

## Backward Differentiation Approximations of Nonlinear Differential/Algebraic Systems

By Kathryn E. Brenan and Bjorn E. Engquist\*

**Abstract.** Finite difference approximations of dynamical systems modelled by nonlinear, semiexplicit, differential/algebraic equations are analyzed. Convergence for the backward differentiation method is proved for index two and index three problems when the numerical initial values obey certain constraints. The appropriate asymptotic convergence rates and the leading error terms are determined.

**1. Introduction.** The most general systems of differential/algebraic equations (DAE's) arise in the *fully implicit* form,

$$(1.1) \quad R(t, x, x') = 0,$$

where  $x = x(t)$  and  $R$  are vectors of dimension  $r$ . The Jacobian matrix  $\partial R/\partial x'$  is assumed to be singular. If it is nonsingular, system (1.1) is an implicit set of ordinary differential equations (ODE's) and can be, at least theoretically, reformulated as  $x' = f(t, x)$ .

The  $k$ -step backward differentiation formula (BDF) was introduced by Gear [8] in 1971 for the numerical integration of DAE systems. This approach has been universally accepted in industry as well as in academia as the standard numerical method for solving DAE systems. To apply this method to (1.1), replace the derivatives  $x'$  by their BDF approximation with constant stepsize  $h$  and evaluate the equations at  $t_n$ , where  $t_n = t_0 + nh$ . The resulting system of difference equations,

$$(1.2) \quad R \left( t_n, x_n, h^{-1} \sum_{i=0}^k \alpha_i x_{n-i} \right) = 0$$

is then solved for the numerical approximation  $x_n$  of  $x(t_n)$ .

The analytical and numerical theory for linear  $R$  with constant coefficients is well understood [5], [7], [25], [27]. The variable coefficient and nonlinear cases are much harder. In contrast to the theory for the numerical approximation of ODE's, the convergence results for linear, constant coefficient systems do not extend to all linear, variable coefficient or nonlinear problems. In this paper we will restrict our interest to the class of *solvable* DAE's. A solvable DAE is one for which solutions exist, and for which solutions having the same initial values are identical [9]. Now, while the  $k$ -step BDF have been shown to converge on solvable, linear constant

---

Received April 11, 1986; revised November 9, 1987.

1980 *Mathematics Subject Classification* (1985 Revision). Primary 65L05.

\*This work was supported in part by NSF Grant No. DMS85-03294, ARO Grant No. DAAG29-85-K-0190, DARPA Grant in the ACMP Program, ONR Grant N00014-86-K-0691, and NASA Langley Grant NAG1-270.

coefficient systems [25], they may be unstable for some fully implicit DAE's. For example, in [10] a fully implicit, linear, time-varying solvable DAE system is presented for which the one-step BDF is unstable and thus does not converge. On the other hand, numerical experiments strongly suggest that the BDF are convergent for some classes of DAE's [1]. This example demonstrates that in order to establish a convergence theory for numerical approximations to the solutions of nonlinear (including linear time-varying) DAE's, a subclass of problems of the form (1.1) has to be chosen. It is the purpose of this paper to define an appropriate subclass of nonlinear DAE's which cover many of the applications, and to develop a rigorous convergence theory for that class. The stability results derived in [1] for the linear homogeneous systems corresponding to this subclass of nonlinear systems will be summarized as well. Before the results can be presented, some background material must be given.

The behavior of DAE systems is directly related to a property called the *index* [5], also sometimes referred to as the *degree of nilpotency* [28]. For example, consider the following subclass of system (1.1) which we will refer to as *semiexplicit* DAE's:

$$(1.3) \quad y' = E(t, y, u),$$

$$(1.4) \quad 0 = H(t, y, u),$$

where  $y = y(t)$  and  $E$  are  $l$  vectors, and  $u = u(t)$  and  $H$  are  $m$  vectors. Semiexplicit DAE systems arise in dynamic simulations of mechanical problems [17], fluid flow problems [18], [21], and optimal [4] and trajectory control [2] problems. Essentially, the index of system (1.3), (1.4) is one more than the number of times it is necessary to differentiate the algebraic equations before the algebraic variables  $u$  can be explicitly determined. Hence, if the Jacobian matrix  $\partial H/\partial u$  is nonsingular, then the system has index one [11], since no differentiation is required. If  $\partial H/\partial u$  is singular, then the index is at least two. If the algebraic subsystem (1.4) is not present, or if the DAE system is simply a set of implicit ODE's, the index is zero. A reduction algorithm described in [9] gives a precise definition of the index.

Recall that initial values of ODE's may be specified arbitrarily. However, initial values for a DAE system of index  $\nu$  must in general satisfy consistency relations which may involve up to  $\nu - 1$  derivatives of some of the variables. For example, initial values for semiexplicit DAE systems must at least satisfy the algebraic equations (1.4).

Not only is the index of the system important in the characterization of the system's solutions, but it is also critical to the convergence and stability properties of numerical approximations. For example, the relation of the index to the stability and convergence properties of the BDF is very well understood for linear, constant coefficient systems,

$$(1.5) \quad Ax' + Bx = f(t),$$

where  $A$  and  $B$  are constant  $r \times r$  matrices. If system (1.5) with index  $\nu$  is solved by the  $k$ -step BDF ( $k \leq 6$ ) with constant stepsize  $h$ , the numerical solution converges globally to  $O(h^k)$  accuracy after  $(\nu - 1)k + 1$  steps [25]. If the one-step BDF (i.e., Backward Euler) is used to solve an index three system, the global error is  $O(1)$

after one step [11]. However, after three steps of constant length  $h$ , the numerical solution is  $O(h)$  accurate.

If the system's index is restricted to one, then the numerical solution determined by the  $k$ -step BDF converges to the solution of the general implicitly formulated problem (1.1) [9]. However, as was mentioned earlier, subclasses of (1.1) must be considered when proving convergence for higher index systems. Since numerical solutions of some higher index, semiexplicit systems have been experimentally determined to  $O(h^k)$  accuracy by the  $k$ -step BDF, it is natural to investigate the convergence properties of the BDF on this subclass of (1.1) [1], [2].

The subclasses of DAE's studied in this paper are higher index problems (i.e., index greater than one) which are sometimes labelled *ill-posed* [19], [20] or *algebraically incomplete* [23]. The classification of these higher index systems as ill-posed should be understood in the strict sense that they may not exhibit a continuous dependence in the maximum norm on the inhomogeneous terms. In fact, it will be necessary to impose additional smoothness on the systems and to consider a weaker form of stability where the solution depends continuously on the inhomogeneous terms and some of their derivatives. Solvability of these systems also requires this regularity. It is still possible for the BDF to produce a convergent numerical solution in spite of this weaker form of stability. Furthermore, it is important in practice to understand how numerical methods behave when applied to higher index systems because higher index problems arise in applications much more frequently than first thought. Until recently, many engineers solving DAE's were not cognizant of the index. Except in special cases, it is still difficult in practice to determine what the index of a given DAE is. In any case, it is useful to identify these special cases. In particular, the nonlinear systems we study here were originally motivated by the need to solve trajectory prescribed path control problems of current interest to the aerospace industry. Typically, these problems are index two or three, semiexplicit, nonlinear DAE's. While sometimes their index can be reduced (by differentiation or some other technique), in general that approach has not been very satisfactory in practice.

In this paper, convergence of the BDF is proven for semiexplicit systems of index two and three. In particular, consider the index two system

$$(1.6) \quad y' = E(t, y, u),$$

$$(1.7) \quad 0 = H(t, y),$$

where the  $m \times m$  matrix  $(\partial H/\partial y)(\partial E/\partial u)$ , evaluated at the solution  $(y(t), u(t))$ , is assumed to be nonsingular for all  $t$  in some interval  $I$ , and  $m \leq l$ . Let us also consider the index three system,

$$(1.8) \quad v' = F(t, v, w, u),$$

$$(1.9) \quad w' = G(t, v, w),$$

$$(1.10) \quad 0 = H(t, w),$$

where the  $m \times m$  matrix  $(\partial H/\partial w)(\partial G/\partial v)(\partial F/\partial u)$  is nonsingular along the solution  $(v(t), w(t), u(t))$  for all  $t \in I$ . In system (1.8) (1.10),  $v$  and  $F$  are vectors

of dimension  $p$ , denoted as  $\dim v = \dim F = p$ , while  $\dim w = \dim G = q$  and  $\dim u = \dim H = m$ . To avoid the case of an overdetermined system, we assume  $m \leq \min(p, q)$ . Throughout this paper, the functions  $E, F, G$ , and  $H$  are assumed to be sufficiently smooth functions of all their arguments as required for the convergence analysis contained herein.

Typically, the variables for which there are explicit differential equations are referred to as *state* variables, while those variables appearing only algebraically will be called the *algebraic* or control variables. In system (1.6), (1.7) the  $l$  state variables are  $y$ , while in system (1.8)–(1.10), the state variables are  $y = (v, w)$  with  $\dim y = l = p + q$ . In both systems, the algebraic variables are  $u$ . In [1] it was proven that a linear, index two system and a linear, index three system corresponding to (1.6), (1.7) and (1.8)–(1.10), respectively, are solvable. The proofs are straightforward and involve reducing the systems to sets of explicit ODE's by differentiating the algebraic equations and substituting for the state variables' derivatives from the DAE's. For *consistent* initial values, the solvability of the DAE's is then established by applying existence and uniqueness theorems for ODE's. Note that a consistent set of initial values must satisfy not only the algebraic equations given in the DAE, but also those equations arising in the reduction process. For the nonlinear DAE's (1.6), (1.7) and (1.8)–(1.10), it is technically necessary to assume there exist solutions satisfying the algebraic equations corresponding to the associated index one problems, including (1.7) and (1.10), respectively. Then, solvability of the nonlinear DAE's may be established as in the linear case—namely, differentiate the algebraic equations and apply the implicit function theorem to solve for the algebraic variables, thereby reducing the DAE's to explicit ODE systems. This last step utilizes the assumptions on the nonsingularity of the matrix products  $(\frac{\partial H}{\partial y})(\frac{\partial E}{\partial u})$  and  $(\frac{\partial H}{\partial w})(\frac{\partial G}{\partial v})(\frac{\partial F}{\partial u})$ .

The facts that system (1.6), (1.7) and the linear system corresponding to (1.8)–(1.10) are index two and three, respectively, was first established in [1] by applying the reduction technique described in [9]. Essentially, the index is one more than the number of times it is necessary to differentiate the algebraic constraints before the algebraic variables appear explicitly. The index of the nonlinear system (1.8)–(1.10) follows in a similar way.

In the following section, convergence theorems for the BDF methods approximating the solutions to (1.6), (1.7) and (1.8)–(1.10) are stated. The rate of convergence is equivalent to the corresponding rate for linear, constant coefficient problems. The leading error term in the asymptotic expansion is given. This section also contains necessary definitions and a brief outline of the proofs. In particular, it is pointed out how this analysis differs from the corresponding theory for ODE's. The convergence proof for the index two system is given in Section 3, while the corresponding proof for the index three system is presented in the Supplement.

This paper is an extension of the convergence results derived in [1], and is a condensed version of the earlier report [3] in which the proofs are carried out in somewhat greater detail. While a convergence analysis of the  $k$ -step BDF was done for the nonlinear, index two system (1.6), (1.7) in [1], the proof presented in [1] was valid for only the *linear*, time-varying index three system corresponding to the nonlinear system (1.8)–(1.10) of interest here. The analysis in this paper

also leads to new results concerning the leading error terms. Extensive numerical experiments have been conducted to verify the rates of convergence for both linear and nonlinear index two and three systems of these forms [1]. Numerical results for a particularly nonlinear, index three problem arising in trajectory control have already been presented in [2].

Using a different approach than the one employed here, convergence results for the BDF have been obtained in [18] for index two systems having a form equivalent to (1.6), (1.7) and for a subclass of the index three systems (1.8)–(1.10). The proofs presented there rely on local rates of convergence established in [1]. Throughout this paper, we will point out the key differences between the results given in [1], [18], and this paper.

**2. Notation, Definition, and Statement of Results.** Before stating the convergence results obtained for systems (1.6), (1.7) and (1.8)–(1.10), it is necessary to present some definitions. A set of initial values  $(y_0, u_0)$  at  $t_0$  is said to be *consistent* for a DAE system if there exists at least one solution  $(y(t), u(t))$  assuming those values at  $t_0$ . For  $k \geq 2$  the  $k$ -step BDF requires not only initial values, but the following set of values here called *starting values*:

$$(2.1) \quad \mathbf{y}_{k-1} = [y_{k-1}^T, y_{k-2}^T, \dots, y_0^T]^T,$$

$$(2.2) \quad \mathbf{u}_{k-1} = [u_{k-1}^T, u_{k-2}^T, \dots, u_0^T]^T,$$

where the superscript  $T$  denotes the transpose. The  $l$  vector  $y_j$  and the  $m$  vector  $u_j$  denote the numerical solution at  $t_j = t_0 + jh$  for  $j = 0, 1, \dots, k - 1$  and  $t_j \in J = [t_0, t_0 + T] \subseteq I$ . We will say these starting values are *numerically consistent to order  $k + 1$*  if there exists a solution to the index  $\nu$  system such that

$$(2.3) \quad \|y_j - y(t_j)\| \leq K_1 h^{k+1},$$

$$(2.4) \quad \|H(t_j, y_j)\| \leq K_2 h^{k+\nu-1}$$

for some constants  $K_1, K_2$ , and  $j = 0, 1, \dots, k - 1$ . In general, we could use any vector norm and corresponding consistent matrix norm during the analysis, but it is convenient to restrict the analysis to the maximum norm. We shall show that the starting values for  $u$  are not critical in the convergence analysis for systems (1.6), (1.7) and (1.8)–(1.10).

Let the BDF difference approximation (1.2) of the semiexplicit problems (1.6), (1.7) and (1.8)–(1.10) have the form,

$$(2.5) \quad \sum_{i=0}^k \alpha_i y_{n-i} = hE(t_n, y_n, u_n),$$

$$(2.6) \quad 0 = H(t_n, y_n),$$

where  $E = (F, G)$  and  $y = (v, w)$  for the index three system (1.8) (1.10). Given starting values satisfying (2.3), (2.4), a numerical solution  $(y_n, u_n)$  of (2.5), (2.6) *converges globally with  $k$ th order accuracy* to a solution of the DAE system (1.6), (1.7) or (1.8)–(1.10) if

$$\|y_n - y(t_n)\| \leq K_3 h^k, \quad \|u_n - u(t_n)\| \leq K_4 h^k$$

for  $t_n \in [t_\theta, t_0 + T] \subseteq J$  where the constants  $K_3$ ,  $K_4$ , and  $\theta$  are independent of  $h$ .

Since instability of the BDF is a concern for fully implicit, higher index DAE's, we summarize some stability results derived in [1] for the linear, homogeneous systems corresponding to the nonlinear systems (1.6), (1.7) and (1.8)–(1.10) of interest here. For *arbitrary* starting values  $\mathbf{y}_{k-1}$ , the  $k$ -step BDF is stable (or *weakly stable* [6]) for these semiexplicit systems because there exist constants  $\mathbf{K}^*$ ,  $\mathbf{K}^{**}$ , and  $\mathcal{E}$  such that the following conditions hold uniformly in  $n$  and  $h$  for  $t_{\bar{n}} \leq t_n \leq t_0 + T$ ,

$$\|\mathbf{y}_n\| \leq \mathbf{K}^* e^{nh\mathcal{E}} \|\mathbf{y}_{k-1}\|, \quad \|\mathbf{u}_n\| \leq \mathbf{K}^{**} \|\mathbf{y}_{k-1}\|$$

where  $\bar{n} = 3k - 1$  (i.e., after  $k + 1$  steps) for index two systems and  $\bar{n} = 4k - 1$  (i.e., after  $2k + 1$  steps) for index three systems. In fact, the state variables  $\mathbf{y}$  are always computed stably in an index two system, but only after  $k + 1$  steps in an index three problem. Moreover, the BDF produces a numerical solution for the algebraic variables having a *boundary layer* of instability of length  $(k + 1)$  steps for index two systems and of length  $(2k + 1)$  steps for index three systems. These weak stability properties have since been studied for more general linear, semiexplicit, solvable, index two systems in [6]. In practice, weak stability describes how the BDF are sensitive to roundoff errors and errors in the starting values.

While it is true that the BDF methods are weakly stable for these semiexplicit systems, it is also clear that the initial steps may contain 'large' errors. In spite of that, we show in this paper that if the starting data is sufficiently accurate, the BDF does produce a numerical solution which converges to a solution of the *nonlinear* system with the expected  $O(h^k)$  rate of convergence. It is precisely one point of this paper to determine how accurate the initial data must be in order to insure convergence (in particular, for index three problems). Not only is this a theoretical concern, but it is also of practical interest as evidenced by the application of the BDF methods to real problems (see numerical results given in [1] or [2]). It is a fact that the Backward Euler method fails to converge at the end of the first integration step when applied to an index three problem. Consequently, serious difficulties arise in practice when using state of the art software implementing the BDF methods (e.g., see [22]) on index three systems. It has even been observed that it is possible for the numerical solution to converge to a solution of the given DAE which is inconsistent with the given initial values for the algebraic variables [2]. This difficulty is due not only to the behavior of the Backward Euler method, but also to the nature of the nonlinear system being solved. In this particular trajectory problem, for a given set of initial state values, there is a nonunique solution for the algebraic variable. A one-step method has no memory, so the  $O(1)$  error introduced into the algebraic variable during the first step may cause the numerical solution to jump to a different solution curve. If the initial values for the state variables are sufficiently accurate, this difficulty does not occur.

Now we can state the convergence results obtained for the  $k$ -step BDF when applied to the semiexplicit DAE systems (1.6), (1.7) or (1.8)–(1.10).

**THEOREM 1.** *There exists a numerical solution of the index two system (1.6), (1.7) by the  $k$ -step BDF with constant stepsize  $h$  for  $k < 7$  which converges globally with  $k$ th order accuracy to a solution of (1.6), (1.7) if the starting values are numerically consistent to order  $k + 1$ .*

**THEOREM 2.** *There exists a numerical solution of the index three system (1.8)–(1.10) by the  $k$ -step BDF with constant stepsize  $h$  for  $k < 7$  which converges globally with  $k$ th order accuracy to a solution of (1.8)–(1.10) after  $k + 1$  steps if the starting values are numerically consistent to order  $k + 1$ .*

The definition of a numerically consistent set of starting values is natural in the following way. For a system of explicit ODE's, a unique solution exists corresponding to any set of *arbitrary* initial values. The numerical solution corresponding to a convergent difference approximation of order  $k$  applied to an ODE will generally require starting values of accuracy  $O(h^k)$ . For a solvable DAE, solutions exist and are uniquely specified when a *consistent set* of initial values are given. It is therefore natural to require that the starting values for the numerical method also satisfy some consistency conditions as well. That is, we cannot expect convergence of the difference method for any arbitrary set of starting values located in a circle of radius  $O(h^k)$  about a consistent set of initial values. Specifically, they must be chosen to be sufficiently close to the manifold containing the solution to the DAE.

In [1] and [18] it appears that convergence of the BDF has been proven for these systems when the starting values are accurate only to  $O(h^k)$ . However, if these starting values are not *numerically consistent*, then  $O(h^k)$  convergence is not obtained until an additional  $k$  steps are taken, where the algebraic equations are satisfied either exactly as in [1] or at least to some sufficiently small tolerance [18] (namely, to  $O(h^{k+1})$  accuracy for index two systems and to  $O(h^{k+2})$  for index three problems). Hence, in [1] or [18] it is proven that the  $k$ -step BDF, given arbitrary starting values of accuracy  $O(h^k)$ , converges with  $O(h^k)$  accuracy to the analytic solution after  $k + 1$  steps for the index two system and after  $2k + 1$  steps for the (linear) index three system. During the initial steps, reduced rates of convergence, namely  $O(h^{k-1})$ , may be observed in the algebraic variables [1]. If the starting values are *numerically consistent*, then  $O(h^k)$  convergence is achieved immediately for the index two system and after  $k + 1$  steps for the index three system.

In [18] it is assumed the algebraic equations are solved at each step to  $O(h^{k+1})$  accuracy for the index two system and to  $O(h^{k+2})$  accuracy for the index three system. These requirements are equivalent in our analysis to the restriction (2.4) for the starting values. In our proofs, we will assume that the algebraic equations are satisfied exactly at each step, but we could instead relax this requirement by enforcing (2.4) on each step. Specifically, the  $k$ -step BDF will converge with  $O(h^k)$  accuracy if the algebraic equations are satisfied to  $O(h^{k+1})$  accuracy for index two problems and to  $O(h^{k+2})$  accuracy ( $k = 1$  requires  $O(h^{k+3})$ ) for index three systems. If the starting values do not satisfy this requirement (i.e., are not numerically consistent), then an additional  $k$  steps will be required before  $O(h^k)$  convergence is attained in all variables.

*Remark.* The convergence results given in Theorems 1 and 2 are also valid, without lengthening the boundary layer, when the consistency relations for the starting values are relaxed somewhat:

1. For the index two system, (2.4) follows immediately from (2.3), and hence is not a restriction.

2. In the convergence analysis for the index three system, we could relax (2.3) to

$$\|w_j - w(t_j)\| = O(h^{k+1}) \quad \text{and} \quad \|v_j - v(t_j)\| = O(h^k).$$

3. It is possible to prove Theorems 1 and 2 given starting values with error  $O(h^k)$  providing they have the following special form: ( $i = 0, 1, \dots, k - 1$ )

$$\begin{aligned} y_i &= y(t_i) + h^k c_{1,i} + O(h^{k+1}) \quad (\text{index two}), \\ w_i &= w(t_i) + h^k c_{1,i} + O(h^{k+1}) \quad (\text{index three}), \end{aligned}$$

where  $c_{1,i}$  is a bounded vector in the null space of the Jacobian matrix  $H_y$  or  $H_w$ , respectively.

Conditions 2 and 3 are less restrictive than the assumed consistency conditions (2.3), (2.4), but since the convergence proof follows in the same way as given here (except for straightforward technical complications arising in the initial step of the proof), we omit the details.

Before rigorously proving the theorems in Section 3 and the Supplement section, we outline some of the key elements of the analysis, common to the proofs for both the index two and index three systems. We also introduce further notation and state two corollaries concerning the form of the principal error terms.

In the convergence analysis of linear multistep methods applied to explicit ODE's, it is common to derive the following one-step evolution equation for the numerical error  $h^k v_n$  [24]:

$$\omega_n = S_n \omega_{n-1} + h \tilde{f}_n$$

for  $\omega_n = (v_n, v_{n-1}, \dots, v_{n-k+1})$  and  $x_n = x(t_n) + h^k v_n$ . The amplification matrix  $S_n$  is bounded. The vector  $\tilde{f}_n$  contains the local truncation error of the BDF method and the remaining nonlinear terms of  $v_n$  after linearization.

We shall also use a one-step evolution equation for the numerical error, but when such an equation is derived for the BDF methods applied to DAE's, neither the amplification matrix nor the vector  $\tilde{f}_n$  is bounded independent of  $h$ . However, it is still possible to prove convergence for semiexplicit DAE's by utilizing the natural structure of these systems. In particular, it is necessary to bound certain matrix products such as  $S_n S_{n-1} \cdots S_1$  for  $n$  sufficiently large and establish cancellation between different terms in the expression for the evolution of the error. It is of particular importance in the index three case for which even smoothness of the leading error term is needed to cancel terms from different time steps. We shall therefore express the error as a truncated asymptotic error expansion

$$(2.7) \quad y_n - y(t_n) = h^k d(t_n) + h^{k+1} r_n,$$

$$(2.8) \quad u_n - u(t_n) = h^k e(t_n) + h^{k+1} s_n,$$

where for the index three system (1.8) (1.10) we let  $d(t) = (c(t), b(t))$  and  $r_n = (q_n, p_n)$ . In the proofs of both theorems, these expansions simplify the analysis (compare [26]). A similar approach is utilized by Henrici in [15] to determine the asymptotic behavior of the discretization error for general linear multistep

methods applied to explicit ODE's. When expressions (2.7), (2.8) are substituted into Eqs. (2.5), (2.6) and expanded by Taylor series around the analytic solution  $(y(t_n), u(t_n))$ , we obtain the following linear, index two DAE system for the leading error terms of the numerical solution to the index two problem (1.6), (1.7):

$$(2.9) \quad d'(t) = A(t) d(t) + B(t)e(t) + \frac{1}{(k+1)}y^{(k+1)}(t),$$

$$(2.10) \quad 0 = C(t) d(t),$$

where  $A(t) = E_y(\ast)$ ,  $B(t) = E_u(\ast)$ ,  $C(t) = H_y(\ast)$ , and the  $\ast$  denotes evaluation at the analytic solution  $(y(t), u(t))$ . Similarly, for the index three problem (1.8)–(1.10), we get the following linear, index three DAE system for the leading error terms:

$$(2.11) \quad c'(t) = A_{11}(t)c(t) + A_{12}(t)b(t) + A_{13}(t)e(t) + \frac{1}{(k+1)}v^{(k+1)}(t),$$

$$(2.12) \quad b'(t) = A_{21}(t)c(t) + A_{22}(t)b(t) + \frac{1}{(k+1)}w^{(k+1)}(t),$$

$$(2.13) \quad 0 = A_{32}(t)b(t),$$

where  $A_{11}(t) = F_v(\ast)$ ,  $A_{12}(t) = F_w(\ast)$ ,  $A_{13}(t) = F_u(\ast)$ ,  $A_{21}(t) = G_v(\ast)$ ,  $A_{22}(t) = G_w(\ast)$ , and  $A_{32}(t) = H_w(\ast)$ . There exists a unique solution to the systems (2.9), (2.10) and (2.11)–(2.13) above for each set of consistent initial values [1]. The functions  $(d(t), e(t))$  are smooth for smooth  $E$  and  $H$ . Note that the asymptotic error expansion is not only useful for the convergence proof. Such expansions are the basis for automatic error controls, initialization and extrapolation techniques [14].

To prove that the numerical solution converges to the true solution as  $h \rightarrow 0$ , we will show that  $(r_n, hs_n)$ , and for technical reasons sometimes that  $(r_n, s_n)$ , are uniformly bounded for all  $n \geq \theta$  where  $\theta$  is independent of  $n$  and  $h$ . After substituting expressions (2.7), (2.8) into the BDF difference equations (2.5), (2.6) and using the fact that  $(d(t), e(t))$  is the solution to a linear DAE system, we write the remaining difference equations for  $(r_n, s_n)$  in the one-step form. An induction argument for  $i = k, \dots, n - 1$  to  $i = n$  on the remainders  $(r_i, s_i)$  is applied. Since the relations derived during the analysis are implicit functions of the remainders, it is also necessary to construct a fixed-point iteration for  $(r_n, hs_n)$  or  $(r_n, s_n)$ . The implicit function theorem is applied. The particular restrictions on the initial conditions (i.e., (2.3), (2.4)) are essential in this context.

Throughout the analysis, certain special properties of the matrices arising must be utilized. For example, factors of order  $O(1/h)$  are annihilated by projection matrices which appear naturally. Frequently, it is necessary to bound the product of time-dependent matrices having a distinctive block structure. For convenience, we now introduce a notation used to represent these block companion matrices.

*Notation.* Let  $I_s$  be an  $s \times s$  identity matrix,  $0_s$  be an  $s \times s$  zero matrix, and  $Y$  be any  $s \times s$  matrix. Let  $\eta_\mu$ ,  $\mu = 1, 2, \dots, k$ , be scalars. Then we define the  $sk \times sk$  block companion matrix

$$\text{CPM}(\eta_\mu Y) = \begin{pmatrix} \eta_1 Y & \eta_2 Y & \dots & \eta_{k-1} Y & \eta_k Y \\ I_s & 0_s & \dots & 0_s & 0_s \\ 0_s & I_s & \cdot & 0_s & 0_s \\ \vdots & \vdots & \cdot & \vdots & \vdots \\ 0_s & 0_s & \dots & 0_s & 0_s \\ 0_s & 0_s & \dots & I_s & 0_s \end{pmatrix}.$$

We shall let  $\prod_{i=1}^n Y_i$  denote the product of time-dependent matrices  $Y_i = Y(t_i)$  given in descending order such as

$$\prod_{i=1}^n Y_i = Y_n Y_{n-1} \cdots Y_1.$$

To prove convergence, it is necessary only to show that  $(hr_n, hs_n)$  is uniformly bounded for all  $n \geq \theta$ . However, we will in fact prove that  $(r_n, hs_n)$  is uniformly bounded. The convergence proof directly establishes the fact that the principal leading error term in  $y_n$  is smooth, but shows only that the principal leading error term in  $u_n$  is  $h^k(e(t_n) + hs_n)$ . The smoothness of the leading error term in  $u_n$  is established in the following corollaries, the proofs of which are given in Section 3 and the Supplement section, respectively:

**COROLLARY 1.** *If the conditions in Theorem 1 are valid and the algebraic equations are satisfied to  $O(h^{k+2})$  accuracy for the starting values  $y_{k-1}$  and at the end of each integration step, then  $s_n$  is uniformly bounded for all  $n \geq k$  and the principal leading error term in  $u_n$  is  $h^k e(t_n)$ .*

**COROLLARY 2.** *If the conditions in Theorem 2 are valid, then  $s_n$  is uniformly bounded for all  $n \geq 3k$  and the principal leading error term in  $u_n$  is  $h^k e(t_n)$ .*

For the state variables, the leading error terms are smooth immediately ( $n \geq k$ ) if the starting values are numerically consistent and if the algebraic equations are satisfied as in (2.4) at the end of each integration step. However, to derive the leading error terms corresponding to the algebraic variables, it is necessary in general to require the algebraic equations to be satisfied even more accurately than required for convergence. This additional regularity should not be too surprising since even for explicit ODE's, increased accuracy of the starting values is needed in the derivation of the smooth leading error terms for general linear multistep methods [15].

In the analysis given in Section 3 and the Supplement section, we assume for technical simplicity that the difference equations are solved exactly on each integration step. However, the convergence analysis could be modified to include residual terms representing the effect of terminating the iteration process. We have already explained how to relax the requirement of satisfying the algebraic equations exactly. Residuals could also be added to the difference equations corresponding to

the ODE's. A straightforward generalization of the analysis is possible by simply forcing these residuals to be sufficiently small so as to not affect the behavior of the derived error estimates (e.g., in [18] it is assumed that the residuals corresponding to the ODE part are  $O(h^k)$  for the index two system). Slightly stronger conditions may have to be imposed to obtain the results derived here, in particular for those results concerning the leading error terms.

Convergence has recently been proved for variable step meshes for nonlinear, semiexplicit, index two systems in [12]. However, the rate of convergence is not established. We cannot extend the convergence analysis for the index three systems to hold for variable stepsize meshes without loss of accuracy because each time the stepsize is changed, a new boundary layer of reduced convergence rates would be initiated. In particular, the first-order BDF would fail to converge at the end of the first step, following every change in the stepsize.

**3. Convergence Proof for the Index Two System (Theorem 1).** We have already stated that there exists a unique solution  $(d(t), e(t))$  to the linear, index two system (2.9), (2.10) for each set of consistent initial values, but we must show one can always find a set of consistent initial values for any set of starting values satisfying (2.3), (2.4). Clearly, from (2.10), we must select  $d(t_{k-1})$  in the nullspace of  $C(t_{k-1})$ , denoted  $\mathcal{N}(C(t_{k-1}))$ . By assumption (2.3), it follows that  $y_{k-1} = y(t_{k-1}) + h^{k+1}c_{1,k-1}$  for some bounded vector  $c_{1,k-1}$ . Therefore, since we want (2.7) to be consistent with our initial values at  $t_{k-1}$ , we select  $d(t_{k-1}) = 0$  and  $r_{k-1} = c_{1,k-1}$ . We then select  $e(t_{k-1})$  to satisfy the first derivative of (2.10), namely

$$(3.1) \quad e(t) = -[C(t)B(t)]^{-1}C(t)y^{(k+1)}(t)/(k + 1)$$

for  $t = t_{k-1}$ . Since  $(d(t_{k-1}), e(t_{k-1}))$  is a consistent set of initial values for (2.9), (2.10), we can define  $(d(t), e(t))$  for  $t \geq t_{k-1}$  to be the solution. In the more general case that  $d(t_{k-1}) \in \mathcal{N}(C(t_{k-1}))$  but is not zero, as discussed in an earlier remark, one can also always find a consistent set of initial values for the DAE system (2.9), (2.10).

Since the  $k$ -step BDF requires a numerically consistent starting vector  $\mathbf{y}_{k-1}$ , we will define  $d(t_i) = 0$  and  $r_i = c_{1,i}$  for  $i = 0, 1, \dots, k - 2$ , where  $c_{1,i}$  are bounded vectors such that  $y_i = y(t_i) + h^{k+1}c_{1,i}$ . Let  $e(t_i)$  be defined as in (3.1) for  $t = t_i$ ,  $i = 0, 1, \dots, k - 2$ . Then, there exists a constant  $\eta_0$  such that  $\|\mathbf{r}_{k-1}\| \leq \eta_0$  where  $\mathbf{r}_{k-1} = [r_{k-1}^T, \dots, r_0^T]^T$ .

After expanding by Taylor series about the analytic solution  $(y(t_n), u(t_n))$ , and using the fact that  $(d(t), e(t))$  is the solution to (2.9), (2.10) with associated initial conditions, Eqs. (2.5), (2.6) can be rewritten in the form,

$$(3.2) \quad S(t_n) \begin{bmatrix} r_n \\ h s_n \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^k \gamma_i r_{n-i} \\ 0 \end{bmatrix} + \begin{bmatrix} \psi_1(t_n) \\ \psi_2(t_n) \end{bmatrix},$$

where

$$S(t_n) = \begin{bmatrix} I_l - \frac{h}{\alpha_0} A(t_n) & -\frac{1}{\alpha_0} B(t_n) \\ C(t_n) & 0_m \end{bmatrix},$$

$$\gamma_i = -\alpha_i/\alpha_0 \quad \text{where } \alpha_i \text{ are the BDF coefficients,}$$

$$(3.3) \quad \psi_1(t) = -\frac{h}{\alpha_0} R_1(\xi) + \frac{h^k}{\alpha_0(k+1)} d^{(k+1)}(t) - \frac{h^{k+1}}{\alpha_0} R_{12}(\xi) + h^k Q_1(t) \\ + h^{k+1}(W^{(1)}(r) + W^{(2)}(s)) + h^{k+2} g_1(r, s) + O(h^{2k}),$$

$$(3.4) \quad \psi_2(t) = -h^{k-1} Q_2(t) - h^k W^{(3)}(r) - h^{k+1} g_2(r) + O(h^{2k}),$$

$$Q_1(t) = \frac{1}{2\alpha_0} \{d(t)^T [E_{yy} d(t) + E_{uy} e(t)] \\ + e(t)^T [E_{yu} d(t) + E_{uu} e(t)]\} + O(h^k),$$

$$Q_2(t) = \frac{1}{2} d(t)^T H_{yy} d(t) + O(h^k),$$

$$W^{(1)}(r) = \frac{1}{2\alpha_0} \{r^T [E_{yy} d(t) + E_{uy} e(t)] + [d(t)^T E_{yy} + e(t)^T E_{yu}] r\},$$

$$W^{(2)}(s) = \frac{1}{2\alpha_0} \{s^T [E_{yu} d(t) + E_{uu} e(t)] + [d(t)^T E_{uy} + e(t)^T E_{uu}] s\},$$

$$W^{(3)}(r) = \frac{1}{2} \{r^T H_{yy} d(t) + d(t)^T H_{yy} r\},$$

$$g_1(r, s) = \frac{1}{2\alpha_0} \{r^T (E_{yy} r + E_{uy} s) + s^T (E_{yu} r + E_{uu} s)\},$$

$$g_2(r) = \frac{1}{2} r^T H_{yy} r$$

for  $t = t_n$ ,  $r = r_n$ ,  $s = s_n$  and  $\xi = \xi_n$ ,  $t_{n-k} \leq \xi_n \leq t_n$ , and where all the partial derivatives of  $E$  and  $H$  are evaluated at the analytic solution  $(y(t_n), u(t_n))$ . The functions  $R_1(\xi)$  and  $R_{12}(\xi)$  are the remainders from the Taylor series expansions, and hence are bounded functions of  $y^{(k+2)}(\xi)$  and  $d^{(k+2)}(\xi)$ , respectively.

The functions  $Q_1(t)$  and  $Q_2(t)$  contain all the inhomogeneous terms which are functions only of  $d(t)$ ,  $e(t)$ , and partial derivatives of  $E$  and  $H$ . We will carefully analyze the leading terms in the Taylor series part of  $\psi_1$  and  $\psi_2$  (through the quadratic terms in  $r$  and  $s$ ). The inclusion of higher-order terms from the Taylor series will only introduce terms of order higher than those already present. These higher-order terms have no influence on the order of the estimates derived in this analysis, and are presented by the final terms  $O(h^{2k})$  in the definition of  $\psi_1$  and  $\psi_2$ .

Since  $W^{(i)}$  ( $i = 1, 2, 3$ ) are linear, bounded functions of their arguments, there exist bounded matrix operators  $W_i$  such that

$$W^{(i)}(\alpha) = W_i \alpha, \quad \|W_i\| \leq \omega_i \quad \text{for } i = 1, 2, 3.$$

Define  $\mathbf{s}_n = [s_n^T, s_{n-1}^T, \dots, s_{n-k+1}^T]^T$ . Then, if

$$S^{-1}(t_n) = \begin{pmatrix} X_n & Y_n \\ U_n & V_n \end{pmatrix}$$

we can rewrite (3.2) in the one-step form,

$$(3.5) \quad \mathbf{r}_n = F_n \mathbf{r}_{n-1} + f_{1,n},$$

$$(3.6) \quad h\mathbf{s}_n = G_n \mathbf{r}_{n-1} + f_{2,n} + hH\mathbf{s}_{n-1},$$

where

$$\begin{aligned} F_n &= \text{CPM}(\gamma_\mu X_n) \quad (\mu = 1, 2, \dots, k), \\ f_{1,n} &= [(X_n \psi_1(t_n) + Y_n \psi_2(t_n))^T, 0, \dots, 0]^T, \\ G_n &= (I_m, 0, \dots, 0)^T (\gamma_1 U_n, \dots, \gamma_k U_n), \\ f_{2,n} &= [(U_n \psi_1(t_n) + V_n \psi_2(t_n))^T, 0, \dots, 0]^T, \\ H &= \text{CPM}(0_m). \end{aligned}$$

Closed form expressions for the block elements of the matrix inverse of  $S(t_n)$  can easily be derived, but the following approximate relations are sufficient for the analysis ( $t$ -dependence is suppressed):

$$(3.7) \quad X = I_l - B(CB)^{-1}C + O(h),$$

$$(3.8) \quad Y = B(CB)^{-1} + O(h),$$

$$(3.9) \quad U = -\alpha_0(CB)^{-1}C + O(h),$$

$$(3.10) \quad V = \alpha_0(CB)^{-1} + O(h).$$

Since  $C(t)B(t)$  is nonsingular for all  $t \in I$ , it follows that  $S^{-1}(t)$  is as smooth as  $S(t)$ , and hence the block elements of  $S^{-1}(t)$  can be bounded uniformly on  $I$ . Therefore, there exists a constant  $K$  such that

$$\|X(t)\|, \|Y(t)\|, \|U(t)\|, \|V(t)\| \leq K$$

for all  $t \in I$ . In general, we will let the symbol  $K$  denote a generic constant, independent of  $n$  and  $h$ .

Note that Eqs. (3.5), (3.6) are implicit in  $r_n$  and  $hs_n$ . We shall construct a fixed point iteration for their solution  $r_n$  and  $hs_n$ , and an induction argument concerning the past values of  $r_{n-1}$  and  $hs_{n-1}$  will be required. Specifically, let the *induction assumption* be the following:

$$(3.11) \quad \|r_i\| \leq \eta_1 e^{(i-k+1)hL}, \quad \|hs_i\| \leq \eta_2 e^{(i-k+1)hL}$$

for  $i = k, k + 1, \dots, n - 1$ . We will prove there exist such constants  $\eta_1, \eta_2$  and  $L$  and the equations (3.11) are satisfied for  $i = n$ . It follows then from the induction assumption and the definitions of  $\psi_1(t)$  and  $\psi_2(t)$  given in (3.3), (3.4) that for  $i = k, k + 1, \dots, n - 1$

$$(3.12) \quad \|\psi_1(t_i)\| \leq \sigma_1 h, \quad \|\psi_2(t_i)\| \leq \sigma_2 h^{k-1},$$

where to  $O(h)$  accuracy  $\sigma_1$  is independent of  $\eta_1, \eta_2$  and  $L$  if  $k \geq 2$ ,  $\sigma_1$  is dependent on  $\eta_2$  and  $L$  if  $k = 1$  (i.e.,  $\sigma_1 = \tilde{\sigma}_1(\eta_2, L) + O(h)$ ), and  $\sigma_2$  is independent of  $\eta_1, \eta_2$  and  $L$  for all  $k$ .

Define the iterates as

$$\begin{aligned} \mathbf{r}_n^{(\nu)} &= [(r_n^{(\nu)})^T, r_{n-1}^T, \dots, r_{n-k+1}^T]^T, \\ \mathbf{s}_n^{(\nu)} &= [(s_n^{(\nu)})^T, s_{n-1}^T, \dots, s_{n-k+1}^T]^T, \end{aligned}$$

and the fixed-point iteration as

$$(3.13) \quad \mathbf{r}_n^{(\nu+1)} = \mathbf{r}_n^{(0)} + [Z_1(r_n^{(\nu)}, hs_n^{(\nu)})^T, 0, \dots, 0]^T,$$

$$(3.14) \quad h\mathbf{s}_n^{(\nu+1)} = h\mathbf{s}_n^{(0)} + [Z_2(r_n^{(\nu)}, h\mathbf{s}_n^{(\nu)})^T, 0, \dots, 0]^T,$$

where

$$\begin{aligned} Z_1(r, h\mathbf{s}) &= X_n[h^{k+1}(W_1r + W_2s) + h^{k+2}g_1(r, s)] \\ &\quad + Y_n[-h^k W_3r - h^{k+1}g_2(r)] + O(h^{2k}), \\ Z_2(r, h\mathbf{s}) &= U_n[h^{k+1}(W_1r + W_2s) + h^{k+2}g_1(r, s)] \\ &\quad + V_n[-h^k W_3r - h^{k+1}g_2(r)] + O(h^{2k}) \end{aligned}$$

for  $r = r_n^{(\nu)}$  and  $s = s_n^{(\nu)}$ . For starting iterates select

$$(3.15) \quad \mathbf{r}_n^{(0)} = F_n \mathbf{r}_{n-1} + \tilde{f}_{1,n},$$

$$(3.16) \quad h\mathbf{s}_n^{(0)} = G_n \mathbf{r}_{n-1} + \tilde{f}_{2,n} + hH\mathbf{s}_{n-1},$$

where

$$\begin{aligned} \tilde{f}_{1,n} &= [(X_n \tilde{\psi}_1(t_n) + Y_n \tilde{\psi}_2(t_n))^T, 0, \dots, 0]^T, \\ \tilde{f}_{2,n} &= [(U_n \tilde{\psi}_1(t_n) + V_n \tilde{\psi}_2(t_n))^T, 0, \dots, 0]^T, \\ \tilde{\psi}_1(t_n) &= -\frac{h}{\alpha_0} R_1(\xi_n) + \frac{h^k}{\alpha_0(k+1)} d^{k+1}(t_n) - \frac{h^{k+1}}{\alpha_0} R_{12}(\xi_n) + h^k Q_1(t_n), \\ \tilde{\psi}_2(t_n) &= -h^{k-1} Q_2(t_n). \end{aligned}$$

The following three conditions corresponding to the implicit function theorem are sufficient to prove the fixed-point iteration defined in (3.13), (3.14) converges to a solution:

$$(3.17) \quad \|\mathbf{r}_n^{(0)}\| \leq \eta_1 e^{(n-k+1)hL}, \quad \|h\mathbf{s}_n^{(0)}\| \leq \eta_2 e^{(n-k+1)hL},$$

$$(3.18) \quad \|Z(r_n^{(0)}, h\mathbf{s}_n^{(0)})\| \leq \delta/2, \quad \delta > 0, \quad Z = (Z_1^T, Z_2^T)^T;$$

$$(3.19) \quad \|J\| \leq \frac{1}{2} \text{ for any } r \text{ and } h\mathbf{s} \text{ such that } \|r - r_n^{(0)}\| \leq \delta \text{ and } \|h\mathbf{s} - h\mathbf{s}_n^{(0)}\| \leq \delta,$$

where  $J$  is the Jacobian matrix of  $Z(r, h\mathbf{s})$  with respect to  $r$  and  $h\mathbf{s}$ .

If  $r$  and  $h\mathbf{s}$  are bounded, conditions (3.18) and (3.19) are relatively straightforward to verify for the iteration defined in (3.13), (3.14). The expressions for  $Z_1(r, h\mathbf{s})$  and  $Z_2(r, h\mathbf{s})$  imply that there exists a  $\delta, \delta = Kh^k$ , which satisfies (3.18) providing  $r = r_n^{(0)}$  and  $h\mathbf{s} = h\mathbf{s}_n^{(0)}$  are bounded. The Jacobian matrix  $J$  will be linear in  $r$  and  $h\mathbf{s}$ , so condition (3.19) will be satisfied for sufficiently small  $h$  and for all  $r$  and  $h\mathbf{s}$  such that  $\|r - r_n^{(0)}\| \leq \delta$  and  $\|h\mathbf{s} - h\mathbf{s}_n^{(0)}\| \leq \delta$ . Hence, it only remains to prove condition (3.17).

It is quite easy to see the starting iterates for the fixed point iteration as defined by (3.15), (3.16) are bounded for  $n = k$ . The remaining conditions for convergence of the fixed point iteration are also satisfied, so the induction assumption is satisfied at the first step.

Now we shall assume the induction assumption (3.11) is true and prove the initial guesses  $\mathbf{r}_n^{(0)}$  and  $h\mathbf{s}_n^{(0)}$  are bounded as in (3.17). We will bound  $h\mathbf{s}_n^{(0)}$  first because when  $k = 1$  the constant  $\eta_1$  will be chosen dependent on  $\eta_2$  and  $L$ . This

dependency is a result of our decision to bound only  $hs_n^{(0)}$ , and not  $s_n^{(0)}$ . To bound  $s_n^{(0)}$  requires much more algebraic manipulation of the structure of the system [1]. This result is obtained in another way in the corollary.

To verify (3.17), we must utilize important cancellation properties in certain matrix products involving the  $O(1)$  matrices  $X, Y$ , and  $U$  defined in (3.7)–(3.9):

$$(3.20) \quad \|X(t_i)Y(t_j)\|, \|U(t_i)X(t_j)\| \leq Kh$$

for any  $|t_i - t_j| = O(h)$  and some generic constant  $K$  independent of  $h$ . In addition, the proof requires that the matrix product  $\prod_{j=i+1}^n F_j$  be bounded uniformly in  $n$  for  $i = 0, 1, \dots, n - 1$ . Since  $F_j = \text{CPM}(\gamma_\mu X_j)$  where  $X_j$  is a projection matrix to  $O(h)$  accuracy, it follows from a result given by Kreiss [16] and Strang [26] that there exist constants  $K^*$  and  $\mathcal{E}$  such that

$$\left\| \prod_{j=i+1}^n F_j \right\| \leq K^* e^{(n-i)h\mathcal{E}}$$

for all  $n$  and  $i = 0, 1, \dots, n - 1$ . Moreover, using the structure of the block companion matrix  $F_j$  and (3.20), it follows that

$$\left\| G_{n-l} \prod_{j=i+1}^{n-l-1} F_j \right\| \leq \tilde{N}h$$

for all  $n - l - i - 1 \geq k$  and some constant  $\tilde{N}$ . In bounding the starting iterates, we will also need the fact that there exist constants  $\rho_1$  and  $\rho_2$  such that

$$\|\tilde{\psi}_1(t_n)\| \leq \rho_1 h, \quad \|\tilde{\psi}_2(t_n)\| \leq \rho_2 h^{k-1}$$

for any  $t_n \in [t_k, t_0 + T]$ . Now rewrite Eq. (3.16) for the initial guess  $hs_n^{(0)}$  as

$$\begin{aligned} hs_n^{(0)} &= \sum_{l=0}^{k-1} H^l G_{n-l} \left( \prod_{i=k}^{n-l-1} F_i \right) \mathbf{r}_{k-1} + \sum_{l=1}^{k-1} H^l f_{2,n-l} + \tilde{f}_{2,n} \\ &\quad + \sum_{l=0}^{k-1} H^l G_{n-l} \left( \sum_{i=k}^{n-l-2} \left( \prod_{j=i+1}^{n-l-1} F_j \right) f_{1,i} + f_{1,n-l-1} \right). \end{aligned}$$

It is then possible to bound  $hs_n^{(0)}$  using these results in the following way:

$$\begin{aligned} \|hs_n^{(0)}\| &\leq kK \max_{1 \leq i \leq k} |\gamma_i| K^* e^{(n-k)h\mathcal{E}} \eta_0 + K(\rho_1 h + \rho_2 h^{k-1}) \\ &\quad + K(\sigma_1 h + \sigma_2 h^{k-1})[(n - 2k)hk\tilde{N} + k^2 K K^* e^{(k-1)h\mathcal{E}} + (k - 1)] \end{aligned}$$

for all  $n \geq k$ . This bound for  $hs_n^{(0)}$  is independent of  $\eta_1, \eta_2$ , and  $L$  to  $O(h)$  accuracy for all  $k$ . Therefore, there exists an  $\eta_2$  such that for sufficiently small  $h$  and  $L = \mathcal{E}$ , we have

$$\|hs_n^{(0)}\| < \eta_2 e^{(n-k+1)h\mathcal{E}}.$$

This result is valid for  $n \rightarrow \infty, h \rightarrow 0$  and  $t_n$  such that  $t_n \in [t_k, t_0 + T] \subseteq I$ .

Next we bound  $\mathbf{r}_n^{(0)}$  defined by (3.15), which can be rewritten as

$$\mathbf{r}_n^{(0)} = \left( \prod_{j=k}^n F_j \right) \mathbf{r}_{k-1} + \sum_{i=k}^{n-1} \left( \prod_{j=i+1}^n F_j \right) f_{1,i} + \tilde{f}_{1,n}.$$

Using relation (3.20), the structure of  $F_j$ , and the bounds (3.12), it is possible to show that there exists a constant  $\Delta_1$  dependent on  $\eta_2$  and  $L$  if  $k = 1$ , such that  $\|(\prod_{j=i+1}^n F_j) f_{1,i}\| \leq \Delta_1 h$  for  $n - i \geq k$ . Then it follows from this and earlier bounds that

$$\begin{aligned} \|\mathbf{r}_n^{(0)}\| &\leq K^* e^{(n-k+1)h\mathcal{E}} \eta_0 + (n - 2k + 1)h\Delta_1 \\ &\quad + (k - 1)K(\sigma_1 h + \sigma_2 h^{k-1})K^* e^{(k-1)h\mathcal{E}} + K(\rho_1 h + \rho_2 h^{k-1}) \end{aligned}$$

for all  $n \geq k$ . If  $k = 1$  the bound depends on  $\eta_2$  and  $\mathcal{E}$  (through  $\Delta_1$ ), and hence a constant  $\eta_1$  must be chosen dependent on  $\eta_2$  and  $\mathcal{E}$  such that

$$\|\mathbf{r}_n^{(0)}\| < \eta_1 e^{(n-k+1)h\mathcal{E}}.$$

If  $k \geq 2$ ,  $\eta_1$  may be chosen independent of  $\eta_2$  and  $\mathcal{E}$ .

Since both  $hs_n^{(0)}$  and  $\mathbf{r}_n^{(0)}$  are uniformly bounded as in (3.17), it follows immediately that conditions (3.18) and (3.19) of the implicit function theorem are satisfied. Therefore, there exists a solution  $\mathbf{r}_n$  and  $hs_n$  satisfying Eqs. (3.5) and (3.6), which can be computed by straightforward iteration. Moreover, there exists a  $\delta$ ,  $\delta = O(h^k)$ , such that  $\|\mathbf{r}_n - \mathbf{r}_n^{(0)}\| = \|Z_1(\mathbf{r}_n, hs_n)\| \leq \delta$  and  $\|hs_n - hs_n^{(0)}\| = \|Z_2(\mathbf{r}_n, hs_n)\| \leq \delta$ . For sufficiently small  $h$ , it follows that

$$\begin{aligned} \|\mathbf{r}_n\| &\leq \|\mathbf{r}_n^{(0)}\| + \delta \leq \eta_1 e^{(n-k+1)h\mathcal{E}}, \\ \|hs_n\| &\leq \|hs_n^{(0)}\| + \delta \leq \eta_2 e^{(n-k+1)h\mathcal{E}}. \end{aligned}$$

This concludes the induction argument. The global rate of convergence is  $O(h^k)$ , since

$$\begin{aligned} \|y_n - y(t_n)\| &\leq h^k \|d(t_n)\| + h^{k+1} \|\mathbf{r}_n\|, \\ \|u_n - u(t_n)\| &\leq h^k (\|e(t_n)\| + \|hs_n\|) \end{aligned}$$

for all  $n \geq k$ .

*Proof of Corollary 1.* Since the algebraic equations are satisfied to  $O(h^{k+2})$  accuracy, both at the starting values and at the end of each integration step, it follows that  $C(t_i)r_i = \psi_2(t_i) + O(h)$  for  $i = 0, 1, \dots$ . From the difference equations (3.2) corresponding to the ODE part, we have

$$\left[ I_l - \frac{h}{\alpha_0} A(t_n) \right] r_n - \frac{h}{\alpha_0} B(t_n) s_n = \sum_{i=1}^k \gamma_i r_{n-i} + \psi_1(t_n).$$

Let these equations also be solved to accuracy  $O(h^{k+2})$ . Then the residual error may be included as a term of  $O(h)$  accuracy in  $\psi_1$ . Multiply by  $C(t_n)$ , an  $m \times n$  matrix of rank  $m$ , and invert  $C(t_n)B(t_n)$  to obtain a relation for  $s_n$ :

$$\begin{aligned} (3.21) \quad s_n &= -\frac{\alpha_0}{h} (C(t_n)B(t_n))^{-1} C(t_n) \left( \sum_{i=1}^k \gamma_i r_{n-i} + \psi_1(t_n) \right. \\ &\quad \left. - \left[ I_l - \frac{h}{\alpha_0} A(t_n) \right] r_n \right). \end{aligned}$$

In the proof of Theorem 1, it is shown that  $r_n$  and  $hs_n$  are uniformly bounded for  $n \geq k$ . Then,  $\|\psi_1(t_n)\| \leq Kh$  for some constant  $K$  independent of  $n$  and  $h$  for all  $n \geq k$ . For  $k \geq 2$  we have

$$\|C(t_n)r_{n-i}\| \leq \|C(t_{n-i})r_{n-i}\| + O(h) \leq \|\psi_2(t_{n-i})\| + O(h) = O(h)$$

for all  $n \geq k$  and  $i = 0, 1, 2, \dots, k-1$ . Thus, for  $k \geq 2$  it follows that the  $1/h$  factor in (3.21) is annihilated, leaving only terms of order  $O(1)$ . If  $k = 1$  we must be a bit more careful and utilize a cancellation property. The expression (3.21) simplifies to

$$s_n = -\frac{1}{h}(C(t_n)B(t_n))^{-1}C(t_n) \left( r_{n-1} + \psi_1(t_n) - \left[ I_l - \frac{h}{\alpha_0}A(t_n) \right] r_n \right).$$

Using the smoothness of  $Q_2(t)$ , the lowest order (i.e.,  $O(1)$ ) term in  $\psi_2(t)$ , it follows that

$$\begin{aligned} -C(t_n)r_{n-1} + C(t_n)r_n &= -C(t_{n-1})r_{n-1} + C(t_n)r_n + O(h) \\ &= -\psi_2(t_{n-1}) + \psi_2(t_n) + O(h) = O(h). \end{aligned}$$

As for  $k \geq 2$  then, it follows that  $s_n$  is bounded for all  $n \geq k$ , and these bounds can be chosen uniformly since they depend on the uniform bounds for  $r_n$  and  $hs_n$ . Therefore, Eqs. (2.7), (2.8) accurately reflect the relation of the numerical solution to the true solution for all  $n \geq k$ . The functions  $d(t)$  and  $e(t)$  in the principal leading error terms are the unique solutions to the linear, index two DAE system with consistent initial conditions described earlier.

*Remark.* Note that the algebraic equations must be satisfied more closely to prove that  $s_n$  is bounded, than required in the proof of convergence. If they are satisfied only to  $O(h^{k+1})$  accuracy, convergence is still obtained but the error is proportional to  $e(t_n) + hs_n$ .

The Aerospace Corporation  
2350 E. El Segundo Boulevard  
El Segundo, California 90245  
E-mail: brenan@aerospace.aero.org

Department of Mathematics  
University of California  
Los Angeles, California 90024

1. K. E. BRENNAN, *Stability and Convergence of Difference Approximations for Higher Index Differential-Algebraic Systems with Applications in Trajectory Control*, Ph.D. Thesis, University of California at Los Angeles, 1983.
2. K. E. BRENNAN, "Numerical simulation of trajectory prescribed path control problems by the backward differentiation formulas," *IEEE Trans. Automat. Control*, v. AC-31, 1986, pp. 266-269.
3. K. E. BRENNAN & B. E. ENGQUIST, *Backward Difference Approximations of Nonlinear Differential-Algebraic Equations*, ATR-85(9990)-5, The Aerospace Corporation, El Segundo, Ca., 1985 or Dept. of Computer Science, Uppsala University Report #101, Uppsala, Sweden.
4. A. E. BRYSON & Y. C. HO, *Applied Optimal Control*, Hemisphere Publishing Corp., Washington, D.C., 1975.
5. S. L. CAMPBELL, *Singular Systems of Differential Equations I-II*, Pitman, Marshfield, Mass., 1980, 1982.
6. K. D. CLARK, "The numerical solution of some higher index time varying semistate systems by difference methods," *Circuits Systems Signal Process.* v. 6, 1987, pp. 61-75.
7. F. R. GANTMACHER, *The Theory of Matrices*, Vol. I-II, Chelsea, New York, 1959.
8. C. W. GEAR, "The simultaneous numerical solution of differential-algebraic equations," *IEEE Trans. Circuit Theory*, v. CT-18, 1971, pp. 89-95.
9. C. W. GEAR & L. R. PETZOLD, "ODE methods for the solution of differential/algebraic systems," *SIAM J. Numer. Anal.*, v. 21, 1984, pp. 367-384.
10. C. W. GEAR & L. R. PETZOLD, "Differential/algebraic systems and matrix pencils," in *Matrix Pencils* (B. Kågström and A. Ruhe, eds.), Lecture Notes in Math., vol. 973, 1983, pp. 75-89.

11. C. W. GEAR, H. H. HSU & L. PETZOLD, *Differential-Algebraic Equations Revisited*, Proc. Conference on Matrix Pencils, Piteå, Sweden, 1982.
12. C. W. GEAR, B. LEIMKUHLE & G. K. GUPTA, "Automatic integration of Euler-Lagrange equations with constraints," *J. Comput. Appl. Math.*, v. 12 and 13, 1985, pp. 77-90.
13. H. GINGOLD, "A method of global blockdiagonalization for matrix-valued functions," *SIAM J. Numer. Anal.*, v. 9, 1978, pp. 1076-1082.
14. P. HENRICI, *Discrete Variable Methods in Ordinary Differential Equations*, Wiley, New York, 1962.
15. P. HENRICI, *Error Propagation for Difference Methods*, Wiley, New York, 1963.
16. H. O. KREISS, "Über implizite Differenzmethoden für partielle Differentialgleichungen," *Numer. Math.*, v. 5, 1963, pp. 24-47.
17. P. LÖTSTEDT, "Mechanical systems of rigid bodies subject to unilateral constraints," *SIAM J. Appl. Math.*, v. 42, 1982, pp. 281-296.
18. P. LÖTSTEDT & L. R. PETZOLD, "Numerical solution of nonlinear differential equations with algebraic constraints I: Convergence results for backward differentiation formulas," *Math. Comp.*, v. 46, 1986, pp. 491-516.
19. R. MÄRZ, *Multistep Methods for Initial Value Problems in Implicit Differential-Algebraic Equations*, Preprint 22, Humboldt-Univ., Berlin, Sektion Mathematik, 1981.
20. R. MÄRZ, "On initial value problems in differential-algebraic equations and their numerical treatment," *Computing*, v. 35, 1985, pp. 13-37.
21. J. F. PAINTER, *Solving the Navier-Stokes Equations with LSODI and the Method of Lines*, Report UCID-19262, Lawrence Livermore National Laboratory, Livermore, Ca., 1981.
22. L. R. PETZOLD, "Differential/algebraic equations are not ODE's," *SIAM J. Sci. Statist. Comput.*, v. 3, 1982, pp. 367-384.
23. W. C. RHEINBOLDT, "Differential-algebraic systems as differential equations on manifolds," *Math. Comp.*, v. 43, 1984, pp. 473-482.
24. R. D. RICHTMYER & R. W. MORTON, *Differential Methods for Initial-Value Problems*, Interscience, New York, 1967.
25. R. F. SINCOVEC, A. M. ERISMAN, E. L. YIP & M. A. EPTON, "Analysis of descriptor systems using numerical algorithms," *IEEE Trans. Automat. Control*, v. AC-26, 1981, pp. 139-147.
26. G. STRANG, "Accurate partial difference methods II-Nonlinear problems," *Numer. Math.*, v. 6, 1964, pp. 37-46.
27. J. H. WILKINSON, "Linear differential equations and Kronecker's canonical form" in *Recent Applications in Numerical Analysis* (C. deBoor and G. H. Golub, eds.), Academic Press, New York, 1978, pp. 231-241.
28. E. L. YIP & R. F. SINCOVEC, "Solvability, controllability, and observability of continuous descriptor systems," *IEEE Trans. Automat. Control*, v. AC-26, 1981, pp. 702-707.

## Supplement to Backward Differentiation Approximations of Nonlinear Differential/Algebraic Systems

By Kathryn E. Brenan and Bjorn E. Engquist

### A . Convergence Proof for the Index-Three System (Theorem 2).

As we did for the index-two system, we must prove there always exists a set of consistent initial values  $(c(t_{k-1}), b(t_{k-1}), e(t_{k-1}))$  for the linear, index-three DAE system (2.11)-(2.13) corresponding to any set of initial values  $(w_{k-1}, u_{k-1}, v_{k-1})$  satisfying (2.3),(2.4). From (2.13) it is clear we should choose  $b(t_{k-1}) \in \mathcal{N}(A_{32}(t_{k-1}))$ . The consistency condition (2.3) implies we can write  $w_{k-1} = w(t_{k-1}) + h^{k+1}c_{1,k-1}$  and  $v_{k-1} = v(t_{k-1}) + h^{k+1}c_{2,k-1}$  for some bounded vectors  $c_{1,k-1}$  and  $c_{2,k-1}$ . Since we want our initial values to be consistent with the asymptotic error expansions (2.7),(2.8), we select for  $t = t_{k-1}$  and  $i = k - 1$ ,

$$b(t) = 0, \quad p_i = c_{1,i}. \tag{A.1}$$

We must then choose  $c(t_{k-1})$  to satisfy the first derivative of the algebraic equation (2.13) in the linear, index-three system: ( $t = t_{k-1}$ )

$$A_{32}(t)A_{21}(t)c(t) = -\frac{1}{(k+1)}A_{32}(t)w^{(k+1)}(t). \tag{A.2}$$

Since the range of  $A_{32}(t)$ , denoted  $\mathbf{R}(A_{32}(t))$ , equals  $\mathbf{R}(A_{32}(t)A_{21}(t))$  for all  $t \in I$ , there is always a solution  $c(t_{k-1})$  to (A.2). Next we must choose  $q_{k-1}$  so that (2.7) is valid for  $n = k - 1$  (here  $t = t_{k-1}$  and  $i = k - 1$ ):

$$q_i = -\frac{1}{h}c(t) + c_{2,i}. \tag{A.3}$$

In general,  $c(t_{k-1}) \neq 0$  so  $q_{k-1}$  will have order  $O(1/h)$ . The convergence analysis will require us to pick  $q_{k-1}$  such that  $P_{11}(t_{k-1})q_{k-1}$  is bounded, where  $P_{11}(t)$  is a projection, namely  $P_{11}(t) = I_p - A_{13}(t)\Gamma(t)A_{32}(t)A_{21}(t)$  and  $\Gamma(t) = (A_{32}(t)A_{21}(t)A_{13}(t))^{-1}$ . Therefore, if we require  $P_{11}(t_{k-1})c(t_{k-1}) = 0$ , so that  $P_{11}(t_{k-1})q_{k-1}$  is bounded, then from (A.2) it follows that for  $t = t_{k-1}$

$$c(t) = -\frac{1}{(k+1)}A_{13}(t)\Gamma(t)A_{32}(t)w^{(k+1)}(t). \tag{A.4}$$

The consistent initial value for  $e(t_{k-1})$  is completely determined from the equation obtained by differentiating (2.13) twice and substituting for  $c'(t)$  and  $b'(t)$ :

$$e(t) = -\Gamma(t) \left( \left( 2A'_{32}A_{21} + A_{32}A'_{21} + A_{32}(A_{21}A_{11} + A_{22}A_{21}) \right) c(t) + \frac{1}{(k+1)} \left( A_{32}A_{21}v^{(k+1)}(t) + (2A'_{32} + A_{32}A_{22})w^{(k+1)}(t) + A_{32}w^{(k+2)}(t) \right) \right), \tag{A.5}$$

where all the  $A_{ij}$  matrices are evaluated at  $t = t_{k-1}$ .

Since  $P_{11}(t)c(t) = 0$ , it follows from (A.1) and (A.3) that there exists a constant  $\eta_0$  such that

$$\|P_{11}(t)c(t)\| \leq \eta_0, \quad (\text{A.6})$$

$$\|P_{11}(t)g_i\| = \|P_{11}(t)c_{2i}\| \leq \eta_0 \quad (\text{A.7})$$

for  $t = t_{k-1}$  and  $i = k-1$ . Then the solution  $(c(t), b(t), e(t))$  to the linear, index-three system (2.11)-(2.13) satisfying the consistency conditions just given for  $(c(t_{k-1}), b(t_{k-1}), e(t_{k-1}))$  is unique and smooth for  $t \geq t_{k-1}$ . Since the  $k$ -step BDF requires a numerically consistent starting vector  $\mathcal{Y}_{k-1}$ , we will define  $(c(t_k), b(t_k), e(t_k))$  using equations (A.1)-(A.5) when  $t = t_k$ , and  $i = 0, 1, \dots, k-2$ . Let us redefine  $\eta_0$  so that (A.6) and (A.7) are satisfied for all  $i = 0, 1, \dots, k-1$ . If the consistency relations are relaxed as stated in the remark following the statement of Theorem 2, consistent initial values  $(c(t_{k-1}), b(t_{k-1}), e(t_{k-1}))$  can be determined similarly, although in general equations (A.1)-(A.5) are more complicated (because  $b(t_{k-1})$  may not be zero). Then, after substituting equations (2.7)-(2.8) for the numerical solution into the BDF equations (2.5)-(2.6) and after expanding by Taylor series about the analytic solution at  $t_n$ , it follows from (2.11)-(2.13) that the remainders  $(g_n, p_n, s_n)$  must satisfy the system

$$S(t_n) \begin{bmatrix} g_n \\ p_n \\ h^2 s_n \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^k \gamma_i g_{n-i} \\ \sum_{i=1}^k \gamma_i p_{n-i} \\ 0 \end{bmatrix} + \begin{bmatrix} \psi_1(t_n) \\ \psi_2(t_n) \\ \psi_3(t_n) \end{bmatrix}, \quad (\text{A.8})$$

where

$$S(t_n) = \begin{bmatrix} I_p - \frac{\Delta}{\alpha_0} A_{11}(t_n) & -\frac{\Delta}{\alpha_0} A_{12}(t_n) & -\frac{\Delta}{\alpha_0 h} A_{13}(t_n) \\ -\frac{\Delta}{\alpha_0} A_{21}(t_n) & I_q - \frac{\Delta}{\alpha_0} A_{22}(t_n) & 0 \\ 0 & A_{31}(t_n) & 0 \end{bmatrix},$$

$$\gamma_i = -\alpha_i/\alpha_0 \text{ where } \alpha_i \ (i = 0, 1, \dots, k) \text{ are the BDF coefficients,}$$

$$\psi_1(t_n) = \left( -hR_1(\xi_n) + h^k c^{(k+1)}(t_n)/(k+1) - h^{k+1}R_{12}(\xi_n) + h^k Q_1(t_n) + h^{k+1}W_1(g_n, p_n, s_n) + h^{k+2}g_1(g_n, p_n, s_n) \right) / \alpha_0, \quad (\text{A.9})$$

$$\psi_2(t_n) = \left( -hc^* w^{(k+2)}(t_n) + h^2 R_2(\xi_n) + h^k b^{(k+1)}(t_n)/(k+1) - h^{k+1}R_{21}(\xi_n) + h^k Q_2(t_n) + h^{k+1}W_2(g_n, p_n) + h^{k+2}g_2(g_n, p_n) \right) / \alpha_0, \quad (\text{A.10})$$

$$\psi_3(t_n) = -h^k W_3(p_n) - h^{k+1}g_3(p_n), \quad (\text{A.11})$$

and  $\alpha^* = \frac{(-1)^k}{(k+2)!} \sum_{i=1}^k i^{k+2} \alpha_i$ .

The functions  $Q_i(t)$  ( $i = 1, 2, 3$ ) are functions of  $c(t)$ ,  $b(t)$ ,  $e(t)$  and partial derivatives of  $F$ ,  $G$ , and  $H$ , and are not dependent on  $g$ ,  $p$ , or  $s$ . Since  $(c(t), b(t), e(t))$  is the solution to the DAE system (2.11)-(2.13), and since the given functions  $F$ ,  $G$ , and  $H$  may be assumed to be arbitrarily smooth, it follows that there exists a constant  $Q$  such that  $\|Q_i(t)\| \leq Q$  ( $i = 1, 2, 3$ ) for all  $t \in [t_{k-1}, t_k + T]$ . For  $k = 1$  the proof of convergence requires not only this boundedness property, but also the smoothness of the  $O(h)$  terms in  $\psi_2(t)$  and the  $O(1)$  term in  $\psi_3(t)$ , namely  $Q_3(t)$ . These conditions are valid providing the given functions  $F$ ,  $G$ , and  $H$  are sufficiently smooth. Arguments of all their arguments.

The functions  $W_i$  ( $i = 1, 2, 3$ ) are linear functions of their arguments, and for our purposes may be considered as bounded matrix operators. The functions  $R_1(\xi_n)$ ,  $R_2(\xi_n)$ ,  $R_{12}(\xi_n)$ , and  $R_{21}(\xi_n)$  are bounded functions of  $\xi_n$  (not necessarily the same  $\xi_n$ ) where  $t_{n-k} \leq \xi_n \leq t_n$ . They are the remainders from the  $k$ -step BDF formula - i.e., dependent on  $J^{(k+3)}(\xi_n)$ ,  $w^{(k+2)}(\xi_n)$ , and  $\beta^{(k+2)}(\xi_n)$ , respectively.

All of the remaining terms in the Taylor series expansion were contained in the  $g_i$  ( $i = 1, 2, 3$ ) functions, with the leading terms being quadratic functions of the arguments. For example, the leading term in  $g_3(p)$  is  $p^T H_{\omega\omega}(\ast)p/2$ . All remaining terms in  $g_3(p)$  are at least of order  $O(h^4)$ . Similarly, the leading terms in  $g_2(g, p)$  are quadratic functions of  $g$  and  $p$ , and all remaining terms are  $O(h^4)$  at least. These statements are true providing  $g$  and  $p$  are bounded (which will always be the case). The behavior of  $g_1(g, p, s)$  depends on  $s$ , since  $s_n$  is not bounded until  $n \geq 4k-1$ . For  $k \geq 2$ , the leading terms in  $g_1(g, p, s)$  are bounded, quadratic functions of  $g$ ,  $p$ , and  $s$ , and any remaining terms in  $g_1(g, p, s)$  will contribute at worst an  $O(h)$  amount to  $\psi_1(t)$  in the case  $h^2 s$  is bounded, but  $hs$  or  $s$  is not. For  $k = 1$ , we will need to bound  $\psi_1(t)$  to  $O(h)$  when either  $\|hs\|$  or  $\|s\|$  is bounded. Since

$$\|h^2 g_1(g, p, s)\| = \left\| \frac{h}{2} \left[ (hs)^T E_{\omega\omega}(\ast) (hs)^T \right] + O(h^2) \right\|$$

when  $\|hs\|$  is bounded, and  $\|h^2 g_1(g, p, s)\| = O(h^2)$  when  $\|s\|$  is bounded, it follows that  $\psi_1(t)$  is  $O(h)$ . In general, we will need to consider only the leading quadratic terms in each  $g_i$  ( $i = 1, 2, 3$ ) for the analysis which follows.

Let us denote the elements of the  $3 \times 3$  matrix  $S^{-1}(t_n)$  by  $X_{ij}(t_n)$  ( $i, j = 1, 2, 3$ ), where the elements are defined below ( $t$ -dependence is suppressed):

$$\Gamma = (A_{32}A_{31}A_{13})^{-1}, \quad (\text{A.12})$$

$$P_{11} = I_p - A_{13}\Gamma A_{32}A_{31}, \quad (\text{A.13})$$

$$P_{22} = I_q - A_{21}A_{13}\Gamma A_{32}, \quad (\text{A.14})$$

$$X_{11} = P_{11} + \frac{h}{\alpha_0} (A_{11} - A_{13}\Gamma A_{32}A_{32}A_{31} + A_{21}A_{11}) P_{11} + O(h^2), \quad (\text{A.15})$$

$$X_{12} = \frac{-\alpha_0}{h} A_{13}\Gamma A_{32} - A_{13}\Gamma A_{32}A_{32}P_{22} - P_{11}A_{11}A_{13}\Gamma A_{32} + O(h), \quad (\text{A.16})$$

$$X_{13} = \frac{\alpha_0}{h} A_{13}\Gamma + P_{11}A_{11}A_{13}\Gamma - A_{13}\Gamma A_{32}A_{32}A_{21}A_{13}\Gamma + O(h), \quad (\text{A.17})$$

$$X_{21} = \frac{h}{\alpha_0} P_{22}A_{32}A_{31} + \frac{h^2}{\alpha_0^2} P_{22}(A_{32}A_{31} + A_{21}A_{11})P_{11} + O(h^2), \quad (\text{A.18})$$

$$X_{22} = P_{22} + \frac{h}{\alpha_0} P_{22}(A_{32}P_{22} + A_{21}A_{11}A_{13}\Gamma A_{32}) + O(h^2), \quad (\text{A.19})$$

$$X_{23} = A_{21}A_{13}\Gamma + \frac{h}{\alpha_0} P_{22}(A_{32}A_{31} + A_{21}A_{11})A_{13}\Gamma + O(h^2), \quad (\text{A.20})$$

$$X_{31} = -hc^* \Gamma A_{32}A_{31} - h^2 \Gamma A_{32}(A_{32}A_{31} + A_{21}A_{11})P_{11} + O(h^2), \quad (\text{A.21})$$

$$X_{32} = -\alpha_0^2 \Gamma A_{32} + hc^* \Gamma A_{32}(A_{32}A_{31} + A_{21}A_{11})A_{13}\Gamma A_{32} - \alpha_0 \Gamma A_{32}A_{32} + O(h^2), \quad (\text{A.22})$$

$$X_{33} = \alpha_0^2 \Gamma - hc^* \Gamma A_{32}(A_{32}A_{31} + A_{21}A_{11})A_{13}\Gamma + O(h^2). \quad (\text{A.23})$$



We will show these conditions hold for  $\delta = O(h)$  and  $n \geq 4k - 1$ . If condition (A.42) holds, it is straightforward to show conditions (A.43) and (A.44) hold by utilizing the information known about  $X_{ij}$  ( $i, j = 1, 2, 3$ ),  $W_i$  and  $g_i$  ( $i = 1, 2, 3$ ). Note that  $\|Z_1\| = O(h^{k-1})$ ,  $\|Z_2\| = O(h^k)$ , and  $\|Z_3\| = O(h^{k-1})$ , so for  $k \geq 2$  we can select  $\delta = O(h)$ . Here we have used the facts that  $h_{k,n}$  is involved only in the  $W_1$  and  $g_1$  terms, that  $W_1$  is linear, and that  $g_1$  is quadratic in  $h_n$ , to  $O(h)$  accuracy. To see that condition (A.44) holds for bounded  $g, p$ , and  $h_n$ , one need only observe that taking the partial derivatives of  $Z$  with respect to  $(g_n, p_n, h_n)$  does not alter the powers of  $h$ .

The proof that the starting iterates are bounded as in (A.42) involves straightforward algebraic manipulations. The details in the algebra are omitted here. However, the essential steps are given below. The proof for  $k = 2$  requires more exploitation of the structure of the block matrices than when  $k \geq 3$ . In both cases, the proof relies on some projection matrices to annihilate the effects of  $1/h$  factors, as is demonstrated repeatedly in certain matrix product relations.

Since equations (A.25) and (A.26) hold when  $n$  is replaced by  $i$  and  $i = k, k + 1, \dots, n - 1$ , we can rewrite equations (A.39)-(A.40) as

$$r_n^{(c)} = \left( \prod_{i=k}^{n-1} F_i \right) r_{k-1} + \sum_{i=k}^{n-1} \left( \prod_{j=i+1}^{n-1} F_j \right) f_{i,i} + \tilde{f}_{1,n}, \quad (\text{A.45})$$

$$h_n^{(c)} = \sum_{\ell=0}^{k-1} \frac{1}{h} H^\ell G_{n-\ell} \left( \prod_{i=k}^{n-\ell-1} F_i \right) r_{k-1} + \sum_{i=k}^{n-\ell-1} \left( \prod_{j=i+1}^{n-\ell-1} F_j \right) f_{i,i} + f_{i,n-\ell-1} \\ + \frac{1}{h} \sum_{\ell=1}^{k-1} H^\ell f_{j,n-\ell} + \frac{1}{h} \tilde{f}_{j,n}. \quad (\text{A.46})$$

First we bound  $r_n^{(c)}$ , starting with the first term in (A.45) involving  $r_{k-1}$ . Recall that  $F_i$  is a  $lk \times lk$  block companion matrix, with blocks  $X_i$  containing a subblock  $X_{12}(t_i)$  of order  $O(1/h)$ . Since  $n \geq 4k - 1$ , the product  $\prod_{i=k}^{n-1} F_i$  involves  $2k$  or more factors of  $F_i$ , so all blocks are sums of terms with at least two factors of  $X_i$  ( $i = k, \dots, n$ ). We will need the following lemma proven in [1]:

**Lemma 1.** *There exist constants  $K^*$  ( $K^* \geq 1$ ) and  $\mathcal{E}$  such that*

$$\left\| \prod_{j=i+1}^n F_j \right\| \leq K^* e^{(n-i)\mathcal{E}}$$

for all  $n - i \geq 2k$  where  $F_j$  is defined in (A.27).

The proof of this lemma is technical but straightforward, involving the application of a stability theorem for matrix products [16], [26] as well as a theorem verifying the existence of smooth similarity transformations [13]. One trick in the proof involves partitioning the product into groups of  $2k$  factors of  $F_j$  because such groups can be shown to be bounded while individual factors  $F_j$  are not. The eigenvalues and eigenvectors of a group must be analyzed in detail. The details of this analysis can be found in [1] (see Lemma 6.14). Therefore, the product  $\prod_{i=k}^n F_i$  in (A.45) is

$$Z_2(g, p, h_n) = \frac{h^{k+1}}{\alpha_0} X_{31} [W_1(g, p, \delta) + h g_1(g, p, \delta) + h g_2(g, p) + h g_3(g, p)] \\ - h^k X_{32} [W_3(p) + h g_3(p)],$$

$$Z_3(g, p, h_n) = \frac{h^k}{\alpha_0} X_{31} [W_1(g, p, \delta) + h g_1(g, p, \delta)] + \frac{h^k}{\alpha_0} X_{32} [W_2(g, p) + h g_2(g, p)] \\ - h^{k-1} X_{33} [W_3(p) + h g_3(p)].$$

For starting approximations select

$$r_n^{(c)} = F_n r_{n-1} + \tilde{f}_{1,n}, \quad (\text{A.39})$$

$$h_n^{(c)} = \frac{1}{h} G_n r_{n-1} + \frac{1}{h} \tilde{f}_{j,n} + h E r_{n-1}, \quad (\text{A.40})$$

where

$$\tilde{f}_{1,n} = \begin{bmatrix} (\tilde{f}_{1,n}^{(1)})^T, (\tilde{f}_{1,n}^{(2)})^T, 0, \dots, 0 \end{bmatrix}^T, \\ \tilde{f}_{1,n}^{(i)} = \begin{bmatrix} X_{11} \tilde{\psi}_1 + X_{12} \tilde{\psi}_2 + X_{13} \tilde{\psi}_3 \\ X_{21} \tilde{\psi}_1 + X_{22} \tilde{\psi}_2 + X_{23} \tilde{\psi}_3 \end{bmatrix} \Big|_{h_n}, \quad i = 1, 2, \\ \tilde{f}_{j,n} = \begin{bmatrix} (X_{31} \tilde{\psi}_1 + X_{32} \tilde{\psi}_2 + X_{33} \tilde{\psi}_3)^T, 0, \dots, 0 \end{bmatrix}^T, \\ \tilde{\psi}_1(t_n) = \begin{bmatrix} -h R_1(\xi_n) + h^k e^{h(k+1)(t_n)} / (k+1) - h^{k+1} R_{12}(\xi_n) + h^k Q_1(t_n) \end{bmatrix} / \alpha_0, \\ \tilde{\psi}_2(t_n) = \begin{bmatrix} -h \alpha^k \omega^{h(k+2)}(t_n) + h^k \theta^{h(k+1)}(t_n) / (k+1) + h^k R_2(\xi_n) \\ -h^{k+1} R_{22}(\xi_n) + h^k Q_2(t_n) \end{bmatrix} / \alpha_0, \\ \tilde{\psi}_3(t_n) = -h^{k-1} Q_3(t_n).$$

Since  $\tilde{\psi}_i(t_n)$  ( $i = 1, 2, 3$ ) are independent of  $(g_n, p_n, \delta_n)$ , it follows that

$$\|\tilde{\psi}_1(t)\|, \|\tilde{\psi}_2(t)\|, \|\tilde{\psi}_3(t)\| \leq \rho h \quad (\text{A.41})$$

for some constant  $\rho$  and  $t \in [h, t_0 + T]$ .

It is important to note we are iterating for  $h_n$ , not  $h^2 \delta_n$  or  $\delta_n$ . For the implicit function argument, the following three conditions must be satisfied for sufficiently large  $n$  (i.e.,  $n \geq 4k - 1$ ):

I.

$$\|g_n^{(c)}\| < \eta_1 e^{(n-3k+1)\mathcal{E}L} \\ \|p_n^{(c)}\| < \eta_2 e^{(n-3k+1)\mathcal{E}L} \\ \|h_n^{(c)}\| < \eta_3 e^{(n-3k+1)\mathcal{E}L} \quad (\text{A.42})$$

II.

$$\|Z(g_n^{(c)}, p_n^{(c)}, h_n^{(c)})\| \leq \delta/2, \quad \delta > 0, \quad Z = (Z_1^T, Z_2^T, Z_3^T)^T \quad (\text{A.43})$$

III.

$$\|J\| \leq 1/2 \text{ for any } g, p, \text{ and } h_n \text{ such that } \|g - g_n^{(c)}\| \leq \delta, \|p_n - p_n^{(c)}\| \leq \delta, \text{ and} \\ \|h_n - h_n^{(c)}\| \leq \delta \text{ where } J \text{ is the Jacobian matrix of } Z \text{ with respect to } g, p, \text{ and } h_n. \quad (\text{A.44})$$

$O(1)$ ), so we can summarize this result for  $k \geq 2$  as  $\Lambda \leq \Delta_1^*$  for some constant  $\Delta_1^*$  independent of  $\eta, \eta_2, \eta_3$  and  $L$  to  $O(h)$  accuracy.

Next, we bound the terms in (A.48) involving  $\psi_2(t_i)$ . The case when  $k \geq 3$  is fairly straightforward, requiring the use of (A.36) and some information concerning the structure of the matrices (i.e., the location of  $O(1/h)$  elements). If  $k = 2$ , the terms must be broken into two parts:

$$\sum_{i=k}^{n-2k} \dots + \sum_{i=n-3k+1}^{n-2k} \dots$$

The finite number of terms is easily bounded by noting any  $O(1/h)$  terms in  $\prod_{j=i+1}^n F_j$  will multiply the  $O(h)$  term  $X_{23}(t_i)\psi_2(t_i)$ . The remaining terms ( $i = k, \dots, n-2k$ ) involve the product  $\prod_{j=i+1}^n F_j$  containing  $2k$  or more factors of  $F_j$ . Therefore, these terms will contain products involving two or more factors of  $X_j$ , such as

$$X_{i+j} X_{i+j} \begin{bmatrix} X_{13i} \\ X_{23i} \end{bmatrix} = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix},$$

where  $A_1$  is  $O(1)$  and  $A_2$  is  $O(h)$  when  $\ell \neq j$ . This last fact follows by direct computation using relations (A.12)-(A.23). Similar relations hold for  $(\prod_{j \in \sigma} X_j) \begin{bmatrix} X_{13i} \\ X_{23i} \end{bmatrix}$ , where  $\sigma$  contains two or more indices. Then, as for  $k \geq 3$ , these terms can be bounded, and in summary, there exists a  $\Delta_2^*$  such that for  $k \geq 2$  and  $n \geq 4k - 1$

$$\left\| \prod_{j=i+1}^{n-1} F_j \begin{bmatrix} X_{13}(t_i)\psi_2(t_i) \\ X_{23}(t_i)\psi_2(t_i) \end{bmatrix}^T, 0, \dots, 0 \right\|^T \leq \Delta_2^*,$$

where  $\Delta_2^*$  is independent of  $\eta, \eta_2, \eta_3$  and  $L$  to  $O(h)$  accuracy.

The remaining terms  $\tilde{f}_{i,n}$  in  $r_n^{(c)}$  can be bounded easily using the relations (A.41):

$$\|\tilde{f}_{i,n}\| \leq K \max(\rho h + 2\rho, O(h)) \leq \Delta_3^*$$

for some constant  $\Delta_3^*$ . Finally, for  $n \geq 4k - 1$  we have shown

$$\|r_n^{(c)}\| \leq \eta_0^c K^c e^{(n-3k+1)\mathcal{L}} + \Delta_1^* + \Delta_2^* + \Delta_3^*.$$

Since this bound is independent of  $\eta, \eta_2, \eta_3$  and  $L$ , at least to  $O(h)$  accuracy, there exists a constant  $\eta$  and  $L = \mathcal{L}$  such that

$$\|r_n^{(c)}\| < \eta e^{(n-3k+1)\mathcal{L}} \quad \text{for } n \geq 4k - 1.$$

Let  $\eta_1 = \eta$  and  $\eta_2 = \eta$ . Then, for  $n \geq 4k - 1$

$$\|d_n^{(c)}\| \leq \|r_n^{(c)}\| < \eta_1 e^{(n-3k+1)\mathcal{L}}$$

and

$$\|p_n^{(c)}\| \leq \|r_n^{(c)}\| < \eta_2 e^{(n-3k+1)\mathcal{L}}.$$

bounded, but the starting vector  $r_{k-1}$  contains  $O(1/h)$  terms  $g_j, j = 0, 1, \dots, k-1$ . Recall the starting values were chosen to satisfy (A.6),(A.7). Now using the structure of the product  $\prod_{i=k}^n F_i$ , one can show the  $g_j$  terms always appear as  $F_{11}(t_i)g_j$  which are  $O(1)$ . All the terms involving  $F_j$  are also clearly  $O(1)$ . Therefore, there exist constants  $\eta_0^{(j)}, j = 1, 2, \dots, 2k$ , such that

$$\left\| \left( \prod_{i=k}^{2k+\ell-1} F_i \right) r_{k-1} \right\| \leq \eta_0^{(j)}. \tag{A.47}$$

For  $n \geq 5k - 1$ , it follows that

$$\left\| \left( \prod_{i=k}^n F_i \right) r_{k-1} \right\| \leq K^* e^{(n-3k+1)\mathcal{L}} \eta_0^{(j)},$$

since  $\|\prod_{i=k}^{n-3k} F_i\| \leq K^* e^{(n-3k+1)\mathcal{L}}$  from Lemma 1. We summarize this bound as

$$\left\| \left( \prod_{i=k}^n F_i \right) r_{k-1} \right\| \leq \eta_0^* K^* e^{(n-3k+1)\mathcal{L}}$$

for  $n \geq 3k - 1$  and  $\eta_0^* = \max_{i=1, \dots, 2k} \eta_0^{(i)}$ . Next, in the bounding of  $r_n^{(c)}$ , we break the terms  $\sum_{i=k}^{n-1} (\prod_{j=i+1}^n F_j) f_{i,1}$  into two parts:

$$\sum_{i=k}^{n-1} \left( \prod_{j=i+1}^n F_j \right) \begin{bmatrix} R_i [\psi_1(t_i)^T, \psi_2(t_i)^T, 0, \dots, 0]^T + [X_{13}(t_i)\psi_2(t_i)]^T, (X_{23}(t_i)\psi_2(t_i))^T, 0, \dots, 0 \end{bmatrix}^T, \tag{A.48}$$

where  $R_i = \text{diag}(X_{11}, 0, \dots, 0)$ . In bounding both parts, it will be necessary to consider the case  $k = 2$  separately from  $k \geq 3$ . Consider the terms involving  $\psi_1$  and  $\psi_2$ . It is possible to show that there exists a constant  $N_1$  such that (see Lemma 6.15 in [1])

$$\left\| \left( \prod_{j=i+1}^n F_j \right) R_i \right\| \leq N_1 K^* e^{n\mathcal{L}} \quad \text{for } n - i \geq k.$$

For  $n - i < k$ , any  $O(1/h)$  terms in  $(\prod_{j=i+1}^n F_j) R_i$  will multiply the  $\psi_2(t_i)$  term, which is always  $O(h)$  by (A.37). Thus, there exists a constant  $\Delta_1$  independent of  $\eta, \eta_2, \eta_3$  or  $L$  such that

$$\left\| \left( \prod_{j=i+1}^{n-1} F_j \right) R_i [\psi_1(t_i)^T, \psi_2(t_i)^T, 0, \dots, 0]^T \right\| \leq \Delta_1$$

for  $i = n - k + 1, \dots, n - 1$ . Then, using relations (A.35)-(A.38) it follows that

$$\begin{aligned} \Lambda &= \left\| \sum_{i=k}^{n-1} \left( \prod_{j=i+1}^n F_j \right) R_i [\psi_1(t_i)^T, \psi_2(t_i)^T, 0, \dots, 0]^T \right\| \\ &\leq [(2k-1)h + (n-4k+2)h] N_1 K^* e^{n\mathcal{L}} \max(\sigma_{11}, \sigma_{12}, \sigma_2) + (k-1)\Delta_1 \end{aligned}$$

providing  $n \geq 4k - 1$  and  $k \geq 3$ . If  $k = 2$ , the terms  $\sum_{i=k}^{2k-2} \dots$  must be bounded more carefully since  $\|\psi_1(t_i)\| \leq \sigma_{11}$  for  $i = k, k+1, \dots, 3k-2$ . Still, there are a finite number of such terms (each

for some constant  $\Delta_3^*$  independent of  $\eta_1, \eta_2, \eta_3$  and  $L$  to  $O(h)$  accuracy.

Next, using relations (A.24), (A.35)-(A.38) and (A.41), we have

$$\left\| \sum_{\ell=1}^{k-1} \frac{1}{h} H^\ell f_{n-\ell} + \frac{1}{h} f_{n,0} \right\| \leq \Delta_3^*$$

for some constant  $\Delta_3^*$  and  $k \geq 2$ . Finally, it is possible to prove that there exists a constant  $\Delta_3^*$  such that

$$\left\| \sum_{\ell=0}^{k-1} \frac{1}{h} H^\ell G_{n-\ell} f_{1,n-\ell-1} \right\| \leq \Delta_3^* \quad \text{for } k \geq 2.$$

This bound can be obtained by examining matrix products such as  $X_{31}, X_{13}$  and  $X_{32}, X_{23}$ , and by applying the bounds (A.35)-(A.38). To summarise the final bound for  $h_{n,0}^{(o)}$ , we have

$$\|h_{n,0}^{(o)}\| \leq \sum_{i=4}^6 \Delta_i^* \quad \text{for } n \geq 4k-1,$$

where the constants  $\Delta_i^*$  are independent of  $\eta_1, \eta_2, \eta_3$  and  $L$  to  $O(h)$  accuracy. Hence, we can select a constant  $\eta_3$  such that for all  $n \geq 4k-1$

$$\|h_{n,0}^{(o)}\| \leq \|h_{n,0}^{(o)}\| \leq \sum_{i=4}^6 \Delta_i^* < \eta_3 e^{(n-3k+1)h\mathcal{L}}.$$

We have proven there exist constants  $\eta_1, \eta_2, \eta_3$  and  $\mathcal{L}$  such that the starting iterates are bounded uniformly providing  $n \geq 4k-1$ , as in conditions (A.42). It follows immediately that conditions (A.43) and (A.44) are satisfied. Hence, the fixed point iteration converges to a solution  $(\hat{r}_n, \hat{p}_n, \hat{h}_{n,0})$  of equations (A.25) and (A.26). In addition,

$$\begin{aligned} \|\hat{q}_n - q_n^{(o)}\| &= \|Z_1(\hat{q}_n, \hat{p}_n, \hat{h}_{n,0})\| \leq \delta, \\ \|\hat{p}_n - p_n^{(o)}\| &= \|Z_2(\hat{q}_n, \hat{p}_n, \hat{h}_{n,0})\| \leq \delta, \\ \|\hat{h}_{n,0} - h_{n,0}^{(o)}\| &= \|Z_3(\hat{q}_n, \hat{p}_n, \hat{h}_{n,0})\| \leq \delta, \end{aligned}$$

where  $\delta = O(h)$  for  $k \geq 2$ . Therefore, for sufficiently small  $h$ , we have for  $n \geq 4k-1$

$$\begin{aligned} \|q_n\| &\leq \|q_n^{(o)}\| + \delta \leq \eta_1 e^{(n-3k+1)h\mathcal{L}}, \\ \|p_n\| &\leq \|p_n^{(o)}\| + \delta \leq \eta_2 e^{(n-3k+1)h\mathcal{L}}, \\ \|h_{n,0}\| &\leq \|h_{n,0}^{(o)}\| + \delta \leq \eta_3 e^{(n-3k+1)h\mathcal{L}}. \end{aligned}$$

To complete the induction argument, we must still verify the starting assumptions (A.28) - (A.30). However, presuming for the moment that they do hold, then  $\hat{q}_n, \hat{p}_n$  and  $\hat{h}_{n,0}$  are uniformly bounded for  $n \geq 4k-1$ . Hence, the  $k$ -step BDF method produces a numerical solution which converges to  $O(h^k)$  accuracy to a solution of system (1.8)-(1.10) for  $n \geq 4k-1$ .

*Verification of Induction Assumption for  $i \leq 4k-2$ .*

The proof of convergence of the  $k$ -step BDF method depends on the induction assumption, and

The procedure for bounding  $h_{n,0}^{(o)}$ , defined in (A.46), is very similar to the previous analysis for  $r_n^{(o)}$ , although additional properties of matrix products are utilized [3]. We will give only a summary of the intermediate results required for the estimate. The terms involving the initial remainder  $r_{k-1}$  can be bounded with the use of (A.47) and the following properties:

1.  $\|G_{n-\ell} \prod_{i=2k}^{n-\ell-1} F_i\| = O(h)$  for all  $n-\ell \geq 5k$
2.  $\|G_{n-\ell} \left( \prod_{i=k}^{n-\ell-1} F_i \right) r_{k-1}\| = O(h)$  for  $n-\ell \geq 3k$ .

Properties (1) and (2) can be easily proven by noting the following fact:

$$\|X_{31}, X_{32}\|_{k+i, k} X_{1+i, j} X_i = [A_1, A_2],$$

where  $A_1$  is  $O(h^2)$  and  $A_2$  is  $O(h)$  for  $s \neq j$ . This fact is proven by direct computation of the matrix products using relations (A.12)-(A.23) and Taylor series expansion about a common  $t$  value. It follows then that

$$\left\| \sum_{\ell=0}^{k-1} \frac{1}{h} H^\ell G_{n-\ell} \left( \prod_{i=k}^{n-\ell-1} F_i \right) r_{k-1} \right\| \leq \Delta_4^*$$

for  $n \geq 4k-1$  and for some constant  $\Delta_4^*$ .

As we did earlier for  $r_n^{(o)}$ , the terms involving  $f_{1,i}$  are split into two pieces, and the cases  $k=2$  and  $k \geq 3$  are done separately. Using the structure of the matrices and relations (A.35) (A.38), it can be shown that the first part is bounded for  $k \geq 2$  as

$$\left\| \sum_{\ell=0}^{k-1} \frac{1}{h} H^\ell \sum_{i=k}^{n-\ell-2} G_{n-\ell} \left( \prod_{j=i+1}^{n-\ell-1} F_j \right) R_i \left[ \psi_1(t_i)^T, \psi_2(t_i)^T, 0, \dots, 0 \right]^T \right\| \leq \Delta_5^*,$$

where  $\Delta_5^*$  is independent of  $\eta_1, \eta_2, \eta_3$  and  $L$  to  $O(h)$  accuracy. Bounding the second part involving  $\psi_3(t_i)$  is straightforward when  $k \geq 3$ , but when  $k=2$  we will need the fact that there exists a constant  $N$  such that

$$\|H^\ell G_{n-\ell} \prod_{j=i+1}^{n-\ell-1} F_j\| \leq N h^2$$

for  $n-\ell-i-1 \geq 2k$  (see Lemma 6.12 in [1]). Then, we split this part further,

$$\sum_{i=k}^{n-\ell-2k-1} \dots \text{and } \sum_{i=n-\ell-2k}^{n-\ell-2} \dots$$

where the first sum is bounded as

$$\left\| \sum_{\ell=0}^{k-1} \frac{1}{h} H^\ell \sum_{i=k}^{n-\ell-2k-1} G_{n-\ell} \prod_{j=i+1}^{n-\ell-1} F_j \left[ X_{13}(t_i)^T \psi_3(t_i)^T, (X_{32}(t_i) \psi_3(t_i))^T, 0, \dots, 0 \right]^T \right\| \leq k(n-3k)hN K \sigma_3 h^{k-2}.$$

Each term in the finite sum  $\sum_{i=n-\ell-2k}^{n-\ell-2}$  involves products of  $X_{ij}$  which can be bounded to  $O(1)$  by direct computation. Hence, for  $k=2$  as well as for  $k \geq 3$ , we have

$$\left\| \sum_{\ell=0}^{k-1} \frac{1}{h} H^\ell G_{n-\ell} \sum_{i=k}^{n-\ell-2} \left( \prod_{j=i+1}^{n-\ell-1} F_j \right) \left[ X_{13}(t_i)^T, X_{32}(t_i)^T, 0, \dots, 0 \right]^T \psi_3(t_i) \right\| \leq \Delta_5^*$$

in particular on the behavior of  $g_i$ ,  $p_i$  and  $h^2 s_i$  during the first steps. The proof that the induction assumptions (A.28) and (A.29) hold is very similar to the proof just given above. That is, it is necessary to construct a fixed point iteration using equations (A.25) and (A.26), but now we iterate for  $g_i$ ,  $p_i$  and  $h^2 s_i$ , and not for  $h s_i$ .

The starting iterates selected in (A.39) (A.40) may be used except  $g_i^{(0)}$  must be altered slightly to include the  $O(1)$  term  $(h^2 g_i^{(0)})^\top F_{uu}(\bullet)(h^2 s_i^{(0)})/2$ , at least when  $k = 2$ . This term is added because the behavior of the  $h^{k+2} g_i(q, p, s)$  term in  $\psi_i(t)$  will have order  $O(h^{k-2})$  when  $h^2 s$  is bounded, but  $h s$  is not.

The assumptions (A.28) and (A.29) are first shown to hold for  $i = k$ , and then the result follows for later steps ( $i = k + 1, \dots$ ) in the same fashion. In bounding the initial iterates for the fixed point iteration, we will use the same properties of the starting values and projection matrices used in the above induction argument. In fact, the consistency relation (2.4) for the starting values is vital in the proof that  $g_i^{(0)}$  is bounded. Bounding the initial iterates is much simpler here because there are only finite sums involved. This argument verifies that  $g_n$ ,  $p_n$ , and  $h^2 s_n$  are bounded for any finite number of steps.

The induction assumption (A.30) must still be verified before the convergence proof is complete. From (A.28) we know  $g_n$  and  $p_n$  are bounded for  $n = k, k+1, \dots, 4k-2$ . Therefore, we can construct a fixed point iteration just for  $h s_n$ ,  $n = 3k-1, \dots, 4k-2$ , using the difference equation (A.26):

$$h s_n^{(0)} = \frac{1}{h} X_{31}(t_n) \left( \sum_{i=1}^k \gamma_i g_{n-i} + \psi_1(t_n) \right) + \frac{1}{h} X_{32}(t_n) \left( \sum_{i=1}^k \gamma_i p_{n-i} + \psi_2(t_n) \right) + \frac{1}{h} X_{33}(t_n) \psi_3(t_n),$$

where  $\psi_i(t_n) = \psi_i(t_n) -$  terms involving  $s_n$  and

$$h s_n^{(v+1)} = h s_n^{(v)} + Z(h s_n^{(v)}),$$

where

$$Z(h s_n^{(v)}) = \frac{1}{h} X_{31}(t_n) \left( \frac{h^{k+1}}{c_0} \bar{W}_1(s_n^{(v)}) + \frac{h^{k+2}}{c_0} \bar{g}_1(g_n, p_n, s_n^{(v)}) \right),$$

and where  $\bar{W}_1$  and  $\bar{g}_1$  are the terms in  $W_1$  and  $g_1$  involving  $s_n$ , respectively. The proof that this fixed point iteration converges is straightforward. The key to bounding  $h s_n^{(v)}$  lies in the fact that since  $n \geq 3k-1$ , the terms involving  $g_{n-i}$  ( $i = 1, 2, \dots, k$ ) do not reference any of the starting  $O(1/h)$  values of  $g$ . The other important step is noting that  $A_{32}(t) p_2 = \psi_2(t)$  for  $j = k, k+1, \dots$  and  $X_{31}(t_n) = -\alpha_2^2 \Gamma(t_n) A_{32}(t_n) + O(h)$ . Therefore, since  $\|\psi_2(t)\| = O(h^{k-1})$ , it follows that  $\|\sum_{i=1}^k \gamma_i X_{32}(t_n) p_{n-i}/h\| = O(h^{k-2})$ . The induction assumption (A.30) is therefore true when  $n = 3k-1$ , and it follows when  $n = 3k, \dots, 4k-2$  in a similar fashion. Now that the induction assumption has been verified for the first few steps, the induction and fixed point argument presented earlier for  $n \geq 4k-1$  is complete.

**Remark.** We have shown the numerical solution  $u_n = (u_n, u_{n-1}, \dots, u_{n-k+1})$  is globally  $O(h^k)$  accurate for all  $n \geq 3k-1$ . Equivalently, after  $k+1$  steps from  $t_{k-1}$ , the BDF method produces a

numerical solution  $u_{3k}$  at  $t_{3k}$  which has accuracy  $O(h^k)$ . The numerical solution expressed in long vector notation first has  $O(h^k)$  accuracy in all components when  $n = 3k-1$ .

*Convergence Proof for  $k = 1$  (Theorem 2).*

While the convergence proof for  $k = 1$  is very similar to the one for  $k \geq 2$  just presented and thus will not be given in detail, there are some additional difficulties that deserve to be noted. In some respect, these difficulties tend to make the convergence proof for the Backward Euler formula even more interesting. The primary difficulty starts with the fact that  $\|h^2 s_n\|$  may not be bounded, so  $\psi_1(t_1)$  may not be bounded as  $h \rightarrow 0$ . Therefore, the proof must be modified to avoid the use of the induction assumption (A.29). Instead, the proof can be altered to use the facts that  $\|X_{11}(t_1) \psi_1(t_1)\| \leq \Delta_{11}$  and  $\|X_{31}(t_1) \psi_1(t_1)\| \leq \Delta_{31} h$  for some constants  $\Delta_{11}$  and  $\Delta_{31}$  which depend only on  $\eta_0$  and  $\eta_1$  to  $O(h)$  accuracy. The fixed point argument must also be modified to iterate for  $(p_n, g_n, s_n)$  instead of  $(g_n, p_n, h s_n)$ . As a result, the smoothness of the leading error terms follows immediately from the convergence proof for  $k = 1$  and  $n \geq 3$ . The starting iterates for  $g_n^{(0)}$  and  $s_n^{(0)}$  in the fixed point iteration must be altered to include some terms which are now  $O(1)$ :

$$g_n^{(0)} = X_{11}(t_n) (g_{n-1} + \bar{\psi}_1(t_n)) + X_{12}(t_n) (p_{n-1} + \bar{\psi}_2(t_n)) + X_{13}(t_n) (-Q_3(t_n) - h W_3(p_n^{(0)})),$$

$$s_n^{(0)} = \frac{1}{h^2} \left[ X_{31}(t_n) (g_{n-1} + \bar{\psi}_1(t_n)) + X_{32}(t_n) (p_{n-1} + \bar{\psi}_2(t_n)) + h^2 W_3(g_n^{(0)}, p_n^{(0)}) \right] + X_{33}(t_n) (-Q_3(t_n) - h W_3(p_n^{(0)}) - h^2 g_3(p_n^{(0)})).$$

Note that the initial iterates must be chosen in the order  $p_n^{(0)}, g_n^{(0)}$ , and then  $s_n^{(0)}$ , and the definitions of  $Z_1, Z_2, Z_3$  must be modified accordingly. As before, the main difficulty in proving the fixed point iteration converges is to show the initial iterates are uniformly bounded. This proof requires numerous cancellation properties of matrix products and sums of the  $X_{ij}$  matrices, most of which can be verified by straightforward (although tedious) algebraic manipulation and expansion by Taylor series.

As in the proof for  $k \geq 2$ , the expressions for the initial iterates are rewritten as follows:

$$\begin{bmatrix} g_n^{(0)} \\ p_n^{(0)} \end{bmatrix} = \prod_{i=1}^n X_i \begin{bmatrix} g_0 \\ p_0 \end{bmatrix} + \sum_{i=1}^{n-1} \left( \prod_{j=i+1}^n X_j \right) \begin{bmatrix} \psi_1(t_i) \\ \psi_2(t_i) \end{bmatrix} + \prod_{j=i+1}^n X_j \begin{bmatrix} X_{12}(t_i) \\ X_{22}(t_i) \end{bmatrix} \psi_3(t_i) \\ + X_n \begin{bmatrix} \bar{\psi}_1(t_n) \\ \bar{\psi}_2(t_n) \end{bmatrix} - \begin{bmatrix} X_{12}(t_n) \\ X_{22}(t_n) \end{bmatrix} Q_3(t_n) - \begin{bmatrix} 0 \\ h X_{12}(t_n) W_3(p_n^{(0)}) \end{bmatrix},$$

$$s_n^{(0)} = \frac{1}{h^2} \left[ X_{31}(t_n), X_{32}(t_n) \right] \left( \prod_{i=1}^{n-1} X_i \begin{bmatrix} g_0 \\ p_0 \end{bmatrix} + \sum_{i=1}^{n-1} \prod_{j=i+1}^{n-1} X_j \begin{bmatrix} \psi_1(t_i) \\ \psi_2(t_i) \end{bmatrix} \right) \\ + \sum_{i=1}^{n-1} \prod_{j=i+1}^{n-1} X_j \begin{bmatrix} X_{12}(t_i) \\ X_{22}(t_i) \end{bmatrix} \psi_3(t_i) \\ + \begin{bmatrix} X_{12}(t_{n-1}) \\ X_{22}(t_{n-1}) \end{bmatrix} \psi_3(t_{n-1}) + \begin{bmatrix} \bar{\psi}_1(t_n) \\ \bar{\psi}_2(t_n) + h^2 W_3(g_n^{(0)}, p_n^{(0)}) \end{bmatrix} \\ + \frac{1}{h^2} X_{33}(t_n) (-Q_3(t_n) - h W_3(p_n^{(0)}) - h X_{33}(t_n) W_3(p_n^{(0)})) - h^2 g_3(p_n^{(0)}) \quad (\text{A.49})$$

These starting iterates can then be bounded uniformly for all  $n \geq 3$ . This bounding process requires a careful, but straightforward examination of all the terms, in particular those terms involving  $q_0$  and  $\psi_1(t_1)$ . Terms involving  $q_0$  can generally be bounded with the help of (A.7). Note that the terms with  $\psi_1(t_1)$  are actually either  $X_{11}(t_1)\psi_1(t_1)$  or  $X_{31}(t_1)\psi_1(t_1)$ , which can be bounded. In the bounding of  $q_0^{(i)}$  it is necessary to use the fact that the  $O(1/h)$  terms in

$$X_{12}(t_n) \left( Q_3(t_n) + hW_3(p_n^{(i)}) \right) - \left( X_{11}(t_n)X_{12}(t_{n-1}) + X_{12}(t_n)X_{23}(t_{n-1}) \right) \psi_3(t_{n-1})$$

cancel because  $Q_3(t)$  is smooth and

$$X_{11}(t_n)X_{12}(t_{n-1}) + X_{12}(t_n)X_{23}(t_{n-1}) = \frac{\alpha_0}{h}A_{12}(t_n)\Gamma(t_n) + O(1).$$

The bounding of  $q_0^{(i)}$  is even more complicated, requiring smoothness of  $\psi_2(t_n)$  and  $\psi_3(t_n)$  and matrix product properties such as

$$\begin{aligned} \|X_{31}(t_n), X_{32}(t_n)\|_{X_{n-1}X_{n-2}} &\leq Nh^2, \\ \|X_{31}(t_n), X_{32}(t_n)\|_{X_{n-1}X_{n-2}X_{n-3}} &\leq Nh^2 \end{aligned}$$

for some constant  $N$ , and

$$\|X_{31}(t_2), X_{32}(t_2)\|_{X_2X_1} \begin{bmatrix} q_0 \\ p_0 \end{bmatrix} \leq \beta_1 h^2$$

for some constant  $\beta_1$  dependent on  $\eta_0$ . The smoothness of  $\psi_2(t)$  can be used to obtain a cancellation of the  $O(1/h)$  terms between the  $\psi_2(t_{n-1})$  and the  $\psi_2(t_n)$  term in (A.49). In particular, there exists a constant  $\beta_2$  which depends on  $\eta_1^*$ ,  $\eta_2^*$ ,  $\eta_1$ ,  $\eta_2$  and  $\mathcal{E}$ , but not on  $\eta_0$  to  $O(h)$  accuracy such that for  $n \geq 2$

$$\| (X_{31}(t_n)X_{12}(t_{n-1}) + X_{32}(t_n)X_{21}(t_{n-1})) \psi_2(t_{n-1}) + X_{32}(t_n) (\psi_2(t_n) + h^2 W_2(q_n^{(i)}, p_n^{(i)})) \| \leq \beta_2 h^2.$$

The bounding of  $q_0^{(i)}$  will also require a cancellation of terms involving  $\psi_3(t_i)$  at  $i = n, n-1$ , and  $n-2$ :

$$\begin{aligned} &\| \frac{1}{h^2} X_{32}(t_n) \left[ -Q_3(t_n) - hW_3 \left( p_n^{(i)} - hX_{23}(t_n)W_3(p_n^{(i)}) \right) - h\eta_3(p_n^{(i)}) \right] \\ &+ \frac{1}{h^2} X_{31}(t_n), X_{32}(t_n) \left[ \begin{pmatrix} X_{12}(t_{n-1}) \\ X_{23}(t_{n-1}) \end{pmatrix} \psi_3(t_{n-1}) + X_{n-1} \begin{bmatrix} X_{13}(t_{n-2}) \\ X_{23}(t_{n-2}) \end{bmatrix} \psi_3(t_{n-2}) \right] \| < \infty. \end{aligned}$$

Matrix product properties such as

$$\|X_{31}(t_n), X_{32}(t_n)\|_{X_{n-1}X_{n-2}} \begin{bmatrix} X_{13}(t_{n-3}) \\ X_{23}(t_{n-3}) \end{bmatrix} \| = O(h^2)$$

will also be utilized. For more details of the bounding of the starting iterates, see [3]. It should be clear from the examples given above that the convergence proof for  $k = 1$  requires additional

matrix product and cancellation of terms properties in order to resolve the difficulties associated with the behavior of  $s_1$ .

The induction assumptions (A.28) and (A.30) must also be verified before the convergence argument is complete. Again a different approach must be taken than the one used when  $k \geq 2$ . In particular, the argument must not require any knowledge of the behavior of  $h^2 s_1$ , but instead must rely on the behavior of

$$\begin{aligned} X_{11}(t_i)\psi_1(t_i) &= X_{11}(t_i)(I_p - hA_{11}(t_i))q_i - hX_{11}(t_i)A_{12}(t_i)p_i - X_{11}(t_i)q_{i-1}, \\ X_{31}(t_i)\psi_1(t_i) &= X_{31}(t_i)(I_p - hA_{11}(t_i))q_i - hX_{31}(t_i)A_{12}(t_i)p_i - X_{31}(t_i)q_{i-1}. \end{aligned}$$

These difficulties can be resolved fully, allowing the starting induction assumptions (A.28) and (A.30) to be verified. Hence the Backward Euler solution does converge to a solution of the index-three system as  $n \rightarrow \infty$ ,  $h \rightarrow 0$  for  $t_n \in [t_1, t_0 + T]$ .

Note we have shown  $\|h s_n\| \leq \eta_0^*$  for  $n = 2$  and  $\|s_n\| \leq \eta_0 e^{n\alpha h}$  for all  $n \geq 3$ . Therefore, convergence of the numerical solution for the algebraic variables  $u$  is not obtained until the second Backward Euler step (i.e.,  $n = 2$ ). The fact that we can not bound either  $s_n$  or  $h s_n$  on the first step is not a deficiency in the proof; convergence is not in general obtained on the first step, and an  $O(1)$  error in the numerical solution will be apparent in the algebraic variables. This  $O(1)$  error will be present even if the initial values used to start the numerical method are consistent (i.e., contain no errors) with a true solution of the system.

While we have proven the numerical solution converges to a true solution of the system, we have not proven it converges to a solution consistent with the given initial values for the algebraic variables. Because of the  $O(1)$  error in the algebraic variables at the end of the first step, and since the convergence analysis does not depend on past values of  $u$ , but only on the current time level (i.e., on  $u_n$  at  $t_n$ ), it is possible for the numerical solution to jump to a different solution curve. This behavior has been observed in practice [2]. It will not occur if the algebraic variables appear only linearly in the system.

*Remark.* If the initial values do not satisfy (2.4), then convergence is not obtained until  $n \geq 3$ . At the end of the first step, the consistency condition (2.4) is satisfied for  $t = t_1$ . The method now requires two more steps, starting from consistent initial values, before convergence is obtained in the algebraic variables.

*Proof of Corollary 2.*

For technical simplicity, we first prove this corollary under the assumption that the difference equations are all solved exactly at each step. We know from the difference equations (A.8) that

$$A_{32}(t_n)p_n = \psi_3(t_n) \quad \text{for } n \geq k. \quad (\text{A.50})$$

Next we use the difference equations (A.8) to derive a relation for  $s_n$ . Namely, the first block

By definition of  $\psi_2(t)$  in (A.10),  $\psi_2(t_n) = h\mathfrak{P}(t_n) + O(h^2)$  where  $\mathfrak{P}(t_n)$  is a smooth function, namely  $\mathfrak{P}(t) = -\alpha^* \mathfrak{w}^{(k+2)}(t)/\alpha_c$ . Therefore,

$$\begin{aligned} \psi_2(t_n) - \sum_{i=1}^k \gamma_i \psi_2(t_{n-i}) &= h \left( \mathfrak{P}(t_n) - \sum_{i=1}^k \gamma_i \mathfrak{P}(t_{n-i}) \right) + O(h^2) \\ &= h \mathfrak{P}(t_n) \left( 1 - \sum_{i=1}^k \gamma_i \right) + O(h^2) = O(h^2), \end{aligned} \quad (\text{A.57})$$

since  $\sum_{i=1}^k \gamma_i = 1$ . We are trying to prove the terms given in (A.53) are  $O(1)$ , or equivalently those given by (A.55) are  $O(h)$ . After applying (A.56) and (A.57), there are still some terms left in (A.55) which must be shown to be  $O(h)$ , namely

$$\frac{\alpha_c}{h} A_{32}(t_n) \left( p_n - 2 \sum_{i=1}^k \gamma_i p_{n-i} + \sum_{j=1}^k \gamma_j \sum_{i=1}^k \gamma_i p_{n-i-j} \right). \quad (\text{A.58})$$

Using the consistency relations (A.50), the terms in (A.58) can be rewritten as

$$\begin{aligned} \frac{\alpha_c}{h} \left( \psi_2(t_n) - 2 \sum_{i=1}^k \gamma_i \psi_2(t_{n-i}) + \sum_{j=1}^k \gamma_j \sum_{i=1}^k \gamma_i \psi_2(t_{n-i-j}) \right) - \\ h A_{32}(t_n) \left( 2 \sum_{i=1}^k \gamma_i \gamma_j p_{n-i} - \sum_{j=1}^k \gamma_j \sum_{i=1}^k \gamma_i (\delta_{ij} + j) p_{n-i-j} \right) + O(h). \end{aligned} \quad (\text{A.59})$$

Note

$$\psi_2(t_n) - \sum_{i=1}^k \gamma_i \psi_2(t_{n-i}) = h^{k-1} \left( -Q_2(t_n) + \sum_{i=1}^k \gamma_i Q_2(t_{n-i}) \right) + O(h^k)$$

and

$$- \sum_{i=1}^k \gamma_i \psi_2(t_{n-i}) + \sum_{j=1}^k \gamma_j \sum_{i=1}^k \gamma_i \psi_2(t_{n-i-j}) = - \sum_{i=1}^k \gamma_i \left( \psi_2(t_{n-i}) - \sum_{j=1}^k \gamma_j \psi_2(t_{n-i-j}) \right) = O(h^k),$$

since  $Q_2(t)$  is a smooth function of  $t$ . Except for the terms involving  $A_{32}(t)$  in (A.59), we have shown all the terms to be  $O(h)$ :

$$\begin{aligned} - \alpha_c A_{32}(t_n) \left( 2 \sum_{i=1}^k \gamma_i p_{n-i} - \sum_{j=1}^k \gamma_j \sum_{i=1}^k \gamma_i (\delta_{ij} + j) p_{n-i-j} \right) = \\ - \alpha_c A_{32}(t_n) \left( \sum_{i=1}^k \gamma_i (p_{n-i} - \sum_{j=1}^k \gamma_j p_{n-i-j}) + \sum_{j=1}^k \gamma_j \sum_{i=1}^k \gamma_i p_{n-i-j} \right) = O(h) \end{aligned}$$

from (A.56) providing  $n \geq 3k$ . Hence, we have shown the terms in (A.55) are  $O(h)$ , and equivalently those in (A.53) are  $O(1)$ . It follows from (A.52) that  $\mathfrak{e}_n$  is bounded. Since the bounds for  $\mathfrak{q}_n$ ,  $p_n$ , and  $h_n$  are uniform for sufficiently large  $n$ , we can bound  $\mathfrak{e}_n$  uniformly. Therefore, the principal leading error term in  $\mathfrak{u}_n$  is  $h^k \mathfrak{e}(t_n)$  for  $n \geq 3k$ .

difference equation is

$$\left( I_p - \frac{h}{\alpha_c} A_{11}(t_n) \right) \mathfrak{q}_n - \frac{h}{\alpha_c} A_{12}(t_n) p_n - \frac{h}{\alpha_c} A_{13}(t_n) h_n = \sum_{i=1}^k \gamma_i \mathfrak{q}_{n-i} + \psi_1(t_n). \quad (\text{A.51})$$

Multiply (A.51) by  $A_{32}(t_n) A_{31}(t_n)$  and invert the coefficient matrix of  $\mathfrak{e}_n$  to obtain

$$\mathfrak{e}_n = - \frac{\alpha_c \Gamma(t_n) A_{32}(t_n) A_{31}(t_n)}{h} \left( \sum_{i=1}^k \gamma_i \mathfrak{q}_{n-i} + \psi_1(t_n) \right) - \left( I_p - \frac{h}{\alpha_c} A_{11}(t_n) \right) \mathfrak{q}_n + \frac{h}{\alpha_c} A_{12}(t_n) p_n. \quad (\text{A.52})$$

For bounded  $\mathfrak{q}_n$ ,  $p_n$ , and  $h_n$ , we know  $\psi_1(t_n)$  is order  $O(h)$ . Thus, in order to prove  $\mathfrak{e}_n$  is bounded, we need only to prove the terms

$$- \frac{\alpha_c \Gamma(t_n) A_{32}(t_n) A_{31}(t_n)}{h} \left( \sum_{i=1}^k \gamma_i \mathfrak{q}_{n-i} - \mathfrak{q}_n \right) \quad (\text{A.53})$$

are  $O(1)$ . The second block difference equation in (A.8) is

$$- \frac{h}{\alpha_c} A_{21}(t_n) \mathfrak{q}_n + \left( I_q - \frac{h}{\alpha_c} A_{22}(t_n) \right) p_n = \sum_{i=1}^k \gamma_i p_{n-i} + \psi_2(t_n).$$

Then,

$$A_{21}(t_{n-i}) \mathfrak{q}_{n-i} = \frac{\alpha_c}{h} \left( \left( I_q - \frac{h}{\alpha_c} A_{22}(t_{n-i}) \right) p_{n-i} - \sum_{j=1}^k \gamma_j p_{n-i-j} - \psi_2(t_{n-i}) \right)$$

for  $i = 1, 2, \dots, k$  and  $n \geq 2k$ , and hence

$$\begin{aligned} A_{21}(t_n) \left( \mathfrak{q}_n - \sum_{i=1}^k \gamma_i \mathfrak{q}_{n-i} \right) &= \left( \frac{\alpha_c}{h} I_q - A_{22}(t_n) \right) \left( p_n - \sum_{i=1}^k \gamma_i p_{n-i} \right) - \\ \frac{\alpha_c}{h} \left( \sum_{j=1}^k \gamma_j p_{n-j} - \sum_{i=1}^k \gamma_i p_{n-i} \right) &+ \left( \psi_2(t_n) - \sum_{i=1}^k \gamma_i \psi_2(t_{n-i}) \right) + O(h) \end{aligned} \quad (\text{A.54})$$

for  $n \geq 2k$ . Multiply (A.54) by  $A_{32}(t_n)$  and rearrange the terms to obtain

$$\begin{aligned} A_{32}(t_n) A_{21}(t_n) \left( \mathfrak{q}_n - \sum_{i=1}^k \gamma_i \mathfrak{q}_{n-i} \right) &= \frac{\alpha_c}{h} A_{32}(t_n) \left( p_n - 2 \sum_{i=1}^k \gamma_i p_{n-i} + \sum_{j=1}^k \gamma_j \sum_{i=1}^k \gamma_i p_{n-i-j} \right) \\ - A_{32}(t_n) \left[ A_{22}(t_n) \left( p_n - \sum_{i=1}^k \gamma_i p_{n-i} \right) + \frac{\alpha_c}{h} \left( \psi_2(t_n) - \sum_{i=1}^k \gamma_i \psi_2(t_{n-i}) \right) \right] &+ O(h). \end{aligned} \quad (\text{A.55})$$

From equation (A.25) it follows that  $p_n = \sum_{i=1}^k \gamma_i X_{32}(t_n) p_{n-i} + O(h)$  for  $n \geq 2k$  and for  $\mathfrak{q}_n$ ,  $p_n$ ,  $h_n$ , bounded ( $k \geq 2$ ). Then,

$$\begin{aligned} p_n - \sum_{i=1}^k \gamma_i p_{n-i} &= \sum_{i=1}^k -\gamma_i A_{31}(t_n) A_{13}(t_n) \Gamma(t_n) A_{32}(t_n) p_{n-i} + O(h) \\ &= \sum_{i=1}^k -\gamma_i A_{31}(t_n) A_{13}(t_n) \Gamma(t_n) \psi_2(t_{n-i}) + O(h) \\ &= O(h^{k-1}) + O(h). \end{aligned} \quad (\text{A.56})$$

This proof can be generalized in a straightforward way to include residual terms representing the effect of not solving the difference equations exactly. Specifically, if one introduces residual terms of order  $O(h^{k+2})$ ,  $O(h^{k+3})$ , and  $O(h^{k+3})$ , respectively in the difference equations corresponding to  $v'$ ,  $w'$ , and the algebraic equations, then one can modify the definitions of  $\psi_1$ ,  $\psi_2$ , and  $\psi_3$  to include corresponding terms of order  $O(h)$ ,  $O(h^2)$ , and  $O(h^3)$ . The proof then follows as before. Moreover, if the starting values  $y_{k-1}$  also satisfy the algebraic equations to  $O(h^{k+3})$  accuracy, then the smoothness of the leading error term  $h^k e(t)$  for  $u$  follows for  $n \geq 2k$  and  $k \geq 2$ . To establish this result for  $k = 1$ , a more careful analysis utilising the smoothness of  $\psi_2(t)$  is required.

*Remark.* If the result in Corollary 2 were to hold for the numerical solution expressed in long vector notation, one would have to require that  $n \geq 4k - 1$ . Specifically, the principal leading error term in  $u_n$  is  $h^k(e(t_n), e(t_{n-1}), \dots, e(t_{n-k+1}))$  for all  $n \geq 4k - 1$ .