

THE NUMERICAL SOLUTION OF NONLINEAR VOLTERRA INTEGRAL EQUATIONS OF THE SECOND KIND BY COLLOCATION AND ITERATED COLLOCATION METHODS*

J. G. BLOM† AND H. BRUNNER‡

Abstract. The subject of this paper is a variable stepsize one-step method of collocation type for solving general nonlinear second kind Volterra integral equations. We extend the iterated collocation method corresponding to polynomial spline collocation to nonlinear Volterra integral equations of the second kind. The resulting superconvergence properties of either the collocation approximation or the iterated collocation approximation are used to obtain (local and global) error estimates which in turn form the basis of a variable stepsize code. The performance of this code is illustrated by means of numerous test problems.

Key words. nonlinear Volterra integral equation, polynomial splines, collocation, iterated collocation, superconvergence, error estimates, variable stepsize

AMS(MOS) subject classifications. Primary 65R20; secondary 45D05, 45L10

0. Introduction. Most of the work on the numerical solution of Volterra integral equations of the second kind done to date deals with fixed-step methods. Recently, some automatic codes have been published. Some of these employ fixed stepsizes and make use of extrapolation tableaux to control the global error (VOLCON [13], IVRXX(C) [19]). Others use variable stepsizes based on local error controlling techniques (VOLTEX [14], RKVIEP [12], INTSOL [16], ORION [3], [17]).

The variable stepsize method we propose offers the user the choice to control the global or the local error, with the possibility to control also the uniform error, as suggested by Arndt [2].

It is believed by the authors that iterated collocation is a very robust and not too expensive method to control the global error. Since there exists a class of problems to which this method cannot be applied, some other combinations of collocation methods are investigated. The performance of these methods was considerably improved by using uniform error control.

The variable stepsize collocation methods described below have been implemented in a FORTRAN code COLVI2 (see [4]).

1. Polynomial spline collocation.

1.1. Exact collocation equations. Consider the nonlinear second kind Volterra integral equation,

$$(1.1.1) \quad y(t) = g(t) + \int_0^t k(t, s, y(s)) ds, \quad t \in I := [0, T],$$

where g and k are given continuous functions on I and on $S \times \mathbb{R}$ (with $S := \{(t, s) : 0 \leq s \leq t \leq T\}$), respectively, and where k is such that (1.1.1) possesses a unique solution $y \in C(I)$ (compare, e.g., Miller [18], or Brunner and van der Houwen [10]).

In order to discretize (1.1.1) let

$$\Pi_N: 0 = t_0 < t_1 < \cdots < t_N = T \quad (N \geq 1)$$

* Received by the editors October 30, 1985; accepted for publication (in revised form) May 29, 1986.

† Centre for Mathematics and Computer Science, Kruislaan 413, 1098 SJ Amsterdam, the Netherlands.

‡ Department of Mathematics and Statistics, Memorial University of Newfoundland, St. John's, Newfoundland, Canada ALC 557.

denote a partition (or: mesh) for the given interval I , and set

$$\begin{aligned} h_n &:= t_{n+1} - t_n \quad (n=0, \dots, N-1), \\ \sigma_0 &:= [t_0, t_1], \quad \sigma_n := (t_n, t_{n+1}] \quad (n=1, \dots, N-1), \\ Z_N &:= \{t_n: n=1, \dots, N-1\} \quad (\text{the set of interior mesh points}), \end{aligned}$$

and

$$\bar{Z}_N := Z_N \cup \{T\}.$$

In the following analysis it will always be assumed that the mesh sequence $\{\Pi_N\}$ ($N \in \mathbb{N}$) is quasi-uniform; i.e.,

$$\frac{\max_{(n)}(h_n)}{\min_{(n)}(h_n)} \leq \gamma < \infty$$

uniformly for $N \in \mathbb{N}$ (for ease of notation we have suppressed the superscript N in $t_n^{(N)}$, $h_n^{(N)}$, indicating the dependence of these quantities on N). In the following $h := \max_{(n)}(h_n)$ denotes the diameter of the mesh Π_N .

The exact solution y of (1.1.1) will be approximated in the polynomial spline space

$$(1.1.2) \quad S_{m-1}^{(-1)}(Z_N) := \{u: u|_{\sigma_n} =: u_n \in \pi_{m-1}, n=0, \dots, N-1 \ (m \geq 1)\},$$

whose elements reduce, on each subinterval σ_n , to (real) polynomials of degree not exceeding $m-1$ and which, in general, possess (finite) discontinuities at their knots Z_N . Here, the value $m=1$ yields, of course, the space of step functions having knots Z_N .

Let $\{c_j\}$, with $0 \leq c_1 < \dots < c_m \leq 1$, be a given set of parameters, and define the sets

$$(1.1.3a) \quad X_n := \{t_{n,j} := t_n + c_j h_n: j=1, \dots, m\} \quad (n=0, \dots, N-1),$$

and

$$(1.1.3b) \quad X(N) := \bigcup_{n=0}^{N-1} X_n.$$

The desired approximation $u \in S_{m-1}^{(-1)}(Z_N)$ will be determined by requiring that u satisfy the integral equation (1.1.1) on the finite set $X(N)$ ("collocation" on $X(N)$):

$$u(t) = g(t) + \int_0^t k(t, s, u(s)) ds \quad \text{for all } t \in X(N).$$

This *collocation equation* may also be written as

$$(1.1.4) \quad u_n(t) = g(t) + \int_{t_n}^t k(t, s, u_n(s)) ds + \sum_{i=0}^{n-1} \int_{t_i}^{t_{i+1}} k(t, s, u_i(s)) ds, \quad t \in X_n \quad (n=0, \dots, N-1),$$

thus explicitly exhibiting its recursive character. Note that if the collocation parameters $\{c_j\}$ are such that $c_1=0$ and $c_m=1$, then the approximation defined by (1.1.4) will be continuous on I ; i.e., we then have

$$u \in S_{m-1}^{(0)}(Z_N) := S_{m-1}^{(-1)}(Z_N) \cap C(I).$$

For the subsequent analysis we introduce the integrals

$$(1.1.5) \quad \Phi_{n,i}^{(j)}[u_i] := \begin{cases} \int_0^1 k(t_{n,j}, t_i + \tau h_i, u_i(t_i + \tau h_i)) d\tau & \text{if } 0 \leq i \leq n-1, \\ \int_0^{c_j} k(t_{n,j}, t_n + \tau h_n, u_n(t_n + \tau h_n)) d\tau & \text{if } i = n \quad (j=1, \dots, m). \end{cases}$$

With this notation the collocation equation (1.1.4) becomes

$$(1.1.6) \quad u_n(t_{n,j}) = g(t_{n,j}) + h_n \Phi_{n,n}^{(j)}[u_n] + \sum_{i=0}^{n-1} h_i \Phi_{n,i}^{(j)}[u_i], \quad j = 1, \dots, m \quad (n = 0, \dots, N-1).$$

Setting

$$(1.1.7) \quad u_n(t_n + \tau h_n) = \sum_{l=1}^m L_l(\tau) Y_{n,l}, \quad t_n + \tau h_n \in \sigma_n,$$

with $Y_{n,l} := u_n(t_{n,l})$, and where the polynomials

$$L_l(\tau) := \prod_{\substack{\kappa=1 \\ \kappa \neq l}}^m \frac{(\tau - c_\kappa)}{(c_l - c_\kappa)}$$

represent the Lagrange canonical polynomials for the collocation parameters $\{c_j\}$, we see that (1.1.6) represents, for each $n = 0, \dots, N-1$, a nonlinear system in \mathbb{R}^m for the vector $Y_n := (Y_{n,1}, \dots, Y_{n,m})^T$. Once Y_n has been determined, the approximation u on the subinterval σ_n is given by (1.1.7).

1.2. Discretization of the collocation equations. In most applications the integrals (1.1.5) occurring in the collocation equation (1.1.6) cannot be evaluated analytically, and one is forced to resort to employing suitable quadrature formulas for their approximation. In the following we shall use *m-point interpolatory quadrature formulas* of the form

$$(1.2.1) \quad \hat{\Phi}_{n,i}^{(j)}[u_i] := \sum_{l=1}^m w_l k(t_{n,j}, t_{i,l}, u_i(t_{i,l})) \quad \text{if } 0 \leq i \leq n-1;$$

for the approximation of $\Phi_{n,n}^{(j)}[u_n]$ we shall employ either

$$(1.2.2a) \quad \tilde{\Phi}_{n,n}^{(j)}[u_n] := \sum_{l=1}^m w_l k(t_{n,j}, t_{n,l}, u_n(t_{n,l})),$$

or

$$(1.2.2b) \quad \hat{\Phi}_{n,n}^{(j)}[u_n] := \sum_{l=1}^m w_{j,l} k(t_{n,j}, t_n + c_j c_l h_n, u_n(t_n + c_j c_l h_n)).$$

Here, the quadrature weights are given by

$$w_l := \int_0^1 L_l(\tau) d\tau,$$

and by

$$w_{j,l} := c_j w_l \quad (j, l = 1, \dots, m).$$

Note that when using the quadrature approximation (1.2.2a) we need kernel values $k(t, s, \cdot)$ at points (t, s) no longer contained in S (since $t_{n,l} > t_{n,j}$ for $l > j$), while (1.2.2b) is based on the abscissas $t_n + c_j c_l h_n \leq t_{n,j}$ (i.e., the affine images of $t_{n,l}$ in $[t_n, t_{n,j}]$).

The resulting fully discretized collocation equations are thus given by, respectively,

$$(1.2.3a) \quad \tilde{u}_n(t_{n,j}) = g(t_{n,j}) + h_n \tilde{\Phi}_{n,n}^{(j)}[\tilde{u}_n] + \sum_{i=0}^{n-1} h_i \hat{\Phi}_{n,i}^{(j)}[\tilde{u}_i], \quad j = 1, \dots, m \quad (n = 0, \dots, N-1);$$

and

$$(1.2.3b) \quad \hat{u}_n(t_{n,j}) = g(t_{n,j}) + h_n \hat{\Phi}_{n,n}^{(j)}[\hat{u}_n] + \sum_{i=0}^{n-1} h_i \hat{\Phi}_{n,i}^{(j)}[\hat{u}_i],$$

$$j = 1, \dots, m \quad (n = 0, \dots, N-1);$$

they yield approximations \tilde{u} and \hat{u} in $S_{m-1}^{(-1)}(Z_N)$ (whose restrictions to the subinterval σ_n are given by interpolation formulas analogous to (1.1.7)). These approximations will, in general, be different from the "exact" collocation approximation u defined by (1.1.6); due to the differing quadrature errors associated with (1.2.2a) and (1.2.2b) we have also $\tilde{u} \neq \hat{u}$. However, the approximations u , \tilde{u} , and \hat{u} all exhibit the same order of (global and local) convergence; in particular, we have the following results (see, e.g., Brunner and van der Houwen [10]; compare also de Hoog and Weiss [15], Brunner [7] and Brunner and Nørsett [11]).

THEOREM 1.1. *Let g and k be smooth functions, such that (1.1.1) has a unique solution y which is sufficiently many times continuously differentiable on I , and let u , \tilde{u} , $\hat{u} \in S_{m-1}^{(-1)}(Z_N)$ denote, respectively, the collocation approximations determined by (1.1.6), (1.2.3a), and (1.2.3b).*

Moreover, let p , \tilde{p} and \hat{p} be the largest integers satisfying

$$\max_{t_n \in \bar{Z}_N} |y(t_n) - u(t_n)| = O(N^{-p}),$$

$$\max_{t_n \in \bar{Z}_N} |y(t_n) - \tilde{u}(t_n)| = O(N^{-\tilde{p}})$$

and

$$\max_{t_n \in \bar{Z}_N} |y(t_n) - \hat{u}(t_n)| = O(N^{-\hat{p}}) \quad (\text{as } N \rightarrow \infty, \text{ with } Nh \leq \gamma T \text{ fixed}).$$

(a) *If the collocation parameters $\{c_j\}$ are given by the zeros of the Legendre polynomials $P_m(2s-1)$ (i.e., the Gauss points for $(0, 1)$), then*

$$(1.2.4a) \quad p = \tilde{p} = \hat{p} = m.$$

(b) *If the $\{c_j\}$ are the zeros of $P_{m-1}(2s-1) - P_m(2s-1)$ (i.e., the Radau II points for $(0, 1]$), then*

$$(1.2.4b) \quad p = \tilde{p} = \hat{p} = 2m - 1.$$

(c) *If the $\{c_j\}$ are the zeros of $s(s-1)P'_{m-1}(2s-1)$ (i.e., the Lobatto points for $[0, 1]$), then*

$$(1.2.4c) \quad p = \tilde{p} = \hat{p} = 2m - 2.$$

Note that \tilde{p} attains the indicated values only if the kernel $k(t, s, \cdot)$ is smooth not only for $(t, s) \in S$ but also for $(t, s) \in S' := \{(t, s) : 0 \leq s \leq t + \delta\} \cap I \times I$, for some $\delta > 0$. This is, of course, a consequence of the special form of the quadrature approximation (1.2.2a).

We also recall that the order of global convergence on I of these three spline approximations equals m . Hence, while collocation at the Radau II points or at the Lobatto points leads to a higher order of convergence on \bar{Z}_N ("local superconvergence"), collocation using the Gauss points does not have this local superconvergence property (in contrast to polynomial spline collocation applied to initial-value problems for ordinary differential equations).

In order to prepare the ground for the convergence analysis for iterated collocation (§ 2) we consider a somewhat different discretization of the collocation equation (1.1.6). To be precise, let $z \in S_m^{(-1)}(Z_N)$ be determined by the (exact) collocation equation

$$(1.2.5) \quad z_n(t_{n,j}) = g(t_{n,j}) + h_n \Phi_{n,n}^{(j)}[z_n] + \sum_{i=0}^{n-1} h_i \Phi_{n,i}^{(j)}[z_i],$$

where now $j = 1, \dots, m+1$, $t_{n,j} := t_n + c_j h_n$, with $0 \leq c_1 < \dots < c_m < c_{m+1} = 1$. Here, the integrals $\Phi_{n,i}^{(j)}$ are as in (1.1.5) (except that now j varies between 1 and $m+1$). These integrals are approximated by m -point (instead of $(m+1)$ -point!) *interpolatory quadrature formulas* using abscissas based on the first m of the $m+1$ collocation parameters; i.e., we employ again the quadrature approximations (1.2.1), (1.2.2a), (1.2.2b), with z replacing u , and with j ranging from 1 to $m+1$. Accordingly, the two possible discretizations of (1.2.5) are given by

$$(1.2.6a) \quad \tilde{z}_n(t_{n,j}) = g(t_{n,j}) + h_n \tilde{\Phi}_{n,n}^{(j)}[\tilde{z}_n] + \sum_{i=0}^{n-1} h_i \hat{\Phi}_{n,i}^{(j)}[\tilde{z}_i],$$

$$j = 1, \dots, m+1 \quad (n = 0, \dots, N-1),$$

and

$$(1.2.6b) \quad \hat{z}_n(t_{n,j}) = g(t_{n,j}) + h_n \hat{\Phi}_{n,n}^{(j)}[\hat{z}_n] + \sum_{i=0}^{n-1} h_i \hat{\Phi}_{n,i}^{(j)}[\hat{z}_i],$$

$$j = 1, \dots, m+1 \quad (n = 0, \dots, N-1);$$

these equations yield, respectively, approximations \tilde{z} and \hat{z} in the space $S_m^{(-1)}(Z_N)$. In analogy to (1.1.7) we set, e.g.,

$$(1.2.7) \quad \tilde{z}_n(t_n + \tau h_n) = \sum_{l=1}^{m+1} L_l^*(\tau) \tilde{Z}_{n,l}, \quad t_n + \tau h_n \in \sigma_n,$$

with $\tilde{Z}_{n,l} := \tilde{z}_n(t_{n,l})$, and with $L_l^*(\tau)$ denoting the l th Lagrange canonical polynomial (of degree m) with respect to c_1, \dots, c_{m+1} .

Consider first (1.2.6a): it follows from the quadrature approximation (1.2.2a) and from (1.2.7) that

$$\tilde{\Phi}_{n,n}^{(j)}[\tilde{z}_n] = \sum_{l=1}^m w_l k(t_{n,j}, t_{n,l}, \tilde{Z}_{n,l}), \quad j = 1, \dots, m+1$$

(note that $L_{m+1}^*(c_l) = 0$ for $l = 1, \dots, m$). Hence, the first m equations of (1.2.6a) (corresponding to $j = 1, \dots, m$) constitute a nonlinear system in \mathbb{R}^m for $\tilde{Z}_{n,1}, \dots, \tilde{Z}_{n,m}$; the last equation ($j = m+1$) then determines, in an explicit way, the approximation $\tilde{Z}_{n,m+1}$ which, by $c_{m+1} = 1$, equals $\tilde{z}_n(t_{n+1})$ and thus approximates $y(t_{n+1})$.

Turning to the discretization (1.2.6b), we observe that now, by (1.2.2b),

$$\hat{\Phi}_{n,n}^{(j)}[\hat{z}_n] = \sum_{l=1}^m w_l k\left(t_{n,j}, t_n + c_j h_n, \sum_{s=1}^{m+1} L_s^*(c_j c_l) \hat{Z}_{n,s}\right), \quad j = 1, \dots, m+1,$$

depends not only on $\hat{Z}_{n,1}, \dots, \hat{Z}_{n,m}$, but also on $\hat{Z}_{n,m+1}$; hence, (1.2.6b) represents a nonlinear system in \mathbb{R}^{m+1} , with $\hat{Z}_{n,m+1}$ approximating $y(t_{n+1})$.

The following theorem shows that the approximations generated by the discretizations (1.2.6a) and (1.2.6b) exhibit the same order of convergence.

THEOREM 1.2. Let g and k be smooth functions, such that the integral equation (1.1.1) possesses a unique solution y which is sufficiently many times continuously differentiable on I , and let $z, \tilde{z}, \hat{z} \in S_m^{(-1)}(Z_N)$ denote, respectively, the collocation approximations determined by (1.2.5), (1.2.6a), and (1.2.6b). If p, \tilde{p} and \hat{p} are the largest integers satisfying

$$\max_{t_n \in \bar{Z}_N} |y(t_n) - z(t_n)| = O(N^{-p}),$$

$$\max_{t_n \in \bar{Z}_N} |y(t_n) - \tilde{z}(t_n)| = O(N^{-\tilde{p}})$$

and

$$\max_{t_n \in \bar{Z}_N} |y(t_n) - \hat{z}(t_n)| = O(N^{-\hat{p}}),$$

and if the first m collocation parameters c_1, \dots, c_m are the zeros of $P_m(2s-1)$ (Gauss points for $(0, 1)$), while $c_{m+1} = 1$, then

$$(1.2.8) \quad p = \tilde{p} = \hat{p} = 2m.$$

The assertion $\tilde{p} = 2m$ will hold only if the kernel $k(t, s, \cdot)$ can be extended smoothly to the domain $S' := \{(t, s) : 0 \leq s \leq t + \delta\} \cap I \times I$, where $\delta > 0$ is suitably chosen.

A proof of this result may be found in Brunner [7], or in Brunner and van der Houwen [10].

We mention in passing that the discretized collocation equation (1.2.6a) represents an m -stage implicit Volterra-Runge-Kutta method of Pouzet type (cf. Brunner, Hairer and Nørsett [9]). It will be shown in the next section (Theorem 2.1) that this method is equivalent to a certain discretized version of an iterated collocation method, in the sense that they both produce identical approximations to y on \bar{Z}_N .

2. Iterated collocation approximations.

2.1. Exact iterated collocation. Suppose that the collocation approximation $u \in S_{m-1}^{(-1)}(Z_N)$ (defined by (1.1.6)) has been computed. The (exact) iterated collocation approximation u^I corresponding to u is then defined by

$$(2.1.1) \quad u^I(t) := g(t) + \int_0^t k(t, s, u(s)) ds, \quad t \in I.$$

(Compare also Brunner [8] and the references cited there for some background and history of iterated approximations for (linear) Fredholm and Volterra integral equations.) It is readily verified that u^I has the following properties:

- (i) $u^I \in C(I)$ (while u is, in general, not continuous on I);
- (ii) $u^I(t) = u(t)$ for all $t \in X(N)$.

This second property implies, in particular, that if $c_m = 1$ (i.e., if \bar{Z}_N is a subset of the set $X(N)$ of collocation points), then $u^I(t_n) = u(t_n)$ for all $t_n \in \bar{Z}_N$.

In the following we shall be interested in the values $u^I(t_n)$, $t_n \in \bar{Z}_N$, in the case when $c_m < 1$. Setting

$$(2.1.2) \quad \Phi_{n,i}[u_i] := \int_0^1 k(t_n, t_i + \tau h_i, u_i(t_i + \tau h_i)) d\tau, \quad 0 \leq i \leq n-1 \leq N-1,$$

we may write (2.1.1) in the form

$$(2.1.3) \quad u^I(t_n) = g(t_n) + \sum_{i=0}^{n-1} h_i \Phi_{n,i}[u_i], \quad t_n \in \bar{Z}_N.$$

2.2. Discretized iterated collocation. The integrals (2.1.2) will, in general, have to be approximated by appropriate numerical quadrature; as in § 1.2 we choose m -point interpolatory quadrature formulas of the form

$$(2.2.1) \quad \hat{\Phi}_{n,i}[u_i] := \sum_{l=1}^m w_l k(t_n, t_{i,l}, u_i(t_{i,l})),$$

with quadrature weights $\{w_l\}$ as in (1.2.1). Moreover, as indicated in § 1.2, we shall usually not be able to compute u itself but only either $\tilde{u} \in S_{m-1}^{(-1)}(\mathcal{Z}_N)$, or $\hat{u} \in S_{m-1}^{(-1)}(\mathcal{Z}_N)$ (i.e., the solutions of the discretized collocation equations (1.2.3a), (1.2.3b)). Accordingly, the *discretization* of (2.1.3) will be given either by

$$(2.2.2a) \quad \tilde{u}^I(t_n) := g(t_n) + \sum_{i=0}^{n-1} h_i \hat{\Phi}_{n,i}[\tilde{u}_i], \quad t_n \in \bar{\mathcal{Z}}_N,$$

with $\tilde{u} \in S_{m-1}^{(-1)}(\mathcal{Z}_N)$ defined by (1.2.3a), or by

$$(2.2.2b) \quad \hat{u}^I(t_n) := g(t_n) + \sum_{i=0}^{n-1} h_i \hat{\Phi}_{n,i}[\hat{u}_i], \quad t_n \in \bar{\mathcal{Z}}_N,$$

with $\hat{u} \in S_{m-1}^{(-1)}(\mathcal{Z}_N)$ given by (1.2.3b).

In the convergence analysis for $y(t_n) - \tilde{u}^I(t_n)$ and $y(t_n) - \hat{u}^I(t_n)$ the following result will play an important role.

LEMMA 2.1. *Let the iterated approximations $\tilde{u}^I(t_n)$ and $\hat{u}^I(t_n)$ be given by (2.2.2a) and (2.2.2b), and assume that g and k in (1.1.1) satisfy the customary smoothness hypotheses. If the collocation parameters $\{c_j\}$ are the zeros of $P_m(2s-1)$ (Gauss points for $(0, 1)$), then*

$$(2.2.3) \quad \max_{t_n \in \bar{\mathcal{Z}}_N} |\tilde{u}^I(t_n) - \hat{u}^I(t_n)| = O(N^{-2m}) \quad (\text{as } N \rightarrow \infty, Nh \leq \gamma T).$$

Proof. Assuming that $k(t, s, y)$ has a bounded partial derivative $\partial k / \partial y$, with L denoting an upper bound for its absolute value, we obtain, using (2.2.2a) and (2.2.2b) together with (2.2.1),

$$\begin{aligned} |\tilde{u}^I(t_n) - \hat{u}^I(t_n)| &\leq \sum_{i=0}^{n-1} h_i |\hat{\Phi}_{n,i}[\tilde{u}_i] - \hat{\Phi}_{n,i}[\hat{u}_i]| \\ &\leq h \sum_{i=0}^{n-1} \sum_{l=1}^m |k(t_n, t_{i,l}, \tilde{u}_i(t_{i,l})) - k(t_n, t_{i,l}, \hat{u}_i(t_{i,l}))| \\ &\leq hL \sum_{i=0}^{n-1} \sum_{l=1}^m |\tilde{u}_i(t_{i,l}) - \hat{u}_i(t_{i,l})|, \quad t_n \in \bar{\mathcal{Z}}_N. \end{aligned}$$

Here, we have used the fact that weights $\{w_l\}$ of the Gaussian quadrature formula are all positive and hence bounded by one (since $\sum_{l=1}^m w_l = 1$). An argument analogous to the one used in Brunner [8, pp. 1138–1139], involving the Gauss quadrature errors and a discrete Gronwall inequality, yields, observing (1.2.3a) and (1.2.3b),

$$|\tilde{u}_i(t_{i,l}) - \hat{u}_i(t_{i,l})| \leq Ch^{2m}, \quad l = 1, \dots, m \quad (i = 0, \dots, N-1),$$

with some constant C not depending on N . The estimate (2.2.3) now follows since $Nh \leq \gamma T$ for all $N \in \mathbb{N}$. \square

THEOREM 2.1. *Let $0 \leq c_1 < \dots < c_m < 1$. Assume that $\tilde{u} \in S_{m-1}^{(-1)}(\mathcal{Z}_N)$ is the solution of the discretized collocation equation (1.2.3a), and let $\tilde{z} \in S_m^{(-1)}(\mathcal{Z}_N)$ be determined by the discretized collocation equation (1.2.6a) (m -stage implicit Volterra–Runge–Kutta*

method of Pouzet type), where the collocation parameters are $0 \leq c_1 < \dots < c_m < c_{m+1} = 1$. It then follows that

$$(2.2.4) \quad \tilde{u}^I(t_n) = \tilde{z}(t_n), \quad t_n \in \bar{Z}_N;$$

i.e., the values of the discretized iterated collocation approximation \tilde{u}^I determined by (2.2.2a) are identical to the approximations $\tilde{z}(t_n)$ furnished by the m -stage implicit VRK method of Pouzet type (1.2.6a).

Proof. Using

$$\tilde{u}_i(t_i + \tau h_i) = \sum_{l=1}^m L_l(\tau) \tilde{Y}_{i,l}, \quad t_i + \tau h_i \in \sigma_i,$$

with $\tilde{Y}_{i,l} := \tilde{u}_i(t_{i,l})$, we may write (1.2.3a) as

$$(2.2.5a) \quad \tilde{Y}_{n,j} = g(t_{n,j}) + \sum_{i=0}^n h_i \sum_{l=1}^m w_l k(t_{n,j}, t_{i,l}, \tilde{Y}_{i,l}),$$

$$j = 1, \dots, m \quad (n = 0, \dots, N-1),$$

and (2.2.2a) becomes

$$(2.2.5b) \quad \tilde{u}^I(t_n) = g(t_n) + \sum_{i=0}^{n-1} h_i \sum_{l=1}^m w_l k(t_n, t_{i,l}, \tilde{Y}_{i,l}), \quad t_n \in \bar{Z}_N.$$

On the other hand, it follows from (1.2.7) that the first m equations of (1.2.6a) ($j = 1, \dots, m$) assume the form

$$(2.2.6a) \quad \tilde{Z}_{n,j} = g(t_{n,j}) + \sum_{i=0}^n h_i \sum_{l=1}^m w_l k(t_{n,j}, t_{i,l}, \tilde{Z}_{i,l}),$$

$$j = 1, \dots, m \quad (n = 0, \dots, N-1).$$

This nonlinear system (2.2.6a) is identical with the nonlinear system (2.2.5a); hence, they possess, for all sufficiently small values of the mesh diameter h , the same solution: $\tilde{Y}_{n,j} = \tilde{Z}_{n,j}$, $j = 1, \dots, m$ ($n = 0, \dots, N-1$).

For $j = m+1$ in (1.2.6a) we find, using $c_{m+1} = 1$ (and hence $t_{n,m+1} = t_{n+1}$),

$$(2.2.6b) \quad \tilde{z}(t_{n+1}) = \tilde{Z}_{n,m+1} = g(t_{n+1}) + \sum_{i=0}^n h_i \sum_{l=1}^m w_l k(t_{n+1}, t_{i,l}, \tilde{Z}_{i,l}) \quad (n = 0, \dots, N-1).$$

Since $\tilde{Z}_{i,l} = \tilde{Y}_{i,l}$ ($l = 1, \dots, m$; $i = 0, \dots, N-1$), (2.2.6b) coincides with (2.2.5b) (with n replaced by $n+1$), thus verifying (2.2.4). \square

COROLLARY 2.1. Let g and k in (1.1.1) be subject to the smoothness hypotheses of Theorem 1.2, and let c_1, \dots, c_m be the zeros of $P_m(2s-1)$. If $\tilde{u}^I(t_n)$ is defined by (2.2.2a), with $\tilde{u} \in S_{m-1}^{(-1)}(Z_N)$ denoting the solution of the discretized collocation equation (1.2.3a), then we have

$$(2.2.7) \quad \max_{t_n \in \bar{Z}_N} |y(t_n) - \tilde{u}^I(t_n)| = O(N^{-2m})$$

(as $N \rightarrow \infty$, with $Nh \leq \gamma T$).

Proof. The above local superconvergence result follows readily from Theorem 1.2 and from (2.2.4). \square

In most applications one will not work with the discretized collocation equation (1.2.3a) and the corresponding discretized iterated collocation approximation (2.2.2a), since it may not be possible to extend the given kernel smoothly to the domain S' . Instead, the pair of discretized equations (1.2.3b), (2.2.2b) will be used: as (1.2.3b)

employs only kernel values $k(t, s, \cdot)$ with $(t, s) \in S$, (1.2.3b) and (2.2.2b) represent the natural discretizations of exact collocation and iterated collocation for (1.1.1). The following local superconvergence result holds.

THEOREM 2.2. *Let g and k in (1.1.1) be subject to the smoothness hypotheses of Theorem 1.2, and let c_1, \dots, c_m be the zeros of $P_m(2s-1)$. If $\hat{u}^I(t_n)$ is determined by (2.2.2b), where $\hat{u} \in S_{m-1}^{(-1)}(Z_N)$ is the solution of the discretized collocation equation (1.2.3b), then we have*

$$(2.2.8) \quad \max_{t_n \in \bar{Z}_N} |y(t_n) - \hat{u}^I(t_n)| = O(N^{-2m})$$

(as $N \rightarrow \infty$, $Nh \leq \gamma T$).

Proof. The result of Theorem 2.2 can be established in complete analogy to (2.2.7) and (2.2.4), using (1.2.6b), (1.2.3b), and (2.2.2b).

An alternative way to verify (2.2.8) consists of writing

$$|y(t_n) - \hat{u}^I(t_n)| \leq |y(t_n) - \tilde{u}^I(t_n)| + |\tilde{u}^I(t_n) - \hat{u}^I(t_n)|, \quad t_n \in \bar{Z}_N.$$

Since the collocation parameters c_1, \dots, c_m are the Gauss points for $(0, 1)$, (2.2.8) follows by Corollary 2.1 and by Lemma 2.1. \square

We note in passing that the "exact" iterated collocation approximation u^I defined by (2.1.1), with $u \in S_{m-1}^{(-1)}(Z_N)$ being the solution of the exact collocation equation (1.1.6), satisfies also

$$(2.2.9) \quad \max_{t_n \in \bar{Z}_N} |y(t_n) - u^I(t_n)| = O(N^{-2m})$$

(as $N \rightarrow \infty$, $Nh \leq \gamma T$).

This is an immediate consequence of

$$|y(t_n) - u^I(t_n)| \leq |y(t_n) - \hat{u}^I(t_n)| + |\hat{u}^I(t_n) - u^I(t_n)|$$

and of the form of the Gauss quadrature errors determining the order of $\hat{u}^I(t_n) - u^I(t_n)$. Hence, (2.2.9) and (2.2.8) generalize the results of Brunner [8] to nonlinear Volterra integral equations.

As an illustration for the above local superconvergence results, we mention the case of *linear* polynomial spline collocation ($m = 2$): if $\hat{u} \in S_1^{(-1)}(Z_N)$ is determined by (1.2.3b), with $c_1 = (3 - \sqrt{3})/6$, $c_2 = (3 + \sqrt{3})/6$, then the corresponding values $\hat{u}^I(t_n)$ given by (2.2.2b) are convergent of order $\hat{p} = 4$: $\max_{t_n \in \bar{Z}_N} |y(t_n) - \hat{u}^I(t_n)| = O(N^{-4})$.

2.3. Polynomial solutions to (1.1.1). Suppose that the exact solution to (1.1.1) is a polynomial, $y \in \pi_{m-1}$. It is then clear that the (exact) collocation approximation $u \in S_{m-1}^{(-1)}(Z_N)$ defined by (1.1.6) satisfies $u(t) = y(t)$ for all $t \in I$, since π_{m-1} is a subspace of $S_{m-1}^{(-1)}(Z_N)$. This no longer holds, however, for the collocation approximation $\hat{u} \in S_{m-1}^{(-1)}(Z_N)$ determined by the *discretized* collocation equation (1.2.3b). To see this, note first that

$$\begin{aligned} y(t) - \hat{u}(t) &= (y(t) - u(t)) + (u(t) - \hat{u}(t)) \\ &= u(t) - \hat{u}(t) = \sum_{l=1}^m L_l(\tau)(Y_{n,l} - \hat{Y}_{n,l}), \quad t = t_n + \tau h_n \in \sigma_n. \end{aligned}$$

Let

$$E_{n,i}^{(j)}[u_i] := \Phi_{n,i}^{(j)}[u_i] - \hat{\Phi}_{n,i}^{(j)}[u_i] \quad (i \leq n)$$

denote the quadrature errors induced by the approximations (1.2.1), (1.2.2b) to the integrals (1.1.5). It follows from (1.1.6) and (1.2.3b), invoking a standard argument

based on a discrete Gronwall inequality, that the order of the differences $Y_{n,l} - \hat{Y}_{n,l}$ is given by the order of the quadrature errors $E_{n,l}^{(j)}[u_i]$. According to the above we thus obtain, e.g.,

$$\max_{t_n \in \bar{Z}_N} |y(t_n) - \hat{u}(t_n)| = \begin{cases} O(N^{-2m}) & \text{if } \{c_j\} \text{ are the Gauss points;} \\ O(N^{-(2m-1)}) & \text{if } \{c_j\} \text{ are the Radau II points;} \\ O(N^{-(2m-2)}) & \text{if } \{c_j\} \text{ are the Lobatto points.} \end{cases}$$

Note in particular that if the collocation parameters $\{c_j\}$ are the zeros of $P_m(2s-1)$ (Gauss points), then the errors $y(t_n) - \hat{u}(t_n)$ and $y(t_n) - \hat{u}^I(t_n)$ (cf. (2.2.8)) exhibit the same order of local superconvergence:

$$\max_{t_n \in \bar{Z}_N} |y(t_n) - \hat{u}(t_n)| = O(N^{-2m}), \quad \max_{t_n \in \bar{Z}_N} |y(t_n) - \hat{u}^I(t_n)| = O(N^{-2m})$$

(as $N \rightarrow \infty$, $Nh \leq \gamma T$).

This result will have certain implications in connection with global error estimation based on collocation at the Gauss points and corresponding iterated collocation (see § 3.2.1).

3. Variable stepsize method. A one-step method with variable stepsize can be separated into three major components:

- a basic fixed-step method
- an error estimator
- a stepsize strategy.

In this chapter these three components will be discussed.

3.1. Fixed-step method. The underlying fixed-step method will be a collocation method as described in the previous chapters. The integrals will be approximated by interpolatory quadrature formulas. For the sake of convenience we give a summary of the formulas used to approximate the solution at the point t_{n+1} (cf. 1.2.3b):

$$(3.1.1) \quad \hat{Y}_{n,j} = g(t_{n,j}) + \hat{F}_n(t_{n,j}) + h_n c_j \sum_{l=1}^s w_l k \left(t_{n,j}, t_n + c_j c_l h_n, \sum_{\kappa=1}^m L_\kappa(c_j c_l) \hat{Y}_{n,\kappa} \right),$$

$$\hat{F}_n(t) := \sum_{i=0}^{n-1} h_i \sum_{l=1}^s w_l k(t, t_{i,l}, \hat{Y}_{i,l}), \quad j = 1, \dots, m; \quad s = m \text{ or } m-1.$$

If $c_m = 1$, the approximate solution in t_{n+1} , \hat{u}_{n+1} equals $\hat{Y}_{n,m}$. Otherwise, \hat{u}_{n+1} is either computed by Lagrange interpolation:

$$(3.1.2) \quad \hat{u}_{n+1} = \sum_{\kappa=1}^m L_\kappa(1) \hat{Y}_{n,\kappa},$$

which is the discrete analogue of (1.1.7), or by iterated collocation (cf. 2.2.2b):

$$(3.1.3) \quad \hat{u}_{n+1}^I = g(t_n + h_n) + \hat{F}_n(t_n + h_n) + h_n \sum_{l=1}^s w_l k(t_n + h_n, t_{n,l}, \hat{Y}_{n,l}).$$

The free collocation parameters can be used to get a method that combines high accuracy with a minimum amount of work involved. The order p of the method is assumed to give an indication of the accuracy of the solution. The amount of computational work involved is normally expressed by the number of kernel evaluations needed.

For various choices of collocation parameters these values are (see Theorem 1.1, Theorem 1.2 and Theorem 2.2):

(a) Gauss-Legendre points c_1, \dots, c_m ; $s = m$,

$$p = m,$$

$$n_{\text{keV}} = \frac{N(N-1)}{2} m^2 + N n_I m^2;$$

(b) Gauss-Legendre points $c_1, \dots, c_{m-1}, c_m = 1$; $s = m - 1$,

$$p = 2m - 2,$$

$$n_{\text{keV}} = \frac{N(N-1)}{2} (m^2 - m) + N n_I (m^2 - m);$$

(c) Gauss-Legendre points c_1, \dots, c_m + Iterated collocation; $s = m$,

$$p = 2m,$$

$$n_{\text{keV}} = \frac{N(N-1)}{2} (m^2 + m) + Nm + N n_I m^2;$$

(d) Lobatto points c_1, \dots, c_m ; $s = m$,

$$p = 2m - 2,$$

$$n_{\text{keV}} = \frac{N(N-1)}{2} (m-1)^2 + 2N(m-1) + N n_I (m-1)^2;$$

(e) Radau II points c_1, \dots, c_m ; $s = m$,

$$p = 2m - 1,$$

$$n_{\text{keV}} = \frac{N(N-1)}{2} m^2 + N n_I m^2.$$

Here, N denotes the number of steps taken and n_I the number of iterations needed to solve the (non)linear system for $\hat{Y}_{n,j}$; $n_I = 1$ if the integral equation has a linear kernel, and n_I is bounded by a constant for nonlinear kernels.

Notice the low number of kernel evaluations in the case of Lobatto points as compared to the amount needed when using (the same number of) Radau II points. This decrease in the number of kernel evaluations is caused by the fact that one can make use of the coincidence of the points $t_{n-1} + c_m h_{n-1}$ and $t_n + c_1 h_n$ and reformulate (3.1.1) into:

$$(3.1.1') \quad \begin{aligned} \hat{Y}_{n,1} &= \hat{u}_n, \\ \hat{Y}_{n,j} &= g(t_{n,j}) + \hat{F}_n(t_{n,j}) + h_n c_j w_1 k(t_{n,j}, t_n, \hat{u}_n) \\ &\quad + h_n c_j \sum_{l=2}^m w_l k\left(t_{n,j}, t_n + c_l h_n, \sum_{\kappa=1}^m L_\kappa(c_l c_1) \hat{Y}_{n,\kappa}\right), \quad j = 2, \dots, m, \\ \hat{F}_n(t) &:= h_0 w_1 k(t, t_0, y(t_0)) \\ &\quad + \sum_{i=0}^{n-2} \left[h_i \sum_{l=2}^{m-1} w_l k(t, t_{i,l}, \hat{Y}_{i,l}) + (h_i + h_{i+1}) w_m k(t, t_i + h_i, \hat{Y}_{i,m}) \right] \\ &\quad + h_{n-1} \sum_{l=2}^m w_l k(t, t_{n-1,l}, \hat{Y}_{n-1,l}). \end{aligned}$$

From the above, it seems that the method with the best accuracy/computational work ratio will be obtained by using the Lobatto points.

3.2. Error estimation. To estimate the error in the approximate solution \hat{u}_{n+1} in the point t_{n+1} one has the choice between global and local error estimation. Both will be treated below.

3.2.1. Global error estimation. The global error in the point t_{n+1} is defined by $|y(t_{n+1}) - \hat{u}_{n+1}|$. It is normally approximated by the difference between \hat{u}_{n+1} and a reference solution that is assumed to have a higher accuracy. We will compute the reference solution with a higher-order method using the same stepsize.

If we consider the possible combinations of methods to approximate the solution and to compute the reference solution two of these appear to have the preference with respect to reliability and efficiency: Gauss–Legendre for the computation of \hat{u}_{n+1} combined with iterated collocation, using the same collocation points, for the reference solution, and Lobatto combined with a Lobatto method of the next higher order.

If the error in \hat{u}_{n+1} is $O(N^{-2\kappa})$, then the number of kernel evaluations for these two variants is approximately:

$$\begin{aligned} \text{Gauss}_{2\kappa} + \text{iterated collocation:} & \quad (2\kappa^2 + \kappa)N^2, \\ \text{Lobatto}_{\kappa+1} + \text{Lobatto}_{\kappa+2:} & \quad \frac{2\kappa^2 + 2\kappa + 1}{2} N^2. \end{aligned}$$

For each method the index denotes the number of collocation points needed in order to obtain the desired order.

Note, that Gauss + ($c_m = 1$), respectively Radau II combinations, need roughly the same number of kernel evaluations as the Lobatto method, and can be expected to be equally robust.

Although the Gauss method seems obviously less efficient, it has some additional advantages. Its implementation is straightforward and it needs considerably less overhead. Besides, the order of accuracy not only of the reference solution but also of the quadrature part of the approximation doubles the order of $\hat{u}(t)$ in the steppoints, which makes it more robust for a variable stepsize method (see also § 3.3). However, as mentioned in § 2.3, it has one drawback. If the solution of the integral equation is a polynomial of degree $< m$, $\hat{u}(t)$ has an error of $O(N^{-2m})$ not only in the collocation points but also in the steppoints. Since the reference solution is equally accurate this results in a very unreliable error estimate. It is however possible to detect this failure automatically and switch to another method whenever this undesired superconvergence in the steppoints occurs (see § 4.2).

3.2.2. Local error estimation. In analogy to what is common practice in the theory of initial-value problems for ordinary differential equations, we define the local error as the error resulting from a single step of the approximating method, not taking into account the errors inherited from previous steps. So, the local error in the point t_{n+1} is defined by:

$$(3.2.2.1) \quad LE_{n+1} := \left| \int_{t_n}^{t_{n+1}} k(t_{n+1}, s, y_n(s)) ds - h_n \sum_{l=1}^s w_l k(t_{n+1}, t_{n,l}, \hat{Y}_{n,l}) \right|,$$

with

$$(3.2.2.2) \quad y_n(t) := g(t) + \hat{F}_n(t) + \int_{t_n}^t k(t, s, y_n(s)) ds$$

(cf. Arndt [1]).

An estimation of the local error is obtained by approximating the integral by a higher-order method using the same stepsize, i.e.,

$$(3.2.2.3) \quad \int_{t_n}^{t_{n+1}} k(t_{n+1}, s, y_n(s)) ds \approx h_n \sum_{l=1}^{\bar{s}} \tilde{w}_l k(t_{n+1}, t_n + \tilde{c}_l h_n, \tilde{Y}_{n,l}),$$

where

$$(3.2.2.4) \quad \begin{aligned} \tilde{Y}_{n,j} &= g(t_n + \tilde{c}_j h_n) + \hat{F}_n(t_n + \tilde{c}_j h_n) \\ &+ h_n \tilde{c}_j \sum_{l=1}^{\bar{s}} \tilde{w}_l k\left(t_n + \tilde{c}_j h_n, t_n + \tilde{c}_j \tilde{c}_l h_n, \sum_{\kappa=1}^{\bar{m}} \tilde{L}_\kappa(\tilde{c}_j \tilde{c}_l) \tilde{Y}_{n,\kappa}\right). \end{aligned}$$

Note that it is assumed that $\tilde{c}_{\bar{m}} = 1$.

3.3. Stepsize strategy. Having an estimation of the error, we will now discuss how to determine the stepsize to limit this error to a user specified tolerance.

The stepsize strategy consists of advancing the solution with a trial stepsize, then accept or reject the result and repeating the process with as trial stepsize a modification of the previous stepsize. This modification should account for the (lack of) accuracy in the accepted or rejected trial solution. It should result in as large a stepsize as possible while still providing an accuracy that satisfies the specified tolerance.

The control of the global error in a steppoint, say t_{n+1} , is a straightforward process if a global error estimation is available. A trial solution will then be accepted if the global error estimation is less than the specified error tolerance, i.e.,

$$(3.3.1) \quad GEE_{n+1} \leq TOL.$$

The modification of the stepsize reads:

$$(3.3.2) \quad h_{new} = h_n \left(\frac{TOL}{GEE_{n+1}} \right)^{1/p}.$$

Here, p is the order of the underlying method, and h_{new} stands for h_{n+1} if the trial solution has been accepted, and otherwise for h_n .

Note that, once a certain point t_{n+1} has been reached, we cannot control the error in the lag term part any more, since the approximating solution for that part of the interval has already been accepted. This means that it is defensible to control only the local error and assume that the global error in the endpoint will be more or less equal to the sum of the local errors.

In this case one accepts a trial solution only if

$$(3.3.3) \quad LEE_{n+1} \leq \frac{TOL \cdot h_n}{t_{n+1}},$$

and modifies the stepsize by:

$$(3.3.4) \quad h_{new} = h_n \left(\frac{TOL \cdot h_n}{t_{n+1} \cdot LEE_{n+1}} \right)^{1/q},$$

where q is the order of the quadrature formula to approximate the integral.

So far, the stepsize strategy is completely analogous to the strategy used in variable stepsize methods for ordinary differential equations. However, the intention is to limit the global error over the whole interval to the user-defined tolerance, i.e., to have

$$|y(t) - \hat{u}(t)| \leq TOL \quad \text{for } 0 \leq t \leq T.$$

The disadvantage of proceeding in the same way as in the case of ODEs is, as Arndt ([1], [2]) points out, that the consequence of accepting the approximated solution to future quadrature computations is not taken into account. This means that the assumption that the sum of the local errors equals the global error in the endpoint is not necessarily valid.

To obviate the problem of the possibly large errors in future lag term computations, once t_{n+1} has been reached, one has to demand that:

$$|y(t) - \hat{u}(t)| \leq TOL \quad \text{for all } t \in [t_{n+1}, T].$$

This reads in discretized form:

$$(3.3.5) \quad \left| \sum_{i=0}^{l-1} \left(\int_{t_i}^{t_{i+1}} k(t_i, s, y(s)) ds - h_i \sum_{j=1}^s w_j k(t_i, t_{i,j}, \hat{Y}_{i,j}) \right) \right| \leq TOL,$$

for $t_i = T, T-H, \dots, t_{n+1}$, with H some fixed number.

Let

$$(3.3.6) \quad EE_{i+1}^l \approx \int_{t_i}^{t_{i+1}} k(t_i, s, y(s)) ds - h_i \sum_{j=1}^s w_j k(t_i, t_{i,j}, \hat{Y}_{i,j}).$$

Then (3.3.5) is satisfied when

$$\sum_{i=n}^{l-1} |EE_{i+1}^l| + \left| \sum_{i=0}^{n-1} EE_{i+1}^l \right| \leq TOL.$$

In this way one gets, as Arndt suggests, a modified stepsize strategy that strives for a uniform control of the error, i.e., let

$$(3.3.7) \quad EE_{n+1} = \max_{t_{n+1} \leq t_i \leq T} \left((t_i - t_n) \cdot \frac{|EE_{n+1}^l|}{h_n} + \left| \sum_{i=0}^{n-1} EE_{i+1}^l \right| \right).$$

Accept the trial solution in t_{n+1} if

$$(3.3.8) \quad EE_{n+1} \leq TOL.$$

The modification of the stepsize is

$$(3.3.9) \quad h_{new} = h_n \left(\frac{TOL}{EE_{n+1}} \right)^{1/q},$$

with q the order of the quadrature formula.

Note that although global error control does not lead to an erroneous error estimate in the endpoint, it is certainly possible that by allowing too large a stepsize at the start, it will be impossible to meet the tolerance as the integration advances, due to an unacceptably large error in the lag term. Therefore, to avoid problems later on, it seems advisable to use the global error control without the uniform error control only if $p \ll q$, which is the case for the Gauss collocation method combined with iterated collocation to estimate the global error.

4. Implementation. For the miscellaneous programming matters and heuristics, e.g., stepsize constraints, the solution of nonlinear systems, etc., we refer to Blom and Brunner [4]. Here, we consider two more general topics.

4.1. Approximation of the lag term. The lion's share of the computational work in all methods is the approximation of the lag term by some quadrature rule. Moreover, for all methods except the iterated collocation method the additional storage space needed for all $\hat{Y}_{n,j}$ is substantial.

However, the number of kernel evaluations needed to compute the reference solution can be significantly reduced by approximating the lag term by Lagrangian interpolation instead of using a quadrature formula. Interpolation over m points yields

$$(4.1.1) \quad \hat{F}_n^{(1)}(t_n + \tau h_n) = \sum_{\kappa=1}^m \hat{F}_n(t_{n,\kappa})$$

and adds an additional interpolation error of $O(h_n^m)$, unless the kernel satisfies $k \in \pi_{m-2}$. So, in general, this interpolation gives rise to a reference solution of $O(h_n^m)$, which is not adequate.

However, one can also interpolate over $2m$ points

$$(4.1.2) \quad \hat{F}_n^{(2)}(t_n + \tau h_n) = \sum_{\kappa=1}^{\bar{m}} \bar{L}_\kappa \left(\frac{h_{n-1} + \tau h_n}{h_{n-1} + h_n} \right) \hat{F}_n(t_{n-1} + \bar{c}_\kappa (h_{n-1} + h_n)),$$

where

$$\bar{m} = 2m$$

and

$$\begin{aligned} \bar{c}_j &= c_j \frac{h_{n-1}}{h_{n-1} + h_n}, & j = 1, \dots, m, \\ \bar{c}_j &= \frac{h_{n-1} + c_l h_n}{h_{n-1} + h_n}, & j = \bar{m} - m + 1, \dots, \bar{m}; \quad l = j + m - \bar{m}, \end{aligned}$$

and where the \bar{L}_κ denote the Lagrange canonical polynomials corresponding to the \bar{c}_j . This interpolation introduces an error of $O((h_{n-1} + h_n)^{\bar{m}})$, unless $k \in \pi_{\bar{m}-2}$, which is small enough to preserve the legitimacy of the reference solution in case the approximating method is based on Gauss collocation.

For all other methods that are mentioned in § 3.1 this kind of approximation of the lag term cannot be used when computing the global error estimation, since the error in the \hat{F}_n will always be a spoiling factor, no matter how accurate the interpolation will be. However, since local error estimation needs only the values of \hat{F}_n in some intermediate points, one can use the interpolation formula (4.1.2) which is sufficiently accurate for all methods described.

Note that for Lobatto collocation the value of \bar{m} is $2m - 1$ instead of $2m$, due to the coincidence of $t_{n-1} + c_m h_{n-1}$ and $t_n + c_1 h_n$.

4.2. Detection of a polynomial solution when using Gauss collocation. As mentioned in § 2.3 the Gauss method yields an error of $O(N^{-2m})$ in the steppoints if the solution is an element of π_{m-1} . If the reference solution is equally or less accurate, then this results in an unreliable global error estimate. One can detect this error by computing the approximations $\hat{Y}_{n,j}$ to $y(t_n + c_j h_n/2)$ and use Lagrangian interpolation to compute the solution at some point of the interval $[t_n, t_n + h_n/2]$, say t_n , both with these approximations and with the $\hat{Y}_{n,j}$. If $y \notin \pi_{m-1}$, then $e_1 = |y(t_n) - \hat{u}_n| \approx O(h_n^m)$ and $e_2 = |y(t_n) - \hat{u}_n| \approx O((h_n/2)^m)$. So, if $e_1/e_2 \approx 2^m$, the order of the Gauss method is as expected.

To avoid the computation of the quadrature of the lag term, we use Lagrange interpolation over the points 0.0, c_j , 1.0. No additional kernel evaluations are needed, since the $\hat{F}_n(t_n)$ and the $\hat{F}_n(t_{n+1})$ already have been calculated during the computation of the reference solution in t_n , respectively t_{n+1} . Now, if $e_1/e_2 \ll 2^m$, then it is almost certain that $y \in \pi_{m-1}$.

Note that if both $y \in \pi_{m-1}$ and $\int_0^t k(t, s, y(s)) ds \in \pi_{m+1}$, then it is not clear what the value of the ratio e_1/e_2 will be. But in this case the approximation of both the lag term and the solution is "exact", which means that the only problem that can arise is an unjustifiable rejection of the current combination of methods for approximating the solution and the reference solution.

5. Numerical performance. The ideas presented in the preceding sections have resulted in a FORTRAN code named COLV12; compare Blom and Brunner [4], where a description of this code and its usage is given.

To show some properties of the code, we tested it on the six problems that were used by Schlichte [19] to compare his method, IVRKX, with the codes ORION of Bader and Kunkel [3], VOLTEX of Hock [14], and VOLCON of Hairer, Lubich and Schlichte [13]. The performance data of these codes for the six problems are listed in the Appendix. We also added the results of VE1 of Bownds and Appelbaum [6], [5] and of INTSOL of Jones and McKee [16], insofar as we could find these results in the literature. (See Tables A1 and A2.)

The tests were performed on a CDC CYBER-750, in single precision (a machine precision of approximately 14 digits).

The problems are:

(1) (problem 4 from Schlichte)

$$y(t) = \frac{1}{2} t^2 e^{-t} + \int_0^t \frac{(t-s)^2}{2} e^{s-t} \cdot y(s) ds \quad \text{for } t \in [0, 5]$$

with solution:

$$y(t) = \frac{1}{3} \left(1 - e^{-\frac{1}{2}t} \left[\cos \left(\frac{1}{2} t \sqrt{3} \right) + \sqrt{3} \sin \left(\frac{1}{2} t \sqrt{3} \right) \right] \right);$$

(2) (problem 15 from Schlichte)

$$y(t) = 1 + (\sin(t))^2 - \int_0^t 3 \sin(t-s) \cdot y^2(s) ds \quad \text{for } t \in [0, 5]$$

with solution:

$$y(t) = \cos(t);$$

(3) (problem 19 from Schlichte)

$$y(t) = \cos(t) - \int_0^t \frac{2}{(t-s+2)^2} \cdot (y(s) + y^3(s)) ds \quad \text{for } t \in [0, 40]$$

with solution:

$$y(40) \approx -0.65013110133344;$$

(4) (problem 20 from Schlichte)

$$y(t) = 1 + \int_0^t (t-s)^3 (4-t+s) e^{s-t} \cdot \frac{y^4(s)}{1+2y^2(s)+2y^4(s)} ds \quad \text{for } t \in [0, 10]$$

with solution:

$$y(10) = 1.2599558233723;$$

(5) (problem 16 from Schlichte)

$$y(t) = e^{-t} + \int_0^t e^{s-t} [y(s) + e^{-y(s)}] ds \quad \text{for } t \in [0, 40]$$

with solution:

$$y(t) = \ln(t + e);$$

(6) (problem 21 from Schlichte)

$$y(t) = t - 1 + (1 + t^2) e^{-t^2} + \int_0^t t^2 e^{-ts} \cdot y(s) ds \quad \text{for } t \in [0, 5]$$

with solution:

$$y(t) = t.$$

The following combinations of methods for the approximation of the solution, respectively the reference solution, have been used:

Gm + IG	m points Gauss + iterated collocation,
1Gm + 1G(m + 1)	m - 1 Gauss points + (c _m = 1), together with m Gauss points + (c _{m+1} = 1),
Lm + L(m + 1)	m, combined with m + 1 points Lobatto collocation,
Rm + R(m + 1)	m, together with m + 1 points Radau collocation.

As stepsize strategies we used global error estimate/global error control (G), and, if applicable, local error estimate/local error control (L). In the latter the lag term needed for the reference solution has been approximated by interpolation (cf. 4.1.2).

The meaning of the headers in the tables is:

sd	number of significant digits: $-^{10}\log y(T) - \hat{u}(T) / \max(1.0, y(T))$;
EE/GE	ratio between the global error estimate in T and the real global error. Note: if the local error is uniformly controlled, the global error in T is estimated by: $\sum_{i=0}^{N-1} EE_{i+1}^l$, with $t_i = T$ (cf. 3.3.6);
N	number of intervals used;
#F tot	number of failed steps in total;
#F 0	number of failed steps at the start in $t = 0.0$;
#kev	total number of kernel evaluations needed.

Table 5.1 gives the results that are obtained when no special care has been taken. For example, for the (iterated) Gauss method, no testing on polynomial solutions has been done, and for the other methods, the error is not uniformly controlled.

For all experiments the initial stepsize was 1.0; the maximal stepsize allowed was 5.0, which was actually used in problem 5.

All failures in Table 5.1 were due to a rejection of a step with the minimum allowed stepsize, which was 5 E-3.

As was expected, problem 6, with solution $y = t$, resulted in a bad performance for the Iterated Gauss method; but problem 5, with solution $y = \ln(t + e)$, also seemed to suffer from a bad error estimate. Both problems were rerun with the polynomial test on and with an automatic escape to a local error estimate and uniform local error control Gauss_m, Gauss_{m+1} + (c_{m+2} = 1) combination. The escape was made after two consecutive occurrences of polynomial behavior, and the last three steps were discarded. As can be seen in Table 5.2, in most cases an escape was needed.

In Table 5.3 the results are listed that were obtained by using a uniform error control, or rather an error control in discrete points spaced 1.0. Here, if a step was rejected which already had a minimum size, the tolerance was relaxed by a factor 4 and the computation was resumed.

Recall that when using iterated Gauss collocation to compute the reference solution there is no need to employ a uniform error control (see § 3.3).

TABLE 5.1
Results for method combinations, without special features.

Prob. 1		Requested tolerance 1 E-4					Requested tolerance 1 E-7						
Method	sd	EE/GE	N	#F		#kev	sd	EE/GE	N	#F		#kev	
				tot	0					tot	0		
G	G4+IG	5.2	1	8	2	2	760	9.3	1	33	5	3	12100
G	1G4+1G5	4.6	1	4	0	0	320	7.1	1	12	2	1	2592
G	L4+L5	4.5	1	4	0	0	278	7.1	1	12	2	1	2123
G	R3+R4	4.5	1	5	1	0	425	---	---	---	---	---	---
G	G8+IG	7.4	1	3	0	0	432	9.8	1	4	0	0	720
G	1G6+1G7	8.2	1	3	0	0	432	8.2	1	3	0	0	432
G	L6+L7	8.5	1	3	0	0	399	8.1	1	4	0	0	654
G	R5+R6	6.8	1	3	0	0	366	8.0	1	4	0	0	610
L	1G4+1G5	4.6		4	0	0	276	7.1		11	2	1	1311
L	L4+L5	4.5		4	0	0	282	7.1		11	2	1	1213
L	R3+R4	5.6		6	2	0	472	4.9		20	2	0	2612
L	1G6+1G7	8.2		3	0	0	429	8.2		3	0	0	429
L	L6+L7	8.5		3	0	0	438	8.1		4	0	0	634
L	R5+R6	6.8		3	0	0	366	8.0		4	0	0	538
Prob. 2		Requested tolerance 1 E-4					Requested tolerance 1 E-7						
Method	sd	EE/GE	N	#F		#kev	sd	EE/GE	N	#F		#kev	
				tot	0					tot	0		
G	G4+IG	6.0	5E-1	9	3	1	1596	8.4	1	43	5	2	21988
G	1G4+1G5	---	---	---	---	---	---	---	---	---	---	---	---
G	L4+L5	---	---	---	---	---	---	---	---	---	---	---	---
G	R3+R4	---	---	---	---	---	---	---	---	---	---	---	---
G	G8+IG	6.2	1	3	0	0	1136	8.4	1E-1	4	1	0	2152
G	1G6+1G7	5.2	1E-1	3	0	0	1224	---	---	---	---	---	---
G	L6+L7	4.1	1	3	0	0	1070	---	---	---	---	---	---
G	R5+R6	---	---	---	---	---	---	---	---	---	---	---	---
L	1G4+1G5	1.5		7	1	0	1396	4.9		22	5	1	6429
L	L4+L5	1.3		8	2	0	1535	4.6		23	4	1	5425
L	R3+R4	2.8		12	4	1	2251	5.6		45	8	2	14730
L	1G6+1G7	5.2		3	0	0	1221	6.3		5	1	0	2490
L	L6+L7	4.1		3	0	0	1109	5.9		5	1	0	2246
L	R5+R6	2.9		4	1	0	1672	5.4		7	1	1	2740
Prob. 3		Requested tolerance 1 E-4					Requested tolerance 1 E-7						
Method	sd	EE/GE	N	#F		#kev	sd	EE/GE	N	#F		#kev	
				tot	0					tot	0		
G	G4+IG	4.2	1	70	31	2	75696	8.7	1	366	72	4	1615392
G	1G4+1G5	5.6	2	34	14	0	31552	7.9	1	98	41	2	231828
G	L4+L5	4.7	1	43	20	0	38350	8.6	1	104	50	2	214553
G	R3+R4	4.5	1	44	18	0	39759	7.3	1	175	52	2	513900
G	G8+IG	6.2	1	27	13	0	48824	7.1	1	54	19	2	154616
G	1G6+1G7	5.3	1	15	2	0	14832	7.2	1	33	9	0	61242
G	L6+L7	5.3	1	17	1	0	13679	7.9	1	39	12	1	73395
G	R5+R6	5.0	1	28	7	0	37235	7.8	1	45	10	1	88511
L	1G4+1G5	6.1		54	21	0	34093	9.5		176	45	2	262811
L	L4+L5	7.2		57	13	0	26104	9.9		190	47	2	226241
L	R3+R4	6.3		90	30	0	58567	9.2		358	68	2	714682
L	1G6+1G7	7.3		26	5	0	21167	9.2		50	18	0	72508
L	L6+L7	6.8		28	5	0	21113	9.3		50	11	1	55681
L	R5+R6	6.8		28	3	0	18805	9.8		61	24	1	87408

TABLE 5.1. (Continued)

Prob. 4		Requested tolerance 1 E-4					Requested tolerance 1 E-7						
Method	sd	EE/GE	N	#F		#kev	sd	EE/GE	N	#F		#kev	
				tot	0					tot	0		
G	G4+IG	4.5	1	11	5	1	2224	7.7	1	50	9	3	29588
G	1G4+1G5	4.9	1	10	0	0	2368	7.5	1	32	3	2	19168
G	L4+L5	4.8	1	10	0	0	1920	7.4	1	33	3	2	16102
G	R3+R4	5.2	1	15	2	1	3841	7.3	1	60	2	2	48541
G	G8+IG	7.5	1	4	0	0	1360	8.8	1	7	1	1	3240
G	1G6+1G7	5.7	1	5	0	0	1872	7.5	1	7	0	0	3240
G	L6+L7	5.6	1	5	0	0	1641	7.1	1	7	0	0	2761
G	R5+R6	6.0	1	5	0	0	1525	8.0	1	9	0	0	3843
L	1G4+1G5	6.1		12	1	0	2255	8.1		42	3	2	15303
L	L4+L5	6.3		13	1	0	2168	8.1		44	3	2	13338
L	R3+R4	4.7		21	2	1	3749	7.6		90	3	2	42831
L	1G6+1G7	5.7		5	0	0	1657	8.3		8	0	0	2896
L	L6+L7	5.6		5	0	0	1526	8.4		8	1	0	3001
L	R5+R6	5.5		6	0	0	1689	7.8		10	1	0	3709
Prob. 5		Requested tolerance 1 E-4					Requested tolerance 1 E-7						
Method	sd	EE/GE	N	#F		#kev	sd	EE/GE	N	#F		#kev	
				tot	0					tot	0		
G	G4+IG	3.2	-5E-3	11	0	0	1672	5.8	-6E-3	39	2	2	16984
G	1G4+1G5	4.3	1	29	0	0	15776	7.3	1	93	1	1	146048
G	L4+L5	4.3	1	30	0	0	13335	7.3	1	98	1	1	127043
G	R3+R4	4.2	1	51	1	0	35325	7.3	1	227	3	2	658650
G	G8+IG	9.2	5E-5	10	0	0	5432	9.3	4E-4	11	0	0	6672
G	1G6+1G7	5.0	1	10	0	0	5616	7.5	1	17	0	0	13968
G	L6+L7	4.9	1	10	0	0	4868	7.5	1	18	0	0	13252
G	R5+R6	4.3	1	11	0	0	5551	7.5	1	24	0	0	21594
L	1G4+1G5	4.2		27	0	0	7349	7.3		92	1	1	61043
L	L4+L5	4.6		33	1	0	8367	7.4		99	1	1	54443
L	R3+R4	4.1		47	0	0	13581	7.2		208	2	2	212869
L	1G6+1G7	5.3		11	0	0	4693	7.0		15	3	0	8856
L	L6+L7	5.5		11	1	0	4538	8.0		19	0	0	9258
L	R5+R6	4.9		12	1	0	4791	7.0		20	1	0	9826
Prob. 6		Requested tolerance 1 E-4					Requested tolerance 1 E-7						
Method	sd	EE/GE	N	#F		#kev	sd	EE/GE	N	#F		#kev	
				tot	0					tot	0		
G	G4+IG	3.9	-7E-3	3	0	0	120	3.9	-3E-5	5	2	0	400
G	1G4+1G5	---	-	--	-	-	----	---	-	--	-	-	----
G	L4+L5	---	-	--	-	-	----	---	-	--	-	-	----
G	R3+R4	---	-	--	-	-	----	---	-	--	-	-	----
G	G8+IG	10.4	-4E-4	3	0	0	432	10.4	-4E-4	3	0	0	432
G	1G6+1G7	5.5	1	3	0	0	432	---	-	--	-	-	----
G	L6+L7	5.4	1	3	0	0	399	---	-	--	-	-	----
G	R5+R6	4.7	1	3	0	0	366	---	-	--	-	-	----
L	1G4+1G5	2.5		4	1	0	327	3.8		9	3	1	1044
L	L4+L5	2.4		4	1	0	332	3.8		9	3	1	987
L	R3+R4	1.9		5	1	0	336	3.3		14	5	1	1679
L	1G6+1G7	5.5		3	0	0	429	5.5		4	1	0	745
L	L6+L7	5.4		3	0	0	438	5.5		4	1	0	744
L	R5+R6	4.7		3	0	0	366	4.7		4	1	0	635

TABLE 5.2
Results for Iterated Gauss; esc. to Gauss_m + Gauss_{m+1} + (c_{m+2} = 1), local error control.

Prob. 5		Requested tolerance 1 E-4					Requested tolerance 1 E-7						
Method		sd	EE/GE	N	#F tot	#kev 0	sd	EE/GE	N	#F tot	#kev 0		
G	G4+IG		escape	at	t =	2.3		escape	at	t =	3.2		
L	G4+1G6	4.5	6E-2	14	6	0	8604	6.2	1E-3	33	9	2	21541
G	G8+IG		escape	at	t =	21.0		escape	at	t =	2.7		
L	G8+1G10	9.2	-1E-1	10	3	0	11537	9.3	-3E-2	11	2	0	16360
Prob. 6		Requested tolerance 1 E-4					Requested tolerance 1 E-7						
Method		sd	EE/GE	N	#F tot	#kev 0	sd	EE/GE	N	#F tot	#kev 0		
G	G4+IG		escape	at	t =	2.8		escape	at	t =	1.8		
L	G4+1G6	7.3	1	7	4	2	1417	9.6	1	11	5	2	2436
G	G8+IG	10.4	-4E-4	3	0	0	624	10.4	-4E-4	3	0	0	624
L	G8+1G10												

The failures in Table 5.3 were caused by a lack of working storage, which occurred, depending on the method used, after about 200–350 successful steps.

Furthermore, we tried for one problem what the influence is if only the stepsize is adjusted according to the estimated error, but if no steps are rejected. Although in most cases this results in fewer kernel evaluations (if the code reaches the endpoint), it is by far not as reliable.

To present a survey of the performance of the code on these six problems a weighted sum has been computed over all (available) “sd”- and “#kev”-values. The weights are inversely proportional to the square of the length of the interval on which the solution has been computed.

Table 5.4 contains the figures of the average performance of the different methods. For the iterated collocation method the results have been taken from Table 5.1 and 5.2 for problems 1–4, respectively 5 and 6. For the other methods the “sd”- and “#kev”-values of Table 5.3 have been used. The results in the first column are the average results computed over all problems. The second contains the average values when the “bad ones” are left out, i.e., for the Gauss+Iterated collocation method problems 5 and 6, and for all other methods problem 2.

To illustrate how the code performs in extreme cases we tried an example with a discontinuous solution. We chose for this test the problem:

$$y(t) = g(t) + \int_0^t y(s) ds \quad \text{for } t \in [0, 5]$$

with:

$$y(t) = \begin{cases} t & t \leq 2.5 \\ 1/t & t > 2.5 \end{cases} \quad \text{and} \quad g(t) = \begin{cases} t - t^2/2, & t \leq 2.5, \\ 1/t - 3.125 - \ln(t) + \ln(2.5), & t > 2.5. \end{cases}$$

In Table 5.5 the results are listed that were obtained using the same features as in Table 5.4; i.e., in case of the Gauss method the solution was tested on polynomial behavior and in all other methods the error was controlled uniformly in discrete points spaced 1.0.

In all tests it was allowed to relax the tolerance. This has been done repeatedly in the last point before the discontinuity for all entries in the table where the number of failures in 2.5 is greater than 10; for the G8+IG combination until the tolerance reached a value greater than 0.4 and for all other methods until it was greater than 1.6.

TABLE 5.3
Results with uniform error control and tolerance relaxing.

Prob. 1		Requested tolerance 1 E-4					Requested tolerance 1 E-7						
Method	sd	EE/GE	N	#F		# kev	sd	EE/GE	N	#F		# kev	
				tot	0					tot	0		
G	1G4+1G5	4.8	1	4	0	0	390	7.2	1	10	1	1	2016
G	L4+L5	4.6	1	4	0	0	368	7.2	1	11	1	1	2056
G	R3+R4	4.6	1	5	1	0	548	7.2	1	18	2	2	4752
G	1G6+1G7	8.2	1	3	0	0	531	9.0	1	4	0	0	830
L	L6+L7	8.5	1	3	0	0	516	8.9	1	4	0	0	784
G	R5+R6	6.8	1	3	0	0	465	8.3	1	4	0	0	731
L	1G4+1G5	4.8	5E-1	4	0	0	346	7.3	3E-1	11	1	1	1498
L	L4+L5	4.7	1	4	0	0	372	7.3	4E-1	11	1	1	1478
L	R3+R4	5.7	-1	6	1	0	536	8.2	2E-1	22	3	2	3631
L	1G6+1G7	8.2	1	3	0	0	528	9.0	1	4	0	0	742
L	L6+L7	8.5	-2	3	0	0	555	8.9	1	4	0	0	764
L	R5+R6	6.8	2	3	0	0	465	8.3	1	4	0	0	659

Prob. 2 ¹		Requested tolerance 1 E-4					Requested tolerance 1 E-7						
Method	sd	EE/GE	N	#F		# kev	sd	EE/GE	N	#F		# kev	
				tot	0					tot	0		
G	1G4+1G5	3.2	1	53	14	1	55521	6.6	1	24	16	2	20762
G	L4+L5	3.1	1	19	16	1	9849	7.4	1	24	16	2	16953
G	R3+R4	4.3	1	70	4	2	65979	6.0	1	54	17	3	50337
G	1G6+1G7	4.8	1	4	0	0	1766	6.4	1	115	13	1	506569
G	L6+L7	4.1	1	4	0	0	1577	---	-	---	-	-	---
G	R5+R6	---	-	---	-	-	---	6.1	1	8	18	1	12305
L	1G4+1G5	3.5	-1E-1	9	2	1	2204	6.3	-5E-2	26	3	2	8058
L	L4+L5	3.8	-1E-1	9	2	1	1943	6.3	-5E-2	27	3	2	7345
L	R3+R4	3.1	2E-2	15	5	2	3446	5.7	8E-3	52	4	3	18608
L	1G6+1G7	4.8	2E-1	4	0	0	1678	6.4	-4E-2	6	1	1	3026
L	L6+L7	4.1	5E-2	4	0	0	1557	6.3	-3E-2	6	1	1	2776
L	R5+R6	3.0	-5E-3	4	0	0	1452	6.3	-1E-2	8	1	1	3481

Prob. 3		Requested tolerance 1 E-4					Requested tolerance 1 E-7						
Method	sd	EE/GE	N	#F		# kev	sd	EE/GE	N	#F		# kev	
				tot	0					tot	0		
G	1G4+1G5	5.0	1	33	6	0	30327	7.4	1	105	26	2	252649
G	L4+L5	4.5	1	38	12	0	37036	8.3	1	113	32	2	242699
G	R3+R4	4.6	1	51	24	0	63912	7.5	1	186	23	2	526699
G	1G6+1G7	4.8	1	16	0	0	17587	8.3	1	37	8	0	82764
G	L6+L7	6.1	1	18	2	0	21046	8.9	1	40	12	1	90095
G	R5+R6	5.3	1	27	5	0	40615	7.8	1	44	7	1	91051
L	1G4+1G5	6.1	1	55	25	0	47213	9.6	1	177	51	2	297461
L	L4+L5	6.5	3E-1	60	20	0	44197	9.9	2	189	44	2	263067
L	R3+R4	6.3	1	89	30	0	73244	9.2	1	359	68	2	769421
L	1G6+1G7	7.3	-1	26	5	0	28064	9.1	1	48	16	0	80400
L	L6+L7	7.7	6	28	4	0	28346	9.2	1	50	12	1	71930
L	R5+R6	6.9	-6	28	3	0	25616	10.0	2	63	22	1	106765

¹ The tolerance has been relaxed two times in all global methods except for the R3 + R4, the 1G6 + 1G7, and the L6 + L7 combinations with tolerance 1 E-4.

The "bad" error estimates for the Gauss method can be explained by the fact that the approximation is "exact" until the discontinuity has been reached and therefore the code gives no warning that the solution is a polynomial. The step that covers the discontinuity results in one warning, but since in the next step the solution is no longer a polynomial no escape will be made.

TABLE 5.3. (Continued)

Prob. 4		Requested tolerance 1 E-4					Requested tolerance 1 E-7						
Method	sd	EE/GE	N	#F		#kev	sd	EE/GE	N	#F		#kev	
				tot	0					tot	0		
G	1G4+1G5	4.9	1	10	0	0	2697	7.4	1	32	2	2	20422
G	L4+L5	4.8	1	10	0	0	2370	7.5	1	34	2	2	18776
G	R3+R4	5.3	1	15	1	1	4320	7.3	1	60	2	2	51008
G	1G6+1G7	5.8	1	5	0	0	2158	8.8	1	8	1	0	4904
G	L6+L7	5.7	1	5	0	0	1979	8.6	1	8	1	0	4401
G	R5+R6	6.0	1	5	0	0	1822	8.8	1	10	1	0	5818
L	1G4+1G5	6.1	-1	12	1	0	2724	8.1	-3E-1	42	2	2	16896
L	L4+L5	6.3	-5	13	1	0	2816	8.1	-2E-1	44	2	2	15506
L	R3+R4	4.8	-8E-2	21	2	1	4603	7.8	-2E-1	90	2	2	46162
L	1G6+1G7	5.7	-2E+1	5	0	0	1901	8.5	2	8	1	0	3837
L	L6+L7	5.6	-5	5	1	0	2124	8.2	-1	8	0	0	3177
L	R5+R6	6.7	-2E+1	6	1	0	2310	7.9	1	11	1	0	4816
Prob. 5		Requested tolerance 1 E-4					Requested tolerance 1 E-7						
Method	sd	EE/GE	N	#F		#kev	sd	EE/GE	N	#F		#kev	
				tot	0					tot	0		
G	1G4+1G5	4.3	1	29	0	0	19535	7.3	1	92	1	1	155629
G	L4+L5	4.3	1	30	0	0	18375	7.3	1	97	1	1	141564
G	R3+R4	4.2	1	51	1	0	42108	7.3	1	225	3	2	678032
G	1G6+1G7	5.0	1	10	0	0	7948	7.5	1	17	0	0	17543
G	L6+L7	4.9	1	10	0	0	7624	7.5	1	18	0	0	17555
G	R5+R6	4.3	1	11	0	0	7960	7.5	1	24	0	0	26478
L	1G4+1G5	4.2	3E-2	27	0	0	10898	7.3	2E-2	92	1	1	73391
L	L4+L5	4.6	3E-2	33	1	0	14145	7.4	3E-2	100	1	1	72694
L	R3+R4	4.1	2E-2	47	0	0	19580	7.2	2E-2	208	2	2	240743
L	1G6+1G7	5.3	-1E-2	11	0	0	7102	7.0	-6E-4	15	3	0	13036
L	L6+L7	5.5	1	11	1	0	7853	8.0	6E-2	19	0	0	14016
L	R5+R6	4.9	-3E-1	12	1	0	7838	7.0	-6E-3	20	1	0	14468
Prob. 6		Requested tolerance 1 E-4					Requested tolerance 1 E-7						
Method	sd	EE/GE	N	#F		#kev	sd	EE/GE	N	#F		#kev	
				tot	0					tot	0		
G	1G4+1G5	7.7	1	9	2	2	1805	8.7	1	21	3	3	8174
G	L4+L5	7.1	1	9	2	2	1648	8.6	1	21	3	3	6918
G	R3+R4	6.3	1	12	2	2	2392	8.4	1	38	4	4	19892
G	1G6+1G7	5.6	1	5	0	0	1223	10.1	1	8	2	2	3154
G	L6+L7	5.5	1	5	0	0	1152	10.1	1	8	2	2	2922
G	R5+R6	5.7	1	6	1	1	1573	10.2	1	9	2	2	3340
L	1G4+1G5	6.8	1	9	2	2	1277	9.2	1	21	3	3	4401
L	L4+L5	6.9	1	9	2	2	1315	9.1	1	22	3	3	4383
L	R3+R4	6.9	1	12	2	2	1549	8.7	1	40	4	4	10090
L	1G6+1G7	5.6	1	5	0	0	1019	10.0	1	8	2	2	2339
L	L6+L7	5.5	1	5	0	0	1050	9.9	1	8	2	2	2345
L	R5+R6	5.7	1	6	1	1	1260	9.7	1	9	2	2	2401

6. Conclusions. It is clear that, in general, higher order results in a better performance.

In most cases local error estimation used in the higher-order methods does not seem to yield a significant decrease in the total number of kernel evaluations. Since it is less reliable than global error estimate/control, which moreover supplies a global error estimate in the endpoint for free, it seems advisable to use global error control in the stepsize strategy.

TABLE 5.4
Average performance.

Method		All problems				Selected problems			
		Tol. 1 E-4		Tol. 1 E-7		Tol. 1 E-4		Tol. 1 E-7	
		sd	#kev	sd	#kev	sd	#kev	sd	#kev
G	G4+IG	5.3	941	8.3	11583	5.0	1023	8.5	16681
G	G8+IG	7.8	579	9.0	1162	6.8	667	8.5	1524
G	1G4+1G5	5.0	9861	7.4	7072	5.3	729	7.6	4334
G	L4+L5	4.7	2220	7.7	6104	5.1	694	7.8	3934
G	R3+R4	4.9	11942	7.3	17759	5.0	1135	7.5	11243
G	1G6+1G7	5.7	743	8.4	85557	5.9	538	8.7	1355
G	L6+L7	5.8	697	8.8	1297	6.1	522	8.8	1297
G	R5+R6	5.6	650	8.1	3277	5.6	650	8.5	1472
L	1G4+1G5	5.3	902	8.0	3995	5.6	642	8.3	3183
L	L4+L5	5.5	874	8.0	3721	5.8	660	8.4	2996
L	R3+R4	5.2	1355	7.8	9942	5.6	937	8.2	8209
L	1G6+1G7	6.2	708	8.3	1421	6.4	514	8.7	1100
L	L6+L7	6.2	709	8.4	1337	6.6	540	8.8	1049
L	R5+R6	5.7	712	8.2	1606	6.2	565	8.6	1231

TABLE 5.5
Results for a problem with a discontinuous solution.

Method		Requested tolerance 1 E-4						Requested tolerance 1 E-7					
		sd	EE/GE	N	#F		#kev	sd	EE/GE	N	#F		#kev
				tot	2.5					tot	2.5		
G	G4+IG	0.1	-6E-5	6	3	2	716	1.8	-1E-3	17	22	22	5912
G	G8+IG	2.2	5E-12	18	16	16	19728	2.2	5E-12	18	21	21	22248
G	1G4+1G5	0.1	-8E-5	5	2	2	780	1.9	2	18	22	22	10405
G	L4+L5	1.8	1	18	17	17	7770	1.8	1	18	22	22	8815
G	R3+R4	0.1	2E-4	7	12	3	2506	3.0	2E+1	18	22	22	8305
G	1G6+1G7	0.1	5E-9	5	2	2	1660	0.1	5E-9	5	2	2	1660
G	L6+L7	1.8	1	20	17	17	20303	1.8	1	20	22	22	22688
G	R5+R6	0.1	2E-7	5	2	2	1440	0.1	1E-7	6	2	2	1817
L	1G4+1G5	0.1	9E-7	6	3	2	841	1.9	2E-1	18	22	22	5701
L	L4+L5	1.8	1E-1	18	17	17	4693	1.8	1E-1	18	22	22	5303
L	R3+R4	0.1	4E-5	7	3	3	801	3.0	2	18	22	22	4587
L	1G6+1G7	0.1	-5E-10	5	2	2	1341	0.1	-5E-10	5	2	2	1341
L	L6+L7	1.8	8E-2	18	17	17	10673	1.8	8E-2	18	22	22	12043
L	R5+R6	0.1	9E-8	5	2	2	1174	0.1	3E-8	6	2	2	1407

If the solution is not polynomial, then the Gauss+Iterated collocation method is the most efficient when both computational time and storage are taken into account, especially for tolerances that are not too tight.

The uniform error control was of great benefit to all other methods. The approximations were more dependable with in most cases only a slight increase in computational work.

Appendix. The results of INTSOL were derived from Kunkel [17], all other results, except those of VEI, from Schlichte [19]. VEIca gives the results of VEI when the kernel has been approximated by a Chebyshev series. The number of terms, m , has been chosen as the minimal number such that the estimated error, as given by VEI, is less than the tolerance. VEIex gives the results when an exact decomposition was provided. In this case no global error estimate can be calculated.

Note that for VOLCON and for IVRKXC, “#kev” stands for the total number of K - plus f -evaluations where $k(t, s, y)$ is written as $K(t-s) \cdot f(s, y)$. For VElex, “#kev” represents the number of evaluations of the right-hand side of the system of ODE’s, “#Fev”, times the dimension of the system, “DIM”, times two. In the case of VEIca, #kev = #Kev + (DIM + 2) · #Fev, where $DIM = 2m + 1$.

TABLE A.1
Results of other published codes.

Prob	Method	Tol. 1 E-4		Tol. 1 E-7		Prob	Method	Tol. 1 E-4		Tol. 1 E-7			
		sd	# kev	sd	# kev			sd	# kev	sd	# kev		
1	IVRKX	7.2	510	9.0	1712	2	IVRKX	5.4	2124	7.7	5208		
	ORION	3.5	923	6.5	1624		ORION	2.0	1442	4.4	2998		
	VOLTEX	5.6	1124	8.7	4110		VOLTEX	2.0	3562	5.6	15092		
	INTSOL	4.0	2722	4.8	127964		INTSOL	1.1	64808				
	VEIca	4.1	1537	7.6	5684		VEIca	5.2	2527	7.6	7550		
	VElex	5.4	414	8.1	630		VElex	1.8	300	4.1	488		
	IVRKXC	7.2	359	9.0	868		IVRKXC	5.4	1041	7.7	1980		
	VOLCON	4.7	238	9.6	762		VOLCON	4.5	1262	8.3	2662		
3	IVRKX	5.0	23763	7.6	54320	4	IVRKX	4.9	1044	7.4	2704		
	ORION	6.3	42325	7.8	182344		ORION	4.7	1805	6.1	6411		
	VOLTEX	5.4	76882	7.6	417014		VOLTEX	5.3	5626	8.2	22918		
	VEIca	2.4	19943	---	----		VEIca	4.5	4157	8.5	12111		
	IVRKXC	5.0	4237	7.6	6358		IVRKXC	4.9	582	7.4	1052		
	VOLCON	6.6	2230	8.5	5370		VOLCON	4.8	366	7.3	1382		
	5	IVRKX	5.2	8625	8.9		21528	6	IVRKX	4.4	1029	9.7	4512
		ORION	5.6	2647	8.1		7114		ORION	6.0	402	8.7	675
VOLTEX		5.2	120752	8.1	160910	VOLTEX	2.5		634	6.3	2648		
INTSOL		3.6	337332			INTSOL	1.7		4312	2.0	9749		
VEIca		3.7	9225	---	----								
VElex		4.3	342	6.8	730								
IVRKXC		5.2	2167	8.9	3864								
VOLCON		6.5	2578	8.2	5170								

TABLE A.2
Average performance.

Method	Tol. 1 E-4		Tol. 1 E-7	
	sd	# kev	sd	# kev
IVRKX	5.4	738	8.4	2216
ORION	4.7	653	6.9	1643
VOLTEX	4.3	1636	7.4	6101
INTSOL	2.6	19278	3.4	68857
VEIca	4.0	1112	7.9	5420
VElex	4.1	232	6.5	373
IVRKXC	5.5	329	8.1	654
VOLCON	5.4	333	8.4	787

REFERENCES

- [1] H. ARNDT (1983), *On step size control for Volterra integral equations*, in Numerical Methods of Approximation Theory, ISNM, 67, L. Collatz, G. Meinardus and H. Werner, eds., Birkhäuser Verlag, Basel-Boston-Stuttgart, pp. 9-17.

- [2] H. ARNDT (1985), *An adaptive step size control for Volterra integral equations*, in Constructive Methods for the Practical Treatment of Integral Equations, ISNM, 73, G. Hammerlin and K.-H. Hoffmann, eds., Birkhäuser Verlag, Basel-Boston-Stuttgart, pp. 44-52.
- [3] G. BADER AND P. KUNKEL, *An adaptive multistep method for the solution of second kind Volterra integral equations*, in preparation.
- [4] J. G. BLOM AND H. BRUNNER, *Algorithm XXX: Discretized collocation and iterated collocation for nonlinear Volterra integral equations of the second kind*, ACM Trans. Math. Software, submitted.
- [5] J. M. BOWNS (1982), *Theory and performance of a subroutine for solving Volterra integral equations*, Computing, 28, pp. 317-332.
- [6] J. M. BOWNS AND L. APPELBAUM (1985), *Algorithm 627: A FORTRAN subroutine for solving Volterra integral equations*, ACM Trans. Math. Software, 11, pp. 58-65.
- [7] H. BRUNNER (1980), *Superconvergence in collocation and implicit Runge-Kutta methods for Volterra-type integral equations of the second kind*, in Numerical Treatment of Integral Equations, ISNM, 53, J. Albrecht and L. Collatz, eds., Birkhäuser Verlag, Basel-Boston-Stuttgart, pp. 54-72.
- [8] ——— (1984), *Iterated collocation methods and their discretization for Volterra integral equations*, SIAM J. Numer. Anal., 21, pp. 1132-1145.
- [9] H. BRUNNER, E. HAIRER AND S. P. NØRSETT (1982), *Runge-Kutta theory for Volterra integral equations of the second kind*, Math. Comp., 39, pp. 147-163.
- [10] H. BRUNNER AND P. J. VAN DER HOUWEN (1986), *The Numerical Solution of Volterra Equations*, North-Holland, Amsterdam.
- [11] H. BRUNNER AND S. P. NØRSETT (1981), *Superconvergence of collocation methods for Volterra and Abel integral equations of the second kind*, Numer. Math., 36, pp. 347-358.
- [12] R. P. DUNCAN (1982), *A Runge-Kutta method using variable stepsizes for Volterra integral equations of the 2nd kind*, Master's thesis, Technical Report 157/82, Department of Computer Science, University of Toronto.
- [13] E. HAIRER, CH. LUBICH AND M. SCHLICHTE (1985), *Fast numerical solution of nonlinear Volterra convolution equations*, this Journal, 6, pp. 532-541.
- [14] W. HOCK (1981), *An extrapolation method with step size control for nonlinear Volterra integral equations*, Numer. Math., 38, pp. 155-178.
- [15] F. DE HOOG AND R. WEISS (1975), *Implicit Runge-Kutta methods for second kind Volterra integral equations*, Numer. Math., 23, pp. 199-213.
- [16] H. M. JONES AND S. MCKEE (1985), *Variable step-size predictor-corrector schemes for second kind Volterra integral equations*, Math. Comp., 44, pp. 391-404.
- [17] P. KUNKEL (1982), *Ein adaptives Verfahren zur Lösung von Volterra'schen Integralgleichungen zweiter Art*, Diplomarbeit, Institut für Angewandte Mathematik, Universität Heidelberg.
- [18] R. K. MILLER (1971), *Nonlinear Volterra Integral Equations*, W. A. Benjamin, Menlo Park, CA.
- [19] M. SCHLICHTE (1984), *Anwendung eines impliziten Runge-Kutta-Verfahrens auf Volterra'sche Integralgleichungen zweiter Art mit Faltungskern*, Diplomarbeit, Institut für Angewandte Mathematik, Universität Heidelberg.