hat{y}_{n,1}^{(K)} + \dots + \varepsilon^\nu \hat{y}_{n,\nu}^{(K)}, \quad \hat{\underline{z}}_n^{(K)} = \hat{z}_{n,0}^{(K)} + \varepsilon \hat{z}_{n,1}^{(K)} + \dots + \varepsilon^\nu \hat{z}_{n,\nu}^{(K)}, \quad (\text{A.56})$$

$$\hat{\underline{Y}}_{ni}^{(K)} = \hat{Y}_{ni,0}^{(K)} + \varepsilon \hat{Y}_{ni,1}^{(k)} + \dots + \varepsilon^\nu \hat{Y}_{ni,\nu}^{(k)}, \quad \hat{\underline{Z}}_{ni}^{(K)} = \hat{Z}_{ni,0}^{(K)} + \varepsilon \hat{Z}_{ni,1}^{(k)} + \dots + \varepsilon^\nu \hat{Z}_{ni,\nu}^{(k)}, \quad (\text{A.57})$$

and

$$\begin{aligned} \Delta \hat{Y}_{ni}^{(K)} &\doteq \hat{Y}_{ni}^{(K)} - \hat{Y}_{ni}^{(k)}, & \Delta \hat{Z}_{ni}^{(K)} &\doteq \hat{Z}_{ni}^{(K)} - \hat{Z}_{ni}^{(k)}. \\ \Delta \hat{\mathcal{K}}_{ni}^{(K-1)} &= \Delta \hat{\mathcal{K}}_{ni,0}^{(K-1)} + \varepsilon \Delta \hat{\mathcal{K}}_{ni,1}^{(K-1)} + \dots + \varepsilon^\nu \Delta \hat{\mathcal{K}}_{ni,\nu}^{(K-1)}, \\ \Delta \hat{\mathcal{L}}_{ni}^{(K-1)} &= \Delta \hat{\mathcal{L}}_{ni,0}^{(K-1)} + \varepsilon \Delta \hat{\mathcal{L}}_{ni,1}^{(K-1)} + \dots + \varepsilon^\nu \Delta \hat{\mathcal{L}}_{ni,\nu}^{(K-1)}, \end{aligned} \tag{A.58}$$

and we use the notation for the remainder

$$\Delta \hat{y}_n^{(K)} \doteq \hat{y}_n^{(K)} - \hat{y}_n^{(k)}, \quad \Delta \hat{z}_n^{(K)} \doteq \hat{z}_n^{(K)} - \hat{z}_n^{(k)}. \tag{A.59}$$

Then using (2.26) and (A.54), (A.55) is then equivalent to

$$\Delta \hat{y}_n^{(K)} = \mathcal{O}(\varepsilon^{\nu+1}/H), \quad \Delta \hat{z}_n^{(K)} = \mathcal{O}(\varepsilon^{\nu+1}/H). \tag{A.60}$$

The proof of (A.60) is similar to the proof of (Thm. 3.8, Chap. VI) in [8]. For the sake of brevity, here we point out the main differences and we give some partial results.

We consider the InDC method (A.9)-(A.10)-(A.11). Inserting (A.56) into (A.2), it gives

$$S_{\hat{f}}^{c_{mi},(k-1)} - S_{\underline{f}}^{c_{mi},(k-1)} = \mathcal{O}(\varepsilon^{\nu+1}), \quad S_{\hat{g}}^{c_{mi},(k-1)} - S_{\underline{g}}^{c_{mi},(k-1)} = \mathcal{O}(\varepsilon^{\nu+1}). \tag{A.61}$$

Similarly inserting the quantities (A.56) (A.57) and (A.58) with  $m$  (index of quadrature nodes) instead of  $n$  into (A.7) and by (A.14), (A.16) and (A.61), we obtain by Lemma (A.7) and  $\nu \leq q^{(0)} + 1$ ,

$$\begin{pmatrix} \hat{Y}_{mi}^{(k)} - \hat{y}_m^{(k)} \\ \varepsilon(\hat{Z}_{mi}^{(k)} - \hat{z}_m^{(k)}) \end{pmatrix} = h \sum_{j=1}^s a_{ij} \begin{pmatrix} \Delta \hat{\mathcal{K}}_{mi}^{(k-1)} \\ \Delta \hat{\mathcal{L}}_{mi}^{(k-1)} \end{pmatrix} + \begin{pmatrix} \mathcal{O}(h\varepsilon^{\nu+1}) \\ \mathcal{O}(\varepsilon^{\nu+1}) \end{pmatrix}. \tag{A.62}$$

This represents the defect when (A.56) (A.57) and (A.58) are inserted into the InDC R-K method (A.9)-(A.10).

From now on the the proof is similar to (Thm. 3.8, in [8]). In fact, applying Theorem 3.6 in [8] to (A.62), it yields

$$\begin{aligned} \|\Delta \hat{Y}_{mi}^{(k)}\| &\leq C(\|\Delta \hat{y}_m^{(k)}\| + \varepsilon \|\Delta \hat{z}_m^{(k)}\|) + \mathcal{O}(\varepsilon^{\nu+1}), \\ \|\Delta \hat{Z}_{mi}^{(k)}\| &\leq C(\|\Delta \hat{y}_m^{(k)}\| + \varepsilon/h \|\Delta \hat{z}_m^{(k)}\|) + \mathcal{O}(\varepsilon^{\nu+1}/h). \end{aligned} \tag{A.63}$$

where here the quantities  $\delta_i$  and  $\theta_i$  in Theorem 3.6 in [8] are given by:  $\delta_i = \mathcal{O}(\varepsilon^{\nu+1})$ ,  $\theta_i = \mathcal{O}(\varepsilon^{\nu+1}/h)$ .

In a similar fashion as the point b) of the proof in Theorem 3.8 in [8], we obtain for the quantities  $\Delta \hat{y}_m^{(k)}$  and  $\Delta \hat{z}_m^{(k)}$  in (A.59) the recursion

$$\begin{pmatrix} \|\Delta \hat{y}_{m+1}^{(k)}\| \\ \|\Delta \hat{z}_{m+1}^{(k)}\| \end{pmatrix} = \begin{pmatrix} 1 + \mathcal{O}(h) & \mathcal{O}(\varepsilon) \\ \mathcal{O}(1) & \alpha + \mathcal{O}(\varepsilon) \end{pmatrix} \begin{pmatrix} \|\Delta \hat{y}_m^{(k)}\| \\ \|\Delta \hat{z}_m^{(k)}\| \end{pmatrix} + \begin{pmatrix} \mathcal{O}(\varepsilon^{\nu+1}) \\ \mathcal{O}(\varepsilon^{\nu+1}/h) \end{pmatrix}, \tag{A.64}$$

with  $\alpha < 1$ . The value of  $\alpha$  is specified in the (Thm. 3.8 in [8]).

Finally applying Lemma 3.9 in [8] to the difference inequalities in (A.64) gives

$$\Delta \hat{y}_m^{(k)} = \mathcal{O}(\varepsilon^{\nu+1}/h), \quad \Delta \hat{z}_m^{(k)} = \mathcal{O}(\varepsilon^{\nu+1}/h). \tag{A.65}$$

for  $mh \leq \text{Const}$ . By  $H = Mh$ , then we get (A.60), *i.e.*, the statement of the theorem. □



**Remark A.10.** We note that from (A.53) for Theorem 3.1 and  $\nu \leq 2$  we get:

$$\begin{aligned} e_n^{(K)} &= \mathcal{O}(H^{\min\{K+1, M\}}) + \mathcal{O}(\varepsilon H) + \mathcal{O}(\varepsilon^2) + \mathcal{O}(\varepsilon^3/H), \\ d_n^{(K)} &= \mathcal{O}(H^{\min\{K+1, M\}}) + \mathcal{O}(\varepsilon H) + \mathcal{O}(\varepsilon^2) + \mathcal{O}(\varepsilon^3/H), \end{aligned} \quad (\text{A.66})$$

and for Theorem 3.2 with  $\nu = q^{(0)} + 1$ :

$$\begin{aligned} e_n^{(K)} &= \mathcal{O}(H^{\min\{s_K, M\}}) + \mathcal{O}(\varepsilon H^{q^{(0)}}) + \dots + \mathcal{O}(\varepsilon^{q^{(0)}+1}) + \mathcal{O}(\varepsilon^{q^{(0)}+2}/H), \\ d_n^{(K)} &= \mathcal{O}(H^{\min\{s_K, M\}}) + \mathcal{O}(\varepsilon H^{q^{(0)}}) + \dots + \mathcal{O}(\varepsilon^{q^{(0)}+1}) + \mathcal{O}(\varepsilon^{q^{(0)}+2}/H). \end{aligned} \quad (\text{A.67})$$

## REFERENCES

- [1] W. Auzinger, H. Hofstätter, W. Kreuzer and E. Weinmüller, Modified defect correction algorithms for ODEs. Part I: General Theory. *Numer. Algorithms* **36** (2004) 135–156.
- [2] K. Böhmer and HJ Stetter, Defect correction methods. Theory and applications (1984).
- [3] A. Christlieb, M. Morton, B. Ong and J.-M. Qiu, Semi-implicit integral deferred correction constructed with high order additive Runge–Kutta methods. *Communications in Mathematical Sciences* (2011).
- [4] A. Christlieb, B. Ong and J.M. Qiu, Comments on high order integrators embedded within integral deferred correction methods. *Commun. Appl. Math. Comput. Sci* **4** (2009) 27–56.
- [5] A. Christlieb, B. Ong and J.M. Qiu, Integral deferred correction methods constructed with high order Runge–Kutta integrators. *Math. Comput.* **79** (2009) 761.
- [6] A. Dutt, L. Greengard and V. Rokhlin, Spectral deferred correction methods for ordinary differential equations. *BIT Numer. Math.* **40** (2000) 241–266.
- [7] C.W. Gear, Differential-algebraic equation index transformations. *SIAM J. Sci. Stat. Comput.* **9** (1988) 39–47.
- [8] E. Hairer and G. Wanner, Solving ordinary differential equations II: stiff and differential algebraic problems, vol. 2. Springer Verlag (1993).
- [9] E. Hairer, C. Lubich and M. Roche, Error of Runge–Kutta methods for stiff problems studied via differential algebraic equations. *BIT Numer. Math.* **28** (1988) 678–700.
- [10] J. Huang, J. Jia and M. Minion, Arbitrary order Krylov deferred correction methods for differential algebraic equations. *J. Comput. Phys.* **221** (2007) 739–760.
- [11] A.T. Layton, On the choice of correctors for semi-implicit picard deferred correction methods. *Appl. Numer. Math.* **58** (2008) 845–858.
- [12] A.T. Layton and M.L. Minion, Implications of the choice of quadrature nodes for picard integral deferred corrections methods for ordinary differential equations. *BIT Numer. Math.* **45** (2005) 341–373.
- [13] A.T. Layton and M.L. Minion, Implications of the choice of predictors for semi-implicit picard integral deferred corrections methods. *Commun. Appl. Math. Comput. Sci.* **1** (2007) 1–34.
- [14] M.L. Minion, Semi-implicit spectral deferred correction methods for ordinary differential equations. *Commun. Math. Sci.* **1** (2003) 471–500.
- [15] R.E. O’Malley Jr, Introduction to singular perturbations, Vol. 14. Applied Mathematics and Mechanics. Technical report, DTIC Document (1974).
- [16] R.D. Skeel, A theoretical framework for proving accuracy results for deferred corrections. *SIAM J. Numer. Anal.* **19** (1982) 171–196.
- [17] A. Tikhonov, B. Vasl’eva and A. Sveshnikov, Differential Equations. Springer Verlag (1985).