

An Attempt to Avoid Exact Jacobian and Nonlinear Equations in the Numerical Solution of Stiff Differential Equations

By Trond Steihaug and Arne Wolfbrandt

Abstract. A class of linear implicit methods for numerical solution of stiff ODE's is presented. These require *only* occasional calculation of the Jacobian matrix while maintaining stability. Especially, an effective second order stable algorithm with automatic stepsize control is designed and tested.

1. Introduction. During the last decade there has been a considerable amount of research on the numerical integration of stiff systems of ODE's. This work indicates that all efficient integration methods for such problems are implicit in character. This is due to the fact that only such methods have the required stability properties. Thus, the practical problem is not the stability restrictions, but the implicitness the need to avoid these give rise to. The relevant question is now, what is the cheapest type of implicitness we have to require.

Mainly, two different approaches to the implicitness can be found in the literature.

The first approach involves the numerical solution of nonlinear algebraic equations by the simplified Newton iteration. The simplification consists of treating the iteration matrix as piecewise constant (which means the use of an approximate Jacobian matrix). Examples of such an approach are semi-implicit Runge-Kutta formulas in Nørsett [4] and the formulas based on backward-differences in Gear [3].

Among recent methods proposed for numerical solution of stiff ODE's are the class of *modified Rosenbrock methods* introduced in Wolfbrandt [6]. When solving the system of equations

$$(1.1) \quad y' = f(y), \quad y(t_0) = y_0,$$

the formulas characterizing these methods are the following:

$$y_1 = y_0 + h \sum_{i=1}^{\nu} b_i k_i,$$
$$k_i = f\left(y_0 + h \sum_{j=1}^{i-1} a_{ij} k_j\right) + h J_0 \sum_{j=1}^i d_{ij} k_j, \quad i = 1, 2, \dots, \nu,$$

where J_0 denotes the Jacobian matrix $\partial f(y_0)/\partial y$.

Received February 22, 1978.

AMS (MOS) subject classifications (1970). Primary 65L05.

© 1979 American Mathematical Society
0025-5718/79/0000-0053/\$04.50

This is an example of the other approach to handle the implicitness. These methods maintain computational efficiency, since if such a method is applied to an m -dimensional system of stiff differential equations, then ν uncoupled systems of m linear equations (with the same matrix if $d_{ii} = d$, all i) will only have to be solved at each integration step. However, they suffer from the practical disadvantage of computing the Jacobian $\partial f/\partial y$ at $y = y_0$.

A natural question is now, is it possible to replace J_0 in the formulas above with an *arbitrary* real square matrix A (usually approximating J_0)? An affirmative answer to this question will be given in this paper. In fact, the W -method introduced in the next section, combines the good things of the two approaches mentioned above.

2. The W -Methods. We consider numerical integration of system (1.1) using a class of methods of the form

$$y_1 = y_0 + h \sum_{i=1}^{\nu} b_i k_i,$$

$$W(h, d_{ii}, A)k_i = f\left(y_0 + h \sum_{j=1}^{i-1} a_{ij}k_j\right) + hA \sum_{j=1}^{i-1} d_{ij}k_j, \quad i = 1, 2, \dots, \nu,$$

where $W(h, d_{ii}, A) = I - hd_{ii}A$, and A is a real square matrix such that $W(h, d_{ii}, A)$ is invertible. The methods above will be called ν -stage W -methods.

Remark. We note that if we choose $A \equiv 0$, or $A \equiv J_0$, then the W -methods reduce to the classical explicit Runge-Kutta methods and the modified Rosenbrock methods, respectively.

3. Order Conditions. In this section we will consider the order conditions for the W -methods.

We confine our attention first to the case $A \equiv 0$. Thus, the set of order conditions for a W -method includes those for a classical Runge-Kutta method. This demonstrates clearly the relation between these methods.

Next we turn to the case $A \equiv J_0$, corresponding to the modified Rosenbrock methods. The order conditions for these methods can originally be found in Wolfbrandt [6].

We combine these two sets of order conditions and further observe that the matrices A and J_0 do not generally commute. The order conditions for the W -methods are then readily obtainable. For illustration, the eight conditions for third order accuracy are given in Table I together with the associated 'elementary differentials' and their orders. To simplify these conditions, we have introduced the following notations:

$$M_i = \sum_{j=1}^{i-1} a_{ij}, \quad i = 2, 3, \dots, \nu,$$

$$N_i = \sum_{j=1}^i d_{ij}, \quad i = 1, 2, \dots, \nu.$$

TABLE I
Equations of conditions for a ν -stage W -method

order	'elementary differential'	condition
1	f	$\sum_{i=1}^{\nu} b_i = 1$
2	$f'f$	$\sum_{i=2}^{\nu} b_i M_i = 1/2$
2	Af	$\sum_{i=1}^{\nu} b_i N_i = 0$
3	$f'f'f$	$\sum_{i=2}^{\nu} b_i \sum_{j=2}^{i-1} a_{ij} M_j = 1/6$
3	$f' Af$	$\sum_{i=2}^{\nu} b_i \sum_{j=1}^{i-1} a_{ij} N_j = 0$
3	$Af'f$	$\sum_{i=2}^{\nu} b_i \sum_{j=2}^i d_{ij} M_j = 0$
3	AAf	$\sum_{i=2}^{\nu} b_i \sum_{j=1}^i d_{ij} N_j = 0$
3	$f''ff$	$\sum_{i=2}^{\nu} b_i M_i^2 = 1/3$

We shall now give an upper bound for the maximum order of a W -method. In preparation for this the following result in Wolfbrandt [7] is necessary. Let $\Pi_r = \{\text{polynomials of degree } \leq r\}$ and

$$R_m^n(u) = N_n(u)/D_m(u), \quad \text{where } N_n \in \Pi_n \text{ and } D_m \in \Pi_m.$$

LEMMA 1. *The maximum attainable order of a rational approximation $R_m^n(u)$ to $\exp(u)$ with only real poles is $m + 1$. \square*

THEOREM 2. *The maximum order for a ν -stage W -method is at most $\nu + 1$.*

Proof. First, it is clear that a W -method produces a rational approximation to $\exp(u)$ with only real poles when applied to the scalar equation $y' = \lambda y$.

Therefore, since the maximum order is attained when $A \equiv J_0$, the theorem follows from Lemma 1. \square

We recall from Section 1 that it is desirable from a practical point of view to choose $d_{ii} = d$, all i . The following surprising result, which is proved in Nørsett and Wolfbrandt [5], also shows that this choice is optimal for $A \equiv J_0$ in the sense of minimizing the absolute value of the error constant for a linear problem.

LEMMA 3. *Let $R_m^n(u)$ be a rational approximation to $\exp(u)$ of order $m + 1$ with only real poles, i.e.*

$$R_m^n(u) - \exp(u) = Cu^{m+2} + O(u^{m+3}).$$

Then the absolute value of the error constant C attains its local minimum values when all the poles are equal. \square

The main conclusion of this discussion is that the ‘best’ we can do is choose $d_{ii} = d$, all i . This is also confirmed in Wolfbrandt [6] by computation of acceptability for the rational approximations mentioned in Lemma 3.

4. Explicit Formulas. Guided by the conclusions in Section 3 we will only consider W -methods with $d_{ii} = d$, all i . Moreover, we temporarily regard d as fixed.

To match all terms up to second order for a 2-stage W -method, the coefficients of the formula (with one free parameter) must satisfy the following set of equations:

$$b_1 + b_2 = 1, \quad b_2 a_{21} = 1/2, \quad b_2 d_{21} = -d.$$

A natural way is to choose the free parameter, so that the method reduces to the *improved Euler* method when $A \equiv 0$, i.e.

$$y_1 = y_0 + \frac{h}{4} (k_1 + 3k_2), \quad W(h, d, A)k_1 = f(y_0),$$

$$W(h, d, A)k_2 = f\left(y_0 + \frac{2}{3}hk_1\right) - \frac{4}{3}hdAk_1,$$

where $W(h, d, A) = I - hdA$. The local truncation error T for this method may be written as

$$T = \frac{h^3}{3} \{f'f'f - 3d(f'Af + Af'f) + 6d^2AAf\}(y_0) + O(h^4).$$

We easily see that it is impossible to construct a third order 3-stage W -method. However, we can obtain a second order 3-stage formula with the local truncation error of the form

$$T = \frac{1}{12} h^3 f'f'f(y_0) + O(h^4).$$

It is interesting to note that the matrix A does not occur in the h^3 -term.

For a 4-stage method, we can of course match up to and including h^3 -terms. In fact, the parameters can be chosen so that the method only requires three function evaluations.

Remark. We observe that a ν -stage W -method with less than ν different function evaluations can be seen as a method with powers of $W(h, d, A)$ in its formula.

We also note that a 4-stage W -method can be based on the classical fourth order Runge-Kutta method. In order to implement a W -method with variable stepsize it is necessary to compute the local truncation error. A device, proposed by England [1], allows us to estimate this error. The basic idea is to add extra stages to the method so that the new method is a more accurate one.

We consider the 2-stage W -method based on the *improved Euler* method. A third order estimate of the local truncation error is then of the form

$$T = \frac{h}{8} (k_1 - 5k_2 + 5k_3 - k_4) + O(h^4),$$

where the extra stages k_3 and k_4 satisfy

$$W(h, d, A)k_3 = f(y_1), \quad W(h, d, A)k_4 = f\left(y_1 + \frac{2}{3}hk_3\right) + hdA\left(\frac{2}{3}k_1 + 6k_2\right).$$

We adopt the notation (2, 4)- W for the above method with built-in error estimate. In a sequence of accepted integration steps of equal size this method will require only two function evaluations per step.

Our final task in this section is to extend (2, 4)- W to solve the *nonautonomous* system $y' = f(y, t)$, $y(t_0) = y_0$. We extend by appending $t_0 + ha_i$ to the function evaluation in each stage k_i , $i = 1, 2, 3$, and 4. The corresponding set of order conditions in Table I has to be supplemented with the following ‘nonautonomous’ one

$$(f_t) \quad b_1a_1 + b_2a_2 = 1/2,$$

$$(f_t) \quad c_1a_1 + c_2a_2 + c_3a_3 + c_4a_4 = 1/2,$$

$$(f_{tt}) \quad c_1a_1^2 + c_2a_2^2 + c_3a_3^2 + c_4a_4^2 = 1/3,$$

$$(f_y f_t) \quad c_2a_{21}a_1 + c_3(a_{31}a_1 + a_{32}a_2) + c_4(a_{41}a_1 + a_{42}a_2 + a_{43}a_3) = 1/6,$$

$$(Af_t) \quad c_2d_{21}a_1 + c_3(d_{31}a_1 + d_{32}a_2) + c_4(d_{41}a_1 + d_{42}a_2 + d_{43}a_3) = -d/2,$$

$$(f_{y_t} f) \quad c_2M_2a_2 + c_3M_3a_3 + c_4M_4a_4 = 1/3,$$

where the associated ‘elementary differentials’ are enclosed in parentheses. Thus, it is convenient to choose $a_1 = 0$ and $a_i = M_i$, $i = 2, 3$, and 4.

5. Stability. $A(\alpha)$ -stability, $L(\alpha)$ -stability (or stiff stability) and its weaker associate strong stability, have become generally accepted as appropriate properties of numerical methods suitable for solving stiff ODE’s.

Definition. A numerical method is said to be $A(\alpha)$ -stable, $\alpha \in (0, \pi/2)$, if all its solutions $\{y_m\}$ tend to zero, as $m \rightarrow \infty$, when the method is applied with fixed positive h to the test equation

$$(1.2) \quad y' = \lambda y, \quad \lambda \in S(\alpha),$$

where $S(\alpha) = \{z \in C: z \neq 0, |\arg(-z)| < \alpha\}$.

A numerical method is said to be *A-stable* (*A(0)-stable*), if it is *A(α)-stable* for all (some) $\alpha \in (0, \pi/2)$.*

Definition. An *A(α)-stable* method, $\alpha \in (0, \pi/2)$, is said to be *strongly stable at infinity* if all its solutions $\{y_m\}$, when the method is applied with fixed positive h to the test equation (1.2), satisfy

$$\lim_{|\lambda| \rightarrow \infty} |y_{m+1}/y_m| \leq c < 1 \quad \text{when } \lambda \in S(\alpha).$$

If the real constant c is equal to zero, then the method is called *L(α)-stable* or *stiffly stable*. A numerical method is said to be *L-stable* (*L(0)-stable*), if it is *L(α)-stable* for all (some) $\alpha \in (0, \pi/2)$.

We shall now discuss stability properties of the *W*-methods. For simplicity, we assume that the matrices A and J_0 commute and that they are diagonalizable. Then the archetypical initial value problem is that in (1.2).

As mentioned before, a *W*-method results in a rational approximation to the exponential function with only real poles when applied to (1.2). Therefore, we first present in Table II some of the stability results for the modified Rosenbrock methods, i.e. the case $A \equiv J_0$, given in Wolfbrandt [6]. We have listed in this table acceptability angles and error constants for rational approximations.

$$R_{\nu}^{\nu-1}(u) \equiv N_{\nu-1}(u)/D_{\nu}(u), \text{ where } D_{\nu}(u) = (u - \gamma)^{\nu},$$

to $\exp(u)$ of maximum order (which means that γ is a zero to the Laguerre polynomial of degree ν and $d_{ii} = 1/\gamma$). The entries in the table represent the pair (acceptability angle; error constant) corresponding to $\gamma_{\nu}^{(k)}$, the k th zero of the Laguerre polynomial of degree ν (the zeros are assumed to be ordered in increasing value).

TABLE II
L-acceptability angles and error constants

$\begin{matrix} k \\ \nu \end{matrix}$	1	2	3	4
1	($\pi/2; 50E-2$)			
2	($\pi/2; 14E-1$)	($\pi/2; 40E-3$)		
3	($1.56; 64E-1$)	($\pi/2; 26E-3$)	($1.32; 39E-4$)	
4	($1.54; 42E+0$)	($\pi/2; 27E-3$)	($1.56; 11E-4$)	($-; 36E-5$)

It would be nice to be able to derive the stability properties also for the case A different from J_0 . We shall content ourselves to discuss the *W*-method based on the *improved Euler* method.

* y_m is an approximation to $y(t_m)$, where $t_m = t_0 + mh, m = 0, 1, \dots$

We consider the scalar equation (1.2). Let the eigenvalue of A be denoted by a . Then the following constraints for this method yield

$$A(0)\text{-stable} \quad \text{if and only if } |da| \geq |\lambda|/4 - 1/(2h),$$

$$L(0)\text{-stable} \quad \text{if and only if } |da| = (1 \pm \sqrt{1/2})|\lambda|.$$

It seems possible to handle these stability requirements by examining the local truncation error. Numerical experience shows that this works well for problems with transient to smooth components, while perhaps a more sophisticated control is necessary for problems with smooth to transient components.

6. A Comparison. The equations defining a *semi-implicit Runge-Kutta method* are

$$y_1 = y_0 + h \sum_{i=1}^{\nu} b_i k_i,$$

$$k_i = f\left(y_0 + h \sum_{j=1}^{i-1} a_{ij} k_j + h d_{ii} k_i\right), \quad i = 1, 2, \dots, \nu.$$

This method, as well as a W -method with $A \equiv J_0$, produces a rational approximation to $\exp(h\lambda)$ with only real poles when applied to the test equation $y' = \lambda y$. In view of this relation a brief comparison between these methods will be done in this section.

The semi-implicit Runge-Kutta methods are implicit in character involving solution of ν uncoupled nonlinear equations of the form

$$F(k_i) \equiv k_i - f\left(y_0 + h \sum_{j=1}^{i-1} a_{ij} k_j + h d_{ii} k_i\right) = 0, \quad i = 1, 2, \dots, \nu,$$

by the *simplified Newton iteration*

$$W(h, d_{ii}, J_0)k_i^{(m)} = f\left(y_0 + h \sum_{j=1}^{i-1} a_{ij} k_j^{(M_j)} + h d_{ii} k_i^{(m-1)}\right) - h d_{ii} J_0 k_i^{(m-1)},$$

$$m = 1, 2, \dots, M_i \text{ and } i = 1, 2, \dots, \nu,$$

where $k_i^{(0)}$ is some suitable starting value and M_i is the maximum number of iterations required.

It is a common practice to employ a linear combination of known iterated values to provide $k_i^{(0)}$ which we write as

$$k_i^{(0)} = \sum_{j=1}^{i-1} \frac{d_{ij}}{d_{ii}} k_j^{(M_j)}.$$

Thus, the modified Rosenbrock method in Wolfbrandt [6] can be regarded as a linearization of a semi-implicit Runge-Kutta method with the above choice of $k_i^{(0)}$. Moreover, replacing J_0 by A in the iteration formula above we obtain a W -method.

The traditional motivation behind the choice of the parameters in the expression for $k_i^{(0)}$ is based on obtaining a good initial approximation to k_i . A somewhat more attractive approach, in our opinion, is to avoid the necessity to iterate. This will be illustrated in the following example.

Example. We consider the *generalized midpoint-rule*, alias the θ -method

$$y_1 = y_0 + hk, \quad k = f(y_0 + h\theta k).$$

After two iterations it yields

$$W(h, \theta, J_0)k_1 = f(y_0), \quad W(h, \theta, J_0)k_2 = f(y_0 + h\theta k_1) - h\theta J_0 k_1,$$

where $k_m \equiv k^{(m)}$, $m = 1$ and 2 ($k_0^{(0)} \equiv 0$). Putting $\tilde{y}_1 = y_0 + hk_2$, we obtain a 2-stage W -method.

We recall from Section 3 that the order of \tilde{y}_1 is less than 2 for $\theta \neq 1/2$, otherwise 2.

First, we will assume that an exact Jacobian matrix J_0 is used in the simplified Newton iteration. Then there is no gain in stability to iterate. Moreover, numerical experiments indicate that only a small number of iterations are required in the θ -method. We easily verify that the dominating local truncation error terms for this method are independent of the number of iterations. In fact, the local truncation error can be expressed as

$$T = h^2 \left(\theta - \frac{1}{2} \right) f'f(y_0) + h^3 \left\{ \left(\theta^2 - \frac{1}{3} \right) f''ff + \left(\theta^2 - \frac{1}{6} \right) f'f'f \right\} (y_0) + O(h^4).$$

This shows again that it becomes no great benefit in making several iterations when an exact Jacobian matrix J_0 is used.

We observe that the W -method based on the *improved Euler* method has a local truncation error of the form

$$T = h^3 \left(\frac{1}{6} - d + d^2 \right) f'f'f(y_0) + O(h^4) \quad \text{when } A \equiv J_0.$$

In the asymptotic region many stiff systems behave as linear systems. Therefore, an approximation to the Jacobian matrix J_0 in the simplified Newton iteration can be used for a large number of integration steps. Accordingly, we next assume that J_0 is replaced by an approximation A and consider the effect of iterations in the θ -method. As noted previously, this method with M iterations may be regarded as an M -stage W -method with $a_{ij} = -d_{ij} = \theta$, for $j = i - 1$, $i = 1, 2, \dots, M$ and $b_i = 1$ for $i = M$, otherwise zero. This method will in the following be called *the degenerated W -method*.

In contrast to the case of exact Jacobian the dominating local truncation error constants are now dependent of the first two iterations, but *only of these*. The error constants together with the associated 'elementary differentials' and the number of iterations are listed in Table III below.

From this important result we conclude that iteration is one way to eliminate the influence on accuracy of an approximate Jacobian matrix. However, the special 3-stage W -method in Section 4 has order 2 for all d and share the above-mentioned properties of the degenerated 3-stage method too.

The changes due to iterations of the $A(0)$ -stability contour for the θ -method is illustrated in Figures 6.1–6.2 for $\theta = 1/2$ (maximum order) and $\theta = 1$ (L -stability), respectively. In these figures the convergence constraint is also drawn. To

interpret the figures we note that the stability regions are to the left of the contours.

TABLE III
Error constants

m	$f'f$	Af	$f'f'f$	$Af'f$	$f'Af$	AAf	$f''ff$
1	$-1/2$	θ	$-1/6$	0	0	θ^2	$\theta^2-1/3$
2	$\theta-1/2$	0	$-1/6$	θ^2	θ^2	$-\theta^2$	$\theta^2-1/3$
3	$\theta-1/2$	0	$\theta^2-1/6$	0	0	0	$\theta^2-1/3$
\cdot	\cdot	\cdot	\cdot	\cdot	\cdot	\cdot	\cdot
∞	$\theta-1/2$	0	$\theta^2-1/6$	0	0	0	$\theta^2-1/3$

From these figures we conclude that 'convergence' is stronger (weaker) than 'stability' for $\theta = 1$ ($1/2$). Further, the 'stability' region for the 2-stage W -method based on the improved Euler method contains the regions in the figures for $d \geq 1/2$.

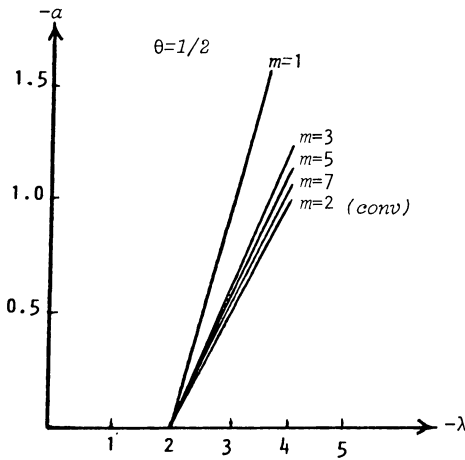


FIGURE 6.1

$A(0)$ -stability contour

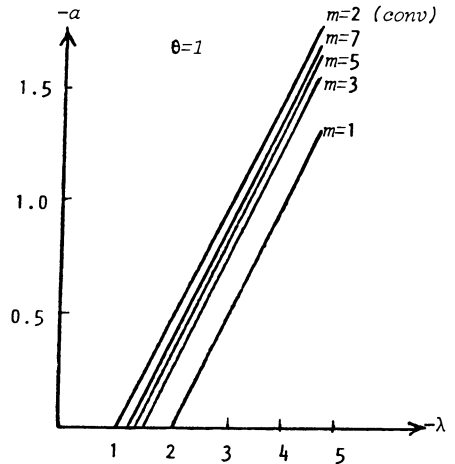


FIGURE 6.2

$A(0)$ -stability contour

Finally, for the θ -method with $\theta = 1$ it is possible to iterate with an approximate Jacobian and yet obtain L -stability, while this is not the case for a W -method.

7. Numerical Examples. The intention of this section is to present an algorithm for solving stiff ODE's. The underlying method is the (2, 4)- W -method with $d = 1 - \sqrt{1/2}$, i.e. $d = 1/\gamma_2^{(2)}$ (see Table II). The design of such an algorithm is a very complex work involving numerous choices in equation solver, error estimation and stepsize control as well as features of detailed programming. We have not attempted to construct a complete algorithm. However, we have concentrated on some of the difficulties associated with the stepsize control and the reevaluation of the matrix A .

The linear systems in the method have the form

$$W(h, d, A)k = f, \quad \text{where } W(h, d, A) \equiv I - dhA.$$

These are solved by Gaussian elimination, i.e. *LU*-decomposition and back-substitution. Therefore, during the computation the *W*-matrix is kept as long as possible and is re-evaluated only after a failure of a specified accuracy requirement for the local truncation error.

The purpose of the control of stepsize is to integrate as efficiently as possible keeping the error under a prescribed level. The optimal stepsize is not available as it depends on the entire solution. Therefore, the stepsize is based on the local truncation error. The algorithm tests that the error-per-step is less than a tolerance requested. The strategy for optimal stepsize is calculated as shown below.

We denote the norm of the error estimate with current stepsize *h* by *E(h)*. Then the optimal stepsize is *ph*, *p* > 0, satisfying *E(ph)* = τ , where τ is the error tolerance. Since the method is of order 2, then $p = (\tau/E(h))^{1/3}$. This strategy is critical to the efficiency, since changes in stepsize will lead to extra *LU*-decompositions, Jacobian evaluations and function evaluations. Therefore, an increase in the stepsize should only be attempted unless $p \geq q > 1$. Also, when changing stepsize upward or downward, the new value would be set at *rh*, for some $r < p$. Finally, before increasing the stepsize it is very economical to retain the old value one step further without making any error estimation.

The strategy we have adopted is based on a combination of heuristics and numerical experiments. The matrix *W* is recalculated only when the stepsize is changed or when the matrix *A* is altered and set equal to the Jacobian matrix. The latter happens only when $E(h) > 0.7\tau$ and $A \neq J_0$. The table below gives the appropriate actions.

TABLE IV
Values of the parameter *r*

Error	$A=J_0$	$A \neq J_0$
$E(h) < \tau$	$\max\{1, 0.9[p]\}$ <i>Accepted step</i>	$\{1, 0.9[p]\}$ <i>Accepted step</i>
$\tau < E(h) < 2\tau$	0.85 <i>p</i> <i>Rejected step</i>	1 <i>Accepted step</i>
$E(h) \geq 2\tau$	0.85 <i>p</i> <i>Rejected step</i>	0.80 <i>p</i> <i>Rejected step</i>

In order to avoid rejection of steps and for stability reasons we have also found it convenient to change the stepsize to $0.9ph$ when $0.8\tau \leq E(h) \leq 2\tau$ in the first step with an 'old' Jacobian matrix. Numerical results for the (2, 4)-*W*-method on four problems are given below. For comparison we have also included results from Nørsett [4], Enright, Hull and Lindberg [2].

Problem 1.

$$y_1' = 0.01 - (1 + (y_1 + 1000)(y_1 + 1))(0.01 + y_1 + y_2), \quad y_1(0) = 0,$$

$$y_2' = 0.01 - (1 + y_2^2)(0.01 + y_1 + y_2), \quad y_2(0) = 0, \quad t \in [0, 100].$$

A relative error measure based on a weighted Euclidean norm is used with a vector of weights made up of the numerically greatest solution points so far encountered for each component.

TABLE V
*Number of function evaluations, Jacobian evaluations
 and steps for Problem 1*

τ	Eh1e1A3	GEAR	IRKSUB4
1E-3	892/50/66	181/20/69	185/12/30
1E-5	1908/97/144	380/36/136	367/25/54

τ	SIRSPN	(2,4)-W
1E-3	250/40/-	92/14/45
1E-5	568/39/-	378/34/181

Problem 2.

$$y' = UBUy + Uw, \quad y(0) = [0, -2, -1, -1]^T,$$

where

$$B = \begin{pmatrix} -10^4 & 10^4 & 0 & 0 \\ -10^4 & -10^4 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -10^{-3} \end{pmatrix},$$

U is the unitary matrix with diagonal elements equal to $-1/2$ and all other elements equal to $1/2$,

$$w = [(z_1^2 - z_2^2)/2, z_1 z_2, z_3^2, z_4^2]^T, \quad z = Uy, \quad t \in [0, 100].$$

The same relative error measure as in Problem 1 is used.

TABLE VI

Number of function evaluations and Jacobian evaluations for Problem 2

τ	GEAR	IMPEX2	IRKSUB4	SIRSPN	(2,4)-W
1E-3	581/40	548/22	479/23	334/26	190/29
1E-5	959/45	1032/24	1095/34	1202/27	694/33

Problem 3.

$$y_1' = -0.04y_1 + 0.01y_2y_3, \quad y_1(0) = 1,$$

$$y_2' = 400y_1 - 100y_2y_3 - 3000y_2^2, \quad y_2(0) = 0,$$

$$y_3' = 30y_2^2, \quad y_3(0) = 0, \quad t \in [0, 40].$$

For this problem the error tolerance $\tau \equiv \tau(t)$ is defined by $\tau(t) = 6he/(t - t_0)$, where ϵ is the global tolerance and $t - t_0$ is the distance from the initial point to the end of the current step (see Enright et al. [2]). Moreover, the maximum norm is used.

TABLE VII

Number of function evaluations, Jacobian evaluations, LU-decompositions, and steps for Problem 3

ϵ	GEAR	SDBASIC	TRAPEX
1E-2	131/14/14/55	241/176/92/28	237/29/29/11

ϵ	GENRK	IMPRK	(2,4)-W
1E-2	99/9/27/9	387/89/89/11	91/15/15/41

Problem 4. This arises from solving parabolic partial differential equations. Consider the parabolic equation

$$\frac{\partial u}{\partial t} - \frac{\partial}{\partial x} \left(a \left(\frac{\partial u}{\partial x} \right) \frac{\partial u}{\partial x} \right) = 0$$

in the region $[0, 1] \times [0, \infty]$ subject to the initial condition $u(x, 0) = u_0(x)$, $x \in (0, 1)$ and the boundary condition $u(0, t) = u(1, t) = 0$, $t \geq 0$. Let $a(z) = \nu z^2 + 1$, $\nu \geq 0$ and $u_0 \equiv 1$. Finite element approximation in the space variable x results in the semidiscrete Galerkin approximation, which is a stiff system of N ODE's of the form

$$B \frac{dU}{dt} + K(U)U = 0, \quad t > 0, \quad U(0) = U_0,$$

where B and K are the mass- and stiffness-finite element matrices, respectively. Here we use linear splines with equidistant knots. An appropriate norm is the weighted Euclidean norm with all weights equal to Δx , where $\Delta x = 1/(N + 1)$.

Numerical results for the (2, 4)- W -method on this problem are shown in Table VIII and Table IX for $\nu = 0$ (which corresponds to a linear problem) and $\nu = 1$, respectively. For illustration we have presented the numerical solution in the following adequate points (x, t) : $x = 0.1, 0.3$ and $t = 0.01, 1, 0$.

Remark. The number of LU -decompositions and the number of Jacobian evaluations is nearly identical for the (2, 4)- W -method on the problems given. In addition, no extra function evaluations are spent on the stepsize control. The only price to be paid for the stepsize control is the cost of one back-substitution.

8. Conclusions. The aim was to introduce a new class of linear implicit methods without exact Jacobian for solving stiff ODE's. We do not claim to have found an ideal algorithm for such problems, but the theoretical and numerical results achieved for the introduced (2, 4)- W -method merit a further investigation of higher order W -methods.

TABLE VIII
Results for (2, 4)- W on Problem 4 with $\nu = 0$

<i>Numerical Solution</i>					
		$t=0.01$		$t=1.0$	
Δx	τ	$x=0.1$	$x=0.3$	$x=0.1$	$x=0.3$
1/10	1E-3	0.4910	0.9646	0.1141E-4	0.2991E-4
1/20	1E-3	0.5100	0.9656	0.1364E-4	0.3570E-4
	1E-5	0.5101	0.9653	0.2163E-4	0.5664E-4
1/40	1E-5	0.5167	0.9655	0.2237E-4	0.5857E-4
	1E-6	0.5185	0.9662	0.2126E-4	0.5565E-4
<i>Exact Solution</i>					
		0.5187	0.9669	0.2035E-4	0.5328E-4
<i>Number of Function Evaluations, Jacobian Evaluation and Steps</i>					
Δx	τ				
1/10	1E-3	18/3/8		61/10/30	
1/20	1E-3	25/4/12		71/11/25	
	1E-5	100/5/49		280/14/139	
1/40	1E-5	136/8/67		318/17/158	
	1E-6	282/8/140		658/17/328	

TABLE IX
Results for (2, 4)-W on Problem 4 with $\nu = 1$

Numerical Solution					
		t=0.01		t=1.0	
Δx	τ	x=0.1	x=0.3	x=0.1	x=0.3
1/10	1E-3	0.3023	0.9508	0.6531E-5	0.1662E-4
1/20	1E-3	0.3034	0.8066	0.4835E-5	0.1267E-4
	1E-5	0.3020	0.8041	0.1503E-4	0.3935E-4
1/40	1E-5	0.3023	0.8041	0.1344E-4	0.3518E-4
	1E-6	0.3026	0.8045	0.1266E-4	0.3315E-4
Number of Function Evaluations, Jacobian evaluations and Steps					
Δx	τ	t=0.01		t=1.0	
1/10	1E-3	38/6/18		85/14/42	
1/20	1E-3	56/10/27		99/17/49	
	1E-5	264/14/130		476/23/236	
1/40	1E-5	380/20/186		610/29/301	
	1E-6	851/17/423		1309/27/652	

School of Organization and Management
Yale University
Box 1A
New Haven, Connecticut 06520

Department of Computer Sciences
Chalmers University of Technology
Fack
S-402 20 Göteborg, Sweden

1. R. ENGLAND, "Error estimates for Runge-Kutta type solutions to systems of ordinary differential equations," *Comput J.*, v. 12, 1969, pp. 166-170.
2. W. H. ENRIGHT, T. E. HULL & B. LINDBERG, "Comparing numerical methods for stiff systems of o.d.e.s," *BIT*, v. 15, 1975, pp. 1-10.
3. C. W. GEAR, *Numerical Initial Value Problems in Ordinary Differential Equations*, Prentice-Hall, Englewood Cliffs, N. J., 1971.
4. S. P. NØRSETT, *Semi-Explicit Runge-Kutta Methods*, Technical Report 6, Dept. of Math., Univ. of Trondheim, 1974.
5. S. P. NØRSETT & A. WOLFBRANDT, "Attainable order of rational approximations to the exponential function with only real poles," *BIT*, v. 17, 1977, pp. 200-208.
6. A. WOLFBRANDT, *A Study of Rosenbrock Processes with Respect to Order Conditions and Stiff Stability*, Thesis, Chalmers Univ. of Technology, Göteborg, Sweden, 1977.
7. A. WOLFBRANDT, "A note on a recent result of rational approximation to the exponential function," *BIT*, v. 17, 1977, pp. 367-368.