# EXTRAPOLATED IMPLICIT-EXPLICIT TIME STEPPING*

### EMIL M. CONSTANTINESCU† AND ADRIAN SANDU‡

**Abstract.** This paper constructs extrapolated implicit-explicit time stepping methods that allow one to efficiently solve problems with both stiff and nonstiff components. The proposed methods are based on Euler steps and can provide very high order discretizations of ODEs, index-1 DAEs, and PDEs in the method-of-lines framework. Implicit-explicit schemes based on extrapolation are simple to construct, easy to implement, and straightforward to parallelize. This work establishes the existence of perturbed asymptotic expansions of global errors, explains the convergence orders of these methods, and studies their linear stability properties. Numerical results with stiff ODE, DAE, and PDE test problems confirm the theoretical findings and illustrate the potential of these methods to solve multiphysics multiscale problems.

**Key words.** extrapolation methods, implicit-explicit methods, stiff problems, differential algebraic systems, ordinary and partial differential equations

**AMS subject classifications.** 34E05, 65L06, 65L80

**DOI.** 10.1137/080732833

**1. Introduction.** Several areas of science and engineering require numerical simulations of multiphysics multiscale systems. Examples include mechanical, aerospace, environmental, and chemical engineering; astrophysics; biology; meteorology; and oceanography [3, 9]. Multiphysics multiscale systems are characterized by multiple simultaneous physical processes evolving at very different time scales. Processes can be informally categorized according to their dynamics into fast (stiff) and slow (nonstiff). For example, consider advection-diffusion-reaction systems where the advection is slow while the diffusion and chemistry are typically fast [14, 27, 30].

The dynamics of a process determine the best numerical solution strategy. Explicit time discretizations are effective for slow processes because their computational cost per step is relatively low. On the other hand, implicit methods are more efficient for fast processes because their step sizes are typically not limited by stability considerations [16, 17]. Time integration of multiscale processes is challenging because neither purely explicit nor purely implicit methods are adequate. Explicit methods require prohibitively small time steps (limited by the fastest time scale in the system). Implicit methods require the solution of (non)linear systems of equations that involve *all* the processes in the model; this is both computationally expensive and difficult to implement [18, 23].

The implicit-explicit (IMEX) approach has been developed to alleviate these difficulties. The IMEX idea is to combine an implicit scheme for the stiff components with an explicit scheme for the nonstiff components such that the overall discretization method has the desired stability and accuracy. IMEX linear multistep (LM) methods have been proposed in [3, 13, 20], and IMEX Runge–Kutta (RK) schemes have been developed in [2, 5, 25, 31]. These methods are generally limited to low orders of consistency (typically, lower than five). High-order IMEX RK methods are difficult to construct because of the large number of order conditions. IMEX LM methods have increasing stability restrictions with increasing order.

In this study we construct a new family of IMEX methods using extrapolation. We are concerned with solving the following problem:

$$(1.1) \qquad y'(x) = F(x,y), \quad F(x,y) = f(x,y) + g(x,y), \quad x \geq x_0, \ y(x_0) = y_0,$$

where $f$ and $g$ represent the nonstiff and the stiff processes, respectively. Our approach is to apply an explicit time discretization to $f$ and an implicit time discretization to $g$ and to achieve high orders of consistency by extrapolation [11, 18, 19]. This paper extends the pioneering work of Deuflhard, Hairer, and Zugck [11, 12] on extrapolated linearly implicit and midpoint rule to extrapolated IMEX methods.

This work brings the following contributions. We propose a new family of IMEX methods that are simple to implement, can attain very high orders of convergence, and are parallelizable. We investigate their linear stability properties and prove the existence of perturbed asymptotic expansions for the global discretization errors. We illustrate the theoretical findings on ODE, DAE, and PDE test problems.

The paper is organized as follows. Section 2 offers a review of the extrapolation methods along with their consistency and linear stability properties. The existence of an asymptotic error expansion for the extrapolated IMEX methods applied to index-1 DAEs is established in section 3 and is illustrated in section 4. The global error expansion results are extended to stiff ODEs in section 5, and in section 6 we show numerical evidence that supports the theory. Section 7 presents numerical results for a PDE system. We discuss practical implementation aspects in section 8 and draw conclusions in section 9.

**2. Extrapolation methods.** Consider a sequence of positive integers $\{n_j\}_{1 \leq j \leq M}$, with $n_j < n_{j+1}$, and define a sequence of step sizes $h_1, h_2, h_3, \ldots$ by $h_j = H/n_j$. Further, consider a "base" numerical method to solve (1.1), and denote by $y_h(x)$ the numerical approximation of $y(x)$ obtained with the step size $h$. Different numerical solutions $\{T_{j,1}\}_{1 \leq j \leq M}$ at $x_0 + H$ are obtained by applying $n_j$ steps of the base method with step size $h_j$:

$$(2.1) \qquad T_{j,1} := y_{h_j}(x_0 + H) \quad 1 \leq j \leq M. \qquad \text{[base method]}$$

Assume that the *global error* of the *p*th-order base method employed in (2.1) has an asymptotic expansion of the form

$$(2.2) \qquad y(x) - y_h(x) = e_p(x)\, h^p + \cdots + e_N(x)\, h^N + E_h(x)\, h^{N+1},$$

where $e_i(x)$ are errors that do not depend on $h$, and $E_h$ is bounded for $x_0 \leq x \leq x_{\text{end}}$. This holds for the methods discussed in this paper (see section 2.1). Using the $M$ approximations (2.1) obtained with different $h_j$'s, one can eliminate the error terms in the global error asymptotic expansion (2.2) by Richardson extrapolation (see [18, Chap. II.9]). High-order numerical approximations of the solution of (1.1) can thus

| (a) $T_{j,k}$ tableau | | | (b) Classical orders | | |
|---|---|---|---|---|---|
| $T_{11}$ | | | $p$ | | |
| $T_{21}$ | $T_{22}$ | | $p$ | $p+1$ | |
| $T_{31}$ | $T_{32}$ | $T_{33}$ | $p$ | $p+1$ | $p+2$ |

be constructed [18, Chap. II, Theorem 9.1]. The most economical approach is given by the Aitken–Neville formula [1, 24]:

$$(2.3) \qquad T_{j,k+1} = T_{j,k} + \frac{T_{j,k} - T_{j-1,k}}{(n_j/n_{j-1}) - 1}, \quad j \le M, \quad k < j.$$

The numerical scheme (2.1)–(2.3) is called the *extrapolation method*. It is customary to represent the solutions $T_{j,k}$ in a tableau; see, for example, Table 1(a). We remark that the extrapolation approach provides a sequence of lower-order embedded methods as illustrated in Table 1(b). This fact can be used for step size ($H$) and order control. The most efficient choice for $n_j$ is the harmonic sequence [10]: $n_j = 1, 2, 3, \ldots$.

**2.1. Base methods.** Typical base methods used to compute (2.1) include the forward Euler method $y^{n+1} = y^n + h\ (f(y^n) + g(y^n))$ and the linearly implicit Euler method [19]:

(2.4a)

$$y^{n+1} = y^n + \left[I - h\ (f + g)'(y^n)\right]^{-1} \left(h\ f(y^n) + h\ g(y^n)\right). \quad \text{[linearly implicit]}$$

Method (2.4a) has been used in [11, 12] as the extrapolation base method for solving stiff ODEs of type (1.1). Explicit Euler and the symmetric base methods are possible but not addressed in this study.

In this paper we extend the work of Deuflhard, Hairer, and Zugck [12] to problems that have both fast and slow components, such as in (1.1). We treat implicitly the fast components and explicitly the slow ones, and we build IMEX extrapolation schemes. To this end we propose three IMEX base methods. *W-IMEX*, *pure-IMEX*, and *split-IMEX* are defined as follows:

(2.4b)

$$y^{n+1} = y^n + \left[I - h\ g'(y^n)\right]^{-1} \left(h\ f(y^n) + h\ g(y^n)\right), \quad \text{[W-IMEX]}$$

(2.4c)

$$y^{n+1} = y^n + h\ f(y^n) + \left[I - h\ g'(y^n)\right]^{-1} \left(h\ g(y^n)\right), \quad \text{[pure-IMEX]}$$

(2.4d)

$$y^{n+1} = y^* + \left[I - h\ g'(y^n)\right]^{-1} \left(h\ g(y^*)\right); \; y^* = y^n + h\ f(y^n). \quad \text{[split-IMEX]}$$

The W-IMEX scheme is essentially the same as the linearly implicit method except for the Jacobian, which is approximated by the Jacobian of the stiff component; this is

typically sufficient for the stability of the numerical algorithm and makes the W-IMEX method computationally cheaper. The pure-IMEX and the split-IMEX schemes use the same approximation of the Jacobian; however, the explicit and implicit parts are treated separately, making them truly IMEX schemes.

The extrapolation of methods (2.4b)–(2.4d) can be shown to be *consistent* for nonstiff problems by following [18, Chap. II.8]; more details can be found in [8]. Their consistency for problems with both stiff and nonstiff components is further analyzed in later sections.

**2.2. Linear stability analysis of the extrapolated IMEX methods.** In this section we investigate the linear stability properties of the proposed extrapolated IMEX methods. A similar analysis can be found in [13]. Consider the linear scalar test problem

$$(2.5) \qquad\qquad y'(t) = \lambda y(t) + \mu y(t) \,,$$

where $\lambda$, $\mu \in \mathbb{C}$ represent the eigenvalues of the nonstiff ($f$) and stiff ($g$) components, respectively. The test problem corresponds to the case where the nonstiff and stiff Jacobians can be simultaneously diagonalized, but it also provides useful insight into the general case where they do not [13].

One step of (2.4) applied to (2.5) gives the solution $y^{n+1} = R(z, w)\, y^n$, where $z = \lambda h$, $w = \mu h$, and $R(z, w)$ is the stability function of the method. The stability functions of the base methods (2.4) are as follows:

$$(2.6a) \qquad\qquad R(z, w) = \frac{1+z}{1-w} \begin{bmatrix} \text{W-IMEX,} \\ \text{split-IMEX} \end{bmatrix} ,$$

$$(2.6b) \qquad\qquad R(z, w) = z + \frac{1}{1-w} \ [\text{pure-IMEX}] \,.$$

The W-IMEX and the split-IMEX methods have the same stability function.

The stability functions of the extrapolated methods are calculated from the extrapolation formula (2.3) as [19, Chap. IV]

$$R_{j,1}(z, w) = R^{n_j}\left(\frac{z}{n_j}, \frac{w}{n_j}\right) , \quad R_{j,k+1}(z, w) = R_{j,k}(z, w) + \frac{R_{j,k}(z, w) - R_{j-1,k}(z, w)}{(n_j/n_{j-k}) - 1} ,$$

where $z = \lambda H$, $w = \mu H$, and $R(z, w)$ is the one-step stability function for the specific base method. The subscripts denote the position in the extrapolation tableau.

The stability region $\mathcal{S}$ of the IMEX method is defined by

$$\mathcal{S} = \{z \in \mathbb{C},\, w \in \mathbb{C} : |R(z, w)| \leq 1\} \subset \mathbb{C} \times \mathbb{C} \,.$$

This definition is of little practical consequence, however, as the set in $\mathbb{C} \times \mathbb{C}$ is difficult to visualize. Therefore, to assess the linear stability, we explore a nonstiff stability region $\mathcal{S}_z \subset \mathbb{C}$ and a stiff stability region $\mathcal{S}_w \subset \mathbb{C}$ such that $\mathcal{S}_z \times \mathcal{S}_w \subset \mathcal{S}$. The method is stable whenever $\lambda h \in \mathcal{S}_z$ for the nonstiff component and $\mu h \in \mathcal{S}_w$ for the stiff one. The choice of $\mathcal{S}_z$ and $\mathcal{S}_w$ is not unique. The two sets can be interpreted as "regular" stability regions of the explicit and of the implicit parts of the IMEX method, respectively.

Desirable stability properties for implicit methods are $A$-stability or $A(\alpha)$-stability [19]. To assess the stability of the IMEX method, we consider the stiff stability regions $\mathcal{S}_w(\alpha) = \{w = \rho e^{i\theta} : \pi - \alpha \leq \theta \leq \pi + \alpha,\ \rho \geq 0\}$ (i.e., the $\alpha$-wedge in the negative half

FIG. 1. *Stiff stability regions $\mathcal{S}_w(\alpha)$ for $\alpha = 90°$ and $30°$ and the corresponding nonstiff stability regions $\mathcal{S}_z$ for several extrapolated IMEX terms with the base methods (2.4). The horizontal and vertical axes represent the real and imaginary components, respectively. We note that a different scaling is used for the right column.*

plane is characteristic for $A(\alpha)$ stability) and determine the maximal nonstiff stability regions $\mathcal{S}_z$ such that $\mathcal{S}_z \times \mathcal{S}_w \subset \mathcal{S}$. To be specific, we consider two stiff stability regions for angle values $\alpha = 90°, 30°$ as shown in Figure 1(left). For each of them the (maximal) nonstiff stability regions are computed for several entries $T_{jk}$ in the extrapolation tableau. These nonstiff stability regions are reported in Figure 1(middle) for the W-IMEX and split-IMEX schemes and in Figure 1(right) for the pure-IMEX scheme. In each case the nonstiff stability regions are nontrivial. They contain a segment of the imaginary axis, which is a desirable property when solving certain PDEs by the method of lines [21]. The explicit stability regions grow for methods $T_{jk}$ farther down the extrapolation tableau (i.e., grow with increasing $j$ and $k$).

Decreasing the stiff stability requirement $\mathcal{S}_w(\alpha)$ by decreasing $\alpha$ leads to an increase in the nonstiff stability region and relaxes the step-size restriction for the entire IMEX method. In many important applications the fast process (diffusion, chemistry) has large eigenvalues $\mu$ on or close to the negative real axis; this property allows relatively large time steps for the entire IMEX method.

Next we turn our attention to the accuracy of the extrapolated IMEX methods.

**3. Global error expansion for extrapolated IMEX methods applied to DAEs.** Consider problems (1.1) with the following special structure. A change of variables exists that splits the solution vector into a purely slow component $y$ (driven

by the slow process $f$) and a purely fast component $z$ (driven by the fast process $g$). We have

$$(3.1) \qquad \begin{pmatrix} y \\ \varepsilon\, z \end{pmatrix}' = \begin{pmatrix} f(y,z) \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ g(y,z) \end{pmatrix} = \begin{pmatrix} f(y,z) \\ g(y,z) \end{pmatrix}.$$

The constant $\varepsilon$ represents the ratio of the fast to slow timescales, determines the stiffness, and provides an appropriate problem scaling. The system (3.1) is a singular perturbation problem (SPP) with the reduced differential algebraic form obtained by taking $\varepsilon \to 0$:

$$(3.2) \qquad \begin{pmatrix} y \\ 0 \end{pmatrix}' = \begin{pmatrix} f(y,z) \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ g(y,z) \end{pmatrix} = \begin{pmatrix} f(y,z) \\ g(y,z) \end{pmatrix}.$$

Under the assumption that $g_z$ is invertible, the system (3.2) is an index-1 DAE.

We first analyze the accuracy of the extrapolated IMEX methods applied to the reduced system (3.2) and then address the discretization of the full problem (3.1). The discussion starts with W-IMEX in section 3.1, continues with pure-IMEX in section 3.2, and focus us on split-IMEX in section 3.3.

**3.1. W-IMEX.** Applying the W-IMEX method (2.4b) to (3.1) yields

$$(3.3) \qquad \begin{pmatrix} I & 0 \\ -hg_y(0) & \varepsilon I - hg_z(0) \end{pmatrix} \begin{pmatrix} y_{i+1} - y_i \\ z_{i+1} - z_i \end{pmatrix} = h \begin{pmatrix} f(y_i, z_i) \\ g(y_i, z_i) \end{pmatrix},$$

where $g_{\{y,z\}}(0) = g_{\{y,z\}}(y_0, z_0)$. Then the reduced form of (3.3) given by $\varepsilon \to 0$ is

$$(3.4) \qquad \begin{pmatrix} I & 0 \\ -hg_y(0) & -hg_z(0) \end{pmatrix} \begin{pmatrix} y_{i+1} - y_i \\ z_{i+1} - z_i \end{pmatrix} = h \begin{pmatrix} f(y_i, z_i) \\ g(y_i, z_i) \end{pmatrix}.$$

To assess the accuracy of the W-IMEX scheme, we first analyze the reduced system (3.4) and then address the full problem (3.3) in section 5.1. Similar work for the extrapolated linearly implicit Euler method can be found in [12] and [19, Chap. VI.5]. We start with the reduced problem and give the following result.

THEOREM 3.1 (global error expansion of the extrapolated W-IMEX method applied to DAEs). *Consider the problem (3.2) with $g_z$ invertible and consistent initial values $(y_0, z_0)$. The global error of the IMEX scheme (3.4) has an asymptotic $h$-expansion of the form*

$$(3.5) \qquad \begin{aligned} y_i - y(x_i) &= \sum_{j=1}^{M} h^j \left( a^{(j)}(x_i) + \alpha_i^{(j)} \right) + \mathcal{O}\left(h^{M+1}\right), \\ z_i - z(x_i) &= \sum_{j=1}^{M} h^j \left( b^{(j)}(x_i) + \beta_i^{(j)} \right) + \mathcal{O}\left(h^{M+1}\right), \end{aligned}$$

*where $a^{(j)}(x)$ and $b^{(j)}(x)$ are smooth functions and the perturbations satisfy*

$$(3.6a) \qquad\qquad \alpha_i^{(1)} = 0, \quad \alpha_i^{(2)} = 0, \quad \beta_i^{(1)} = 0 \quad \forall i \geq 0\,;$$

$$(3.6b) \qquad\qquad \alpha_i^{(3)} = 0, \quad \alpha_i^{(4)} = 0, \quad \beta_i^{(2)} = 0 \quad \forall i \geq 1\,;$$

$$(3.6c) \qquad\qquad \alpha_i^{(j)} = 0 \quad \forall i \geq j - 3, \quad j \geq 5\,;$$

$$(3.6d) \qquad\qquad \beta_i^{(j)} = 0 \quad \forall i \geq j - 2, \quad j \geq 3\,.$$

*The error terms in* (3.5) *are uniformly bounded for* $x_i = ih \leq H$ *if* $H$ *is sufficiently small.*

*Proof.* Following Deuflhard, Hairer, and Zugck [12], the proof consists of two parts: in the first part (a) truncated expansions are constructed, and in the second one (b) an error bound is obtained from a stability estimate.

(a) Consider the truncated expansions of the numerical solution

$$(3.7) \quad \widehat{y}_i = y(x_i) + \sum_{j=1}^{M} h^j \left( a^{(j)}(x_i) + \alpha_i^{(j)} \right) \; ; \; \widehat{z}_i = z(x_i) + \sum_{j=1}^{M} h^j \left( b^{(j)}(x_i) + \beta_i^{(j)} \right) ,$$

such that the defect of $\widehat{y}_i$, $\widehat{z}_i$ inserted in (3.4) is small (see [15]):

$$(3.8) \quad \begin{pmatrix} I & 0 \\ -hg_y(0) & -hg_z(0) \end{pmatrix} \begin{pmatrix} \widehat{y}_{i+1} - \widehat{y}_i \\ \widehat{z}_{i+1} - \widehat{z}_i \end{pmatrix} = h \begin{pmatrix} f(\widehat{y}_i, \widehat{z}_i) \\ g(\widehat{y}_i, \widehat{z}_i) \end{pmatrix} + \mathcal{O}\left(h^{M+2}\right) .$$

The initial values are the exact solution ($\widehat{y}_0 = y_0$, $\widehat{z}_0 = z_0$), and the perturbation terms ($\alpha$, $\beta$) are assumed to satisfy

$$(3.9a) \qquad\qquad a^{(j)}(0) + \alpha_0^{(j)} = 0 , \qquad b^{(j)}(0) + \beta_0^{(j)} = 0 ,$$

$$(3.9b) \qquad\qquad \alpha_i^{(j)} \to 0 , \qquad \beta_i^{(j)} \to 0 \qquad \text{for } i \to \infty .$$

Consider the Taylor expansion for $f(\widehat{y}_i, \widehat{z}_i)$ and $g(\widehat{y}_i, \widehat{z}_i)$ about $(y(x_i), z(x_i))$. Replacing them together with the expansion of their numerical solutions $\widehat{y}_{i+1} - \widehat{y}_i$ and $\widehat{z}_{i+1} - \widehat{z}_i$ in (3.8) and equating the terms in $h$, we get

$$(3.10) \qquad y'(x_i) + \left( \alpha_{i+1}^{(1)} - \alpha_i^{(1)} \right) = f\left(y(x_i), z(x_i)\right) , \quad 0 = g\left(y(x_i), z(x_i)\right) .$$

Using the consistency requirement (3.9b) gives $\alpha_{i+1}^{(1)} = \alpha_i^{(1)}$, which verifies (3.2). Thus one has $\alpha_i^{(1)} = 0 \; \forall i \geq 0$. Next we consider the coefficients of $h^2$ and obtain

$$(3.11a) \qquad \frac{1}{2} y''(x) + \left( a^{(1)} \right)' (x) = f_y(x) a^{(1)}(x) + f_z(x) b^{(1)}(x),$$

$$(3.11b) \qquad - g_y(0) y'(x) - g_z(0) z'(x) = g_y(x) a^{(1)}(x) + g_z(x) b^{(1)}(x),$$

$$(3.11c) \qquad \left( \alpha_{i+1}^{(2)} - \alpha_i^{(2)} \right) = f_z(0) \beta_i^{(1)}, \quad -g_z(0) \left( \beta_{i+1}^{(1)} - \beta_i^{(1)} \right) = g_z(0) \beta_i^{(1)}.$$

System (3.11a)–(3.11b) can be solved by substituting $b^{(1)}(x)$ from (3.11b) in (3.11a), which leads to an ODE in $a^{(1)}$, and together with (3.9a) and $\alpha_0^{(1)} = 0$ gives $a^{(1)}(0) = 0$. Therefore $a^{(1)}(x)$ and $b^{(1)}(x)$ are uniquely determined by (3.11a)–(3.11b).

Relations (3.11c) and $0 = g(y, z)$ for $x = 0$ are used to eliminate the left-hand side of (3.11b): $g_y(0) a^{(1)}(0) + g_z(0) b^{(1)}(0) = 0 \Rightarrow g_z(0) b^{(1)}(0) = 0 \Rightarrow b^{(1)}(0) = 0$.

By (3.9a) one has $\beta_0^{(1)} = 0$. In general $\beta_i^{(1)} = 0 \; \forall i \geq 0$, from (3.11c), which gives $\alpha_i^{(2)} = 0 \; \forall i \geq 0$.

The coefficients of $h^3$ lead to

$$(3.12a) \qquad \left( a^{(2)} \right)' (x) = f_y(x) a^{(2)}(x) + f_z(x) b^{(2)}(x) + r^{(2)}(x) ,$$

$$(3.12b) \qquad 0 = g_y(x) a^{(2)}(x) + g_z(x) b^{(2)}(x) + s^{(2)}(x) ,$$

where $r^{(2)}(x)$ and $s^{(2)}(x)$ are known functions that depend on the derivatives of $y(x)$, $z(x)$, $a^{(1)}(x)$, $b^{(1)}(x)$. The perturbations with the additional cancellations of terms that have coefficients $\alpha_i^{(1)} = 0$ and $\beta_i^{(1)} = 0$ $\forall i$, and using $\alpha_i^{(2)} = 0$ $\forall i$, lead to

$$(3.13a) \qquad \alpha_{i+1}^{(3)} - \alpha_i^{(3)} = f_z(0)\beta_i^{(2)},$$

$$(3.13b) \qquad 0 = g_z(0)\beta_{i+1}^{(2)}.$$

Terms $a^{(2)}(x)$ and $b^{(2)}(x)$ are determined in the same way as $a^{(1)}(x)$ and $b^{(1)}(x)$. One has $a^{(2)}(0) = 0$ from $\alpha_i^{(2)} = 0$. However, $b^{(2)}(0) \neq 0$, and by (3.9a) one has $\beta_0^{(2)} \neq 0$. From (3.13b) one obtains $\beta_i^{(2)} = 0$ $\forall i \geq 1$, and with (3.13a) one has $\alpha_i^{(3)} = 0$ $\forall i \geq 1$.

A recurrence formula can be constructed for the coefficients of $h^{j+1}$ $\forall j \geq 4$:

$$(3.14a) \qquad \left(a^{(j)}\right)'(x) = f_y(x)\,a^{(j)}(x) + f_z(x)\,b^{(j)}(x) + r^{(j)}(x),$$

$$(3.14b) \qquad 0 = g_y(x)\,a^{(j)}(x) + g_z(x)\,b^{(j)}(x) + s^{(j)}(x),$$

$$(3.14c) \qquad \alpha_{i+1}^{(j+1)} - \alpha_i^{(j+1)} = f_z(0)\beta_i^{(j)} + \varrho_i^{(j)},$$

$$(3.14d) \qquad 0 = g_z(0)\beta_{i+1}^{(j)} + \sigma_i^{(j)},$$

where $\varrho_i^{(j)}$ and $\sigma_i^{(j)}$ are linear combinations of expressions that contain as factors $\alpha_{i+1}^{(\ell)}$, $\alpha_{i+1}^{(\ell-1)}$, $\beta_{i+1}^{(\ell-1)}$, $\ell \leq j$. By induction on $j$ with $\varrho_i^{(j)} = 0$ and $\sigma_i^{(j)} = 0$, $i \geq j - 3$, one can show that (3.14d) implies that $\beta_{i+1}^{(j)} = 0$, $i \geq j - 3$. Then relations (3.9b) and (3.14c) give $\alpha_{i+1}^{(j+1)} = 0$, $i \geq j - 3$. This concludes the proof for (3.6c)–(3.6d).

(b) The second part of this proof consists of estimating a bound on the reminder term; that is, differences $\Delta y_i = y_i - \widehat{y}_i$ and $\Delta z_i = z_i - \widehat{z}_i$. Subtracting (3.8) from (3.4) and eliminating $\Delta y_i$ and $\Delta z_i$, we have

$$\begin{pmatrix} \Delta y_{i+1} \\ \Delta z_{i+1} \end{pmatrix} = \begin{pmatrix} \Delta y_i \\ \Delta z_i \end{pmatrix} + \begin{pmatrix} I & 0 \\ -g_y(0) & -g_z(0) \end{pmatrix}^{-1} \begin{pmatrix} h\left(f\left(y_i, z_i\right) - f\left(\widehat{y}_i, \widehat{z}_i\right)\right) \\ g\left(y_i, z_i\right) - g\left(\widehat{y}_i, \widehat{z}_i\right) \end{pmatrix}$$
$$+ \begin{pmatrix} I & 0 \\ -g_y(0) & -g_z(0) \end{pmatrix}^{-1} \begin{pmatrix} \mathcal{O}\left(h^{M+2}\right) \\ \mathcal{O}\left(h^{M+1}\right) \end{pmatrix} = \begin{pmatrix} \Delta y_i \\ \Delta z_i \end{pmatrix}$$
$$+ \begin{pmatrix} I & 0 \\ \mathcal{O}(1) & -g_z(0)^{-1} \end{pmatrix} \begin{pmatrix} h\left(f\left(y_i, z_i\right) - f\left(\widehat{y}_i, \widehat{z}_i\right)\right) \\ g\left(y_i, z_i\right) - g\left(\widehat{y}_i, \widehat{z}_i\right) \end{pmatrix} + \begin{pmatrix} \mathcal{O}\left(h^{M+2}\right) \\ \mathcal{O}\left(h^{M+1}\right) \end{pmatrix}.$$

The application of the Lipschitz condition on $f(y, z)$ and $g(y, z)$, $|\zeta| < 1$ gives

$$(3.15) \qquad \begin{pmatrix} \|\Delta y_{i+1}\| \\ \|\Delta z_{i+1}\| \end{pmatrix} \leq \begin{pmatrix} I & 0 \\ \mathcal{O}(1) & \zeta \end{pmatrix} \begin{pmatrix} \|\Delta y_i\| \\ \|\Delta z_i\| \end{pmatrix} + \begin{pmatrix} \mathcal{O}\left(h^{M+2}\right) \\ \mathcal{O}\left(h^{M+1}\right) \end{pmatrix},$$

with $H$ is sufficiently small. Using [8, Lemma C1] gives $\|\Delta y_i\| + \|\Delta z_i\| = \mathcal{O}\left(h^{M+1}\right)$. $\square$

Next we investigate the extrapolation orders using the base method (3.4) and define

$$(3.16) \qquad Y_{jk} = y_{h_j}\left(x_0 + H\right), \qquad Z_{jk} = z_{h_j}\left(x_0 + H\right)$$

to be the numerical solution of (3.2) after $j$ steps with step size $h_j = H/n_j$, extrapolated with (2.3); in other words, on the $k$th column of the extrapolation tableau.

We note that each extrapolation step (2.3) cancels one smooth term $(\{a,b\}^{(j)})$ from the error expansion (3.5); however, perturbations $\alpha$ and $\beta$ propagate through the extrapolation steps (2.3). Furthermore, we note that the accuracy of the solution on the extrapolation tableau diagonal is affected by terms $\{\alpha, \beta\}_1^{(j)}$, and nonzero smooth terms $a(0)$ and $b(0)$ affect the perturbations $\alpha_0$ and $\beta_0$ through (3.9a); for example, $b^{(2)}(0) \neq 0 \Rightarrow \beta_0^{(2)} \neq 0$.

We prove the following result. Similar approaches are found in [19, Chap. VI, Theorem 5.4] and [12].

THEOREM 3.2 (consistency of the extrapolated W-IMEX applied to DAEs). *For the harmonic sequence $\{1, 2, 3, \ldots\}$ the extrapolated values $Y_{jk}$ and $Z_{jk}$ satisfy*

$$(3.17) \qquad Y_{jk} - y(x_0 + h) = \mathcal{O}\left(H^{r_{jk}}\right), \qquad Z_{jk} - z(x_0 + h) = \mathcal{O}\left(H^{s_{jk}}\right),$$

*where the differential and algebraic orders $r_{jk}$ and $s_{jk}$ are given in Table 2 up to $j = 12$, $k = 12$.*

*Proof.* We use the expansion (3.5). It follows from (3.6a) and (3.9a) that $a(x_0) = 0$ and $b(x_0) = 0$. Since $a^{(j)}(x)$ and $b^{(j)}(x)$ are smooth functions, one obtains $a^{(1)}(x_0 + H) = \mathcal{O}(H)$, and $b^{(1)}(x_0 + H) = \mathcal{O}(H)$, and thus the errors in $Y_{j1}$ and $Z_{j1}$ are of $\mathcal{O}(H^2)$, which gives the first column entries in Table 2 for the W-IMEX scheme. In the same way one can deduce that $a^{(2)}(x_0 + h) = \mathcal{O}(H)$; however, since $\beta_0^{(2)} \neq 0$, by (3.9a) one obtains $b^{(2)}(0) \neq 0$ (in general) and $b^{(2)}(x_0 + h) = \mathcal{O}(1)$. One extrapolation of the numerical method eliminates the terms with $j = 1$ in (3.5). The error is thus $\mathcal{O}(H^3)$ for $Y_{j2}$ and $\mathcal{O}(H^2)$ for $Z_{j2}$. Equivalently, one can expand (3.5) to

$$\begin{cases} y_1 - y(x_1) = h^1 \left(a^{(1)}(x_1) + \alpha_1^{(1)}\right) + h^2 \left(a^{(2)}(x_1) + \alpha_1^{(2)}\right) + \cdots = \mathcal{O}\left(H^2\right), \\ z_1 - z(x_1) = h^1 \left(b^{(1)}(x_1) + \beta_1^{(1)}\right) + h^2 \left(b^{(2)}(x_1) + \beta_1^{(2)}\right) + \cdots = \mathcal{O}\left(H^2\right). \end{cases}$$

However, for $j = 2$ one has $a^{(2)}(x_0 + h) = \mathcal{O}(H)$ and $b^{(2)}(x_0 + h) = \mathcal{O}(1)$, and thus

$$y_1 - y(x_1) = h^2 \left(\mathcal{O}(H)\right) + \cdots = \mathcal{O}\left(H^3\right); \; z_1 - z(x_1) = h^2 \left(\mathcal{O}(1)\right) + \cdots = \mathcal{O}\left(H^2\right).$$

The smooth parts of (3.5) are eliminated one by one; however, the perturbations are not, and the orders are reduced as follows. One order is "lost" on columns $y_{j3}$ and $z_{j2}$ because of $\mathcal{O}(1)$ smooth part expansion; however, thereafter the orders are increasing by using the extrapolation formula (2.3) that cancels the smooth terms. The nonzero perturbation terms affect the orders of the extrapolation method by propagating through (2.3). Specifically, for $y_{jk}$ components one has $\alpha_1^{(5)} \neq 0$, which limits the order on the diagonal for $y_{jj}$, $j \geq 6$ to 4. Using the same argument, one can show that the first subdiagonal $y_{j\,j-1}$, $j \geq 8$, is limited to 5 and the second one $y_{j\,j-2}$, $j \geq 10$, is limited to 6 from $\alpha_2^{(6)} \neq 0$ and $\alpha_3^{(7)} \neq 0$, respectively, and so on. Similarly, for $z_{jk}$ components one has $z_{jj}$, $j \geq 5$ to 3; $z_{j\,j-1}$, $j \geq 7$ to 4; and $z_{j\,j-2}$, $j \geq 9$ to 5, because of $\beta_1^{(4)} \neq 0$, $\beta_2^{(5)} \neq 0$, and $\beta_3^{(6)} \neq 0$, respectively. This process can be continued to find all the entries in Table 2. $\Box$

Of particular interest is the location of the maximum accuracy term in the extrapolation tableau for a given number of steps $j$. A quick inspection of Table 2 reveals that the best choice is $T_{j,j}$ for $j \leq 4$; $T_{j,(j-1)/2+3}$ for $j \geq 5$ and odd; and $T_{j,j/2+2}$ for $j \geq 6$ and even. Boldface fonts are used to identify the location of the most accurate yielding extrapolation tableau term. The theoretical orders for the extrapolated

TABLE 2

*Theoretical local extrapolation orders for linearly implicit, W-IMEX, pure-IMEX, and split-IMEX methods for index-1 DAEs. Boldface fonts represent the "best" or optimal choice for a given number of steps.*

Orders $(r_{jk})$ for component $y_{jk}$ for linearly implicit|W-IMEX|pure-IMEX|split-IMEX

| | $a^{(1)}(\cdot)$ | $a^{(2)}(\cdot)$ | $a^{(3)}(\cdot)$ | $a^{(4)}(\cdot)$ | $a^{(5)}(\cdot)$ | $a^{(6)}(\cdot)$ | $a^{(7)}(\cdot)$ | $a^{(8)}(\cdot)$ | $a^{(9)}(\cdot)$ | $a^{(10)}(\cdot)$ | $a^{(11)}(\cdot)$ | $a^{(12)}(\cdot)$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | **2\|2\|2\|2** | | | | | | | | | | | |
| 2 | 2\|2\|2\|2 | **3\|3\|2\|3** | | | | | | | | | | |
| 3 | 2\|2\|2\|2 | 3\|3\|2\|3 | **4\|3\|3\|3** | | | | | | | | | |
| 4 | 2\|2\|2\|2 | 3\|3\|2\|3 | 4\|3\|3\|3 | **5\|4\|3\|4** | | | | | | | | |
| 5 | 2\|2\|2\|2 | 3\|3\|2\|3 | 4\|3\|3\|3 | 5\|4\|4\|4 | **5\|5\|4\|5** | | | | | | | |
| 6 | 2\|2\|2\|2 | 3\|3\|2\|3 | 4\|3\|3\|3 | 5\|4\|4\|4 | 5\|5\|5\|5 | **6\|6\|4\|5** | | | | | | |
| 7 | 2\|2\|2\|2 | 3\|3\|2\|3 | 4\|3\|3\|3 | 5\|4\|4\|4 | 5\|5\|5\|5 | **6\|6\|5\|6** | 7\|6\|4\|5 | | | | | |
| 8 | 2\|2\|2\|2 | 3\|3\|2\|3 | 4\|3\|3\|3 | 5\|4\|4\|4 | 5\|5\|5\|5 | **6\|6\|6\|6** | **7\|7\|5\|6** | 6\|5\|3\|4 | | | | |
| 9 | 2\|2\|2\|2 | 3\|3\|2\|3 | 4\|3\|3\|3 | 5\|4\|4\|4 | 5\|5\|5\|5 | 6\|6\|6\|6 | **7\|7\|6\|7** | 7\|6\|4\|4 | 6\|5\|3\|4 | | | |
| 10 | 2\|2\|2\|2 | 3\|3\|2\|3 | 4\|3\|3\|3 | 5\|4\|4\|4 | 5\|5\|5\|5 | 6\|6\|6\|6 | **7\|7\|7\|7** | **8\|7\|5\|6** | 7\|6\|4\|5 | 6\|5\|3\|4 | | |
| 11 | 2\|2\|2\|2 | 3\|3\|2\|3 | 4\|3\|3\|3 | 5\|4\|4\|4 | 5\|5\|5\|5 | 6\|6\|6\|6 | 7\|7\|7\|7 | **8\|8\|6\|7** | 8\|7\|5\|6 | 7\|6\|4\|5 | 6\|5\|3\|4 | |
| 12 | 2\|2\|2\|2 | 3\|3\|2\|3 | 4\|3\|3\|3 | 5\|4\|4\|4 | 5\|5\|5\|5 | 6\|6\|6\|6 | 7\|7\|7\|7 | **8\|8\|7\|8** | 9\|8\|6\|7 | 8\|7\|5\|6 | 7\|6\|4\|5 | 6\|5\|3\|4 |
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |

Orders $(s_{jk})$ for component $z_{jk}$ for linearly implicit|W-IMEX|pure-IMEX|split-IMEX

| | $b^{(1)}(\cdot)$ | $b^{(2)}(\cdot)$ | $b^{(3)}(\cdot)$ | $b^{(4)}(\cdot)$ | $b^{(5)}(\cdot)$ | $b^{(6)}(\cdot)$ | $b^{(7)}(\cdot)$ | $b^{(8)}(\cdot)$ | $b^{(9)}(\cdot)$ | $b^{(10)}(\cdot)$ | $b^{(11)}(\cdot)$ | $b^{(12)}(\cdot)$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | **2\|2\|1\|2** | | | | | | | | | | | |
| 2 | 2\|2\|1\|2 | **2\|2\|2\|2** | | | | | | | | | | |
| 3 | 2\|2\|1\|2 | 2\|2\|2\|2 | **3\|3\|2\|3** | | | | | | | | | |
| 4 | 2\|2\|1\|2 | 2\|2\|2\|2 | 3\|3\|3\|3 | **4\|4\|2\|3** | | | | | | | | |
| 5 | 2\|2\|1\|2 | 2\|2\|2\|2 | 3\|3\|3\|3 | 4\|4\|3\|4 | **4\|4\|2\|3** | | | | | | | |
| 6 | 2\|2\|1\|2 | 2\|2\|2\|2 | 3\|3\|3\|3 | 4\|4\|4\|4 | **5\|5\|3\|4** | 4\|4\|2\|3 | | | | | | |
| 7 | 2\|2\|1\|2 | 2\|2\|2\|2 | 3\|3\|3\|3 | 4\|4\|4\|4 | 5\|5\|4\|5 | **5\|5\|3\|4** | 4\|4\|2\|3 | | | | | |
| 8 | 2\|2\|1\|2 | 2\|2\|2\|2 | 3\|3\|3\|3 | 4\|4\|4\|4 | 5\|5\|5\|5 | **6\|6\|4\|5** | 5\|5\|3\|4 | 4\|4\|2\|3 | | | | |
| 9 | 2\|2\|1\|2 | 2\|2\|2\|2 | 3\|3\|3\|3 | 4\|4\|4\|4 | 5\|5\|5\|5 | **6\|6\|5\|6** | **6\|6\|4\|5** | 5\|5\|3\|4 | 4\|4\|2\|3 | | | |
| 10 | 2\|2\|1\|2 | 2\|2\|2\|2 | 3\|3\|3\|3 | 4\|4\|4\|4 | 5\|5\|5\|5 | 6\|6\|6\|6 | **7\|7\|5\|6** | 6\|6\|4\|5 | 5\|5\|3\|4 | 4\|4\|2\|3 | | |
| 11 | 2\|2\|1\|2 | 2\|2\|2\|2 | 3\|3\|3\|3 | 4\|4\|4\|4 | 5\|5\|5\|5 | 6\|6\|6\|6 | 7\|7\|6\|7 | **7\|7\|5\|6** | 6\|6\|4\|5 | 5\|5\|3\|4 | 4\|4\|2\|3 | |
| 12 | 2\|2\|1\|2 | 2\|2\|2\|2 | 3\|3\|3\|3 | 4\|4\|4\|4 | 5\|5\|5\|5 | 6\|6\|6\|6 | 7\|7\|7\|7 | **8\|8\|6\|7** | 7\|7\|5\|6 | 6\|6\|4\|5 | 5\|5\|3\|4 | 4\|4\|2\|3 |
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |

linearly implicit Euler method (2.4a) are also shown [19, 12]. The "best" terms are selected by first identifying the most accurate stiff components and then matching them with the best nonstiff counterparts.

**3.2. Pure-IMEX method.** Applying the pure-IMEX method (2.4c) to (3.1) yields

$$
(3.18) \quad
\begin{pmatrix} I & 0 \\ -hg_y(0) & \varepsilon I - hg_z(0) \end{pmatrix}
\begin{pmatrix} y_{i+1} - y_i \\ z_{i+1} - z_i \end{pmatrix}
= h \begin{pmatrix} f(y_i, z_i) \\ g(y_i, z_i) - hg_y(0)f(y_i, z_i) \end{pmatrix}.
$$

The reduced form given by $\varepsilon \to 0$ is

$$
(3.19) \quad
\begin{pmatrix} I & 0 \\ -hg_y(0) & -hg_z(0) \end{pmatrix}
\begin{pmatrix} y_{i+1} - y_i \\ z_{i+1} - z_i \end{pmatrix}
= h \begin{pmatrix} f(y_i, z_i) \\ g(y_i, z_i) - hg_y(0)f(y_i, z_i) \end{pmatrix}.
$$

We next formulate a similar pair of theorems (i.e., error expansions and extrapolated orders) for the extrapolated pure-IMEX method.

THEOREM 3.3 (global error expansion of the extrapolated pure-IMEX method applied to DAEs). *Consider the problem* (3.2) *with $g_z$ invertible and consistent initial values $(y_0, z_0)$. The global error of the pure-IMEX scheme* (3.19) *has an asymptotic $h$-expansion of the form* (3.5), *where $a^{(j)}(x)$ and $b^{(j)}(x)$ are smooth functions and the perturbations satisfy*

$$
(3.20a) \quad \alpha_i^{(1)} = 0 \quad \forall i \geq 0; \quad \alpha_i^{(2)} = 0, \ \beta_i^{(1)} = 0 \quad \forall i \geq 1; \quad \alpha_i^{(3)} = 0, \ \beta_i^{(2)} = 0 \quad \forall i \geq 2;
$$

$$
(3.20b) \quad \alpha_i^{(j)} = 0 \quad \forall i \geq j - 1, \ j \geq 4; \quad \beta_i^{(j)} = 0 \quad \forall i \geq j, \ j \geq 3.
$$

*The error terms in* (3.5) *are uniformly bounded for $x_i = ih \leq H$ if $H$ is sufficiently small.*

*Proof.* This proof follows the same ideas used in the proof of Theorem 3.1. We begin with part (a) in which the truncated expansions are constructed. The second part follows the same steps as in the W-IMEX method. We focus on the first part only.

We consider the truncated expansions (3.7) with small defects

$$
\begin{pmatrix} I & 0 \\ -hg_y(0) & -hg_z(0) \end{pmatrix}
\begin{pmatrix} \widehat{y}_{i+1} - \widehat{y}_i \\ \widehat{z}_{i+1} - \widehat{z}_i \end{pmatrix}
$$
$$
= h \begin{pmatrix} f(\widehat{y}_i, \widehat{z}_i) \\ g(\widehat{y}_i, \widehat{z}_i) - hg_y(0)f(\widehat{y}_i, \widehat{z}_i) \end{pmatrix} + \mathcal{O}\left(h^{M+1}\right).
$$

The initial values are exact, with perturbation terms satisfying (3.9). By expanding the equation above, one obtains (3.10) for the coefficients of $h$. Using the consistency requirements (3.9b) yields (3.2), and hence $\alpha_i^{(1)} = 0 \ \forall i \geq 0$. The coefficients of $h^2$ give

$$
(3.21a) \quad \frac{1}{2}y''(x) + \left(a^{(1)}\right)'(x) = f_y(x)a^{(1)}(x) + f_z(x)b^{(1)}(x),
$$

$$
(3.21b) \quad -g_y(0)y'(x) - g_z(0)z'(x) + f(x)g_y(0) = g_y(x)a^{(1)}(x) + g_z(x)b^{(1)}(x),
$$

$$
(3.21c) \quad \left(\alpha_{i+1}^{(2)} - \alpha_i^{(2)}\right) = f_z(0)\beta_i^{(1)}, \quad -g_z(0)\left(\beta_{i+1}^{(1)} - \beta_i^{(1)}\right) = g_z(0)\beta_i^{(1)}.
$$

This system can be solved by computing $b^{(1)}(x)$ in (3.21b) and then replacing it in (3.21a) to yield an ODE in $a^{(1)}$. Using (3.9a) and $\alpha_0^{(1)} = 0$, one has $a^{(1)}(0) = 0$. Therefore $a^{(1)}(x)$ and $b^{(1)}(x)$ are uniquely determined by (3.21a)–(3.21b). In contrast

with the W-IMEX method (3.11b), the left-hand side of (3.21b) does not vanish, and hence $b^{(1)}(0) \neq 0$. By (3.9a) one also has $\beta_0^{(1)} \neq 0$. In general $\beta_i^{(1)} = 0 \ \forall i \geq 1$ from (3.21c), and together with (3.9b) one obtains $\alpha_i^{(2)} = 0 \ \forall i \geq 1$.

The coefficients of $h^3$ for the smooth part give

$$(3.22a) \qquad \left(a^{(2)}\right)'(x) = f_y(x)\, a^{(2)}(x) + f_z(x)\, b^{(2)}(x) + r^{(2)}(x)\,,$$

$$(3.22b) \qquad 0 = g_y(x)\, a^{(2)}(x) + g_z(x)\, b^{(2)}(x) + s^{(2)}(x)\,,$$

where $r^{(2)}(x)$ and $s^{(2)}(x)$ are known functions that depend on derivatives of $y(x)$, $z(x)$, $a^{(1)}(x)$, $b^{(1)}(x)$. The perturbations can be expressed as

$$(3.23a) \qquad \alpha_{i+1}^{(3)} - \alpha_i^{(3)} = f_y(0)\alpha_i^{(2)} + \beta_i^{(1)}(\dots) + f_z(0)\beta_i^{(2)}\,,$$

$$(3.23b) \qquad 0 = g_z(0)\beta_{i+1}^{(2)} + \beta_i^{(1)}(\dots) + \alpha_i^{(2)}(\dots)\,.$$

From (3.23), $\beta_i^{(2)} = 0 \ \forall i \geq 2$ and $\alpha_i^{(3)} = 0 \ \forall i \geq 2$. This concludes the proof for hypotheses (3.20a). The general recurrence (3.14) follows. Hypothesis (3.20b) can be easily verified by following the same type of induction on (3.14a)–(3.14b) as in the proof of Theorem 3.1. □

THEOREM 3.4 (consistency of the extrapolated pure-IMEX method applied to DAEs). *For the harmonic sequence $\{1, 2, 3, \dots\}$ the extrapolated values $Y_{jk}$ and $Z_{jk}$ satisfy*

$$(3.24) \qquad Y_{jk} - y(x_0 + h) = \mathcal{O}\left(H^{r_{jk}}\right)\,, \qquad Z_{jk} - z(x_0 + h) = \mathcal{O}\left(H^{s_{jk}}\right)\,,$$

*where the differential and algebraic orders $r_{jk}$ and $s_{jk}$ are given in Table 2.*

*Proof.* The orders in Table 2 for the pure-IMEX method can be recovered by using the same procedure as in the proof of Theorem 3.2 with the error expansion given in Theorem 3.3. The major difference is that now $\alpha_0^{(2)}$ is nonzero, and thus one order is "lost" on the second column of the $y$ component. Then $\alpha_1^{(3)}$ gives the third order on the diagonal. For the $z$ component, $\beta_0^{(1)}$ is nonzero, and hence the first column of the $z$ component is 1. Furthermore, $\beta_1^{(2)}$ does not vanish, and thus $T_{kk}$ has order two for $k \geq 2$. The rest follows from the propagation of the error terms through the extrapolation procedure. □

**3.3. Split-IMEX method.** The split-IMEX method (2.4d) applied to (3.1) yields

$$(3.25)$$
$$\begin{pmatrix} I & 0 \\ -hg_y(0) & \varepsilon I - hg_z(0) \end{pmatrix} \begin{pmatrix} y_{i+1} - y_i \\ z_{i+1} - z_i \end{pmatrix} = h \begin{pmatrix} f(y_i, z_i) \\ g(y_*, z_i) - hg_y(0)f(y_i, z_i) \end{pmatrix},$$

where $y_* = y_i + hf(y_i, z_i)$ and the DAE reduced form given by $\varepsilon \to 0$ is

$$(3.26) \quad \begin{pmatrix} I & 0 \\ -hg_y(0) & -hg_z(0) \end{pmatrix} \begin{pmatrix} y_{i+1} - y_i \\ z_{i+1} - z_i \end{pmatrix} = h \begin{pmatrix} f(y_i, z_i) \\ g(y_*, z_i) - hg_y(0)f(y_i, z_i) \end{pmatrix}.$$

THEOREM 3.5 (global error expansion of the extrapolated split-IMEX method applied to DAEs). *Consider the problem (3.2) with $g_z$ invertible and consistent initial values $(y_0, z_0)$. The global error of the split-IMEX scheme (3.26) has an asymptotic $h$-expansion of the form (3.5), where $a^{(j)}(x)$ and $b^{(j)}(x)$ are smooth functions and the*

*perturbations satisfy*

(3.27a) $$\alpha_i^{(1)} = 0, \quad \alpha_i^{(2)} = 0, \quad \beta_i^{(1)} = 0 \quad \forall i \geq 0;$$

(3.27b) $$\alpha_i^{(3)} = 0, \quad \beta_i^{(2)} = 0 \quad \forall i \geq 1;$$

(3.27c) $$\alpha_i^{(j)} = 0 \quad \forall i \geq j - 2, \quad j \geq 4;$$

(3.27d) $$\beta_i^{(j)} = 0 \quad \forall i \geq j - 1, \quad j \geq 3.$$

*The error terms in* (3.5) *are uniformly bounded for* $x_i = ih \leq H$ *if* $H$ *is sufficiently small.*

*Proof.* This proof follows the same ideas used in Theorem 3.1. We begin with part (a) in which the truncated expansions are constructed. The second part follows the same steps as in the W-IMEX case.

Truncated expansions (3.7) are considered with defects

$$\begin{pmatrix} I & 0 \\ -hg_y(0) & -hg_z(0) \end{pmatrix} \begin{pmatrix} \widehat{y}_{i+1} - \widehat{y}_i \\ \widehat{z}_{i+1} - \widehat{z}_i \end{pmatrix}$$
$$= h \begin{pmatrix} f(\widehat{y}_i, \widehat{z}_i) \\ g(\widehat{y}_*, \widehat{z}_i) - hg_y(0)f(\widehat{y}_i, \widehat{z}_i) \end{pmatrix} + \mathcal{O}\left(h^{M+1}\right),$$

where $\widehat{y}_* = \widehat{y}_i + hf(\widehat{y}_i, \widehat{z}_i)$. The initial values are exact, and the perturbation terms satisfy (3.9). One obtains (3.10) for the coefficients of $h$. Using the consistency requirements (3.9b) gives (3.2), and hence $\alpha_i^{(1)} = 0 \ \forall i \geq 0$. The coefficients of $h^2$ yield

(3.28a)
$$\frac{1}{2}y''(x) + \left(a^{(1)}\right)'(x) = f_y(x)a^{(1)}(x) + f_z(x)b^{(1)}(x),$$

(3.28b)
$$-g_y(0)y'(x) - g_z(0)z'(x) - f(x)(g_y(x) - g_y(0)) = g_y(x)a^{(1)}(x) + g_z(x)b^{(1)}(x),$$

(3.28c)
$$\left(\alpha_{i+1}^{(2)} - \alpha_i^{(2)}\right) = f_z(0)\beta_i^{(1)}, \quad -g_z(0)\left(\beta_{i+1}^{(1)} - \beta_i^{(1)}\right) = g_z(0)\beta_i^{(1)}.$$

The differential equation (3.28a)–(3.28b) can be solved by computing $b^{(1)}(x)$ in (3.28b) and then by replacing it in (3.28a) to yield an ODE in $a^{(1)}$. Using (3.9a) and $\alpha_0^{(1)} = 0$, one has again that $a^{(1)}(0) = 0$. Therefore $a^{(1)}(x)$ and $b^{(1)}(x)$ are uniquely determined by (3.28a)–(3.28b). The left-hand side of (3.28b) at $x = 0$ gives

$$g_y(0)a^{(1)}(0) + g_z(0)b^{(1)}(0) + f(0)g_y(0) - f(0)g_y(0)$$
$$= 0 \Rightarrow g_z(0)b^{(1)}(0) = 0 \Rightarrow b^{(1)}(0) = 0.$$

By (3.9a) and (3.28c) one also has $\beta_0^{(1)} = 0$, and in general $\beta_i^{(1)} = 0 \ \forall i \geq 0$. Further, by using (3.9b) one obtains $\alpha_i^{(2)} = 0 \ \forall i \geq 0$.

The coefficients of $h^3$ give for the smooth part

(3.29a) $$\left(a^{(2)}\right)'(x) = f_y(x)a^{(2)}(x) + f_z(x)b^{(2)}(x) + r^{(2)}(x),$$

(3.29b) $$0 = g_y(x)a^{(2)}(x) + g_z(x)b^{(2)}(x) + s^{(2)}(x),$$

where $r^{(2)}(x)$ and $s^{(2)}(x)$ are known functions that depend on derivatives of $y(x)$, $z(x)$, $a^{(1)}(x)$, $b^{(1)}(x)$. The perturbations can be expressed as $\alpha_{i+1}^{(3)} - \alpha_i^{(3)} = f_z(0)\beta_i^{(2)}$

and $0 = g_z(0)\beta_{i+1}^{(2)}$. Then $\beta_i^{(2)} = 0 \; \forall i \geq 1$, and $\alpha_i^{(3)} = 0 \; \forall i \geq 1$. The coefficients of $h^4$ reveal that the perturbations satisfy

$$(3.30a) \qquad\qquad \alpha_{i+1}^{(4)} - \alpha_i^{(4)} = f_y(0)\alpha_i^{(3)} + f_z(0)\beta_i^{(3)} \,,$$

$$(3.30b) \qquad\qquad 0 = g_z(0)\beta_{i+1}^{(3)} + g_y(0)\alpha_{i+1}^{(3)} + f(0)g_{yz}(0)\beta_i^{(2)} \,.$$

From (3.30) one has $\beta_i^{(3)} = 0 \; \forall i \geq 2$ and $\alpha_i^{(4)} = 0 \; \forall i \geq 2$. The general recurrence formula (3.14) is obtained, and the same procedure as in Theorem 3.1 can be followed. ☐

THEOREM 3.6 (consistency of the extrapolated split-IMEX method applied to DAEs). *For the harmonic sequence $\{1, 2, 3, \dots\}$ the extrapolated values $Y_{jk}$ and $Z_{jk}$ satisfy*

$$(3.31) \qquad Y_{jk} - y(x_0 + h) = \mathcal{O}\left(H^{r_{jk}}\right), \qquad Z_{jk} - z(x_0 + h) = \mathcal{O}\left(H^{s_{jk}}\right),$$

*where the differential and algebraic orders $r_{jk}$ and $s_{jk}$ are given in Table 2.*

*Proof.* The orders in Table 2 for the split-IMEX method can be recovered by using the same procedure as in the proof of Theorem 3.2 with the error expansion given by Theorem 3.5. In contrast with the proof of Theorem 3.4, $\alpha_0^{(3)}$ is nonzero, and thus one order is "lost" on the third column of the $y$ component. Then $\alpha_1^{(4)}$ gives the fourth order on the diagonal. For the $z$ component, $\beta_1^{(2)}$ is nonzero, and hence the second column of the $z$ component is 2. Furthermore, $\beta_1^{(3)}$ does not vanish, and thus the diagonal $T_{kk}$ is 3 for $k \geq 3$. The rest follows from the propagation of the error terms through the extrapolation procedure. ☐

The previous theorem concludes the set of theoretical results for the proposed extrapolation IMEX methods applied to DAEs. The results point to the W-IMEX scheme as being the most accurate; the split-IMEX scheme is computationally cheaper yet remains reasonably accurate.

**4. Numerical results for extrapolated IMEX applied to DAEs.** We illustrate the theoretical findings using two DAE examples: the reduced van der Pol equation and a trigonometric problem for which we have an analytical solution. The reduced van der Pol equation is a typical example for index-1 DAEs. In this case the numerical results using split-IMEX have a slightly higher convergence order than what is predicted by the theory. We explain this phenomenon and propose the trigonometric example to illustrate the theoretical results.

The numerical experiments are implemented in MATLAB using variable-precision arithmetic with 64 digits of accuracy. For van der Pol a numerical reference solution is computed with very high accuracy.

**4.1. Experiments with the reduced van der Pol equation.** The reduced van der Pol equation is

$$(4.1) \qquad\qquad y' = f(y, z) = -z; \quad 0 = g(y, z) = y - \left(\frac{z^3}{3} - z\right).$$

We take the initial conditions $y(0) = -2$ and $z(0) = -2.3553013976081\dots$ that satisfy $g(y(0), z(0)) = 0$. The values of $H$ range from $10^{-1}$ to $10^{-4.5}$.

The orders of the *local errors* for linearly implicit, W-IMEX, and pure-IMEX methods are given in Table 3. These experimental orders match the theoretical ones given in Table 2. The experimental orders for the split-IMEX method, not shown, are

TABLE 3

*Numerical local extrapolation orders for the van der Pol equation using the linearly implicit, W-IMEX, pure-IMEX methods, and for the trigonometric equation using the split-IMEX scheme (based on $L_1$ error norm). These results can be compared with the theoretical ones presented in Table 2.*

Orders component $y_{jk}$ (linearly implicit|W-IMEX|pure-IMEX) for van der Pol and (|split-IMEX) for the trigonometric example

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|
| 2.0\|2.0\|2.0 | | | | | | | | |
| 2.0\|2.0\|1.9 | 3.0\|3.0\|2.0\|3.0 | | | | | | | |
| 2.0\|2.0\|1.9 | 3.0\|3.0\|2.0\|3.0 | 4.0\|3.0\|3.0\|3.0 | | | | | | |
| 2.0\|2.0\|2.0 | 3.0\|3.0\|2.0\|3.0 | 4.0\|3.0\|3.0\|2.9 | 5.0\|4.0\|3.0\|4.0 | | | | | |
| 2.0\|2.0\|2.0 | 3.0\|3.0\|2.0\|3.0 | 4.0\|3.0\|3.0\|3.0 | 5.0\|4.0\|4.0\|4.0 | | | | | |
| 2.0\|2.0\|2.0 | 3.0\|3.0\|2.0\|3.0 | 4.0\|3.0\|3.0\|3.0 | 5.0\|4.0\|4.0\|3.9 | 5.1\|5.0\|3.0\|4.0 | 6.1\|5.1\|3.0\|5.1 | | | |
| 2.0\|2.0\|2.0 | 3.0\|3.0\|2.0\|3.0 | 4.0\|3.0\|3.0\|3.0 | 5.0\|4.0\|4.0\|4.0 | 5.1\|5.0\|4.0\|4.9 | 6.2\|6.0\|4.0\|6.0 | 5.9\|5.0\|3.0\|5.0 | | |
| 2.0\|2.0\|2.0 | 3.0\|3.0\|2.0\|3.0 | 4.0\|3.0\|3.0\|3.0 | 5.0\|4.0\|4.0\|4.0 | 5.1\|5.0\|5.0\|4.9 | 6.2\|6.0\|5.0\|6.0 | 8.4\|5.8\|4.0\|5.8 | 5.9\|5.0\|3.0\|5.0 | |
| 2.0\|2.0\|2.0 | 3.0\|3.0\|2.0\|3.0 | 4.0\|3.0\|3.0\|3.0 | 5.0\|4.0\|4.0\|4.0 | 5.2\|5.0\|5.0\|4.9 | 6.2\|6.0\|6.0\|6.0 | 7.3\|7.0\|5.0\|7.0 | 7.0\|6.0\|4.0\|6.0 | 6\|4.9\|2.9\|4.9 |
| 2.0\|2.0\|2.0 | 3.0\|3.0\|2.0\|3.0 | 4.0\|3.0\|3.0\|3.0 | 5.0\|4.0\|4.0\|4.0 | 5.2\|5.0\|5.0\|4.9 | | | | |

Orders component $z_{jk}$ (linearly implicit|W-IMEX|pure-IMEX) for van der Pol and (|split-IMEX) for the trigonometric example

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|
| 2.0\|2.0\|1.0\|2.0 | | | | | | | | |
| 2.0\|2.0\|1.0\|2.0 | 2.0\|2.0\|2.0\|2.0 | | | | | | | |
| 2.0\|2.0\|1.0\|2.2 | 2.0\|2.0\|2.0\|2.0 | 3.0\|3.0\|2.0\|3.0 | | | | | | |
| 2.0\|2.0\|1.0\|1.8 | 2.0\|2.0\|2.0\|2.0 | 3.0\|3.0\|3.0\|2.9 | 4.0\|4.0\|2.0\|3.0 | | | | | |
| 2.0\|2.0\|1.0\|1.9 | 2.0\|2.0\|2.0\|2.0 | 3.0\|3.0\|3.0\|2.9 | 4.0\|4.0\|3.0\|4.0 | 4.1\|4.1\|2.0\|3.1 | | | | |
| 2.0\|2.0\|1.0\|1.9 | 2.0\|2.0\|2.0\|2.0 | 3.0\|3.0\|3.0\|2.9 | 4.0\|4.0\|4.0\|3.9 | 5.0\|5.0\|3.0\|4.0 | 3.9\|4.0\|2.0\|3.1 | | | |
| 2.0\|2.0\|1.0\|1.9 | 2.0\|2.0\|2.0\|2.0 | 3.0\|3.0\|3.0\|2.9 | 4.0\|4.0\|4.0\|3.9 | 5.0\|5.0\|4.0\|5.0 | 4.6\|4.8\|3.0\|4.0 | 3.9\|4.0\|2.0\|3.2 | | |
| 2.0\|2.0\|1.0\|1.9 | 2.0\|2.0\|2.0\|2.0 | 3.0\|3.0\|3.0\|2.9 | 4.0\|4.0\|4.0\|3.9 | 5.0\|5.0\|5.0\|4.9 | 6.0\|6.0\|4.0\|4.9 | 5.0\|5.0\|3.0\|4.0 | 3.8\|3.9\|1.9\|3 | |
| 2.0\|2.0\|1.0\|1.9 | 2.0\|2.0\|2.0\|2.0 | 3.0\|3.0\|3.0\|2.9 | 4.0\|4.0\|4.0\|3.9 | 5.0\|5.0\|5.0\|4.9 | 6.0\|6.0\|5.0\|6.0 | 6.2\|6.1\|4.0\|5.0 | 4.9\|5.0\|3.0\|4.0 | 4\|4\|2\|3 |
| 2.0\|2.0\|1.0\|2.0 | 2.0\|2.0\|2.0\|2.0 | | | | | | | |

higher than the orders predicted by the theory. We can explain this disagreement by closely inspecting (4.1) and noting that $g_{yz}$ is zero. If we consider this in (3.30), we find that $\beta_2^{(3)}$ is zero and thus $\alpha_2^{(4)} = 0$. This effectively increases the order by one for the diagonal terms corresponding to the $y$ and $z$ components. Next we explore an example where $g_{yz}$ is nonzero in order to illustrate the theoretical findings for the split-IMEX method.

**4.2. Experiments with a trigonometric equation.** We consider the following DAE discretized using the split-IMEX method:

$$(4.2) \qquad y' = \frac{y^2}{z\sqrt{\frac{y^2}{z^2} - 1}}; \quad 0 = z^2 - \frac{1}{1+y^2} - y^2\left(\frac{1}{z^2} - 1\right).$$

The exact solution is $y(t) = \sinh(t)$, $z(t) = \tanh(t)$. We start at $t_0 = 0.5$, where $g_{yz}$ is nonzero. The experimental orders for split-IMEX shown in Table 3 now match the theoretical ones given in Table 2.

**5. Global error expansion for extrapolated IMEX methods applied to stiff ODEs.** In this section we extend the theoretical results for the global error expansion of the proposed methods applied to stiff ODEs. We consider the singular perturbation system (3.1) with initial conditions $(y_0, z_0)$ and $0 < \varepsilon \ll 1$ [16, 4], which is solved by using the W-IMEX (3.3), pure-IMEX (3.18), and split-IMEX (3.25) schemes. The favorable convergence results obtained for DAEs in the previous sections do not extend directly to the stiff ODEs ($0 < \varepsilon \leq H$). In this case the asymptotic expansion of the global error is more complicated, especially for "small" values of $H$. Different convergence regimes can be identified for the numerical approximations in the extrapolation tableau that depend on $H/\varepsilon$.

**5.1. W-IMEX.** We start with the W-IMEX method and consider equations of the following form (in line with (3.12)):

$$(5.1) \qquad a' = f_y(x)a + f_z(x)b + c(x, \varepsilon); \quad \varepsilon b' = g_y(x)a + g_z(x)b + d(x, \varepsilon).$$

The solution described by Lemma 5.5 in [19, Chap. IV] will be the basis for proving the next theorems.

THEOREM 5.1 (global error expansion for the extrapolated W-IMEX applied to stiff ODEs). *Assume that the solution of* (3.1) *is smooth. Under the condition*

$$(5.2) \qquad \left\|\left(I - \gamma g_z(0)\right)^{-1}\right\| \leq \frac{1}{1+\gamma} \quad \text{for} \quad \gamma \geq 1,$$

*the numerical solution of* (3.3) *possesses for* $\varepsilon \leq h$ *a perturbed asymptotic expansion of the form*

$$(5.3a) \qquad y_i = y(x_i) + ha^{(1)}(x_i) + h^2 a^{(2)}(x_i) + \mathcal{O}(h^3)$$
$$- \varepsilon f_z(0)g_z^{-1}(0)\left(I - \frac{h}{\varepsilon}g_z(0)\right)^{-i+1}\left(hb^{(1)}(0) + h^2 b^{(2)}(0)\right),$$

$$(5.3b) \qquad z_i = z(x_i) + hb^{(1)}(x_i) + h^2 b^{(2)}(x_i) + \mathcal{O}(h^3)$$
$$- \left(I - \frac{h}{\varepsilon}g_z(0)\right)^{-i+1}\left(hb^{(1)}(0) + h^2 b^{(2)}(0)\right),$$

*where $x_i = ih \leq H$ with $H$ sufficiently small independent of $\varepsilon$. The smooth functions satisfy $a^{(1)}(0) = \mathcal{O}(\varepsilon h)$, $a^{(2)}(0) = \mathcal{O}(h)$, $b^{(1)}(0) = \mathcal{O}(\varepsilon)$, $b^{(2)}(0) = \mathcal{O}(1)$.*

*Proof.* The proof goes along the lines of Theorem 3.1 and also [19, Theorem 5.6, Chap. VI] and [16]. See also a similar approach for implicit Euler [4]. The following truncated expansions are considered:

$$(5.4) \quad \widehat{y}_i = y(x_i) + \sum_{j=1}^{M} h^j \left( a^{(j)}(x_i) + \alpha_i^{(j)} \right) ; \quad \widehat{z}_i = z(x_i) + \sum_{j=1}^{M} h^j \left( b^{(j)}(x_i) + \beta_i^{(j)} \right) ,$$

such that

$$(5.5) \quad \begin{pmatrix} I & 0 \\ -hg_y(0) & \varepsilon I - hg_z(0) \end{pmatrix} \begin{pmatrix} \widehat{y}_{i+1} - \widehat{y}_i \\ \widehat{z}_{i+1} - \widehat{z}_i \end{pmatrix} = h \begin{pmatrix} f(\widehat{y}_i, \widehat{z}_i) \\ g(\widehat{y}_i, \widehat{z}_i) \end{pmatrix} + \mathcal{O}\left( h^{M+2} \right)$$

is satisfied.

(a) The smooth functions $a(x)$ and $b(x)$ depend on $\varepsilon$ but are independent of $h$. The perturbation terms $\alpha_i^{(j)}$ and $\beta_i^{(j)}$ $\forall i \geq 1$ depend smoothly on $\varepsilon$ and $\varepsilon/h$. Equations (3.9a) and (3.9b) are considered satisfied.

$M = 0$. This case is easily verified.

$M = 1$. Relation (5.4) is inserted in (5.5), and comparing the smooth coefficients of $h^2$ yields

$$(5.6a)$$
$$\left( a^{(1)} \right)'(x) + \frac{1}{2}y''(x) = f_y(x) a^{(1)}(x) + f_z(x) b^{(1)}(x),$$

$$(5.6b)$$
$$\frac{1}{2}\varepsilon z''(x) - g_y(0)y'(x) - g_z(0)z'(x) + \varepsilon \left( b^{(1)} \right)'(x) = g_y(x) a^{(1)}(x) + g_z(x) b^{(1)}(x).$$

By [19, Lemma 5.5, Chap. IV], the initial value $b^{(1)}(0)$ is uniquely determined by $a^{(1)}(0)$. Differentiating $\varepsilon z'(x) = g(y(x), z(x))$ and inserting it in (5.6b) at $x = 0$, we get

$$(5.7)$$
$$g_y(0) a^{(1)}(0) + g_z(0) b^{(1)}(0) = -\frac{1}{2} \left( g_y(0) y'(0) + g_z(0) z'(0) \right) + \varepsilon \left( b^{(1)} \right)'(0) = \mathcal{O}(\varepsilon)$$

with a known right-hand side. The perturbation terms up to $\mathcal{O}(h^2)$ give

$$(5.8a) \qquad \alpha_{i+1}^{(1)} - \alpha_i^{(1)} = h f_y(x_i)\alpha_i^{(1)} + h f_z(x_i)\beta_i^{(1)} ,$$

$$(5.8b) \qquad \varepsilon \left( \beta_{i+1}^{(1)} - \beta_i^{(1)} \right) - h g_y(0) \left( \alpha_{i+1}^{(1)} - \alpha_i^{(1)} \right) - h g_z(0) \left( \beta_{i+1}^{(1)} - \beta_i^{(1)} \right)$$
$$= h g_y(x_i)\alpha_i^{(1)} + h g_z(x_i)\beta_i^{(1)} .$$

Next we eliminate as many terms in (5.8) as possible by replacing $f_y(x_i)$ with $f_y(0)$, $g_y(x_i)$ with $g_y(0)$, and so on. With $x_i = ih$, the following substitution is of order $h$: $f_y(x_i) - f_y(0) = \mathcal{O}(h)$, since $i \leq 1$. Then one is left with

$$(5.9) \qquad \begin{cases} \alpha_{i+1}^{(1)} - \alpha_i^{(1)} = h f_y(0)\alpha_i^{(1)} + h f_z(0)\beta_i^{(1)} + \mathcal{O}(h^2), \\ \varepsilon \left( \beta_{i+1}^{(1)} - \beta_i^{(1)} \right) - h g_y(0)\alpha_{i+1}^{(1)} - h g_z(0)\beta_{i+1}^{(1)} = \mathcal{O}(h^2). \end{cases}$$

In the second expression of (5.9), we note that $\beta_{i+1}^{(1)}$ is multiplied by $\varepsilon$, whereas $\alpha_{i+1}^{(1)}$ is not and thus can be ignored (for $\varepsilon \ll h$). Then one gets

$$(5.10a) \qquad\qquad \alpha_{i+1}^{(1)} - \alpha_i^{(1)} = hf_z(0)\beta_i^{(1)},$$

$$(5.10b) \qquad\qquad \varepsilon\left(\beta_{i+1}^{(1)} - \beta_i^{(1)}\right) = hg_z(0)\beta_{i+1}^{(1)}.$$

The solutions of (5.6), (5.10) when substituted in (5.5) are analyzed next. From (5.10b)

$$(5.11) \qquad\qquad \beta_i^{(1)} = \left(I - \frac{h}{\varepsilon}g_z(0)\right)^{-i}\beta_0^{(1)}.$$

Substituting (5.11) into (5.10a) and using (3.9b), we have

$$(5.12) \qquad\qquad \alpha_i^{(1)} = \varepsilon f_z(0)g_z^{-1}(0)\left(I - \frac{h}{\varepsilon}g_z(0)\right)^{-i+1}\beta_0^{(1)}.$$

Expression (5.12) at $i = 0$ with $\varepsilon \leq h$ yields

$$(5.13) \qquad \alpha_0^{(1)} = \varepsilon f_z(0)g_z^{-1}(0)\left(I - \frac{h}{\varepsilon}g_z(0)\right)\beta_0^{(1)} = \mathcal{O}(h)\,\beta_0^{(1)} = \mathcal{O}(\varepsilon h).$$

In the previous relation we used (5.7) and (3.9a) to bound $\beta_0^{(1)}$. The consistency assumptions (3.9a) with (5.7) and (5.13) and by using Lemma 5.5 in [19, Chap. IV] guarantees that the coefficients $a^{(1)}(0)$, $b^{(1)}(0)$, $\alpha_0^{(1)}$, $\beta_0^{(1)}$ are uniquely determined; moreover, one has $a^{(1)}(0) = \mathcal{O}(\varepsilon h)$ and $b^{(1)}(0) = \mathcal{O}(\varepsilon)$ $(\alpha_i^{(1)} = \mathcal{O}(\varepsilon h)$, $\beta_i^{(1)} = \mathcal{O}(\varepsilon))$. Now the relation (5.5) can be verified for $M = 1$, $\varepsilon \leq h$.

$M = 2$. Relation (5.4) is inserted in (5.5), and comparing the smooth coefficients of $h^3$ we obtain the same form as in (5.1), $\left(a^{(2)}\right)'(x) = a^{(2)}(x)f_y(x) + f_z(x)b^{(2)}(x) + c(x,\varepsilon)$ with known $c(x,\varepsilon)$. Using $\varepsilon z'(x) = g(y(x), z(x))$, and evaluating at $x = 0$, one obtains $\varepsilon\left(b^{(2)}\right)'(0) = g_y(0)a^{(2)}(0) + g_z(0)b^{(2)}(0) + d(0,\varepsilon)$ with known $d(0,\varepsilon)$. It follows again from Lemma 5.5 in [19, Chap. IV] and $d(0,\varepsilon) = \mathcal{O}(1)$ that

$$(5.14) \qquad\qquad g_y(0)a^{(2)}(0) + g_z(0)b^{(2)}(0) = \mathcal{O}(1).$$

Just as in the $M = 1$ case, for the perturbations we require $\alpha_{i+1}^{(2)} - \alpha_i^{(2)} = hf_z(0)\beta_i^{(2)}$ and $\varepsilon(\beta_{i+1}^{(2)} - \beta_i^{(2)}) = hg_z(0)\beta_{i+1}^{(2)}$, and

$(5.15a)$
$$\beta_i^{(2)} = \left(I - \frac{h}{\varepsilon}g_z(0)\right)^{-i}\beta_0^{(2)}, \qquad \alpha_i^{(2)} = \varepsilon f_z(0)g_z^{-1}(0)\left(I - \frac{h}{\varepsilon}g_z(0)\right)^{-i+1}\beta_0^{(2)},$$

$(5.15b)$
$$\alpha_0^{(2)} = \varepsilon f_z(0)g_z^{-1}(0)\left(I - \frac{h}{\varepsilon}g_z(0)\right)\beta_0^{(2)}$$

are obtained just as for (5.11), (5.12), and (5.13), respectively. The values $a^{(1)}(0)$, $b^{(1)}(0)$, $\alpha_0^{(1)}$, $\beta_0^{(1)}$ are uniquely determined by (3.9a), (5.14), and (5.15). By using Lemma 5.5 in [19, Chap. IV], one has that $a^{(2)}(0) = \mathcal{O}(h)$ and $b^{(1)}(0) = \mathcal{O}(1)$;

moreover, by using (3.9a) one obtains $\alpha_i^{(2)} = \mathcal{O}(h)$ for $\varepsilon \leq h$. The verification of (5.5) for $M = 2$ is tedious, but it can be shown to be satisfied in general by using the following remarks. The coefficients of $h^1$ can be ignored since they vanish for large $i$'s. The assumption (5.2) gives $\beta_i^{(1)} = \mathcal{O}\left(\varepsilon 2^{-i}\right)$ and $\beta_i^{(2)} = \mathcal{O}\left(2^{-i}\right)$. These terms can also be neglected; however, in practice, they can give additional convergence regimes that quickly vanish. The convergence $(H \to 0, H/\varepsilon \to \infty)$ will have different slopes that are determined by the ratio of $H$ and $\varepsilon$.

This analysis gets complicated for $M \geq 3$; however, the behavior of the error in practical applications can be understood from the discussion above.

(b) The second part of the proof consists of estimating a bound on the reminder term just as we did for the proof of Theorem 3.1; that is, differences $\Delta y_i = y_i - \widehat{y}_i$ and $\Delta z_i = z_i - \widehat{z}_i$. Subtracting (5.5) from (3.3) and eliminating $\Delta y_i$, $\Delta z_i$, and using (5.2) with $\varepsilon \leq h$, we have

$$(5.16) \qquad \left\| I + \left(\frac{\varepsilon}{h}I - g_z(0)\right)^{-1} g_z(0) \right\| = \left\| \left(I - \frac{h}{\varepsilon}g_z(0)\right)^{-1} \right\| \leq \frac{\varepsilon}{\varepsilon + h} \leq \frac{1}{2}.$$

We therefore obtain (3.15) with $|\zeta| < 1$ and $H$ sufficiently small. Using the same procedure as in the proof of Theorem 3.1, one therefore obtains $\|\Delta y_i\| + \|\Delta z_i\| = \mathcal{O}\left(h^{M+1}\right)$. $\qquad \square$

A close inspection of (5.3) reveals that the global error has different convergence regimes when $\varepsilon \leq h$. We now focus on the global error expansion of the stiff component (5.3), which gives the following leading term:

$$Z_{j1} = \left(I - \frac{h}{\varepsilon}g_z(0)\right)^{-n_j+1} \left(hb^{(1)}(0) + h^2 b^{(2)}(0)\right) = h^2 \left(I - \frac{h}{\varepsilon}g_z(0)\right)^{-n_j+1} b^{(2)}(0).$$

We further consider $g_z(0) \propto -1$. With $H = h/n_j$, one has

$$T_{j1} = \left(\frac{H}{\varepsilon n_j}\right)^2 \left(1 + \frac{H}{\varepsilon n_j}\right)^{-n_j+1} b^{(2)}(0) \quad \text{and} \quad Z_{j1} = \varepsilon^2 T_{j1} b^{(2)}(0).$$

The error propagates through the extrapolation tableau through (2.3). Similar to the behavior of the global error for the linearly implicit method [19, p. 438], the first subdiagonal $(T_{j\,j-1})$ with $n_1 = 1$ gives $T_{j\,j-1} = \text{const.} (H/\varepsilon)^{2-n_2} + \mathcal{O}((H/\varepsilon)^{2-n_2})$, where the constant is determined by (2.3). This suggests a superposition of the convergence slopes predicted for DAEs and a factor $\mathcal{O}\left(\varepsilon^2\right)$ as discussed in [19].

**5.2. Pure-IMEX Method.** We now consider the pure-IMEX method.

THEOREM 5.2 (global error expansion for the extrapolated pure-IMEX method applied to stiff ODEs). *Assume that the solution of (3.1) is smooth. Under the condition (5.2) the numerical solution of (3.18) possesses for $\varepsilon \leq h$ a perturbed asymptotic expansion of form (5.3) with $x_i = ih \leq H$, $H$ sufficiently small independent of $\varepsilon$. The smooth functions satisfy $a^{(1)}(0) = \mathcal{O}(h)$, $a^{(2)}(0) = \mathcal{O}(h)$, $b^{(1)}(0) = \mathcal{O}(1)$, $b^{(2)}(0) = \mathcal{O}(1)$.*

*Proof.* The proof goes along the same lines as for Theorem 5.1. Assumptions (5.1) and (5.4) are considered, and thus (5.5) becomes

$$(5.17) \qquad \begin{pmatrix} I & 0 \\ -hg_y(0) & \varepsilon I - hg_z(0) \end{pmatrix} \begin{pmatrix} \widehat{y}_{i+1} - \widehat{y}_i \\ \widehat{z}_{i+1} - \widehat{z}_i \end{pmatrix}$$
$$= h \begin{pmatrix} f\left(\widehat{y}_i, \widehat{z}_i\right) \\ g\left(\widehat{y}_i, \widehat{z}_i\right) - hg_y(0)f\left(\widehat{y}_i, \widehat{z}_i\right) \end{pmatrix} + \mathcal{O}\left(h^{M+1}\right).$$

For $M = 1$ one obtains $(a^{(1)})'(x) + \frac{1}{2}y''(x) = f_y(x) a^{(1)}(x) + f_z(x) b^{(1)}(x)$ and

$$\frac{1}{2}\varepsilon z''(x) - g_y(0)y'(x) - g_z(0)z'(x) + \varepsilon (b^{(1)})'(x)$$
$$= g_y(x) a^{(1)}(x) + g_z(x) b^{(1)}(x) - f(x)g_y(0).$$

This leads to

$$g_y(0) a^{(1)}(0) + g_z(0) b^{(1)}(0)$$
$$= -\frac{1}{2}(g_y(0) y'(0) + g_z(0) z'(0)) + f(0)g_y(0) + \varepsilon \left(b^{(1)}\right)'(0)$$

with a known right-hand side of $\mathcal{O}(1)$. The perturbation terms up to $\mathcal{O}(h^2)$ give the same expression as in the W-IMEX case (5.8) that yields (5.9) and eventually (5.10). The values for $\alpha_i^{(1)}$ and $\beta_i^{(1)}$ are given by (5.12) and (5.11), respectively. By using the consistency assumptions (3.9a) and (5.11) one obtains

$$(5.19) \qquad \alpha_0^{(1)} = \varepsilon f_z(0)g_z^{-1}(0) \left(I - \frac{h}{\varepsilon}g_z(0)\right) \beta_0^{(1)} = \mathcal{O}(h) \, \beta_0^{(1)} = \mathcal{O}(h),$$

which yields $a^{(1)}(0) = \mathcal{O}(h)$ and $b^{(1)}(0) = \mathcal{O}(1)$ $(\alpha_i^{(1)} = \mathcal{O}(h), \, \beta_i^{(1)} = \mathcal{O}(1))$. With these assumptions (5.17) can be verified.

For $M = 2$ one obtains the same form as (5.1) and again (5.14)–(5.15). Using $b^{(2)}(0) = \mathcal{O}(1)$ yields $a^{(2)}(0) = \mathcal{O}(h)$ and $b^{(1)}(0) = \mathcal{O}(1)$. The rest is similar to Theorem 5.1. $\quad\square$

The convergence behavior of this method is similar to the one for the W-IMEX scheme (section 5.1); however, in this case the superposition of the error has a factor of $\mathcal{O}(\varepsilon)$.

**5.3. Split-IMEX method.** We next consider the split-IMEX method.

THEOREM 5.3 (global error expansion for the extrapolated split-IMEX method applied to stiff ODEs). *Assume that the solution of* (3.1) *is smooth. Under the condition* (5.2) *the numerical solution of* (3.25) *possesses for* $\varepsilon \leq h$ *a perturbed asymptotic expansion of the form* (5.3) *with* $x_i = ih \leq H$, *H sufficiently small independent of* $\varepsilon$. *The smooth functions satisfy* $a^{(1)}(0) = \mathcal{O}(\varepsilon h)$, $a^{(2)}(0) = \mathcal{O}(h)$, $b^{(1)}(0) = \mathcal{O}(\varepsilon)$, $b^{(2)}(0) = \mathcal{O}(1)$.

*Proof.* The proof goes along the same lines as for Theorem 5.1. Assumptions (5.1), (5.4) are considered, and (5.5) (with $\widehat{y}_* = \widehat{y}_i + hf(\widehat{y}_i, \widehat{z}_i)$) becomes

$$\begin{pmatrix} I & 0 \\ -hg_y(0) & \varepsilon I - hg_z(0) \end{pmatrix} \begin{pmatrix} \widehat{y}_{i+1} - \widehat{y}_i \\ \widehat{z}_{i+1} - \widehat{z}_i \end{pmatrix}$$
$$= h \begin{pmatrix} f(\widehat{y}_i, \widehat{z}_i) \\ g(\widehat{y}_*, \widehat{z}_i) - hg_y(0)f(\widehat{y}_i, \widehat{z}_i) \end{pmatrix} + \mathcal{O}\left(h^{M+2}\right).$$

For $M = 1$ one obtains $(a^{(1)})'(x) + \frac{1}{2}y''(x) = f_y(x) a^{(1)}(x) + f_z(x) b^{(1)}(x)$ and

$$\frac{1}{2}\varepsilon z''(x) - g_y(0)y'(x) - g_z(0)z'(x) + \varepsilon (b^{(1)})'(x)$$
$$= g_y(x)a^{(1)}(x) + g_z(x)b^{(1)}(x) + f(x)(g_y(x) - g_y(0)),$$

which leads to $g_y(0)a^{(1)}(0) + g_z(0)b^{(1)}(0) = -\frac{1}{2}(g_y(0) y'(0) + g_z(0)z'(0)) + \varepsilon (b^{(1)})'(0)$ with a known right-hand side of $\mathcal{O}(\varepsilon)$. The perturbation terms up to $\mathcal{O}(h^2)$ give the

same expression as in the W-IMEX case (5.8) that yields (5.9) and eventually (5.10). The values for $\alpha_i^{(1)}$ and $\beta_i^{(1)}$ are given by (5.12) and (5.11), respectively. By using the consistency assumptions (3.9a) and (5.11) one obtains

$$\alpha_i^{(1)} = \varepsilon f_z(0) g_z^{-1}(0) \left( I - \frac{h}{\varepsilon} g_z(0) \right)^{-i+1} \beta_0^{(1)} \quad \text{and} \quad \alpha_0^{(1)} = \mathcal{O}(h) \, \beta_0^{(1)} = \mathcal{O}(\varepsilon h) \, ,$$

which yields $a^{(1)}(0) = \mathcal{O}(\varepsilon h)$ and $b^{(1)}(0) = \mathcal{O}(\varepsilon)$ $(\alpha_i^{(1)} = \mathcal{O}(\varepsilon h), \, \beta_i^{(1)} = \mathcal{O}(\varepsilon))$. With these assumptions, (5.17) can be verified.

For $M = 2$ one obtains the same form as (5.1) and then (5.14)–(5.15). Using $b^{(2)}(0) = \mathcal{O}(1)$ yields $a^{(2)}(0) = \mathcal{O}(h)$ and $b^{(1)}(0) = \mathcal{O}(1)$. The rest is similar to Theorem 5.1.    ☐

The convergence behavior is similar to the W-IMEX scheme (see section 5.1).

**6. Numerical results for extrapolated IMEX applied to stiff ODEs.** We investigate the numerical properties of the extrapolated IMEX methods applied to stiff ODEs. We consider van der Pol's equation, the prototypical stiff ODE example. For comparison we include the Kennedy and Carpenter schemes proposed in [22], denoted here by "ARK(order(embedded order)stages)." All IMEX RK methods require solving a (non)linear system of equations. A simplification occurs if the implicit part is linear and RK methods are of ESDIRK type, which is the case for the ARK methods used here.

The implementation is done in MATLAB using high-precision (64-digit) arithmetic. The experiments consist of integrating the problem with successively smaller steps $H$ and computing the $L_1$ error norm for each step size. We compare the results of the proposed IMEX methods and the above mentioned IMEX RK schemes with a third-order reference solution computed with the stiff solver RODAS-3 [28] and a step size of $10^{-9}$. The nonlinear solver used in the computation of the reference solution and in the IMEX RK methods is based on classical Newton iterations. The process is stopped when the difference between successive iterates is below $10^{-25}$.

We consider van der Pol's equation (see [19, 5])

$$(6.1) \qquad \begin{array}{rcl} y' & = & z \\ \varepsilon \, z' & = & \left(1 - y^2\right) z - y \end{array} = \begin{pmatrix} z \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ \left(1 - y^2\right) z - y \end{pmatrix}$$

with $y(0) = 2$, $z(0) = -\frac{2}{3} + \frac{10}{81}\varepsilon - \frac{292}{2187}\varepsilon^2 - \frac{1814}{19683}\varepsilon^3 + \mathcal{O}\left(\varepsilon^4\right)$, and $\varepsilon = 10^{-5}$ [5]. Figure 2 presents the error for the stiff solution component $(z)$ calculated by using extrapolated linearly implicit and IMEX methods (2.4) with 3, 6, 9, and 12 extrapolation steps. For each extrapolated IMEX scheme the tableau entry $T_{j,k}$ with optimal convergence order is selected from Table 2. The observed convergence rates match the theoretical predictions. The error decreases until it reaches $\mathcal{O}(\varepsilon)$ for pure-IMEX and $\mathcal{O}(\varepsilon^2)$ for the others.

We compare the extrapolated IMEX methods with several IMEX RK methods. Figure 3(a) shows the local errors $(L_1)$ of the stiff component versus the step size for the third- to fifth-order ARK methods. The order reduction phenomenon can be clearly noticed. A detailed explanation of the convergence behavior is given in [5].

The computational cost of the IMEX extrapolation methods increases linearly with each additional extrapolation step. For $T_{jk}$ one needs $j(j+1)/2$ right-hand side evaluations. In contrast, for an $s_i$-implicit, $s_e$-explicit-stage IMEX RK scheme, one needs $\approx [(s_e - s_i) + s_i \times \#$ of Newton iterations] function evaluations. In this study we do not focus on the computational cost, which can change with the implementation/application.
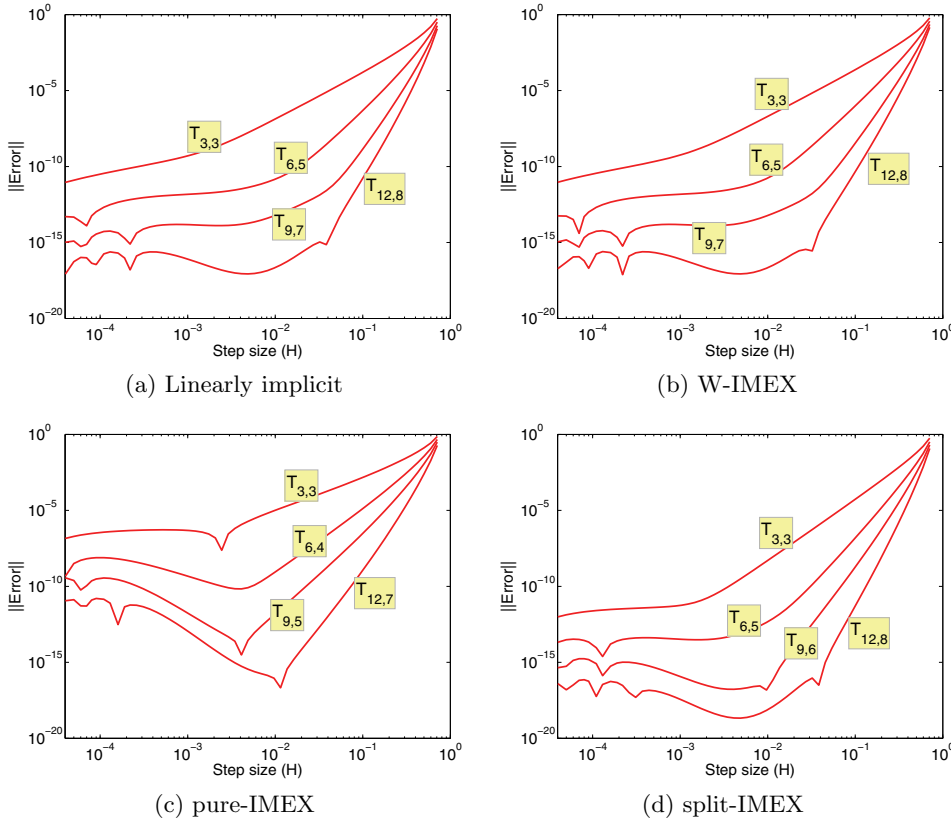
FIG. 2. *Local error versus step size for the stiff solution component of the van der Pol equation using extrapolated linearly implicit and IMEX methods for the optimal convergence rates with 3, 6, 9, and 12 extrapolation steps; that is, the optimal k for each method's $T_{3k}$, $T_{6k}$, $T_{9k}$, $T_{12k}$.*
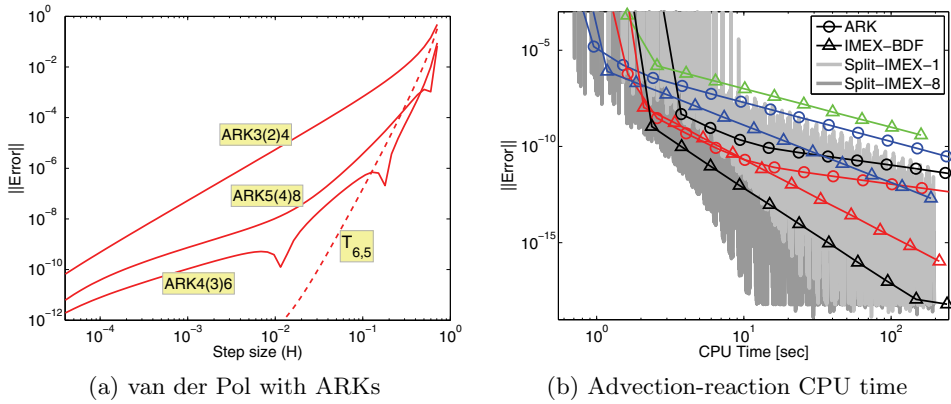


FIG. 3. *(a) Local error versus step size for the stiff solution component of the van der Pol equation using ARK methods and $T_{6,5}$ split-IMEX for comparison. (b) Global error versus CPU time for the advection-reaction equation at $t_{\max} = 1$ solved with IMEX-BDF (orders 2 (green), 3 (blue), 4 (red), 5 (black)), ARK (orders 3 (blue), 4 (red), 5 (black)), and the proposed split-IMEX method for different extrapolation terms (up to order 18) using a sequential (light gray) and a straightforward OpenMP parallel implementation (dark gray) on an 8-core processor.*

**7. Numerical results for PDEs.** We next investigate the time discretization accuracy of the advection-reaction PDE using the extrapolated W-IMEX, pure-IMEX, and split-IMEX schemes. In this section we denote by $x$ the spatial variable and by $t$ the temporal variable. The numerical order of convergence is estimated in the $L_1$ error norm ($\|Err\|_1 = \Delta x/m \sum_{i=1}^{m} |Err_i|$), where $m$ is the total number of variables. The errors at the final time for different step sizes ($H$) are considered.

In [8] we also discuss the order reduction phenomenon due to nonhomogeneous boundary conditions or source terms [6], illustrate it on a typical example [29], and show how to avoid order reduction using a strategy developed for RK methods [7]. This topic needs to be explored further, but it falls outside the scope of this paper.

The advection-reaction PDE is described by the following system with the setting presented in [20]:

$$(7.1) \quad \begin{aligned} y_t + \alpha_1\, y_x &= -k_1 y + k_2 z + s_1, & 0 < x < 1, & \quad \alpha_1 = 1,\, k_1 = 10^6,\, s_1 = 0, \\ z_t + \alpha_2\, z_x &= k_1 y - k_2 z + s_2, & 0 < t \le t_{\max}, & \quad \alpha_2 = 0,\, k_2 = 2k_1,\, s_2 = 1, \end{aligned}$$

with $y(x,0) = 1 + s_2 x$, $z(x,0) = \frac{k_1}{k_2} y(x,0) + \frac{1}{k_2} s_2$, $y(0,t) = 1 - \sin(12t)^4$. The advection term is treated explicitly and the reaction term implicitly because of its numerical stiffness. For the spatial discretization we use the fourth-order central finite difference scheme described in [20]. A uniform grid is considered in space: $x_i = i\Delta x$, $i = 1 \ldots m$ with $\Delta x = 1/m$, $m = 400$. The experimental orders observed at $t_{\max} = 1$ are shown in Table 4. They are generally in agreement with the theoretical predictions. Some components have slightly higher convergence orders, which is expected because of the linearity of this example, which makes W-IMEX and split-IMEX equivalent. We also note that the experimental orders continue to increase with the addition of more terms in the extrapolation tableau.

In Figure 3(b) we show the CPU time versus the global error of the split-IMEX method compared to ARK orders 3–5 [22] and IMEX-BDF orders 2–5 [20] methods. The implementation is done in FORTRAN compiled with quad precision on an 8-core machine, and the resulting linear system is solved by using LAPACK LU factorization. For the split-IMEX scheme we consider orders up to 18, and the errors of all terms in the extrapolation tableau are represented as light gray lines in Figure 3(b). The superposition of various terms gives an apparent oscillatory convergence rate; however, this is just a visual artifact. If the optimal extrapolation term that corresponds to the smallest error is considered (see Table 3), then the split-IMEX method compares well with RK and LM methods on low accuracy and is superior for high-accuracy results. We further considered a straightforward OpenMP parallelization of the extrapolation row calculations. Each row is dynamically allocated on a CPU core. The timing results show that on an 8-core machine, split-IMEX is superior in efficiency to the considered LM and RK methods. No effort has been made to optimize the parallel performance, but additional improvements seem possible by optimizing the code and by employing more CPUs. In comparison, neither the ARK nor the IMEX-BDF can

TABLE 4
*Numerical orders for the advection-reaction PDE with extrapolated W-IMEX|pure-IMEX|split-IMEX schemes ($t_{\max} = 1$, $m = 400$).*

| | | | | |
|---|---|---|---|---|
| 1.0\|1.0\|1.0 | | | | |
| 1.0\|1.0\|1.0 | 2.0\|1.0\|2.0 | | | |
| 1.0\|1.0\|1.0 | 2.0\|1.0\|2.0 | 3.0\|1.9\|3.0 | | |
| 1.2\|1.0\|1.2 | 2.0\|1.0\|2.0 | 3.0\|2.0\|3.0 | 4.0\|2.0\|4.0 | |
| 1.0\|1.0\|1.0 | 2.0\|1.0\|2.0 | 3.0\|2.0\|3.0 | 4.0\|3.0\|4.0 | 5.0\|2.0\|5.0 |

benefit from parallelization. Moreover, LM methods in general, and IMEX-BDF in particular, may become unstable if the eigenvalues of the implicit term are relatively large and close to the imaginary axis, whereas the proposed methods allow for A-stability on the implicit part.

**8. Implementation considerations.** In this section we make several observations regarding the implementation of the proposed IMEX extrapolation methods.

*Construction of extrapolation methods.* The extrapolation procedure provides a set of increasingly accurate results. Lower-order embedded approximations are readily available, and thus a step size $(H)$ control strategy is easy to implement [19]. Because each computational step in the extrapolation procedure is a consistent approximation, one can consider an adaptive-order approach. The implementation consists of coding (2.3) and (2.4). The Jacobian $g'$ is evaluated only at the beginning of the step. Therefore several computational simplifications may occur, especially if $g$ is linear. Very high-order approximations are easily obtained, with no limitation on the theoretical achievable convergence order.

*Computational cost, memory usage, and parallelization.* In the classical setting $(\varepsilon \approx h)$, the extrapolation methods are considered less efficient than the established RK or LM schemes. Depending on the problem type and its stiffness, however, traditional RK or LM schemes require nonlinear solver iterations, whereas extrapolated IMEX schemes may not; they are similar to W-methods but can achieve much higher orders of convergence. Moreover, the extrapolation methods can be easily parallelized [26] as each entry on $T_{j,1}$ can be computed independently. Furthermore, the computational cost is predetermined: cost for $T_{jk} \propto j(j+1)/2$ function evaluations, and thus each entry can be optimally scheduled on multiprocessor or multicore architectures. This strategy could lead to more efficient overall implementations with the total computational cost $\propto j$, as illustrated in Figure 3(b). In contrast, the IMEX RK methods have a cost proportional to the product of the number of implicit stages nonlinear solver iterations and cannot benefit from parallelization. The memory requirements for full-extrapolation tableaux are proportional to $j(j+1)/2$; however, a large number of tableau entries need not be computed, and thus the number of registers required in practice can be reduced.

*Extrapolation methods for stiff systems.* For stiff nonlinear problems, the diagonal entries in the extrapolation tableau are typically not the best approximations. The optimal entries in the extrapolation tableau are emphasized in Table 2. This is equivalent to using a shifted harmonic sequence $n_j = \ell, \ell + 1, \ldots, j = 1, 2, \ldots,$ $\ell \geq 1$ that includes the optimal values (see Table 2). If a sufficiently large number of extrapolation steps is computed, then the diagonal and several other subdiagonal entries are not necessary, and hence cost and memory requirements are alleviated.

**9. Discussion.** In this paper we construct extrapolated IMEX time discretization methods for problems with both stiff and nonstiff components; for example, multiphysics multiscale partial differential equations.

We propose three new extrapolation methods: W-IMEX, pure-IMEX, and split-IMEX. The theoretical study reveals the existence of perturbed global error expansions for each of these base methods. Theoretical predictions of the orders of convergence are made in the DAE and SPP settings. A (scalar) linear stability analysis is performed.

The W-IMEX method resembles the linearly implicit scheme in terms of implementation and performance but is computationally more attractive because it uses only the Jacobian of the stiff component. The closely related pure-IMEX and split-IMEX methods are truly IMEX methods, in that they fully decouple the explicit and the implicit

parts. The split-IMEX method performs one explicit step with the nonstiff component, followed by a linearly implicit step with the stiff component.

Extrapolated IMEX methods have very low implementation costs and can easily deliver very high orders of consistency. Thus they are well suited for high accuracy integration of ODEs, index-1 DAEs, and PDEs via the method of lines. In this study we have not extensively assessed the efficiency of these methods; however, the numerical tests indicate that they compare well with existing IMEX RK and LM methods and are superior when even a straightforward OpenMP parallelization is considered. Additional improvements seem possible by optimizing the extrapolation code and by employing more computational units. IMEX RK and IMEX-BDF cannot benefit from paralellization. Split-IMEX seems a good choice at least for the problems analyzed in this study. It is easier to implement than W-IMEX and has more favorable properties than does the pure-IMEX method.

The proposed IMEX extrapolation methods parallelize well and can take advantage of the emerging multicore hardware architectures. By construction they provide low-order embedded approximations, thus facilitating implementations with variable step size. Moreover, they do not require a predetermined number of stages and thus allow variable-order strategies as well.

## REFERENCES

[1] A. C. AITKEN, *On interpolation by iteration of proportional parts without the use of differences*, in Proc. Edin. Math. Soc., 3 (1932), pp. 56–76.

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

[3] U. M. ASCHER, S. J. RUUTH, AND B. T. R. WETTON, *Implicit-explicit methods for time-dependent partial differential equations*, SIAM J. Numer. Anal., 32 (1995), pp. 797–823.

[4] W. AUZINGER, R. FRANK, AND F. MACSEK, *Asymptotic error expansions for stiff equations: The implicit Euler scheme*, SIAM J. Numer. Anal., 27 (1990), pp. 67–104.

[5] S. BOSCARINO, *Error analysis of IMEX Runge-Kutta methods derived from differential-algebraic systems*, SIAM J. Numer. Anal., 45 (2007), pp. 1600–1621.

[6] P. BRENNER, M. CROUZEIX, AND V. THOMÉE, *Single step methods for inhomogeneous linear differential equations in Banach space*, RAIRO Anal. Numér., 16 (1982), pp. 5–26.

[7] M. H. CARPENTER, D. GOTTLIEB, S. ABARBANEL, AND W.-S. DON, *The theoretical accuracy of Runge-Kutta time discretizations for the initial boundary value problem: A study of the boundary error*, SIAM J. Sci. Comput., 16 (1995), pp. 1241–1252.

[8] E. CONSTANTINESCU AND A. SANDU, *Achieving Very High Order for Implicit Explicit Time Stepping: Extrapolation Methods*, Technical report ANL/MCS-TM-306, Argonne National Laboratory, Mathematics and Computer Science Division Technical Memorandum, Argonne, IL, 2009, available online at http://www.mcs.anl.gov/uploads/cels/papers/TM-306.pdf.

[9] E. M. CONSTANTINESCU, A. SANDU, AND G. R. CARMICHAEL, *Modeling atmospheric chemistry and transport with dynamic adaptive resolution*, Comput. Geosci., 12 (2008), pp. 133–151.

[10] P. DEUFLHARD, *Order and stepsize control in extrapolation methods*, Numer. Math., 41 (1983), pp. 399–422.

[11] P. DEUFLHARD, *Recent progress in extrapolation methods for ordinary differential equations*, SIAM Rev., 27 (1985), pp. 505–535.

[12] P. DEUFLHARD, E. HAIRER, AND J. ZUGCK, *One-step and extrapolation methods for differential-algebraic systems*, Numer. Math., 51 (1987), pp. 501–516.

[13] J. FRANK, W. HUNDSDORFER, AND J. G. VERWER, *On the stability of implicit-explicit linear multistep methods*, Appl. Numer. Math., 25 (1997), pp. 193–205.

[14] L. GEBHARDT, D. FOKIN, T. LUTZ, AND S. WAGNER, *An Implicit-Explicit Dirichlet-Based Field Panel Method for Transonic Aircraft Design*, AIAA-Paper 2002-3145, 2002.

[15] E. HAIRER AND C. LUBICH, *Asymptotic expansions of the global error of fixed-stepsize methods*, Numer. Math., 45 (1984), pp. 345–360.

[16] E. HAIRER AND C. LUBICH, *Extrapolation at stiff differential equations*, Numer. Math., 52 (1988), pp. 377–400.

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

[18] E. HAIRER, S. NØRSETT, AND G. WANNER, *Solving Ordinary Differential Equations* I: *Nonstiff Problems*, 2nd ed., Springer Ser. Comput. Math. 8, Springer-Verlag, Berlin, 1993.

[19] E. HAIRER AND G. WANNER, *Solving Ordinary Differential Equations* II: *Stiff and Differential-Algebraic Problems*, 2nd ed., Springer Ser. Comput. Math. 14, Springer-Verlag, Berlin, 1993.

[20] W. HUNDSDORFER AND S. J. RUUTH, *IMEX extensions of linear multistep methods with general monotonicity and boundedness properties*, J. Comput. Phys., 225 (2007), pp. 2016–2042.

[21] W. HUNDSDORFER AND J. VERWER, *Numerical Solution of Time-Dependent Advection-Diffusion-Reaction Equations*, Springer Ser. Comput. Math. 33, Springer-Verlag, Berlin, 2003.

[22] C. A. KENNEDY AND M. H. CARPENTER, *Additive Runge-Kutta schemes for convection-diffusion-reaction equations*, Appl. Numer. Math., 44 (2003), pp. 139–181.

[23] J. D. LAMBERT, *Numerical Methods for Ordinary Differential Systems: The Initial Value Problem*, John Wiley and Sons, Chichester, UK, 1991.

[24] E. NEVILLE, *Iterative interpolation*, J. Indian Math. Soc., 20 (1934), pp. 87–120.

[25] L. PARESCHI AND G. RUSSO, *Implicit-explicit Runge-Kutta schemes for stiff systems of differential equations*, in Recent Trends in Numerical Analysis, Adv. Theory Comput. Math. 3, Nova Sci. Publ., Huntington, NY, 2001, pp. 269–288.

[26] T. RAUBER AND G. RÜNGER, *Load balancing schemes for extrapolation methods*, Concurrency: Practice Exp., 9 (1997), pp. 181–202.

[27] S. J. RUUTH, *Implicit-explicit methods for reaction-diffusion*, J. Math. Biol., 34 (1995), pp. 148–176.

[28] A. SANDU, J. G. VERWER, J. G. BLOM, E. J. SPEE, G. R. CARMICHAEL, AND F. A. POTRA, *Benchmarking stiff ODE solvers for atmospheric chemistry problems* II: *Rosenbrock methods*, Atmos. Environ., 31 (1997), pp. 3459–3472.

[29] J. M. SANZ-SERNA, J. G. VERWER, AND W. H. HUNDSDORFER, *Convergence and order reduction of Runge-Kutta schemes applied to evolutionary problems in partial differential equations*, Numer. Math., 50 (1987), pp. 405–418.

[30] J. G. VERWER, J. G. BLOM, AND W. HUNDSDORFER, *An implicit-explicit approach for atmospheric transport-chemistry problems*, Appl. Numer. Math., 20 (1996), pp. 191–209.

[31] J. G. VERWER AND B. P. SOMMEIJER, *An implicit-explicit Runge-Kutta-Chebyshev scheme for diffusion-reaction equations*, SIAM J. Sci. Comput., 25 (2004), pp. 1824–1835.