

Error Inhibiting Block One-Step Schemes for Ordinary Differential Equations

A. Ditkowski

School of Mathematical Sciences, Tel Aviv University, Tel Aviv 69978, Israel

S. Gottlieb

Department of Mathematics, University of Massachusetts, Dartmouth, 285 Old Westport Road, North Dartmouth, MA 02747

Abstract

The commonly used one step methods and linear multi-step methods all have a global error that is of the same order as the local truncation error (as defined in [6,13,1,8,15]). In fact, this is true of the entire class of general linear methods. In practice, this means that the order of the method is typically defined solely by order conditions which are derived by studying the local truncation error. In this work we investigate the interplay between the local truncation error and the global error, and develop a methodology which defines the construction of explicit *error inhibiting* block one-step methods (alternatively written as explicit general linear methods [2]). These *error inhibiting schemes* are constructed so that the accumulation of the local truncation error over time is controlled, which results in a global error that is one order higher than the local truncation error. In this work, we delineate how to carefully choose the coefficient matrices so that the growth of the local truncation error is inhibited. We then use this theoretical understanding to construct several methods that have higher order global error than local truncation error, and demonstrate their enhanced order of accuracy on test cases. These methods demonstrate that the error inhibiting concept is realizable. Future work will further develop new error inhibiting methods and will analyze the computational efficiency and linear stability properties of these methods.

Key words: ODE solvers, General linear methods, One-step methods, Global error, local truncation error, Error inhibiting schemes.

Email addresses: `adid@post.tau.ac.il` (A. Ditkowski),
`sgottlieb@umassd.edu` (S. Gottlieb).

1 Introduction

When solving an ordinary differential equation (ODE) of the form

$$\begin{aligned} u_t &= F(t, u), & t \geq 0 \\ u(t=0) &= u_0 \end{aligned} \tag{1}$$

One can evolve the solution forward in time using the first order forward Euler method

$$v_{n+1} = v_n + \Delta t F(t_n, v_n).$$

To obtain a more accurate solution, one can use methods with multiple steps:

$$v_{n+1} = \sum_{j=1}^s a_j v_{n+1-j} + \Delta t \sum_{j=0}^s b_j F(t_{n+1-j}, v_{n+1-j}), \tag{2}$$

known as linear multistep methods [3]. Alternatively, one can use multiple stages, such as Runge–Kutta methods [3]:

$$\begin{aligned} y_i &= F\left(v_n + \sum_{j=1}^m a_{ij} y^{(j)}, t_n + \Delta t \sum_{j=1}^m a_{ij}\right) \quad \text{for } j = 1, \dots, m \\ v_{n+1} &= v_n + \Delta t \sum_{j=1}^m b_j y_j. \end{aligned}$$

The class of general linear methods described in [2,9] combines the use of multiple steps and stages, constructing methods of the form:

$$\begin{aligned} y_i &= \sum_{j=1}^s \tilde{U}_{ij} v_n + \Delta t \sum_{j=1}^m \tilde{A}_{ij} f(y_j) \\ v_{n+1}^i &= \sum_{j=1}^s \tilde{V}_{ij} v_n^i + \Delta t \sum_{j=1}^m \tilde{B}_{ij} f(y_j). \end{aligned} \tag{3}$$

The inclusion of multiple derivatives, such as Taylor series methods [3],

$$v_{n+1} = v_n + \Delta t F(t_n, v_n) + \frac{\Delta t^2}{2} F'(t_n, v_n) + \frac{\Delta t^3}{3!} F''(v^n),$$

is another possibility, and multiple stages and derivatives have also been developed and used successfully [17], [18], [11], [10], [4].

For time-dependent problems the global error, which is the difference between the numerical and exact solution at any given time $t_n = n\Delta t$:

$$E_n = v_n - u(t_n),$$

depends on the local truncation error which, roughly speaking, is the accuracy over one time step. In our case we define the local truncation error as the error

of the method over one time-step, normalized by Δt . For example, the local truncation error for Euler’s method is (following [6,13,1,8,15])

$$\tau = \frac{u(t_{n+1}) - u(t_n) - \Delta t F(t_n, u(t_n))}{\Delta t} \approx O(\Delta t).$$

(To avoid confusion it is important to note that sometimes the truncation error is defined a little differently than we define it above and is not normalized by Δt).

A well known theoretical result, known as the Lax-Richtmeyer equivalence theorem (see e.g. [12], [6], [13]) states that if the numerical scheme is stable then the global error is at least of the same order as the local truncation error. In all the schemes for numerically solving ordinary differential equations (ODEs) that we are familiar with from the literature, the global errors are indeed of the same order as their local truncation errors¹. This is common to other fields in numerical mathematics, such as for finite difference schemes for partial differential equations (PDEs), see e.g. [6,13]. It was recently demonstrated, however, that finite difference schemes for PDEs can be constructed such that their convergence rates, or the order of their global errors, are higher than the order of the truncation errors [5]. In this work we adopt and adapt the ideas presented in [5] to show that it is possible to construct numerical methods for ODEs where the the global error is *one order higher* than the local truncation error. As we discuss below, these schemes achieve this higher order by inhibiting the lowest order term in the local error from accumulating over time, and so we name them *Error Inhibiting Schemes*.

Following an idea in [14], an interesting paper by Shampine and Watt in 1969 [16] describes a class of implicit one-step methods that obtain a block of s new step values at each step. These methods take s initial step values and generate the next s step values, and so on, all in one step. These methods are in fact explicit block one-step methods, and can be written as general linear methods of the form (3) above. Inspired by this form, we construct explicit block one-step methods which are in the form (3), but where the matrix \tilde{U} is an identity matrix, and the matrix \tilde{A} is all zeros; these are known as Type 3 methods in [2]. The major feature of our methods is that in addition to satisfying the appropriate order conditions listed in [2], they have a special structure that mitigates the accumulation of the truncation error, so we obtain a global error that is one order *higher* than predicted by the order conditions in [2], which describe the local truncation error.

In Section 2 we motivate our approach by describing how typical multistep methods can be written and analyzed as block one-step methods: these methods obtain a block of s new step values at each step. We show how this form allows us to precisely describe the growth of the error over the time-evolution.

¹ In the case where the truncation error is defined without the Δt normalization the global error is one order lower than the truncation error.

In Section 3 we then exploit this understanding to develop explicit error inhibiting block one-step methods that produce higher order global errors than possible for typical multistep methods. In Section 4 we present some methods developed according to the theory in Section 3 and we test these methods on several numerical examples to demonstrate that the order of convergence is indeed one order higher than the local truncation error. We also show that, in contrast to our error inhibiting Type 3 method, a typical Type 3 method developed by Butcher in [2] does not satisfy the critical condition for a method to be error inhibiting and therefore produces a global error that is of the same order as the local truncation error. Finally, we present our conclusions in Section 5, and suggest that further investigation of error inhibiting methods shall include the analysis of their linear stability properties, storage implications, and computational efficiency.

2 Motivation

In this section we present the analysis of explicit multistep methods in a block one-step form for a simple linear problem. In this familiar setting we define the local truncation error, the global error, and the solution operator that connects them. We also discuss the stability of a method of this form. We limit our analysis to the linear case so that we can clearly observe the process by which the solution operator interacts with the local truncation error, and results in a global error that is of the same order as the local truncation error. Although we are dealing for the moment with standard multistep methods, this will set the stage for the construction and analysis of error inhibiting block one-step methods.

In order to illustrate the main idea we start with a linear ordinary differential equation (ODE)

$$\begin{aligned} u_t &= f(t) u, & t \geq 0 \\ u(t=0) &= u_0 \end{aligned} \tag{4}$$

where $f(t) < M$, $\forall t \geq 0$ and $f(t)$ is analytic.

An s -step explicit multistep method applied to (4) takes the form

$$v_{n+s} = \sum_{j=0}^{s-1} a_j v_{n+j} + \Delta t \sum_{j=0}^{s-1} b_j F(t_{n+j}, v_{n+j}) = \sum_{j=0}^{s-1} a_j v_{n+j} + \Delta t \sum_{j=0}^{s-1} b_j f(t_{n+j}) v_{n+j} \tag{5}$$

where the time domain is discretized by the sequence $t_n = n \Delta t$, and v_n denotes the numerical approximation of $u(t_n)$. The method (5) is defined by its coefficients $\{a_j\}_{j=0}^{s-1}$ and $\{b_j\}_{j=0}^{s-1}$, which are constant values.

Following [6] we rewrite the method (5) in its block form. To do this, we first

introduce the exact solution vector

$$U_n = (u(t_{n+s-1}), \dots, u(t_n))^T \quad (6)$$

and similarly, the numerical solution vector is

$$V_n = (v_{n+s-1}, \dots, v_n)^T. \quad (7)$$

Now (5) can be written in block form so that it looks like a one step scheme

$$V_{n+1} = Q_n V_n \quad (8)$$

where

$$Q_n = \begin{pmatrix} a_{s-1} + \Delta t b_{s-1} f(t_{n+s-1}) & a_{s-2} + \Delta t b_{s-2} f(t_{n+s-2}) & \dots & a_0 + \Delta t b_0 f(t_n) \\ & I & & \\ & & \ddots & \\ & & & I & 0 \end{pmatrix}. \quad (9)$$

From repeated applications of equation (8) we observe that the numerical solution vector V_n at any time t_n can be related to V_ν for any previous time t_ν

$$V_n = S_{\Delta t}(t_n, t_\nu) V_\nu, \quad \nu \leq n \quad (10)$$

where $S_{\Delta t}$ is the discrete solution operator. This operator can be expressed explicitly by

$$S_{\Delta t}(t_n, t_\nu) = Q_{n-1} \dots Q_{\nu+1} Q_\nu, \quad S_{\Delta t}(t_n, t_n) = I. \quad (11)$$

For simplicity we can define this by

$$\prod_{\mu=\nu}^{n-1} Q_\mu \equiv Q_{n-1} \dots Q_{\nu+1} Q_\nu, \quad \prod_{\mu=n}^{n-1} Q_\mu \equiv I. \quad (12)$$

Note that if each matrix Q_μ is independent of μ (in other words, in the constant coefficient case where f is independent of t), we simply have a product of matrices Q , and the discrete solution operator becomes

$$S_{\Delta t}(t_n, t_\nu) = Q^{n-\nu}. \quad (13)$$

The behavior of a method depends in large part on the accuracy of its solution operator. We begin by defining the local truncation error as the error of the method over one time-step, normalized by Δt :

Definition 1: *The local truncation error τ_n is given by [6,13,1,8,15]*

$$\Delta t \tau_n = U_{n+1} - Q_n U_n \quad (14)$$

Note that in the case of the standard multistep method, where Q_n is given by the matrix (9), the truncation error has only one non-zero element:

$$\boldsymbol{\tau}_n = (\tau_n, 0, \dots, 0)^T. \quad (15)$$

The error that we are most interested in is the difference between the exact error vector and the numerical error vector at time t_n ,

$$E_n = U_n - V_n, \quad (16)$$

known as the global error. At the initial time, we have the error E_0 which is based on the starting values a method of this sort requires: the values v_j , $j = 0, \dots, s-1$ that are prescribed or somehow computed. Typically, v_0 is the initial condition defined in (1) and v_j , $j = 1, \dots, s-1$ are computed to sufficient accuracy using some other numerical scheme. Thus, the value of E_0 is assumed to be as small as needed.

The evolution of the global error (16) depends on the local truncation error defined by (14) and the discrete solution operator given in (8):

$$E_{n+1} = Q_n E_n + \Delta t \boldsymbol{\tau}_n. \quad (17)$$

Unraveling this equality all the way back to E_0 gives

$$E_n = S_{\Delta t}(t_n, 0) E_0 + \Delta t \sum_{\nu=0}^{n-1} S_{\Delta t}(t_n, t_{\nu+1}) \boldsymbol{\tau}_\nu, \quad (18)$$

or, equivalently

$$E_n = \prod_{\mu=0}^{n-1} Q_\mu E_0 + \Delta t \sum_{\nu=0}^{n-1} \left(\prod_{\mu=\nu+1}^{n-1} Q_\mu \right) \boldsymbol{\tau}_\nu. \quad (19)$$

(This formula is obtained from the discrete version of Duhamel's principle, see Lemma 5.1.1 in [6]).

It is clear from (18) that the behavior of the discrete solution operator $S_{\Delta t}(t_n, t_{\nu+1})$ must be controlled for this error to converge. This property defines the stability of the method. Also here we use the stability definition presented in [6], namely:

Definition 2: *The scheme (8) is called stable if there are constants α_s and K_s , independent of Δt , such that for all $0 < \Delta t \leq \Delta t_0$*

$$\|S_{\Delta t}(t_n, t_\nu)\| \leq K_s e^{\alpha_s(t_n - t_\nu)} \quad (20)$$

If the scheme is stable, we can use (20) and (18) to bound the growth of the error:

$$\|E_n\| \leq K_s \left[e^{\alpha_s t_n} \|E_0\| + \max_{0 \leq \nu \leq n-1} \|\tau_\nu\| \phi_h^*(\alpha_s, t_n) \right]. \quad (21)$$

where

$$\phi_{\Delta t}^*(\alpha_s, t_n) = \Delta t \sum_{\nu=0}^{n-1} e^{\alpha_s(t_n-t_{\nu+1})} \approx \int_0^{t_n} e^{\alpha_s(t_n-\zeta)} d\zeta = \begin{cases} \frac{e^{\alpha_s t_n} - 1}{\alpha_s} & \alpha_s \neq 0 \\ t_n & \alpha_s = 0 \end{cases}. \quad (22)$$

Equation (21) means that stability implies convergence:² if the scheme is stable then the global error is controlled by the local truncation error for any given final time. In the formula above it is clear that the global error must have order at least as high as the local truncation error, but the possibility of having a higher order global error is left open.

The first Dahlquist barrier [7,3] states that any explicit s step linear multistep method can be of order p no higher than s . It is the common experience that methods have global error of the same order as the local truncation error. These two together greatly limit the accuracy of the methods we can derive.

Remark 1 *In an Adams-Bashforth scheme the entry in the first row and first column in the term $S_{\Delta t}(t_n, t_\nu) = \prod_{\mu=\nu}^{n-1} Q_\mu$ is equal to $1 + O(\Delta t)$. Therefore the error, due to the accumulation of the contributions from the truncation errors, becomes:*

$$\mathbf{e}_{n+s} = \Delta t \sum_{\nu=0}^{n-1} (1 + O(\Delta t)) \tau_\nu \quad (23)$$

which is approximately the average value of τ_ν over $\nu = 0, \dots, n-1$. This suggests that we may need to look outside the family of linear multistep methods to attain a higher order global error.

The analysis in this section suggests that if the operator Q_n is properly constructed, the growth of the global error described in Equation (19) may be controlled through the properties of the operator Q_n and its relationship with the local truncation error τ_n . However, as implied by the example of the Adams-Bashforth scheme above, we need to construct methods where the operator Q_n is not limited by the structure in this section. In the next section we present the construction of block one-step methods that are error inhibiting. The class of methods described by this block one-step structure is very broad: while all classical multistep methods can be written in this block form, not every such block one-step method can be written as a classical multistep method. Thus, we rely on the discussion in this section with one main change: the structure of the operator Q_n .

² For partial differential equations this result is known as one part of the celebrated Lax-Richtmeyer equivalence theorem. See e.g. [12], [6], [13].

3 An Error Inhibiting Methodology

In Section 2 we rewrote explicit linear multistep methods in a block one-step form, and expressed the relationship between its local and global error. We observed that the growth of the local errors is driven by the behavior of the discrete solution operator Q_n , and in particular its interaction with the local truncation error. Using this insight we show in this section that it is possible to construct such explicit block one-step methods (which are also known as Type 3 DIMSIM methods in [2]) that *inhibit* the growth of the truncation error so that the global error (16) gains an order of accuracy over the local truncation error (14).

We begin in Section 3.1 by describing the construction and analysis of error inhibiting block one-step schemes for the case of linear constant coefficient equations. We then show that this approach yields methods that are also error inhibiting for variable coefficient linear equations in Section 3.2 and nonlinear equations in Section 3.3.

3.1 Error inhibiting schemes for linear constant coefficient equations

Given a linear ordinary differential equation with constant coefficients:

$$\begin{aligned} u_t &= f \cdot u, \quad \text{for } t \geq 0, \\ u(t=0) &= u_0 \end{aligned} \tag{24}$$

where $f \in \mathbb{R}$. We define a vector of length s that contains the exact solution of (24) at times $(t_n + j\Delta t/s)$ for $j = 0, \dots, s-1$

$$U_n = \left(u(t_{n+(s-1)/s}), \dots, u(t_{n+1/s}), u(t_n) \right)^T, \tag{25}$$

and the corresponding vector of numerical approximations

$$V_n = \left(v_{n+(s-1)/s}, \dots, v_{n+1/s}, v_n \right)^T. \tag{26}$$

Note that although we are assuming that the solution u at any given time is a scalar, this entire discussion easily generalizes to the case where u is a vector, with only some cumbersome notation needed. Thus without loss of generality we continue the discussion with scalar notation.

Remark 2 *The notation above emphasizes that this scheme uses s terms for generating the next s terms, unlike the explicit linear multistep methods in the section above which use s terms to generate one term. To match with the notation in Section 2 above, we can replace $\Delta t' = s\Delta t$ thus defining this scheme on integer grid points.*

We define the block one-step method for the linear constant coefficient problem

(24)

$$V_{n+1} = QV_n \quad (27)$$

where

$$Q = A + \Delta t Bf \quad (28)$$

and $A, B \in \mathbb{R}^{s \times s}$. Unlike in the case of classical multistep methods, here we do not restrict the structure of the matrices A and B . Thus, any multistep method of the form (5) can be written in this form (as we saw above), but not every method of the form (28) can be written as a multistep method. In fact, this method is a general linear method of the DIMSIM form (3) with \tilde{A} is all zeroes, \tilde{U} is the identity matrix, $\tilde{V} = A$, and $\tilde{B} = B$. This particular formulation is, as we mentioned above, called a Type 3 DIMSIM in Butcher's 1993 paper [2].

At any time t_n , we know that $u(t_n + \Delta t) = u(t_n) + O(\Delta t)$, so that for the numerical solution V_n to converge to the analytic solution U_n one of the eigenvalues of Q must be equal to $1 + O(\Delta t)$, and its eigenvector must have the form:

$$(1 + O(\Delta t), \dots, 1 + O(\Delta t))^T. \quad (29)$$

The structure of the eigensystem of A , which is the leading part of Q , is critical to the stability of the scheme and the dynamics of the error.

Suppose A is constructed such that:

- (1) $\text{rank}(A) = 1$.
- (2) Its non-zero eigenvalue is equal to one and its corresponding eigenvector is $(1, \dots, 1)^T$
- (3) A can be diagonalized.

Property (2) assures that the method produces the exact solution for the case $f = 0$. Now, since the term $\Delta t Bf$ is only an $O(\Delta t)$ perturbation to A , the matrix Q will have one eigenvalue, $z_1 = 1 + O(\Delta t)$ whose eigenvector has the form

$$\psi_1 = (1 + O(\Delta t), \dots, 1 + O(\Delta t))^T \quad (30)$$

and the rest of the eigenvalues satisfy $z_j = O(\Delta t)$ for $j = 2, \dots, s$.

Since the $\|Q\| = 1 + O(\Delta t)$, we can conclude that there exist constants K_s and α_s such that

$$\|S_{\Delta t}(t_n, t_\nu)\| = \|Q^{n-\nu}\| \leq K_s e^{\alpha_s(t_n - t_\nu)} \quad (31)$$

where $\alpha_s = \|B\| |f|$. Therefore, according to Definition 2, the scheme (27) is stable. By the same argument used above, we can show that the global error will have order that is no less than the order of the local truncation error.

We now turn to the task of investigating the truncation error, τ_n . The definition of the local truncation error in this case is still

$$\Delta t \tau_n = U_{n+1} - Q_n U_n$$

as defined in the previous section in Equation (14).

Remark 3 Since $Q = A + \Delta t B f$ and $u_t = f u$ the local truncation error can be written as

$$\Delta t \boldsymbol{\tau}_n = U_{n+1} - \left(A U_n + \Delta t B \frac{dU_n}{dt} \right).$$

Therefore $\boldsymbol{\tau}_n$ does not explicitly depend on f . This observation is valid for the variable coefficients and the nonlinear case as well.

The definition of the error is

$$E_n = U_n - V_n,$$

as in Equation (16). The evolution of the error is still described by Equation (19)

$$E_n = \prod_{\mu=0}^{n-1} Q_{\mu} E_0 + \Delta t \sum_{\nu=0}^{n-1} \left(\prod_{\mu=\nu+1}^{n-1} Q_{\mu} \right) \boldsymbol{\tau}_{\nu},$$

which in the linear constant coefficient case becomes

$$E_n = Q^n E_0 + \Delta t \sum_{\nu=0}^{n-1} Q^{n-\nu-1} \boldsymbol{\tau}_{\nu}. \quad (32)$$

The main difference between this case and the linear multistep method in Section 2 is that the structure of Q is different, and that unlike (15), in this case all the entries in $\boldsymbol{\tau}_n$ are typically non-zero.

Equation (32) indicates that there are several sources for the error at the time t_n :

- (1) *The initial error E_0 which is the error in the initial condition V_0 :* This error is caused primarily by the numerical scheme used to compute the first $s - 1$ elements in V_0 . We assume these errors can be made arbitrary small. The initial value, which is the final element of V_0 , is taken from the analytic initial condition and is considered to be accurate to machine precision.
- (2) *The term $\Delta t \boldsymbol{\tau}_{n-1}$, which is the last term in the sum in the right hand side of (32):* This term is clearly, by definition, of the size $O(\Delta t) \|\boldsymbol{\tau}_{n-1}\|$.
- (3) *The summation*

$$\Delta t \sum_{\nu=0}^{n-2} Q^{n-\nu-1} \boldsymbol{\tau}_{\nu}, \quad (33)$$

which are all the rest of the terms in the sum in the right hand side of (18): This is the term we need to bound to control the growth of the truncation error.

The terms in the sum (33) are all comprised of the discrete solution operator Q multiplying the local truncation error. This leads us to the major observation that is the key to constructing error inhibiting methods: **if the local truncation error lives in the subspace of eigenvectors that correspond**

to the eigenvalues of $O(\Delta t)$, then the growth of the truncation error will be inhibited, and the global error will be one order higher than the local truncation error.

Recall that Q has one dominant eigenvalue that has the form $1 + O(\Delta t)$ and all the others are $O(\Delta t)$. Correspondingly, two subspaces can be defined

$$\Psi_1 = \text{span} \{ \psi_1 \} \quad \text{and} \quad \Psi_1^c = \text{span} \{ \psi_2, \dots, \psi_s \}$$

where ψ_j is the eigenvector associated with each eigenvalue z_j . As ψ_j can be normalized, we assume that $\|\psi_j\| = O(1)$. It should be noted that while Ψ_1 and Ψ_1^c are linearly independent, they are not orthogonal subspaces. Furthermore, since the matrix A is diagonalizable by construction, its eigenvectors span \mathbb{R}^s . Since $\boldsymbol{\tau}_\nu \in \mathbb{R}^s$, it can be written as

$$\boldsymbol{\tau}_\nu = \gamma_1 \psi_1 + \sum_{j=2}^s \gamma_j \psi_j \quad (34)$$

where $\gamma_1 \psi_1 \in \Psi_1$ and $\sum_{j=2}^s \gamma_j \psi_j \in \Psi_1^c$.

Of course, the truncation error $\boldsymbol{\tau}_\nu$ is determined by the entries of Q . To ensure that the local truncation error is mostly in the space Ψ_1^c of eigenvectors which correspond to the eigenvalues of size $O(\Delta t)$, we choose the entries of Q (i.e. the entries of A and B) such that $\gamma_1 = O(\Delta t)$, which will mean that

$$\|\gamma_1 \psi_1\| = O(\Delta t) \|\boldsymbol{\tau}_\nu\| . \quad (35)$$

Using this, we can bound product of the discrete solution operator and the truncation error,

$$\begin{aligned} \|Q\boldsymbol{\tau}_\nu\| &= \left\| \gamma_1 Q\psi_1 + \sum_{j=2}^s \gamma_j Q\psi_j \right\| \leq \|\gamma_1 Q\psi_1\| + \left\| \sum_{j=2}^s \gamma_j Q\psi_j \right\| \\ &= \|\gamma_1 z_1 \psi_1\| + \left\| \sum_{j=2}^s \gamma_j z_j \psi_j \right\| \leq |z_1| \|\gamma_1 \psi_1\| + \max_{j=2, \dots, s} |z_j| \left\| \sum_{j=2}^s \gamma_j \psi_j \right\| \\ &\leq |z_1| \|\gamma_1 \psi_1\| + \max_{j=2, \dots, s} |z_j| \|\boldsymbol{\tau}_\nu - \gamma_1 \psi_1\| \\ &\leq (1 + O(\Delta t)) O(\Delta t) \|\boldsymbol{\tau}_\nu\| + O(\Delta t) \|\boldsymbol{\tau}_\nu\| = O(\Delta t) \|\boldsymbol{\tau}_\nu\| \end{aligned}$$

where z_j are the eigenvalues of Q . Therefore we have

$$\|Q\boldsymbol{\tau}_\nu\| \leq O(\Delta t) \|\boldsymbol{\tau}_\nu\| . \quad (36)$$

Whenever the condition (36) is satisfied, we can show that the sum (33) above is bounded:

$$\begin{aligned}
\left\| \Delta t \sum_{\nu=0}^{n-2} Q^{n-\nu-1} \boldsymbol{\tau}_\nu \right\| &= \Delta t \left\| \sum_{\nu=0}^{n-2} Q^{n-\nu-1} \boldsymbol{\tau}_\nu \right\| \leq \Delta t \sum_{\nu=0}^{n-2} \|Q^{n-\nu-2}\| \|Q \boldsymbol{\tau}_\nu\| \\
&\leq \Delta t \sum_{\nu=0}^{n-2} \|Q\|^{n-\nu-2} O(\Delta t) \|\boldsymbol{\tau}_\nu\| \\
&\leq \Delta t \left(\max_{\nu=0, \dots, n-2} \|\boldsymbol{\tau}_\nu\| \right) \sum_{\nu=0}^{n-2} (1 + c\Delta t)^{n-\nu-2} O(\Delta t) \\
&\leq \Delta t \left(\max_{\nu=0, \dots, n-2} \|\boldsymbol{\tau}_\nu\| \right) \sum_{\nu=0}^{n-2} \left[e^{c\Delta t} (1 + O(\Delta t^2)) \right]^{n-\nu-2} O(\Delta t) \\
&\leq \Delta t \left(\max_{\nu=0, \dots, n-2} \|\boldsymbol{\tau}_\nu\| \right) \sum_{\nu=0}^{n-2} \left[e^{c(t_{n-2}-t_\nu)} (1 + O(\Delta t)) \right] O(\Delta t) \\
&\leq O(\Delta t) \left(\max_{\nu=0, \dots, n-2} \|\boldsymbol{\tau}_\nu\| \right) \phi_{\Delta t}^*(c, T). \tag{37}
\end{aligned}$$

(Recall (22) for the definition of $\phi_{\Delta t}^*(c, T)$.)

In the final equation, T is the final time, and the term $\phi_{\Delta t}^*(c, T)$ is therefore a constant. Thus we have the bound

$$\left\| \Delta t \sum_{\nu=0}^{n-2} Q^{n-\nu-1} \boldsymbol{\tau}_\nu \right\| \leq O(\Delta t) \max_{\nu=0, \dots, n-2} \|\boldsymbol{\tau}_\nu\|. \tag{38}$$

Putting this all together into (32), we obtain

$$\|E_n\| = O(\Delta t) \max_{\nu=0, \dots, n-1} \|\boldsymbol{\tau}_\nu\|. \tag{39}$$

Thus, if the coefficients of A and B are chosen so that we can control the size of $\|Q \boldsymbol{\tau}_\nu\|$ in (36), we can obtain a scheme that inhibits the growth of the local truncation error, so that the global error is one order more accurate than its truncation error.

3.2 Linear variable-coefficient equations

In the previous section we showed how to construct an error inhibiting method by choosing the coefficients in A and B so that the local truncation error lives (mostly) in the space that is spanned by the eigenvectors corresponding to eigenvalues that are of $O(\Delta t)$. In this section we show that under the same criteria as above, these methods are also error inhibiting when applied to a *variable coefficient* linear ordinary differential equation:

$$\begin{aligned}
u_t &= f(t)u, \quad t \geq 0 \\
u(t=0) &= u_0
\end{aligned} \tag{40}$$

where $f(t)$ assumed to be analytic or as smooth as needed, and bounded. In this case the scheme is given by a time-dependent evolution operator Q_n which

may change each time-step:

$$V_{n+1} = Q_n V_n \quad (41)$$

where

$$Q_n = A + \Delta t B \begin{pmatrix} f(t_{n+(s-1)/s}) & & & \\ & f(t_{n+(s-2)/s}) & & \\ & & \ddots & \\ & & & f(t_n) \end{pmatrix} \quad (42)$$

and the matrices A and B are the same as described above for the constant coefficient scheme.

Since $f(t)$ is an analytic function, Q_n can be written as

$$Q_n = A + \Delta t B f(t_n) + \Delta t^2 B f'(t_n) \begin{pmatrix} ((s-1)/s) & & & \\ & ((s-2)/s) & & \\ & & \ddots & \\ & & & 0 \end{pmatrix} + O(\Delta t^3) \quad (43)$$

We can also say then that

$$Q_n = A + \Delta t B f(t_n) + O(\Delta t^2) B f'(t_n) = \tilde{Q}_n + O(\Delta t^2). \quad (44)$$

Each \tilde{Q}_n has the same structure as Q in the constant coefficient case. In particular

$$\|\tilde{Q}_n\| = (1 + O(\Delta t)) \leq 1 + c\Delta t, \quad \forall n. \quad (45)$$

Furthermore, as was pointed out in Remark 3, since the local truncation error τ_n does not depend explicitly on $f(t)$ at any time t_n , we can write τ_n as a linear combination of the eigenvectors of A that correspond to the zero eigenvalues. Thus, τ_n lives (mostly) in the space that is spanned by the eigenvectors of any matrix \tilde{Q}_n corresponding to eigenvalues that are of $O(\Delta t)$. We can then follow the same analysis as in (35)–(36), to obtain the bound

$$\|\tilde{Q}_{n+1} \tau_n\| = O(\Delta t) \|\tau_n\|, \quad \forall n. \quad (46)$$

In this case, Equation (18) takes the modified form (for $n \geq 1$)

$$\begin{aligned} E_n &= \prod_{\mu=0}^{n-1} Q_\mu E_0 + \Delta t \sum_{\nu=0}^{n-1} \left(\prod_{\mu=\nu+1}^{n-1} Q_\mu \right) \tau_\nu \\ &= \prod_{\mu=0}^{n-1} Q_\mu E_0 + \Delta t \sum_{\nu=0}^{n-2} \prod_{\mu=\nu+1}^{n-1} (\tilde{Q}_\mu + O(\Delta t^2)) \tau_\nu + \Delta t \tau_{n-1} \end{aligned}$$

The first term is negligible because we assume that the initial error can be made arbitrarily small, and the final term is clearly of order $\Delta t \boldsymbol{\tau}_{n-1}$. Using (45), (46) and the same analysis as in (35)–(38) we have

$$\begin{aligned}
\left\| \Delta t \sum_{\nu=0}^{n-2} \left(\prod_{\mu=\nu+1}^{n-1} \tilde{Q}_\mu \right) \boldsymbol{\tau}_\nu \right\| &= \left\| \Delta t \sum_{\nu=0}^{n-2} \left(\prod_{\mu=\nu+2}^{n-1} \tilde{Q}_\mu \right) (\tilde{Q}_{\nu+1} \boldsymbol{\tau}_\nu) \right\| \\
&\leq \Delta t \sum_{\nu=0}^{n-2} \left\| \prod_{\mu=\nu+2}^{n-1} \tilde{Q}_\mu \right\| \left\| \tilde{Q}_{\nu+1} \boldsymbol{\tau}_\nu \right\| \\
&\leq \Delta t \sum_{\nu=0}^{n-2} O(1 + O(\Delta t))^{n-\nu-2} O(\Delta t) \|\boldsymbol{\tau}_\nu\| \\
&\leq O(\Delta t) \max_{\nu=0, \dots, n-2} \|\boldsymbol{\tau}_\nu\|.
\end{aligned}$$

Putting these all together we have

$$\|E_n\| = O(\Delta t) \max_{\nu=0, \dots, n-1} \|\boldsymbol{\tau}_\nu\|. \quad (47)$$

This simple proof shows that even for the variable coefficient case, the schemes constructed as described above have a higher order error than would be expected from the truncation error. In the next subsection we extend this analysis to the general nonlinear case.

3.3 Nonlinear equations

Finally, we analyze the behavior of methods satisfying the assumptions in Section 3.1 when applied to nonlinear problems. Consider the nonlinear equation

$$\begin{aligned}
u_t &= f(u(t), t), \quad t \geq 0 \\
u(t=0) &= u_0
\end{aligned} \quad (48)$$

where $f(u, t)$ assumed to be analytic in u and t . We now use the scheme

$$V_{n+1} = AV_n + \Delta t B \begin{pmatrix} f(v_{n+(s-1)/s}, t_{n+(s-1)/s}) \\ \vdots \\ f(v_n, t_n) \end{pmatrix} \quad (49)$$

where the matrices A and B are as constructed above for the constant coefficients problem.

As defined in (14), the exact solution to (48) and the truncation error are

related by

$$U_{n+1} = AU_n + \Delta t B \begin{pmatrix} f(u_{n+(s-1)/s}, t_{n+(s-1)/s}) \\ \vdots \\ f(u_n, t_n) \end{pmatrix} + \Delta t \tau_n. \quad (50)$$

Note that by Taylor expansion

$$f(v_\nu, t_\nu) = f(u_\nu, t_\nu) + f_u(u_\nu, t_\nu)(v_\nu - u_\nu) + r(v_\nu - u_\nu),$$

where $f_u(u, t) = \partial f(u, t)/\partial u$ and $|r(v_\nu - u_\nu)| \leq c_1|v_\nu - u_\nu|^2$. Subtracting (49) from (50) and assuming that $E_n = U_n - V_n \ll 1$ gives

$$E_{n+1} = AE_n - \Delta t B \begin{pmatrix} f_u(u_{n+(s-1)/s}, t_{n+(s-1)/s}) \\ \ddots \\ f_u(u_n, t_n) \end{pmatrix} E_n + \Delta t \tau_n + \Delta t R(E_n) \quad (51)$$

where $\|R(E_n)\| \leq c_1\|E_n\|^2$. Equation (51) means that as long as $O(E_n^2) \ll O(\tau_n)$, the equation for the error E_n can be analyzed in essentially the same way as for the linear variable coefficient case, and the same estimates hold.

In order to evaluate the time interval in which $O(E_n^2) \ll O(\tau_n)$ we note that although the term $R(E_n)$ in (51) is not a non-homogeneous term but rather a function of E_n , we can still use the approach used in [6, Theorem 5.1.2]) to prove stability for a perturbed solution operator. As in [6, Theorem 5.1.2]), we use the discrete version of Duhamel's principle to obtain

$$E_n = \prod_{\mu=0}^{n-1} \hat{Q}_\mu E_0 + \Delta t \sum_{\nu=0}^{n-1} \left(\prod_{\mu=\nu+1}^{n-1} \hat{Q}_\mu \right) \tau_\nu + \Delta t \sum_{\nu=0}^{n-1} \left(\prod_{\mu=\nu+1}^{n-1} \hat{Q}_\mu \right) R(E_\nu) \quad (52)$$

where

$$\hat{Q}_n = A - \Delta t B \begin{pmatrix} f_u(u_{n+(s-1)/s}, t_{n+(s-1)/s}) \\ \ddots \\ f_u(u_n, t_n) \end{pmatrix}. \quad (53)$$

Taking the norm of (52) and using the triangle inequality we obtain

$$\|E_n\| \leq \left\| \prod_{\mu=0}^{n-1} \hat{Q}_\mu E_0 \right\| + \left\| \Delta t \sum_{\nu=0}^{n-1} \left(\prod_{\mu=\nu+1}^{n-1} \hat{Q}_\mu \right) \tau_\nu \right\| + \left\| \Delta t \sum_{\nu=0}^{n-1} \left(\prod_{\mu=\nu+1}^{n-1} \hat{Q}_\mu \right) R(E_\nu) \right\|. \quad (54)$$

As in the linear case we assume that the initial error, E_0 is arbitrary small, so the first term is negligible. If $\hat{Q}_{\nu+1}$ is constructed such that $\|\hat{Q}_{\nu+1}\boldsymbol{\tau}_\nu\| = \Delta t O(\boldsymbol{\tau}_\nu)$ then using the same analysis as in variable coefficient case the second term in (54) is less or equal to $\Delta t c_0 \phi_h^*(c, t_n) \max_{\nu=0, \dots, n-1} \|\boldsymbol{\tau}_\nu\|$. As to the third term, the same arguments can be used to show that it is bounded by

$$\left\| \Delta t \sum_{\nu=0}^{n-1} \left(\prod_{\mu=\nu+1}^{n-1} \hat{Q}_\mu \right) R(E_\nu) \right\| \leq c_1 \phi_h^*(c, t_n) \|E_n\|^2, \quad (55)$$

so that (54) (with the substitution of (55) for the final term) can be re-arranged to obtain

$$\|E_n\| (1 - c_1 \phi_h^*(c, t_n) \|E_n\|) \leq \Delta t c_0 \phi_h^*(c, t_n) \max_{\nu=0, \dots, n-1} \|\boldsymbol{\tau}_\nu\|. \quad (56)$$

If $c_1 \phi_h^*(c, t_n) \|E_n\| < 1/2$, we obtain

$$\|E_n\| \leq 2 \Delta t c_0 \phi_h^*(c, t_n) \max_{\nu=0, \dots, n-1} \|\boldsymbol{\tau}_\nu\| \quad (57)$$

This estimate holds as long as

$$c_1 \phi_h^*(c, t_n) \|E_n\| \leq 2 \Delta t c_0 c_1 (\phi_h^*(c, t_n))^2 \max_{\nu=0, \dots, n-1} \|\boldsymbol{\tau}_\nu\| \leq \frac{1}{2}, \quad (58)$$

which is satisfied for all times t_n such that $\Delta t \phi_h^*(c, t_n) = O(1)$.

Therefore

$$\|E_n\| = O(\Delta t) \max_{\nu=0, \dots, n-1} \|\boldsymbol{\tau}_\nu\|.$$

for the nonlinear case as well.

4 Some Error Inhibiting Explicit Schemes

In the previous section we define sufficient conditions for methods of the form

$$V_{n+1} = Q V_n \quad (59)$$

where

$$Q = A + \Delta t B f$$

to be error inhibiting. These are

C1. $\text{rank}(A) = 1$.

C2. Its non-zero eigenvalue is equal to 1 and its corresponding eigenvector is

$$(1, \dots, 1)^T.$$

C3. A can be diagonalized.

C4. The matrices A and B are constructed such that when the local truncation error is multiplied by the discrete solution operator we have the bound:

$$\|Q\tau_\nu\| \leq O(\Delta t) \|\tau_\nu\|.$$

This is accomplished by requiring the local truncation error to live in the space of the eigenvectors of A that correspond to the zero eigenvalues.

In this section we present several schemes which were constructed using the approach presented in the previous section. In Section 4.1, we present a block one-step method that evolves two steps (v_n and $v_{n+\frac{1}{2}}$) to obtain the next two steps (v_{n+1} and $v_{n+\frac{3}{2}}$). This method has truncation error (14) that is second order, while its global order (16) is third order. We demonstrate that the expected convergence rate is attained on several sample nonlinear problems. In this section we also show that a typical Type 3 DIMSIM method (derived in [2]) that satisfies the first three conditions above but not the fourth, has truncation error of order two, and its global error is of the same order. This demonstrates the importance of condition **C4**.

Next, in Section 4.2 we present a block one-step method that evolves three steps v_n , $v_{n+\frac{1}{3}}$ and $v_{n+\frac{2}{3}}$ to obtain v_{n+1} , $v_{n+\frac{4}{3}}$ and $v_{n+\frac{5}{3}}$. This method has truncation error (14) that is third order, while its global order (16) is fourth order, as we demonstrate on several sample problems. Finally, to show that the methods in each class are not unique, we present two other methods of this type and show that their global error is of one order higher than the local truncation error on a sample nonlinear system.

4.1 A third order error inhibiting method with $s = 2$.

In this subsection we define an explicit block one-step with $s = 2$ that satisfies the conditions **C1** – **C4** above. This method takes the values of the solution at the times t_n and $t_{n+\frac{1}{2}}$ and obtains the solution at the time-level t_{n+1} and $t_{n+\frac{3}{2}}$. The exact solution vector for this problem is

$$U_n = \left(u(t_{n+1/2}), u(t_n) \right)^T$$

and, similarly, the corresponding vector of numerical approximations is

$$V_n = \left(v_{n+1/2}, v_n \right)^T.$$

The scheme is given by:

$$V_{n+1} = \frac{1}{6} \begin{pmatrix} -1 & 7 \\ -1 & 7 \end{pmatrix} V_n + \frac{\Delta t}{24} \begin{pmatrix} 55 & -17 \\ 25 & 1 \end{pmatrix} \begin{pmatrix} f(v_{n+1/2}, t_{n+1/2}) \\ f(v_n, t_n) \end{pmatrix}, \quad (60)$$

and has truncation error

$$\boldsymbol{\tau}_n = \frac{23}{576} \begin{pmatrix} 7 \\ 1 \end{pmatrix} \frac{d^3}{dt^3} u(t_n) \Delta t^2 + O(\Delta t^3). \quad (61)$$

The matrix A can be diagonalized as follows:

$$A = \frac{1}{6} \begin{pmatrix} -1 & 7 \\ -1 & 7 \end{pmatrix} = \frac{1}{6} \begin{pmatrix} 1 & 7 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & \\ & 0 \end{pmatrix} \begin{pmatrix} -1 & 7 \\ 1 & -1 \end{pmatrix}. \quad (62)$$

Observe that the leading order of the truncation error (61) is in the space of the second eigenvector of A , the one that corresponds to the zero eigenvalue. Also, as was pointed out in Remark 3, $\boldsymbol{\tau}_n$ depends only on this eigenvector of A and a multiple that is not directly dependent on f but only on the third derivative of the solution u . This underscores the analysis in Sections 3.2 and 3.3 that demonstrates that the error inhibiting property carries through for variable coefficient and nonlinear problems.

To study the behavior of the global error we use the fact shown in Section 3.3 that even for a nonlinear equation it is sufficient to analyze the matrix

$$Q = A + \Delta t B f \quad (63)$$

where f is a constant. In this case:

$$Q = \frac{1}{6} \begin{pmatrix} 1 + \frac{f\Delta t}{2} + \frac{f^2\Delta t^2}{8} + O(\Delta t^3) & 7 + 36f\Delta t + 228f^2\Delta t^2 + O(\Delta t^3) \\ 1 & 1 \end{pmatrix} \\ \begin{pmatrix} 1 + f\Delta t + \frac{f^2\Delta t^2}{2} + \frac{f^3\Delta t^3}{6} + O(\Delta t^4) & \\ \frac{4f\Delta t}{3} - \frac{f^2\Delta t^2}{2} - \frac{f^3\Delta t^3}{6} + O(\Delta t^4) & \end{pmatrix} \\ \begin{pmatrix} -1 + \frac{71f\Delta t}{12} + \frac{107f^2\Delta t^2}{36} + O(\Delta t^3) & 7 - \frac{65fk}{12} - \frac{209f^2\Delta t^2}{36} + O(\Delta t^3) \\ 1 - \frac{71f\Delta t}{12} - \frac{107f^2\Delta t^2}{36} + O(\Delta t^3) & -1 + \frac{65fk}{12} + \frac{209f^2\Delta t^2}{36} + O(\Delta t^3) \end{pmatrix} \quad (64)$$

Recall that, neglecting the initial error E_0 , we can say that the global error is (16)

$$E_n = \Delta t \sum_{\nu=0}^{n-1} Q^{n-\nu-1} \boldsymbol{\tau}_\nu$$

Putting together equations (61) and (64) we see that each term $Q \boldsymbol{\tau}_\nu$ contributes to the error in two ways:

- The first contribution is due to the fact that $\boldsymbol{\tau}_\nu$ is almost co-linear with the second eigenvector ψ_2 . The order of this contribution is

$$|z_2| \|\psi_2 \boldsymbol{\tau}_\nu\| = O(\Delta t) \cdot O(\|\Delta t \boldsymbol{\tau}_\nu\|) = O(\Delta t^3)$$

where the term $|z_2|$ is the second eigenvalue which is of order $O(\Delta t)$.

- The second contribution to the error comes from the component of $\boldsymbol{\tau}_\nu$ that is a multiple γ_1 of the first eigenvector ψ_1 ,

$$|z_1| \|\gamma_1 \psi_1 \boldsymbol{\tau}_\nu\| = O(\Delta t) \cdot O(\|\boldsymbol{\tau}_\nu\|) = O(\Delta t^3)$$

the term γ_1 is of $O(\Delta t)$ because $\boldsymbol{\tau}_\nu$ lives mostly in the space of ψ_2 .

While each of the terms in $\Delta t Q \boldsymbol{\tau}_\nu$ has order $O(\Delta t^2) \cdot O(\|\boldsymbol{\tau}_\nu\|) = O(\Delta t^4)$, as the method is evolved forward, the errors accumulate over time, and sum of all contributions from all the times gives us a global error of order $O(\Delta t) \cdot O(\|\boldsymbol{\tau}_n\|) = O(\Delta t^3)$.

Example 1a: To demonstrate that this method indeed performs as designed we study its behavior on a nonlinear scalar equation of the form:

$$\begin{aligned} u_t &= -u^2 = f(u) \ , \quad t \geq 0 \\ u(t=0) &= 1 \ . \end{aligned} \tag{65}$$

We evolve the solution of this equation to time $T = 1$ using the scheme (60). The initial steps are computed exactly. The plots of the errors and the truncation errors are presented in Figure 1(a). Both errors are shown for the first component, $v_{n+1/2}$ (denoted v(1) in the legend) and the second component, v_n (denoted v(2) in the legend). Clearly, although the truncation error is only second order (denoted tr err v(1) and tr err v(2) in the legend), the global error is third order, as predicted by the theory.

Example 1b: It is important that the method will perform as designed on a nonlinear system as well. To demonstrate this, we solve the the van der Pol system

$$\begin{aligned} u_t^{(1)} &= u^{(2)} \\ u_t^{(2)} &= 0.1[1 - (u^{(1)})^2]u^{(2)} - u^{(1)} \end{aligned} \tag{66}$$

using the same scheme (60). As this is a system, it is important that both components are examined. Thus, the vector of the numerical solution has two components for the time level t_n , denoted by v(2), and two components for the time level $t_{n+\frac{1}{2}}$, denoted by v(1). In Figure 1(b) the convergence plot of the components of $u^{(1)}$ and $u^{(2)}$ are presented. Once again, we see that the convergence rate is indeed third order.

Remark 4 *It is important to note that not all Type 3 DIMSIM methods have the EIS property! The property that the local truncation error lives in the space spanned by the eigenvectors of A that correspond to the zero eigenvalues*

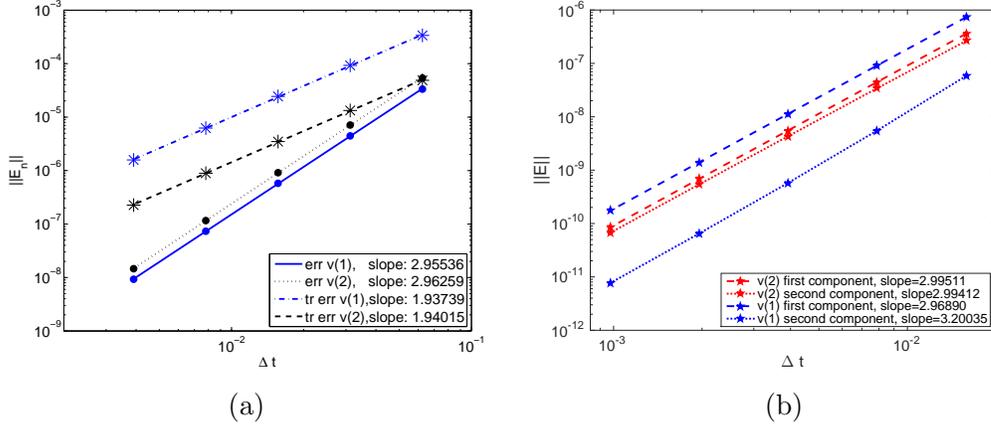


Fig. 1. Convergence plots using the scheme (60). (a) The errors and truncation errors vs. Δt , for several values of Δt , for the numerical solution of (65). (b) The errors vs. Δt for each component of the solution, computed for several values of Δt , for the numerical solution of the van der Pol equation (66).

is needed for the error inhibiting behavior to occur, and this property is not generally satisfied. To observe this, we study the DIMSIM scheme of types 3 presented by J. C. Butcher in [2].

Consider the scheme

$$\begin{pmatrix} v_{n+2} \\ v_{n+1} \end{pmatrix} = \frac{1}{4} \begin{pmatrix} 7 & -3 \\ 7 & -3 \end{pmatrix} \begin{pmatrix} v_{n+1} \\ v_n \end{pmatrix} + \frac{\Delta t}{8} \begin{pmatrix} 9 & -7 \\ -3 & -3 \end{pmatrix} \begin{pmatrix} f(v_{n+1}, t_{n+1}) \\ f(v_n, t_n) \end{pmatrix} \quad (67)$$

given in [2]. This scheme has truncation error

$$\tau_n = \frac{1}{48} \begin{pmatrix} 23 \\ 3 \end{pmatrix} \frac{d^3}{dt^3} u(t_n) \Delta t^2 + O(\Delta t^3). \quad (68)$$

The matrix A can be diagonalized as follows:

$$A = \frac{1}{4} \begin{pmatrix} 7 & -3 \\ 7 & -3 \end{pmatrix} = \begin{pmatrix} 1 & 3/7 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & \\ & 0 \end{pmatrix} \frac{1}{4} \begin{pmatrix} 7 & -3 \\ -7 & 7 \end{pmatrix}. \quad (69)$$

The truncation error τ_n can be written as a linear combination of the two eigenvectors of A as follows:

$$\tau_n = \left[\frac{19}{24} \begin{pmatrix} 1 \\ 1 \end{pmatrix} - \frac{35}{48} \begin{pmatrix} 3/7 \\ 1 \end{pmatrix} \right] \frac{d^3}{dt^3} u(t_n) \Delta t^2 + O(\Delta t^3). \quad (70)$$

Unlike the EIS scheme (60), here the first term in this expansion is of the order of $O(\tau_n) = O(\Delta t^2)$. Therefore a term of the order of $\Delta t O(\tau_n) = O(\Delta t^3)$ is accumulated at each time step, so that the global error is second order.

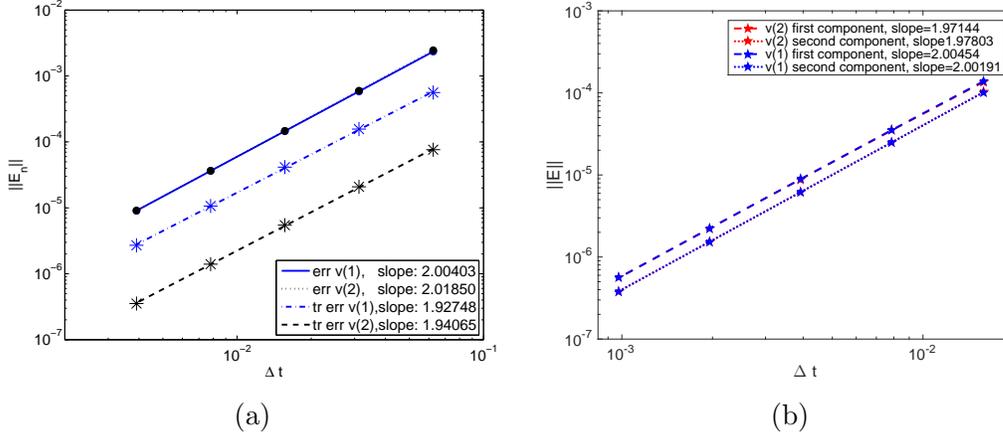


Fig. 2. Convergence plots using Butcher's scheme (67). (a) The errors and truncation errors vs. Δt , for several values of Δt , for the numerical solution of (65). Note that the errors for $v(1)$ and $v(2)$ are virtually identical so these error lines coincide. (b) The errors vs. Δt for each component of the solution, computed for several values of Δt , for the numerical solution of the van der Pol equation (66). Note that for this problem as well the behavior of this method on both components is virtually identical, so the error lines for each component of the solution coincide. Both the local truncation errors and the global errors are second order: this is not an error inhibiting scheme.

We note that both this method (67) and our error inhibiting method (60) satisfy the order conditions in Theorem 3.1 of [2] only up to second order ($p = 2$). However, as we see in Figure 2, when the method (67) is used to simulate the solution of the problems (65) and (66) we have second order accuracy, while the error inhibiting method (60) gave third order accuracy (Figure 1).

4.2 A fourth order error inhibiting method with $s = 3$.

In this subsection we present an error inhibiting method with $s = 3$ that takes the values of the solution at the times t_n , $t_{n+\frac{1}{3}}$, and $t_{n+\frac{2}{3}}$ and uses these three values to obtain the solution at the time-level t_{n+1} , $t_{n+\frac{4}{3}}$, and $t_{n+\frac{5}{3}}$. The exact solution vector is given by

$$U_n = \left(u(t_{n+2/3}), u(t_{n+1/3}), u(t_n) \right)^T,$$

and the corresponding vector of numerical approximations is

$$V_n = \left(v_{n+2/3}, v_{n+1/3}, v_n \right)^T.$$

Consider the error inhibiting scheme

$$\begin{aligned}
V_{n+1} = & \frac{1}{768} \begin{pmatrix} 467 & -1996 & 2297 \\ 467 & -1996 & 2297 \\ 467 & -1996 & 2297 \end{pmatrix} V_n + \\
& \frac{\Delta t}{1152} \begin{pmatrix} 5439 & -6046 & 3058 \\ 2399 & -1694 & 1362 \\ 703 & 354 & 626 \end{pmatrix} \begin{pmatrix} f(v_{n+2/3}, t_{n+2/3}) \\ f(v_{n+1/3}, t_{n+1/3}) \\ f(v_n, t_n) \end{pmatrix}, \quad (71)
\end{aligned}$$

which has a local truncation error of third order,

$$\begin{aligned}
\tau_n = & \frac{1}{373248} \begin{pmatrix} 43699 \\ 12787 \\ 2227 \end{pmatrix} \frac{d^4}{dt^4} u(t_n) \Delta t^3 + O(\Delta t^4) \\
\approx & \begin{pmatrix} 0.117078 \\ 0.0342587 \\ 0.00596654 \end{pmatrix} \frac{d^4}{dt^4} u(t_n) \Delta t^3 + O(\Delta t^4). \quad (72)
\end{aligned}$$

However, it can be verified that for the linear case, the product

$$Q_n \tau_n = O(\Delta t \tau_n) = O(\Delta t^4).$$

Given the analysis in Section 3.3 above, this result will carry over to the nonlinear case, and thus this method will have a fourth order global error, despite the third order truncation error.

To demonstrate this result we revisit the two examples (65) and (66) in the previous subsection and use the scheme (71) to evolve them forward in time. The results, shown in Figure 3, are exactly as we expect: although the truncation errors (seen for the problem (65) in Figure 3(a)) are only third order, the errors are fourth order for both problems (65) and the van der Pol problem (66).

4.2.1 Other fourth order error inhibiting methods with $s = 3$.

The methods above are not unique, in fact other methods can be derived using this approach. In this section we present two additional error inhibiting methods with $s = 3$ that have local truncation error that is third order but demonstrate fourth order global error on a nonlinear system.

The first method is

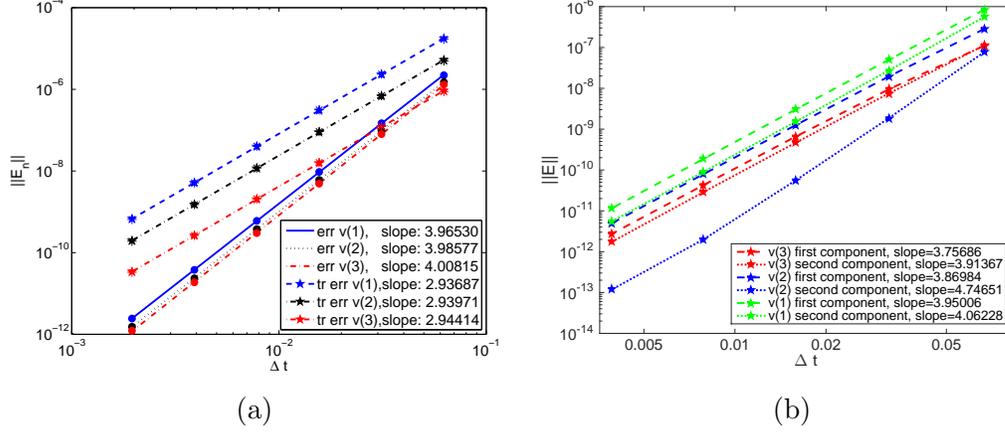


Fig. 3. Convergence plots using the scheme (71). (a) The errors and truncation errors vs. Δt , for several values of Δt , for the numerical solution of (65). (b) The errors vs. Δt for each component of the solution, computed for several values of Δt , for the numerical solution of the van der Pol equation (66). As expected, we observe fourth order accuracy for the errors, although the truncation errors are third order.

$$\begin{aligned}
 V_{n+1} = & \frac{1}{1020} \begin{pmatrix} 449 & -1966 & 2537 \\ 449 & -1966 & 2537 \\ 449 & -1966 & 2537 \end{pmatrix} V_n + \\
 & \frac{\Delta t}{6120} \begin{pmatrix} 29123 & -32576 & 15789 \\ 12973 & -9456 & 6779 \\ 3963 & 1424 & 2869 \end{pmatrix} \begin{pmatrix} f(v_{n+2/3}, t_{n+2/3}) \\ f(v_{n+1/3}, t_{n+1/3}) \\ f(v_n, t_n) \end{pmatrix}, \quad (73)
 \end{aligned}$$

and has a local truncation error of third order,

$$\begin{aligned}
 \tau_n = & \frac{1}{991440} \begin{pmatrix} 115733 \\ 33623 \\ 5573 \end{pmatrix} \frac{d^4}{dt^4} u(t_n) \Delta t^3 + O(\Delta t^4) \\
 \approx & \begin{pmatrix} 0.116732 \\ 0.0339133 \\ 0.00562112 \end{pmatrix} \frac{d^4}{dt^4} u(t_n) \Delta t^3 + O(\Delta t^4). \quad (74)
 \end{aligned}$$

The second method is

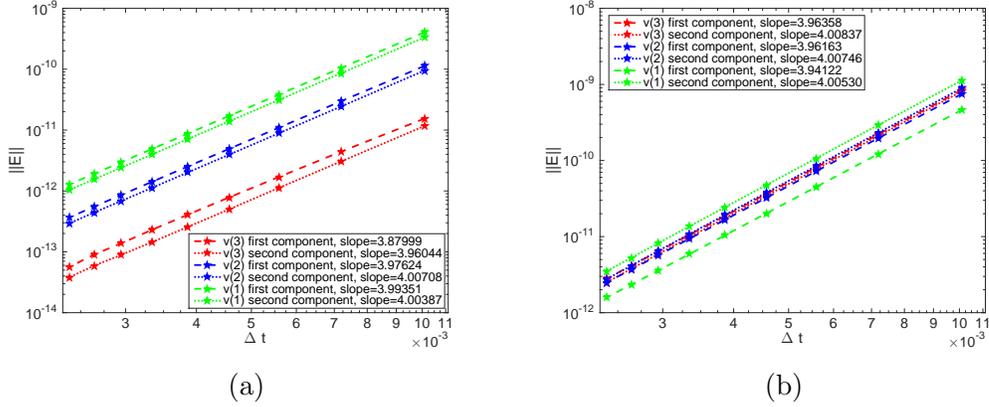


Fig. 4. Convergence plots van der Pol equation (66). The plots show the errors vs. Δt for each component of the solution, computed for several values of Δt for (a) the scheme (73) and (b) the scheme (75). As expected, we observe fourth order accuracy for the errors, although the truncation errors computed above are third order.

$$V_{n+1} = \begin{pmatrix} -\frac{101}{96} & \frac{97}{24} & -\frac{191}{96} \\ -\frac{101}{96} & \frac{97}{24} & -\frac{191}{96} \\ -\frac{101}{96} & \frac{97}{24} & -\frac{191}{96} \end{pmatrix} V_n + \Delta t \begin{pmatrix} \frac{733}{144} & -\frac{431}{72} & \frac{23}{12} \\ \frac{353}{144} & -\frac{53}{24} & \frac{4}{9} \\ \frac{47}{48} & -\frac{31}{72} & -\frac{7}{36} \end{pmatrix} \begin{pmatrix} f(v_{n+2/3}, t_{n+2/3}) \\ f(v_{n+1/3}, t_{n+1/3}) \\ f(v_n, t_n) \end{pmatrix}. \quad (75)$$

The truncation error is also third order

$$\begin{aligned} \tau_n &= \begin{pmatrix} \frac{5303}{46656} \\ \frac{1439}{46656} \\ \frac{119}{46656} \end{pmatrix} \frac{d^4}{dt^4} u(t_n) \Delta t^3 + O(\Delta t^4) \\ &= \begin{pmatrix} 0.113662 \\ 0.0308428 \\ 0.00255058 \end{pmatrix} \frac{d^4}{dt^4} u(t_n) \Delta t^3 + O(\Delta t^4) \end{aligned} \quad (76)$$

Both these methods satisfy

$$Q_n \tau_n = O(\Delta t \tau_n) = O(\Delta t^4)$$

as well. As above, this property results in an error inhibiting mechanism that produced a global error of order four. This can be seen once again in Figure 4, using the nonlinear problem (66) above. The results of method (73) are on the left and of (75) are on the right.

5 Conclusions

While it is generally assumed that the global error will be of the order of the local truncation error, in this work we presented an approach to creating methods that have a global error of higher order than predicted by the local truncation error. To accomplish this, we used the block formulation of a method $V_{n+1} = Q_n V_n$ where the discrete solution operator $Q_n = A + \Delta t B F_n$ is comprised of matrices of coefficients A and B , and the matrix operator F_n . We show that if A is a diagonalizable matrix of rank one, that has only one nonzero eigenvalue, $z_1 = 1$, that corresponds to the eigenvector of all ones, then the error inhibiting property will occur if the leading part of the local truncation error error for the linear constant coefficient case ($F_n = F = a$ constant) is spanned by the eigenvectors corresponding to the zero eigenvalues of A (to the leading order). We show that a method that has these properties will have a global error that has higher order than the local error, on nonlinear problems.

After presenting the concept behind these methods we use the theoretical properties above to develop block one-step methods that are in the family of Type 3 DIMSIM methods presented in [2]. We demonstrate in numerical examples on nonlinear problems (including a nonlinear system) that these methods have global error that is one order higher than the local truncation errors. We also show that this is in contrast to another Type 3 DIMSIM method which has a matrix A that satisfies the first three properties **C1** – **C3**, but does not satisfy the error inhibiting property **C4**, that the local truncation error is in the space spanned by the eigenvectors of A that correspond to the zero eigenvalues, and indeed does not give us a global error that is higher than the local truncation error on nonlinear test problems.

The major development in this work is the concept of an error inhibiting method and the new approach for developing methods that are constructed to control the growth of the local truncation error. While the newly developed methods presented in this work can be used in place of currently standard methods (particularly in place of type 3 DIMSIM methods) to obtain higher order accuracy, it is not yet known how they compare to other methods in terms of other important properties. In future work we intend to the study of the computational efficiency and storage requirements of these methods and the analysis of their linear stability regions. We expect that this will also lead to further development of error inhibiting methods that have other favorable properties.

Acknowledgements: *The authors wish to thank Professor John Butcher for a very helpful discussion, and in particular for his valuable advice on general linear methods, especially the Type 3 DIMSIM methods.*

The work of Sigal Gottlieb was supported by AFOSR grant FA9550-15-1-0235.

References

- [1] Myron B. Allen and Eli L. Isaacson, *Numerical analysis for applied science*, John Wiley & Sons, 1998.
- [2] John C. Butcher, *Diagonally-implicit multi-stage integration method*, Applied Numerical Mathematics **11** (1993), 347–363.
- [3] John C Butcher, *Numerical methods for ordinary differential equations*, John Wiley & Sons, 2008.
- [4] Robert PK Chan and Angela YJ Tsai, *On explicit two-derivative Runge–Kutta methods*, Numerical Algorithms **53** (2010), no. 2-3, 171–194.
- [5] A Ditkowski, *High order finite difference schemes for the heat equation whose convergence rates are higher than their truncation errors*, Spectral and High Order Methods for Partial Differential Equations ICOSAHOM 2014, Springer, 2015, pp. 167–178.
- [6] Bertil Gustafsson, Heinz-Otto Kreiss, and Joseph Oliger, *Time dependent problems and difference methods*, vol. 24, John Wiley & Sons, 1995.
- [7] Ernst Hairer, SP Nørsett, and Gerhard Wanner, *Solving ordinary, differential equations i, nonstiff problems*, 2Ed. Springer-Verlag, 2000, 2000.
- [8] Eugene Isaacson and Herbert Bishop Keller, *Analysis of numerical methods*, Dover Publications, Inc, 1994.
- [9] Zdzislaw Jackiewicz, *General linear methods for ordinary differential equations*, John Wiley & Sons, 2009.
- [10] K Kastlunger and Gerhard Wanner, *On turan type implicit Runge–Kutta methods*, Computing **9** (1972), no. 4, 317–325.
- [11] KH Kastlunger and Gerhard Wanner, *Runge–Kutta processes with multiple nodes*, Computing **9** (1972), no. 1, 9–24.
- [12] Peter D Lax and Robert D Richtmyer, *Survey of the stability of linear finite difference equations*, Communications on pure and applied mathematics **9** (1956), no. 2, 267–293.
- [13] Alfio Quarteroni, Riccardo Sacco, and Fausto Saleri, *Numerical mathematics*, vol. 37, Springer Science & Business Media, 2010.
- [14] J Barkley Rosser, *A Runge-Kutta for all seasons*, SIAM Review **9** (1967), 41717452.
- [15] Granville Sewell, *The numerical solution of ordinary and partial differential equations*, World Scientific, 2015.
- [16] Larry F Shampine and H A Watts, *Block implicit one-step methods*, Math. Comp **23** (1969), 73117740.

- [17] Hisayoshi Shintani et al., *On one-step methods utilizing the second derivative*, Hiroshima Mathematical Journal **1** (1971), no. 2, 349–372.
- [18] ———, *On explicit one-step methods utilizing the second derivative*, Hiroshima Mathematical Journal **2** (1972), no. 2, 353–368.