

5 AUG/80.

CS-80-09

Do NOT REMOVE

-THIS report will appear in:

SIAM J. ON NUMERICAL

ANALYSIS, VOL 17, #5, OCT/80

-CONTACT AUTHORS FOR COPY.

The Extrapolation of First Order Methods for  
Parabolic Partial Differential Equations II

A. R. Gourlay <sup>†</sup>  
and  
J. L. Morris <sup>††</sup>

<sup>†</sup> IBM Scientific Centre, Athelstan House, Winchester, England

<sup>††</sup> Department of Computer Science, University of Waterloo, Ontario, Canada

## 1. Introduction

In a recent paper, Lawson and Morris [5] described second order algorithms for solving parabolic partial differential equations in several space variables. The essential motivation there was to derive  $L_0$ -stable methods of second order which would remove constraints on the discrete time step which are present in A-stable methods like, for example, the Crank Nicolson method. These constraints occur essentially because the parabolic differential equation gives rise to a system of ordinary differential equations which are stiff (see Lambert [4]). In all the cases considered in [5], the basic algorithm used was the first order Backward Euler method. To obtain second order accuracy, a combination of results obtained by applications of the Backward Euler method using time step lengths of different size was proposed so that the  $L_0$ -stability property of the Backward Euler method is retained by the resulting second order method.

In the current paper, the authors consider generalizations of this approach to achieve higher order methods which retain the property of  $L_0$ -stability. The approach adopted is essentially one of extrapolation where a low order method (the first order Backward Euler method or  $\theta$ -method) is applied on a sequence of different sized time steps and a linear combination of the results used to match the required number of terms in the Taylor expansion of the theoretical solution (of the system of ordinary differential equations arising from the spatial discretization). The idea of extrapolation for increased order is, of course, not new. Richardson extrapolation [8] is now classic as is its generalizations to Ordinary Differential Equations (see [4]). Recent alternative approaches to increasing the accuracy of the time

integration associated with Galerkin methods for solving parabolic differential equations include the defect corrections method of Zadunaisky [10] and Stetter [9], Frank [2]; the deferred correction of Fox [1] has received attention in Pereyra [6] and recently Saylor [7] has made a comprehensive study of all these for linear parabolic problems.

x ✓ The present approach differs in detail, not principle, in the manner in which the increased accuracy is achieved<sup>†</sup>. However, one essential requirement that is imposed is that the resulting high order methods be  $L_0$ -stable, a feature not present in the previously described applications.

We restrict our attention to the parabolic equation in one space variable although in principle the idea carries over to many space variables provided an efficient algorithm exists for solving the sparse matrix system which arises, e.g. [3] George. A further assumption made in this paper is that the coefficient of diffusion is constant. If it depends on the space variable(s) no change to the given algorithm is necessary. Time dependent (and nonlinear) coefficients will be considered at a later date.

---

<sup>†</sup> We further believe there are applications in the area of stiff methods for ordinary differential equations and applications to systems of hyperbolic equations; these applications will form the basis of future papers.

## 2. Second Order $L_0$ -Stable Methods

Consider the constant coefficient homogeneous parabolic differential equation

$$(2.1) \quad \frac{\partial u}{\partial t} = Lu \quad (x,t) \in I \times [0 < t \leq T]$$

and  $I = [0,1]$ , say. Equation (2.1) is subject to the initial condition

$$u(x,0) = f(x) \quad x \in I, \quad f \text{ given,}$$

and boundary conditions

$$u(x,t) = 0 \quad x \in \delta I, \text{ the end points of } I.$$

In the usual manner introduce a uniform discretization  $h$  of  $I$  and denote  $x = ih$ ,  $i$  a nonnegative integer. Replace  $t$  by  $m\tau$  where  $\tau$  is a constant time step and  $m$  is a nonnegative integer.

On the resulting set of points we replace in the usual way, the spatial derivatives in  $L$  by divided differences. So, if  $L \equiv \frac{\partial^2}{\partial x^2}$  we might replace  $L$  by  $\frac{\delta x^2}{h^2}$  where  $\delta x$  is the usual central difference operator. Applying the resulting difference scheme at each spatial grid point, equation (2.1) is then replaced by a system of ordinary differential equations. For example, with the above replacement of  $L$ , assuming there are  $N$  grid point interior to  $I$ , then the system of ordinary differential equations would be

$$(2.2) \quad \frac{d\mu}{dt} = A\mu$$

$\mu = \mu(t)$  an  $N$ -vector of unknowns (the approximations to  $u(t)$  at each of the grid points) and

$$A = \frac{1}{h^2} \begin{bmatrix} -2 & 1 & & & 0 \\ 1 & -2 & 1 & & \\ & 1 & -2 & 1 & \\ & & \ddots & \ddots & \ddots \\ 0 & & & \ddots & \ddots & 1 \\ & & & & 1 & -2 \end{bmatrix}$$

(has dimension  $N$ ).

The theoretical solution of (2.2) is of course

$$\underline{\mu}(t+\tau) = \exp(\tau A)\underline{\mu}(t)$$

$$\underline{\mu}(0) = \underline{f} \quad \text{- the vector of initial values.}$$

The first order accurate Backward Euler method then defines an approximation  $\underline{v}(t+\tau)$  to  $\underline{\mu}(t+\tau)$  by

$$(2.3) \quad (I - \tau A)\underline{v}(t+\tau) = \underline{v}(t)$$

$$\underline{v}(0) = \underline{f}.$$

In contrast, the  $\theta$ -method can be used to approximate  $\underline{\mu}(t+\tau)$  and is given by

$$(2.4) \quad [I - \tau(1-\theta)A]\underline{v}(t+\tau) = [I + \tau\theta A]\underline{v}(t)$$

which reproduces (2.3) when  $\theta = 0$ . Equation (2.4) is again first order accurate for all  $\theta (\neq 1/2)$  and attains second order accuracy when  $\theta = 1/2$ , namely for the Crank Nicolson method.

If we apply eq. (2.4) over two steps, we can write (at least) two possible approximations to  $\underline{\mu}(t+2\tau)$ , namely

$$(2.5) \quad [I - 2\tau(1-\theta)A]\underline{v}^{(1)}(t+2\tau) = [I + 2\tau\theta A]\underline{v}(t)$$

or

$$(2.6) \quad [I - \tau(1-\theta)A]^2 \underline{v}^{(2)}(t+2\tau) = [I + \tau\theta A]^2 \underline{v}(t)$$

both of which are first order accurate.

If we now propose a linear combination of  $\underline{v}^{(1)}$  and  $\underline{v}^{(2)}$ ,

namely

$$(2.7) \quad \underline{y}(t+2\tau) = \alpha \underline{y}^{(2)} + (1-\alpha) \underline{y}^{(1)}$$

then the parameters  $\alpha$  and  $\theta$  can be chosen to achieve second order accuracy. To obtain the appropriate expressions for  $\alpha$  and  $\theta$  we progress as follows:

Rewriting eq. (2.5) as

$$\underline{y}^{(1)}(t+2\tau) = [I - 2\tau(1-\theta)A]^{-1} [I + 2\tau\theta A] \underline{y}(t)$$

we can expand the matrix inverse so that

$$(2.8) \quad \underline{y}^{(1)}(t+2\tau) = [I + 2\tau A + 4(1-\theta)(\tau A)^2 + 8(1-\theta)^2(\tau A)^3 + \dots] \underline{y}(t).$$

Similarly

$$(2.9) \quad \underline{y}^{(2)}(t+2\tau) = [I + 2\tau A + (3-2\theta)(\tau A)^2 + 2(1-\theta)(2-\theta)(\tau A)^3 + \dots] \underline{y}(t).$$

substituting expressions (2.8) and (2.9) we obtain

$$(2.10) \quad \underline{y}(t+2\tau) = [I + 2\tau A + \{(3-2\theta)\alpha + 4(1-\alpha)(1-\theta)\} \tau^2 A^2 + \{2\alpha(1-\theta)(2-\theta) + 8(1-\alpha)(1-\theta)^2\} \tau^3 A^3 + \dots] \underline{y}(t)$$

A comparison with the expansion of  $\exp(2\tau A)$  in

$$\underline{y}(t+2\tau) = \exp(2\tau A) \underline{y}(t)$$

indicates second order accuracy is achieved if

$$\alpha(3-2\theta) + 4(1-\alpha)(1-\theta) = 2$$

Namely if  $(\alpha-2)(2\theta-1) = 0$ .

So second order accuracy is achieved

$$(1) \text{ for all } \alpha \text{ if } \theta = 1/2$$

$$(2) \text{ for all } \theta \text{ if } \alpha = 2.$$

Note:  $\theta = 1/2$  produces a linear combination of Crank Nicolson schemes.

$\theta = 0 \quad \alpha = 2$  reproduces the Lawson-Morris scheme [5].

It is interesting to consider the third order terms in the expansion in eq. (2.10). It is seen that third order accuracy is possible if

$$\theta = 1/2 \quad \text{and} \quad \alpha = 4/3.$$

(It is theoretically possible to achieve third order accuracy with  $\alpha = 2$  and  $\theta^2 - \theta + 1/3 = 0$ , but this involves complex arithmetic which is of little interest.)

In  $\theta = 1/2$ ,  $\alpha = 4/3$ , the reader will recognize the familiar generalization of Simpson's rule from the Trapezoidal rule for numerical integration.

In addition to achieving second and higher order accuracy it is necessary that the resulting schemas be stable. Furthermore, for reasons expounded in [5] the stability sought is  $L_0$ -stability.

To consider stability of the second order  $\theta$ -method (2.5, 2.6, 2.7) consider the symbol of the algorithm

$$S(z) = \alpha \left[ \frac{1-\theta z}{1+(1-\theta)z} \right]^2 + (1-\alpha) \left[ \frac{1-2\theta z}{1+2(1-\theta)z} \right]$$

where  $z = \tau\lambda$ ,  $\lambda$  an eigenvalue of  $A$ . For positive definite  $A$ ,  $z > 0$  and hence for  $L_0$ -stability we require that

$$\max_{z \geq 0} |S(z)| \leq 1$$

$$\text{and} \quad \lim_{z \rightarrow \infty} S(z) = 0$$

Now,

$$\lim_{z \rightarrow \infty} S(z) = \frac{\theta(\theta+\alpha-1)}{(1-\theta)^2}$$

Combining the requirement that this tends to 0 as  $z \rightarrow \infty$  with (i) and (ii) above we have

$$(a) \quad \text{For } \theta = 1/2, \quad \lim_{z \rightarrow \infty} S(z) = 2\alpha - 1 = 0 \quad \text{if } \alpha = 1/2$$

and



(b) For  $\alpha = 2$

$$\lim_{z \rightarrow \infty} S(z) = \frac{\theta(\theta+1)}{(1-\theta)^2} = 0 \quad \text{if } \theta = 0 \text{ or } \theta = -1.$$

in this case there are two possible  $L_0$ -stable members:

namely;  $\alpha = 2$ ;  $\theta = 0$  - the Lawson-Morris scheme [5]

and

$\alpha = 2$ ;  $\theta = -1$  - (a novel algorithm).

To ensure that  $|S(z)| \leq 1 \forall z \geq 0$  we graphed the symbol for increasing  $z$ . The respective symbols are depicted in figure 1.

The three methods are indeed  $L_0$ -stable as can be verified from the graphs.

An analysis of the symbol for  $\theta = 1/2$  and  $\alpha = 4/3$  (the third order accurate method) indicates conditional stability. For example if  $L = \frac{\partial^2}{\partial x^2}$  and the second order derivative is replaced by central differences then the resulting algorithm is stable if

$$(2.11) \quad \tau/h^2 \leq 3.23205.$$

Although this is not A-stable, for some applications the restriction imposed by (2.11) may not be serious in practice. However, see sections 3 and 4.

For numerical experiments for the homogeneous problem in this paper we will restrict our attention to solving the one-space dimensional heat equation

$$(2.12) \quad \frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} \quad (0 \leq x \leq 2) \times t \geq 0$$

subject to  $u(x,0) = 1$

and  $u(0,t) = u(1,t) = 0$ .

This problem was used in [5]. The theoretical solution is given by

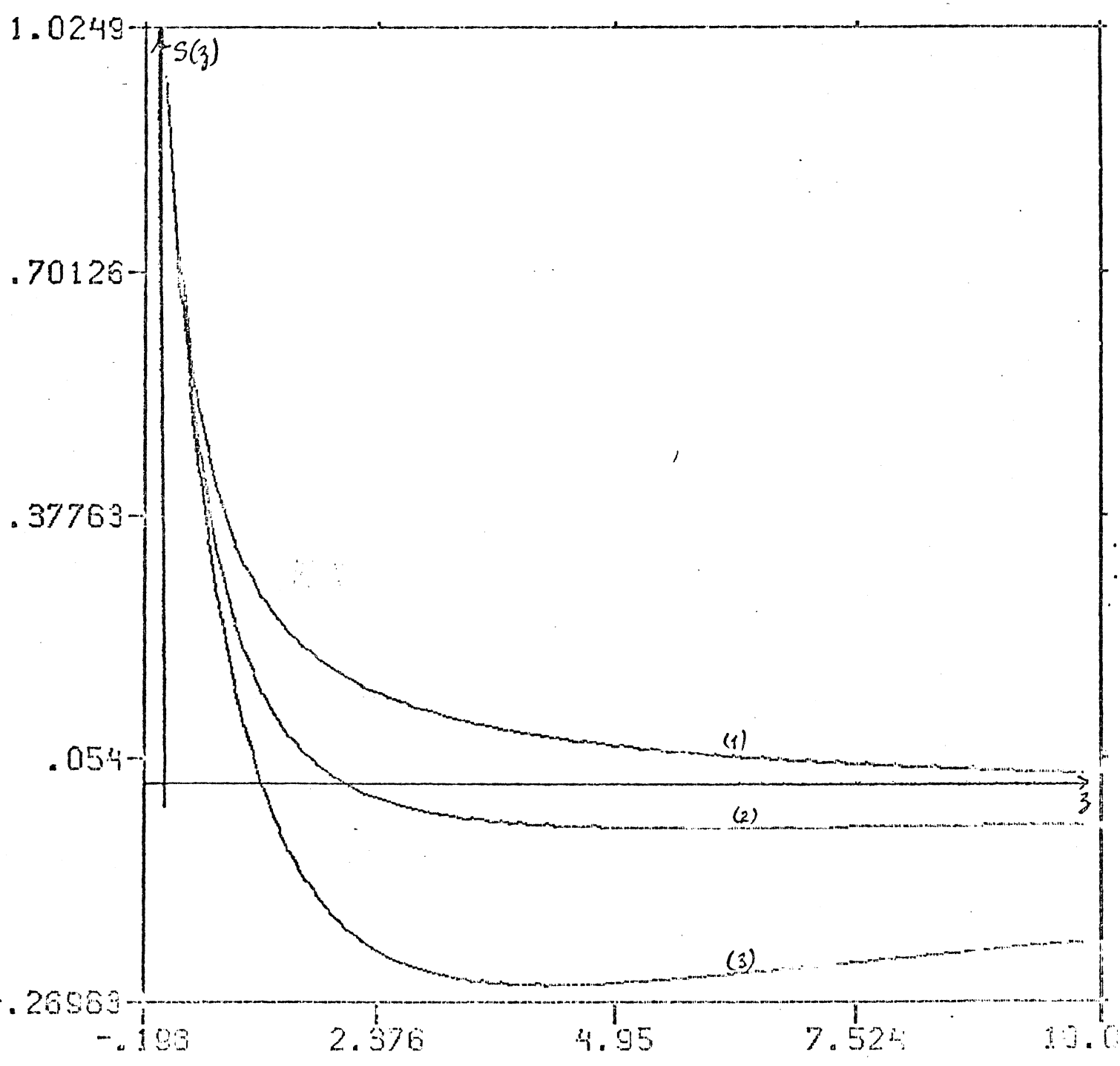


Figure 1 - Symbol for second order schemes: (1)  $\theta = -1; \alpha = 2$   
(2)  $\theta = 0; \alpha = 2$   
(3)  $\theta = 1/2; \alpha = 1/2$

$$u(x,t) = \sum_{n=1}^{\infty} [1-(-1)^n] \frac{2}{n\pi} \sin\left(\frac{n\pi x}{2}\right) \exp\left(-\frac{n^2 \pi^2 t}{4}\right)$$

This simple problem represents a situation where A-stable methods tend to do rather poorly, especially for large values of the mesh ratio  $r = \tau/h^2$  (see [5]). The interval  $[0,2]$  was divided into 40 subintervals thereby defining a value of  $h = 0.05$ . We tested the three algorithms ( $\alpha=1/2, \theta=1/2$ ;  $\alpha=2, \theta=0$ ;  $\alpha=2, \theta=-1$ ) for mesh ratios = 10 and 40 computing the solution at  $t=1.2$ . The maximum error found in each case is given in table 1. In figures 2 and 3 the computed solutions are given together with the theoretical solution.

From the table 1 of results and the figures it would appear that the algorithm  $\theta=0, \alpha=2$  performs more accurately than either of the other algorithms; the disparity being more marked for  $r = 40$ .

r	$\theta=0; \alpha=2$	$\theta=-1; \alpha=2$	$\theta=1/2; \alpha=1/2$
10.0	0.48E-03	0.23E-02	0.31E-02
40.0	0.45E-02	0.20E-01	0.12E-02

Table 1 - maximum error at  $t=1.2$  (second order methods)

The Lawson-Morris algorithm appears to perform best in that the maximum error appears to be somewhat smaller than either of the other two  $L_0$ -stable members. This behaviour is, of course, dependent on the error constant which in turn will depend upon the particular problem being solved. An additional advantage of the case  $\theta=0; \alpha=2$  is that the number of operations required to advance the solution over a time step of length  $2\tau$  is smaller than either of the other two cases.

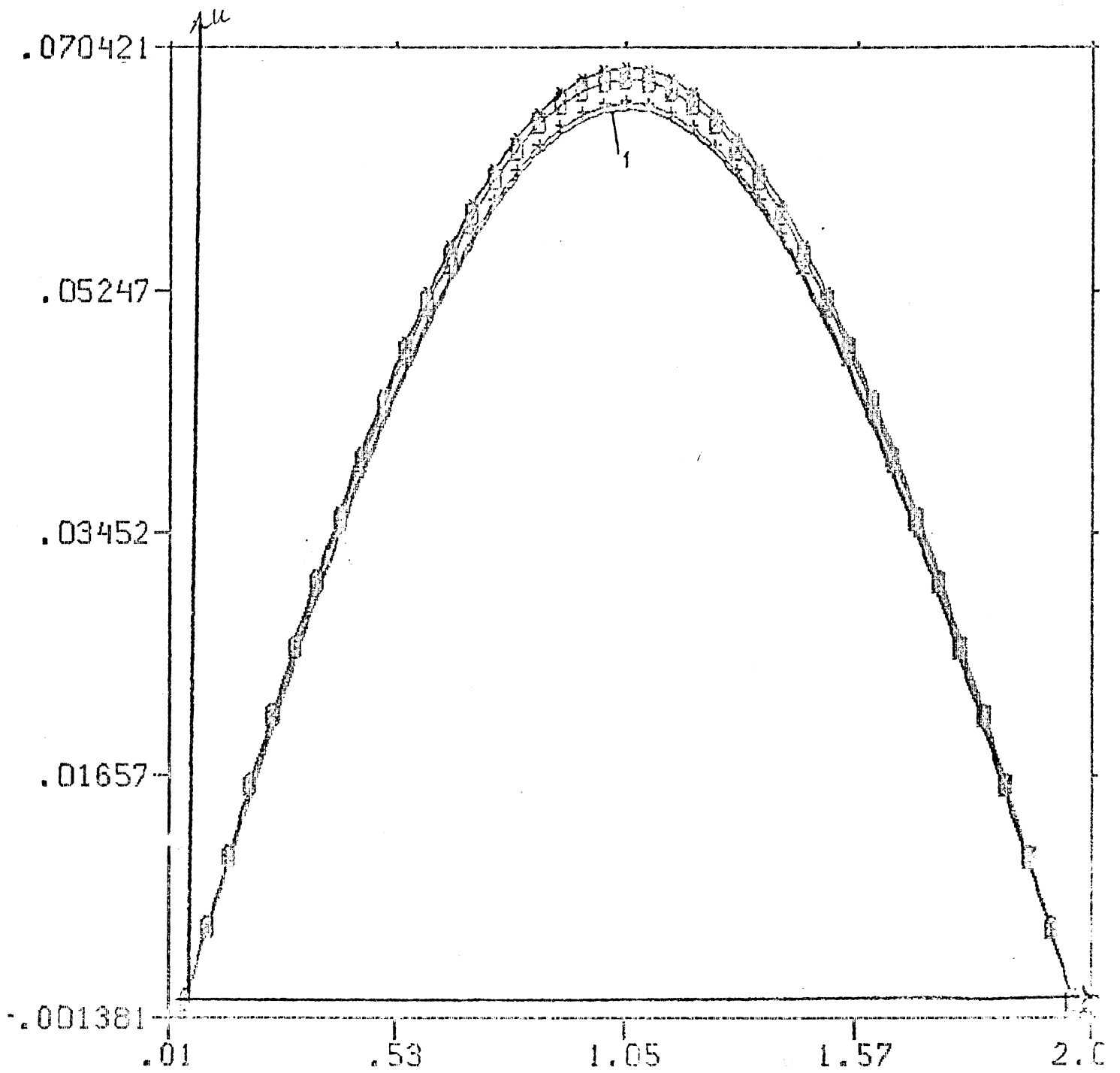


Figure 2 - Solution for problem (2.12) at  $t=1.2$ ;  $h=0.05$ ,  $\tau=0.025$  ( $r=10$ )

Fourier sum (1).  $\theta=0; \alpha=2$  (+).  $\theta=1/2; \alpha=1/2$  (\*).  $\theta=-1; \alpha=2$  ( $\phi$ ).

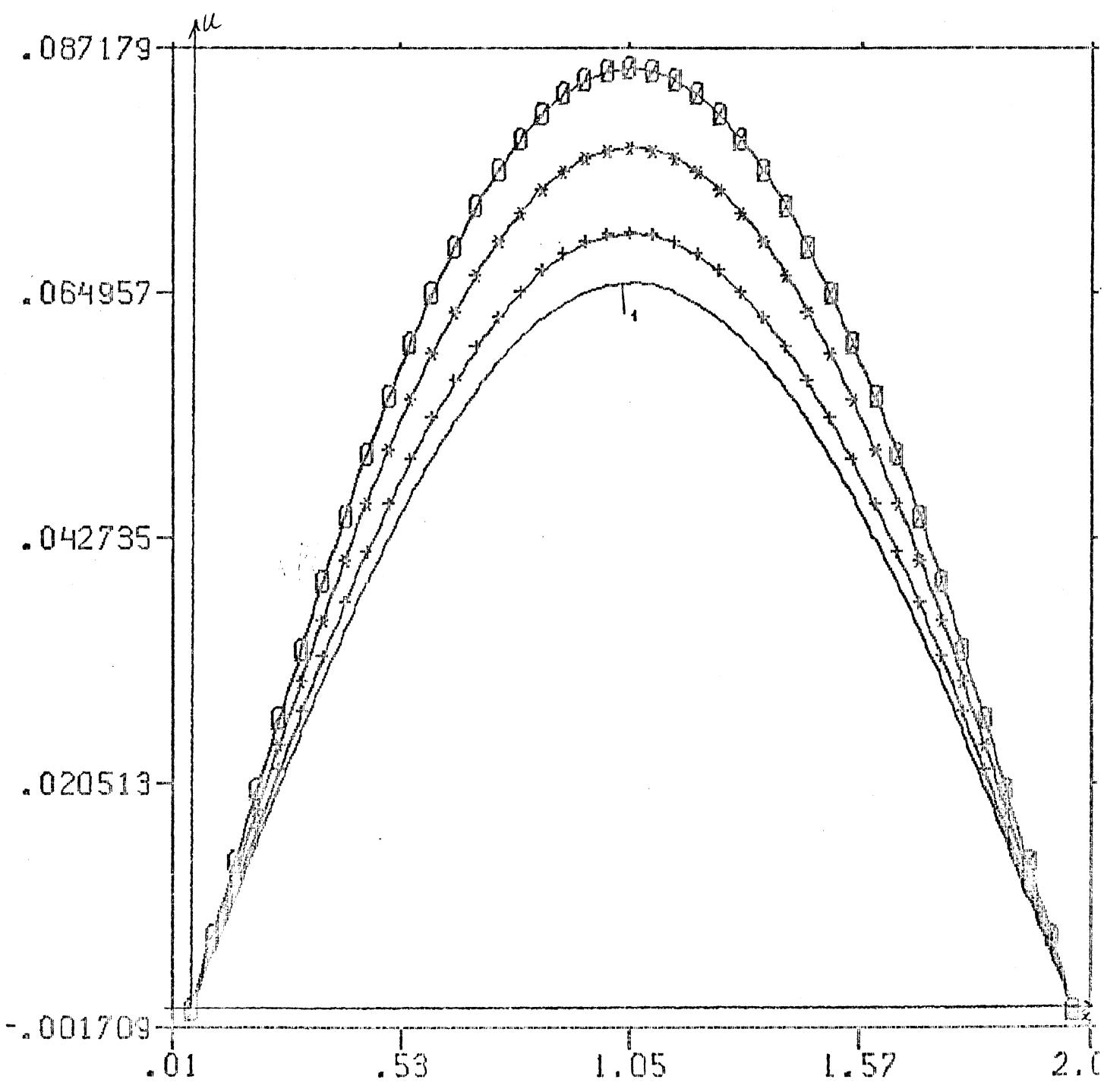


Figure 3 - Solutions for problem (2.11) at  $t=1.2$ ;  $h=0.05$ ,  $\tau=0.1$  ( $r=40$ )  
Fourier sum (1).  $\theta=0; \alpha=2$  (+).  $\theta=1/2; \alpha=1/2$  (\*).  $\theta=-1; \alpha=2$  ( $\phi$ ).

It should be pointed out that 'efficient' versions exist in a manner similar to that described in Lawson and Morris [5], but such reorganization does lead to an increased complexity unnecessary in the Lawson-Morris method ( $\theta=0; \alpha=2$ ).

Additionally it should be noted that for a time step of size  $2\tau$  three tridiagonal systems have to be solved, two of which have the same coefficient matrix which avoids the need for refactoring. In contrast the Crank Nicolson method requires two tridiagonal solutions with the same coefficient matrix for the constant coefficient (in time) partial differential equation. However, the Crank Nicolson method is not  $L_0$ -stable, and to ensure damping of errors requires the time step to be chosen so that

$$\frac{\tau}{h} \leq \frac{2}{\pi} , \text{ approximately.}$$

This restriction, consequently, makes the Crank Nicolson method considerably less efficient than any of the members of the  $L_0$ -stable family described here. (See [5] for further details.)

### 3. Third order $L_0$ -stable algorithms

One of the parametrizations in the previous section produced a third order accurate algorithm which (unfortunately) was not  $L_0$ -stable. In this section we address this question of higher order  $L_0$ -stable methods.

Define  $L_\tau$ , a difference operator, as

$$L_\tau \equiv [I - \tau(1-\theta)A]^{-1} [I + \tau\theta A]$$

Then,  $L_{2\tau} \equiv [I - 2\tau(1-\theta)A]^{-1} [I + 2\tau\theta A]$ , etc.

Consider the sequence of vectors defined by

$$(3.1) \quad \underline{y}^{(1)}(t+3\tau) = L_{\tau}^3 \underline{y}(t)$$

$$(3.2) \quad \underline{y}^{(2)}(t+3\tau) = L_{2\tau} L_{\tau} \underline{y}(t)$$

$$(3.3) \quad \underline{y}^{(3)}(t+3\tau) = L_{3\tau} \underline{y}(t)$$

and then

$$(3.4) \quad \underline{y}(t+3\tau) = \alpha \underline{y}^{(1)} + \beta \underline{y}^{(2)} + (1-\alpha-\beta) \underline{y}^{(3)}$$

$\alpha$ ,  $\beta$ ,  $\theta$  are to be chosen so that eq. (3.4) produces an expansion which agrees with terms up to, at least,  $O(\tau^3)$  in

$$\exp(3\tau A) = I + 3\tau A + \frac{9}{2}\tau^2 A^2 + \frac{9}{2}\tau^3 A^3 + \dots$$

After tedious manipulation it may be shown that ((3.1), (3.2), (3.3) and (3.4)) possesses third order accuracy if

$$\alpha = -\beta = 9/2 \forall \theta .$$

In addition to the order conditions  $L_0$ -stability will impose constraints on the parameters. The symbol for scheme (3.4) is given by

$$S(z) = \alpha \left( \frac{1-\theta z}{1+(1-\theta)z} \right)^3 + \beta \left( \frac{1-\theta z}{1+(1-\theta)z} \right) \left( \frac{1-2\theta z}{1+2(1-\theta)z} \right) \\ + (1-\alpha-\beta) \left( \frac{1-3\theta z}{1+3(1-\theta)z} \right).$$

For  $\alpha = -\beta = 9/2$

$$\lim_{z \rightarrow \infty} S(z) = \frac{-\theta}{(1-\theta)} (2+\theta)(\theta+1/2)$$

so that for  $L_0$ -stability we require one of

$$(3.5) \quad \theta = 0; \theta = -1/2; \theta = -2 .$$

The symbols for the parameters defined by (3.5) are shown in figure 4. As can be seen from the figure each of the algorithms has a symbol whose absolute value is bounded by 1 and by construction has a limit = 0 as  $z \rightarrow \infty$ . Hence, each of these algorithms is  $L_0$ -stable.

It is easy to show that  $\theta = 1/2$  produces a fourth order "three-stage" algorithm which is  $L_0$ -stable when  $\alpha = 3/4$  and  $\beta = 1/2$ . The associated symbol is also plotted in figure 4. It is interesting to note that this latter algorithm produces a symbol with very different characteristics to the third order  $L_0$ -stable algorithms. Although it is not clear that this behaviour of the symbol is significant, the behaviour of this algorithm appears somewhat less satisfactory than the third order  $L_0$ -stable variant  $\theta = 0, \alpha = -\beta = 9/2$  when used to compute the solution of eq. (2.12). The graphs of the computed results are given in figures 5 and 6<sup>†</sup> and a summary of the maximum errors is given in table 2. We note in passing that it is here possible to choose  $\alpha$  or  $\beta$  to reduce the number of stages in the algorithm. For example if  $\beta = 0$  then fourth order accuracy is still retained with  $\theta = 1/2$  if  $\alpha = 9/8$ . It is a simple matter to show that this is  $A_0$ -stable. This two stage method in fact requires four systems of linear equations to be solved, three of which comprise the same coefficient matrix, for each step of size  $3\tau$ . In contrast, the two stage method  $\alpha = 0,$

---

<sup>†</sup> We omit the theoretical solution from figures 5 and 6 for reasons of clarity; to the scale used it is essentially the curve corresponding to  $\theta = 0$ .



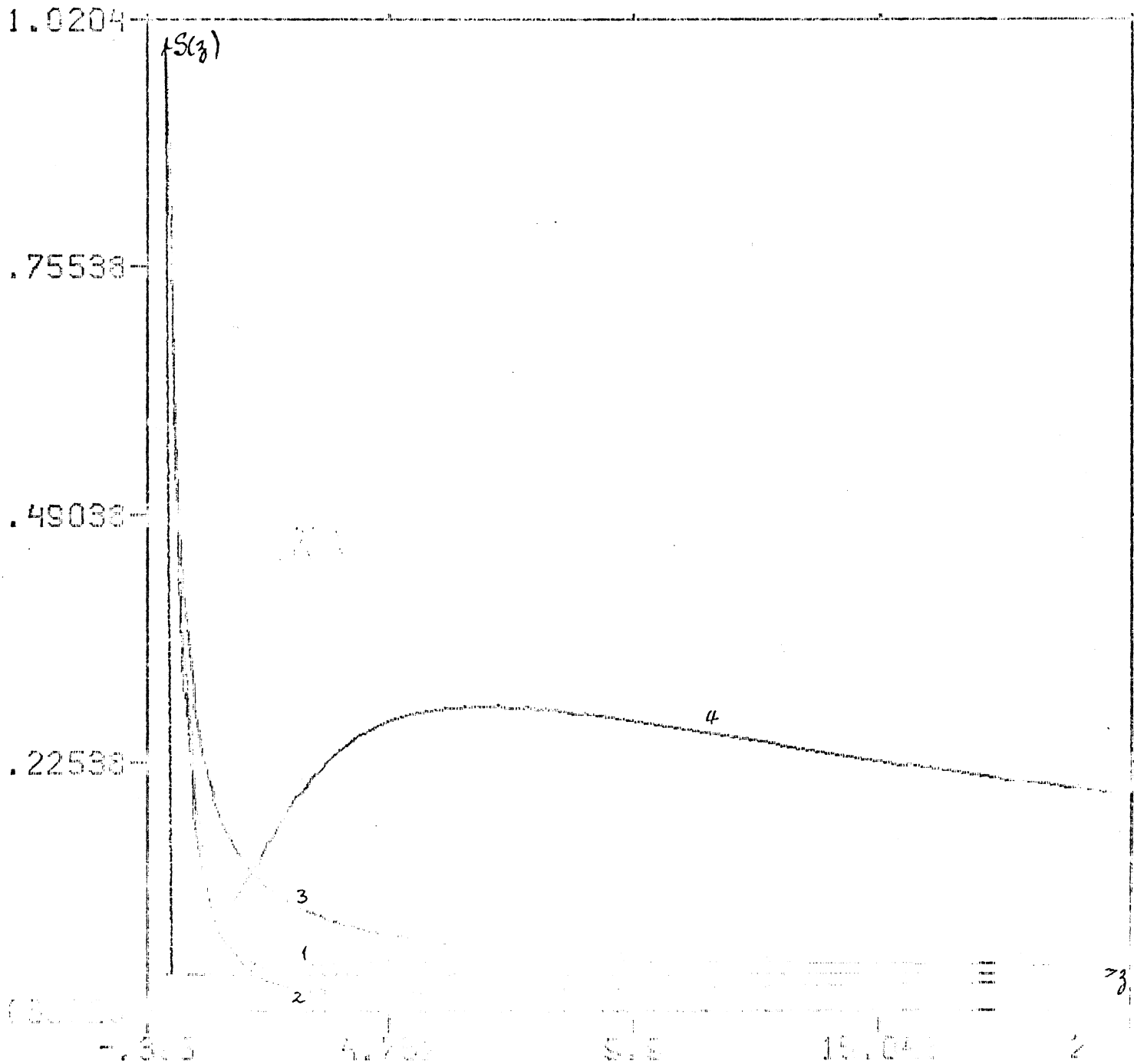


Figure 4 - Symbols for third order schemes:

- (1)  $\theta=0; \alpha=-\beta=9/2$
- (2)  $\theta=-1/2; \alpha=-\beta=9/2$
- (3)  $\theta=-2; \alpha=-\beta=9/2$
- (4)  $\theta=1/2; \alpha=3/4; \beta=1/2$

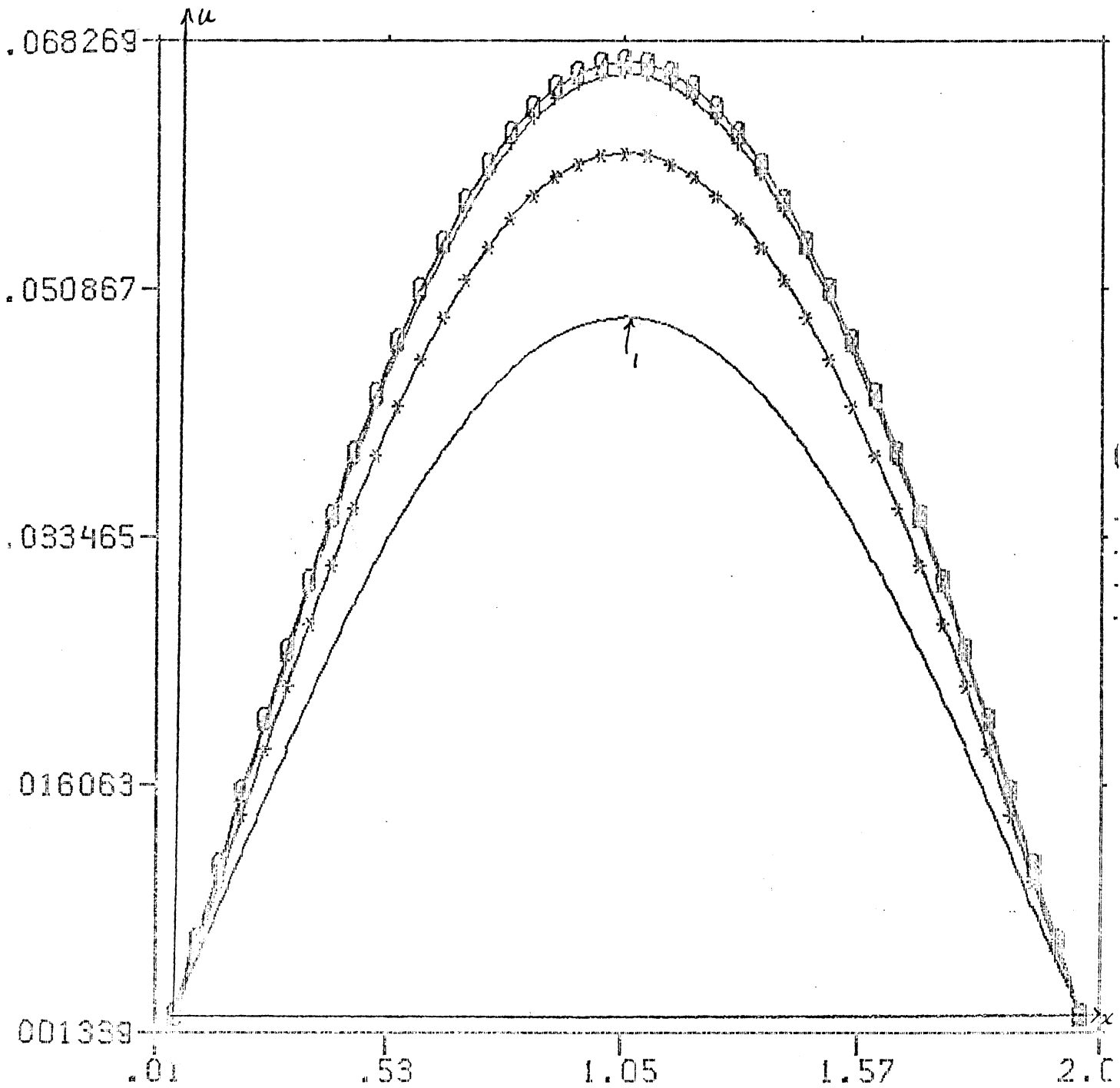


Figure 5 - Solution for problem (2.11) at  $t=1.2$ ;  $h=0.05$ ;  $\tau=0.025$  ( $r=10$ )

$\theta=-2; \alpha=-\beta=4.5$  (1);  $\theta=-1/2; \alpha=-\beta=4.5$  (\*);  $\theta=0; \alpha=-\beta=4.5$  (+);  $\theta=1/2; \alpha=3/4; \beta=1/2$  ( $\phi$ )

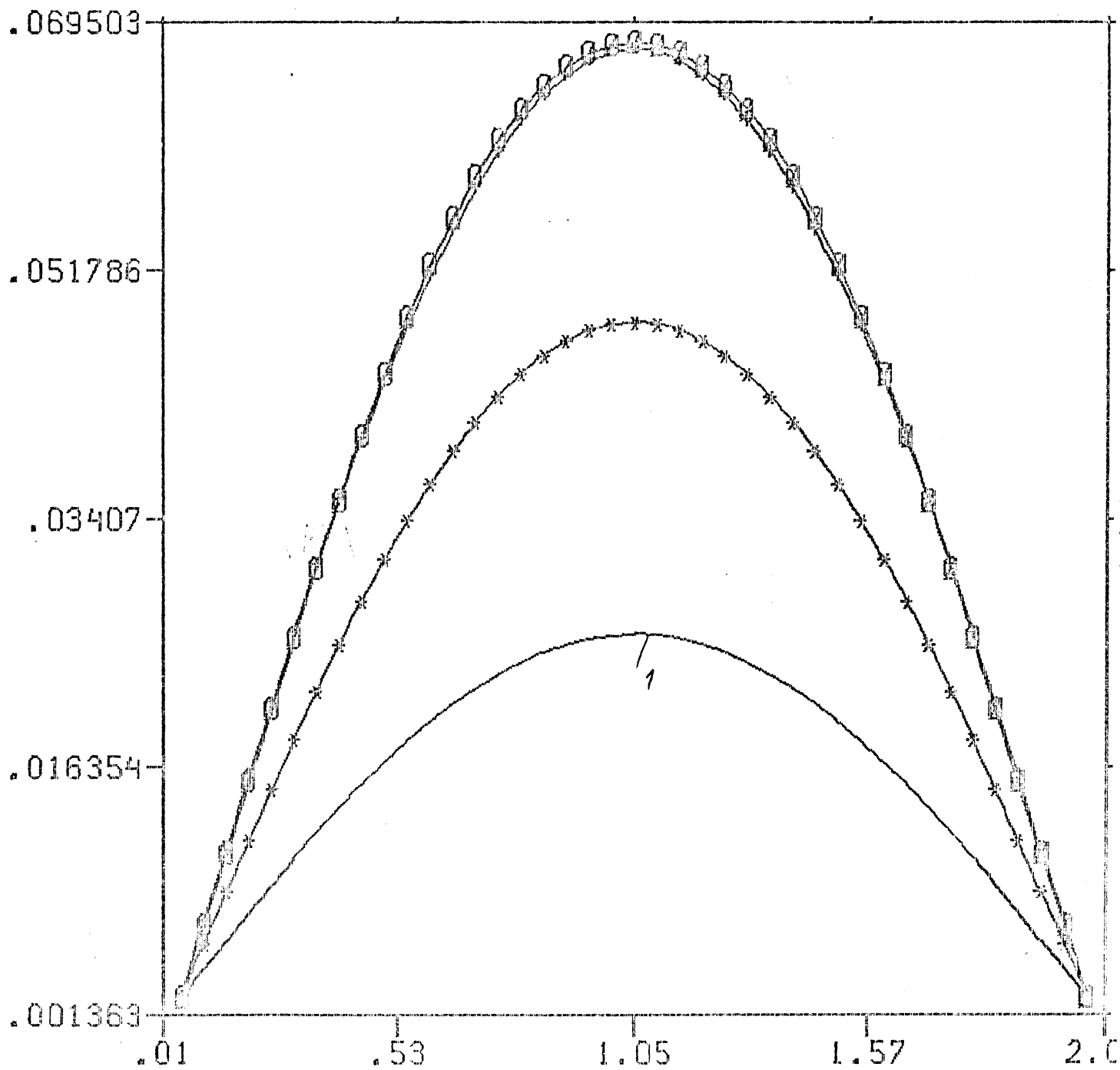


Figure 6 - Solution for problem (2.11) at  $t=1.2$ ;  $h=0.05$ ;  $\tau=0.1$  ( $r=40$ )

$\theta=-2; \alpha=-\beta=4.5$  (1);  $\theta=-1/2; \alpha=-\beta=4.5$  (\*);  $\theta=0; \alpha=-\beta=4.5$  (+);  $\theta=1/2; \alpha=3/4; \beta=1/2$  ( $\phi$ )

$\beta = 3/2$ ,  $\theta = 1/2$  is fourth order accurate but only conditionally stable with a condition on  $r$  given by  $r|\lambda| \leq 6^{2/3}$ ,  $|\lambda|$  the maximum modulus eigenvalue of  $A$ .

r	$\alpha = -\beta = 9/2$			$\alpha = 3/4; \beta = 1/2$
	$\theta = 0;$	$\theta = -1/2;$	$\theta = -2$	$\theta = 1/2$
10.0	0.13E-03	0.54E-02	0.17E-01	0.10E-02
40.0	0.17E-02	0.18E-01	0.40E-01	0.24E-02

Table 2 - Maximum errors at  $t=1.2$  (third order methods)

The most accurate member of the family appears to be the algorithm defined by  $\theta=0, \alpha=9/2, \beta=-9/2$  although the method defined by  $\theta=1/2, \alpha=3/4, \beta=1/2$  is close behind. The  $L_0$ -stability of this latter method proves to be important in the context of third order accuracy algorithms when time dependent source terms are introduced to the partial differential equation.

A comparison between the entries of tables 1 and 2 indicate that the third order method  $\theta=0$  is indeed more accurate than the second order scheme  $\theta=0$ .

#### 4. Fourth order $L_0$ -stable algorithms

In this section we seek to generalize the multistage concept introduced in section 3 to produce alternative  $L_0$ -stable methods of fourth order. To achieve this, introduce the difference operators  $L_\tau, L_{2\tau}, L_{3\tau}, L_{4\tau}$ , in a similar manner to that defined in the previous section. Then, consider the sequence of vectors defined by

$$(4.1) \quad \underline{y}^{(1)}(t+4\tau) = L_\tau^4 \underline{v}(t)$$

$$(4.2) \quad \underline{y}^{(2)}(t+4\tau) = L_{3\tau} L_\tau \underline{v}(t)$$

$$(4.3) \quad \underline{y}^{(3)}(t+4\tau) = L_{2\tau}^2 \underline{v}(t)$$

$$(4.4) \quad \underline{y}^{(4)}(t+4\tau) = L_{2\tau} L_\tau^2 \underline{v}(t)$$

$$(4.5) \quad \underline{y}^{(5)}(t+4\tau) = L_{4\tau} \underline{v}(t)$$

and then the 5 stage algorithm is defined by

$$(4.6) \quad \underline{y}(t+4\tau) = \alpha \underline{y}^{(1)} + \beta \underline{y}^{(2)} + \gamma \underline{y}^{(3)} + \delta \underline{y}^{(4)} + (1-\alpha-\beta-\gamma-\delta) \underline{y}^{(5)}.$$

The parameters  $\alpha, \beta, \gamma, \delta$  and  $\theta$  are to be chosen to produce fourth order accuracy and  $L_0$ -stability.

The order conditions are obtained by considering the expansion of the matrix inverses contained in (4.6) and comparing this expansion with that for  $\exp(4\tau A)$ . After considerable manipulation the following conditions are obtained.

For fourth order accuracy either

$$(4.7) \quad \theta = 1/2 \quad \text{and} \quad 10\alpha + 6\beta + 8\gamma + 9\delta = 32/3.$$

or

$$(4.8) \quad \left\{ \begin{array}{l} 8-6\alpha-3\beta-4\gamma-5\delta=0 \quad \text{and} \quad 2\alpha+\delta=16/3 \\ \text{and} \quad (16/3+\alpha-3\beta)-4(16/3+\alpha-3\beta)\theta+(112/3+4\alpha-15\beta)\theta^2-(32-6\beta)\theta^3=0. \end{array} \right.$$

To these conditions must also be added those arising from the requirement of  $L_0$ -stability. The five stage method produces a symbol given by

$$(4.9) \quad S(z) = \alpha \left[ \frac{1-\theta z}{1+(1-\theta)z} \right]^4 + \beta \left[ \frac{1-3\theta z}{1+3(1-\theta)z} \right] \left[ \frac{1-\theta z}{1+(1-\theta)z} \right] + \gamma \left[ \frac{1-2\theta z}{1+2(1-\theta)z} \right]^2 + \delta \left[ \frac{1-2\theta z}{1+2(1-\theta)z} \right] \left[ \frac{1-\theta z}{1+(1-\theta)z} \right]^2 + (1-\alpha-\beta-\gamma-\delta) \left[ \frac{1-4\theta z}{1+4(1-\theta)z} \right].$$

Hence for  $L_0$ -stability it is necessary that

$$(4.10) \quad \lim_{z \rightarrow \infty} S(z) = \frac{\alpha\theta^4}{(1-\theta)^4} + \frac{\beta\theta^2}{(1-\theta)^2} + \frac{\gamma\theta^2}{(1-\theta)^2} - \frac{\delta\theta^3}{(1-\theta)^3} - (1-\alpha-\beta-\gamma-\delta) \frac{\theta}{1-\theta} = 0.$$

Hence we require to choose  $\alpha, \beta, \gamma, \delta$  and  $\theta$  so that (4.10) and (4.7)

or (4.8) are satisfied.

Imposing these conditions we find the following possibilities.

$\theta$	$\alpha$	$\beta$	$\gamma$	$\delta$
0	8	40/9	0	-32/3
0	0	16/9	-6	16/3
0	-16/3	0	-10	16
0	8/3	8/3	-4	0
0	-20	-44/9	-21	136/3
1/2	0	1/2	0	23/27
1/2	0	-40/12	23/6	0
1/2	23/12	-12/12	0	0
1/2	0	0	1/2	20/27
1/2	1/2	0	0	17/27

The associated symbols are depicted in figures 7-16. By construction each scheme has a symbol whose limit as  $z \rightarrow \infty$  is 0. As can be seen the schemes are all  $L_0$ -stable. However, it is noted that the family of schemes with  $\theta=1/2$  has a distinctly different behaviour to the algorithms with  $\theta=0$ . It would be surmised that such algorithms with  $\theta=1/2$  produce considerably less damping, particularly for components

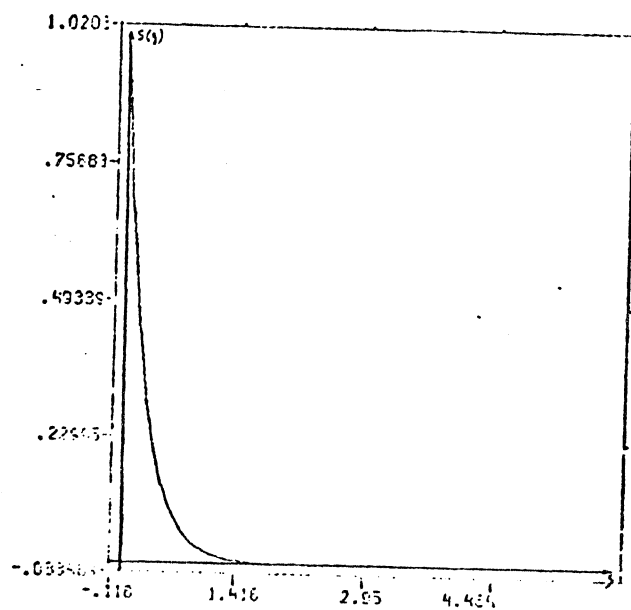


Figure 7 - Symbol for fourth order scheme  $\theta=0; a=8; b=40/9; \gamma=0; \delta=-32/3$

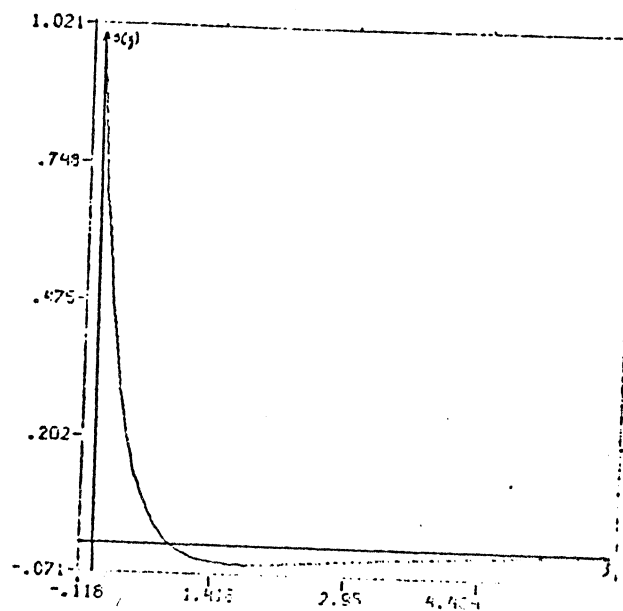


Figure 8 - Symbol for fourth order scheme  $\theta=0; a=0; b=16/9; \gamma=-6; \delta=16/3$

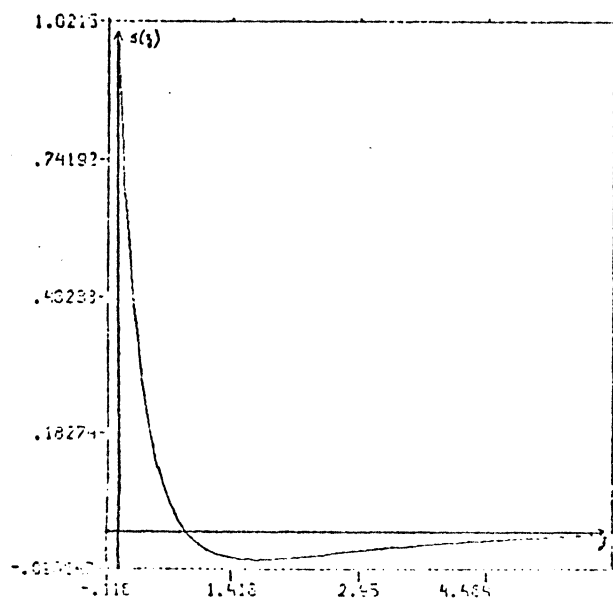


Figure 9 - Symbol for fourth order scheme  $\theta=0; a=-16/3; b=0; \gamma=-10; \delta=16$

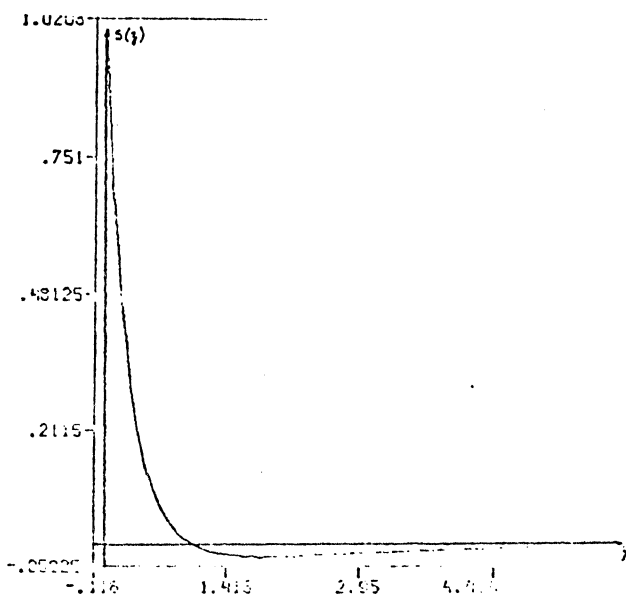


Figure 10 - Symbol for fourth order scheme  $\theta=0; a=9/3; b=8/3; \gamma=-4; \delta=0$

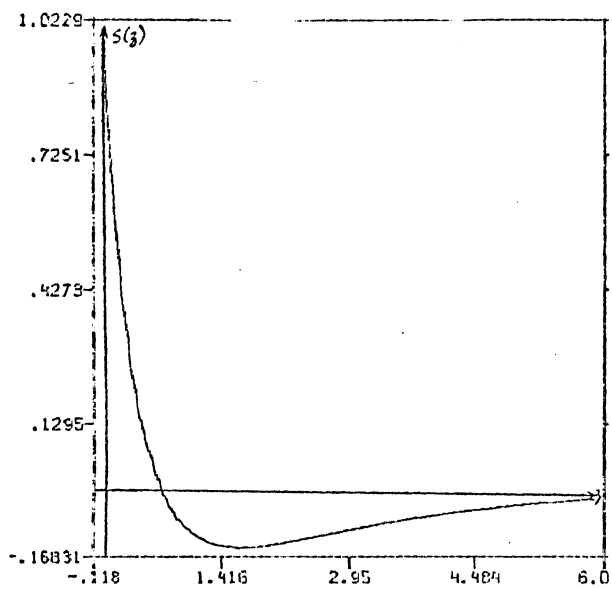


Figure 11 - Symbol for fourth order scheme  $\theta=0; \alpha=-20; \beta=-44/9; \gamma=-21; \delta=136/3$

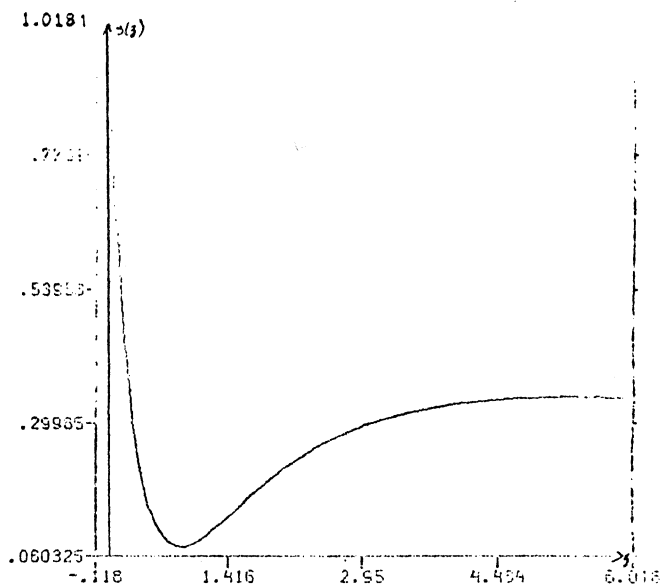


Figure 12 - Symbol for fourth order scheme  $\theta=0.5; \alpha=0; \beta=1/2; \gamma=0; \delta=23/27$

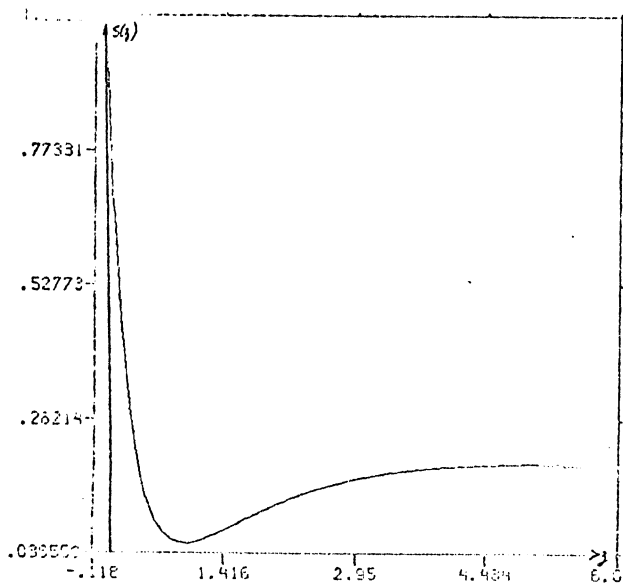


Figure 13 - Symbol for fourth order scheme  $\theta=0.5; \alpha=23/12; \beta=-17/12; \gamma=0; \delta=0$

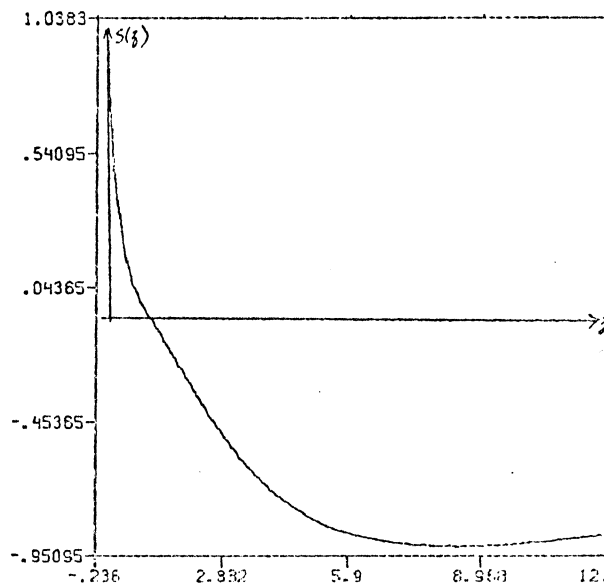


Figure 14 - Symbol for fourth order scheme  $\theta=1/2; \alpha=23/12; \beta=-17/12; \gamma=0; \delta=0$



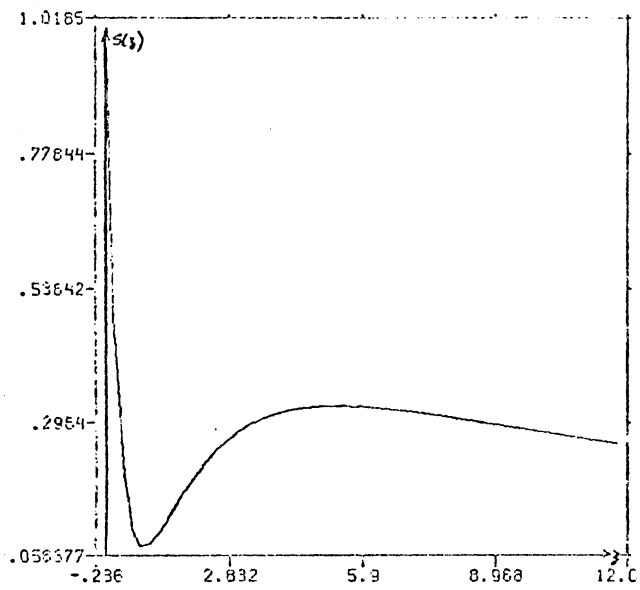


Figure 15 - Symbol for fourth order scheme  $\theta=0.5; \alpha=0; \beta=0; \gamma=1/2; \delta=20/27$

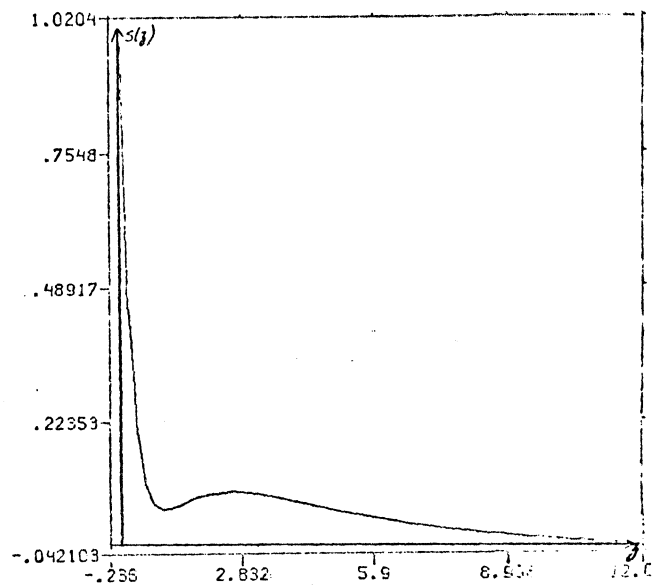


Figure 16 - Symbol for fourth order scheme  $\theta=1/2; \alpha=1/2; \beta=0; \gamma=0; \delta=17/27$

in the solution such that the product of the eigenvalue and time steps lies on the symbol which is substantially different from zero (for example  $z = 6$  in figure 12). For this reason we suspect that solutions obtained from the  $\theta=1/2$  - family are inferior, particularly for problems in which discontinuous initial/boundary values occur.

To test the behaviour of these algorithms we computed the solution of problem (2.11) for  $r=10$  and  $r=40$  in a similar manner to previous sections. The numerical results obtained are depicted in figures 17-20. In figure 17 the computed results, to the thickness of the curve, are identical for  $r=10$ . When  $r$  increases to 40 there is now some discernible difference between the various parameterizations as shown in figure 18. For  $\theta=1/2$  the behaviour of the family is much less satisfactory, and much more unpredictable. For example the results obtained from the algorithms with parameter values given by

$$(\alpha, \beta, \gamma, \delta) = (0, -10/3, 23/6, 0) \text{ and } (23/12, -17/12, 0, 0)$$

are really quite unacceptable. A summary of the respective accuracies is given in table 3 where the maximum absolute errors are shown. From both the table and figures it is seen that the algorithms given by  $\theta=0$  are very accurate, with little to choose between them when  $r=10$  and with the algorithm  $(\alpha, \beta, \gamma, \delta) = (8, 40/9, 0, -32/3)$  being marginally superior for  $r=40$ . The most inaccurate member of the  $\theta=0$  - family appears to be that given by  $(\alpha, \beta, \gamma, \delta) = (-20, -44/9, -21, 136/3)$ . A comparison of the entries in table 3 for  $\theta=0$   $(\alpha, \beta, \gamma, \delta) = (8, 40/9, 0, -32/3)$  with the corresponding entries in tables 1 and 2 indicate that the fourth order algorithm is, indeed, more accurate.

Four of the algorithms, with  $\theta=0$  possess one coefficient = 0, consequently the number of stages present in these cases is four. However, the organization of these stages is such that some four stage

members of the  $\theta=0$  - family are more efficient than others.

	r	$\alpha$	$\beta$	$\gamma$	$\delta$	$(1-\alpha-\beta-\gamma-\delta)$	max. error	# factorization	# solves
$\theta=0$	10	8	$\frac{40}{9}$	0	$-\frac{32}{3}$	$-\frac{7}{9}$	0.18 E-04	4	7
	40						0.39 E-03		
	10	0	$\frac{16}{9}$	-6	$\frac{16}{3}$	$-\frac{1}{9}$	0.13 E-05	4	7
	40						0.84 E-03		
	10	$-\frac{16}{3}$	0	-10	16	$-\frac{31}{3}$	0.13 E-04	3	8
	40						0.16 E-02		
$\theta=1/2$	10	0	$\frac{1}{2}$	0	$\frac{23}{27}$	$-\frac{19}{54}$	0.37 E-03	4	5
	40						0.89 E-02		
	10	0	$-\frac{10}{3}$	$\frac{23}{6}$	0	$\frac{1}{2}$	0.14 E-01	4	5
	40						0.79 E-01		
	10	$\frac{23}{12}$	$-\frac{17}{12}$	0	0	$\frac{1}{2}$	0.15 E-01	3	6
	40						0.91 E-01		
$\theta=1/2$	10	0	0	$\frac{1}{2}$	$\frac{20}{27}$	$-\frac{13}{54}$	0.15 E-02	3	6
	40						0.55 E-02		
	10	$\frac{1}{2}$	0	0	$\frac{17}{27}$	$-\frac{7}{54}$	0.12 E-02	3	6
	40						0.31 E-02		

Table 3 - Errors for the fourth order algorithms

The number of factorizations of the tridiagonal matrices required for the complete computation (the constant coefficient partial differential equation necessitates a single factorization; time dependent coefficients would require this factorization to be performed each cycle) and the number of forward-backward solutions required each interval of time  $4\tau$  are also summarized in table 3. We have also included the associated 'costs' for the  $\theta=1/2$  family. We have not taken into account any difference in costs associated with the fact that  $\theta=0$  produces a

right hand side without matrix/vector multiplication whereas for  $\theta=1/2$ , as formulated in this paper, there is a matrix/vector multiplication for each constituent part of the algorithm. We omit this since there is a simple means of avoiding this matrix/vector multiplication in a manner similar to that conventionally used for the Crank Nicolson method.

Assuming constant coefficients and hence assuming matrix factorizations are of no consequence, it can be seen from table 3 that for  $\theta=0$  the two algorithms  $(\alpha,\beta,\gamma,\delta) = (8,^{40}/9,0,-^{32}/3)$  and  $(0,^{16}/9,-6,^{16}/3)$  are equally efficient. Moreover both methods produce similar accuracies so either would appear an excellent choice. (The former method has a marginally superior performance for larger  $r$  so perhaps it is the favourite.)

For  $\theta=1/2$ , all members require few<sup>e</sup> solves and hence are more efficient. As we have seen each member of this class has substantially larger error for the same values of  $r$ . However, because of the increased efficiency we can take a smaller  $r$ , perform more steps and still obtain the solution with the same effective cost. This was investigated but the numerical results are still inferior to the schemes based on  $\theta=0$ . For brevity we omit the numerical details.

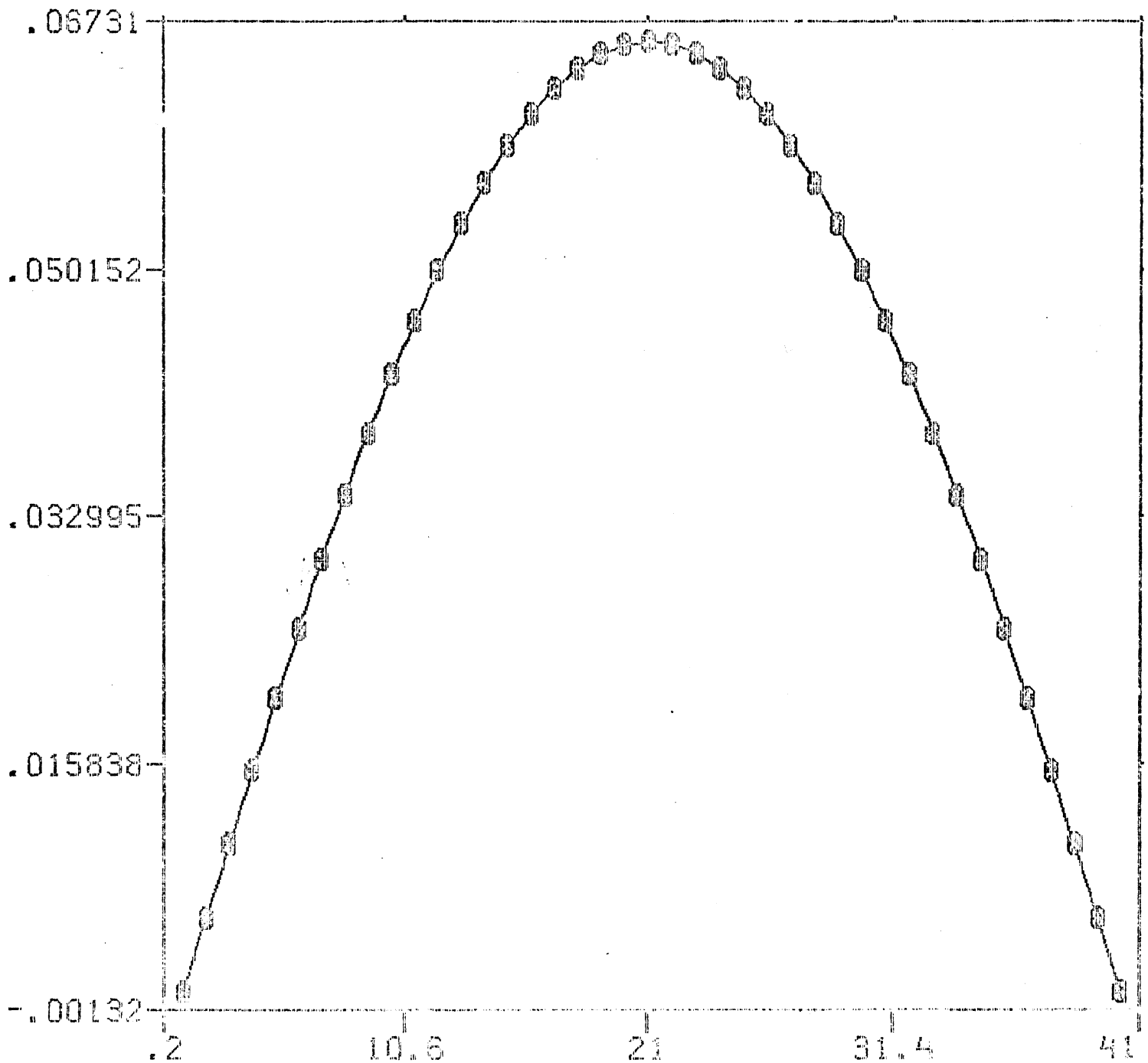


Figure 17 - Solution of problem (2.11) at  $t=1.2$ ;  $h=0.05$ ;  $\tau=0.025$  ( $r=10$ )  
(5 coincident curves)

$$\begin{aligned}
 \theta=0: \quad & \alpha=8, & \beta=40/9, & \gamma=0, & \delta=-32/3 \\
 & \alpha=0, & \beta=16/9, & \gamma=-6, & \delta=16/3 \\
 & \alpha=-16/3, & \beta=0, & \gamma=-10, & \delta=16 \\
 & \alpha=8/3, & \beta=8/2, & \gamma=-4, & \delta=0 \\
 & \alpha=-20, & \beta=-44/9, & \gamma=-21, & \delta=136/3
 \end{aligned}$$

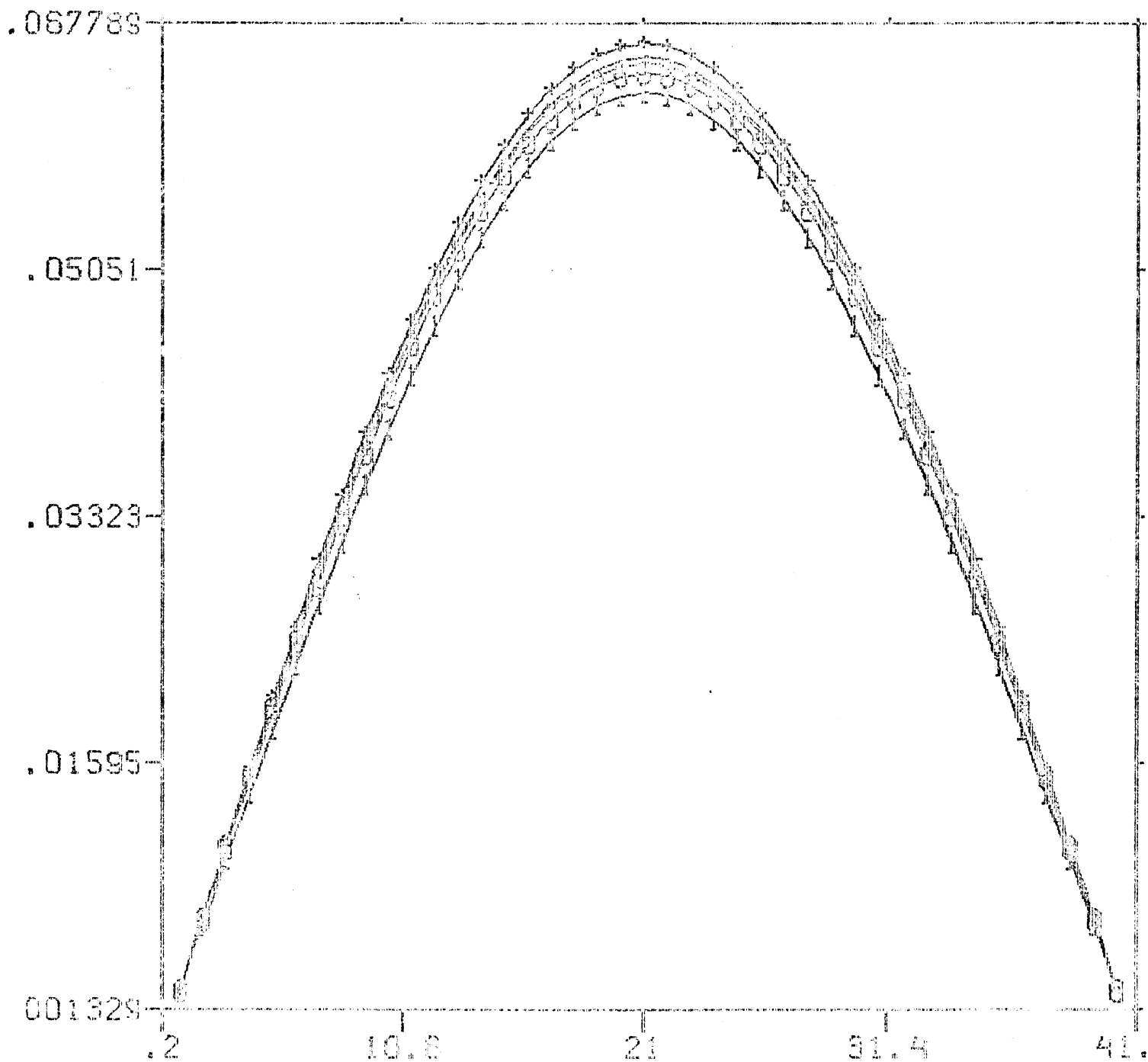


Figure 18 - Solution of problem (2.11) at  $t=1.2$ ;  $h=0.05$ ;  $\tau=0.01$  ( $r=40$ )

$\theta=0$ :	$\alpha=8$ ,	$\beta=40/9$ ,	$\gamma=0$ ,	$\delta=-32/3$	curve (+)
	$\alpha=0$ ,	$\beta=16/9$	$\gamma=-6$ ,	$\delta=16/3$	curve (*)
	$\alpha=-16/3$ ,	$\beta=0$ ,	$\gamma=-10$ ,	$\delta=16$	curve ( $\theta$ )
	$\alpha=8/3$ ,	$\beta=8/3$ ,	$\gamma=-4$ ,	$\delta=0$	curve ( $\bullet$ )
	$\alpha=-20$ ,	$\beta=-44/9$ ,	$\gamma=-21$ ,	$\delta=136/3$	curve (l)

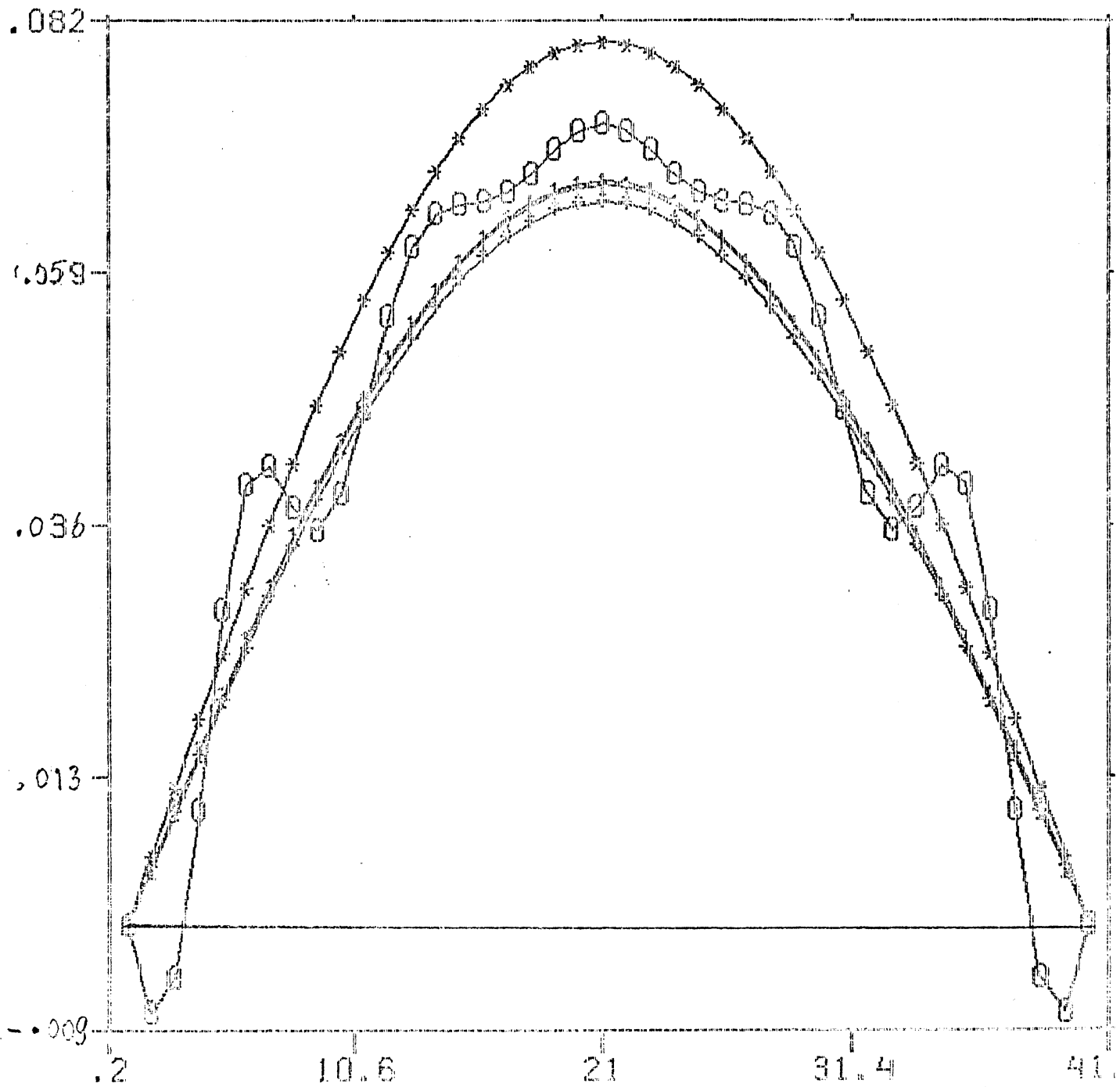


Figure 19 - Solution of problem (2.11) at  $t=1.2$ ;  $h=0.05$ ;  $\tau=0.025$  ( $r=10$ )

$\theta=1/2$ :	$\alpha=0$ ,	$\beta=1/2$ ,	$\gamma=0$ ,	$\delta=23/27$	curve (+)
	$\alpha=0$ ,	$\beta=-10/3$ ,	$\gamma=23/6$ ,	$\delta=0$	curve (*)
	$\alpha=23/12$ ,	$\beta=-17/12$ ,	$\gamma=0$ ,	$\delta=0$	curve ( $\theta$ )
	$\alpha=0$ ,	$\beta=0$ ,	$\gamma=1/2$ ,	$\delta=20/27$	curve ( $\cdot$ )
	$\alpha=1/2$ ,	$\beta=0$ ,	$\gamma=0$ ,	$\delta=17/27$	curve (l)

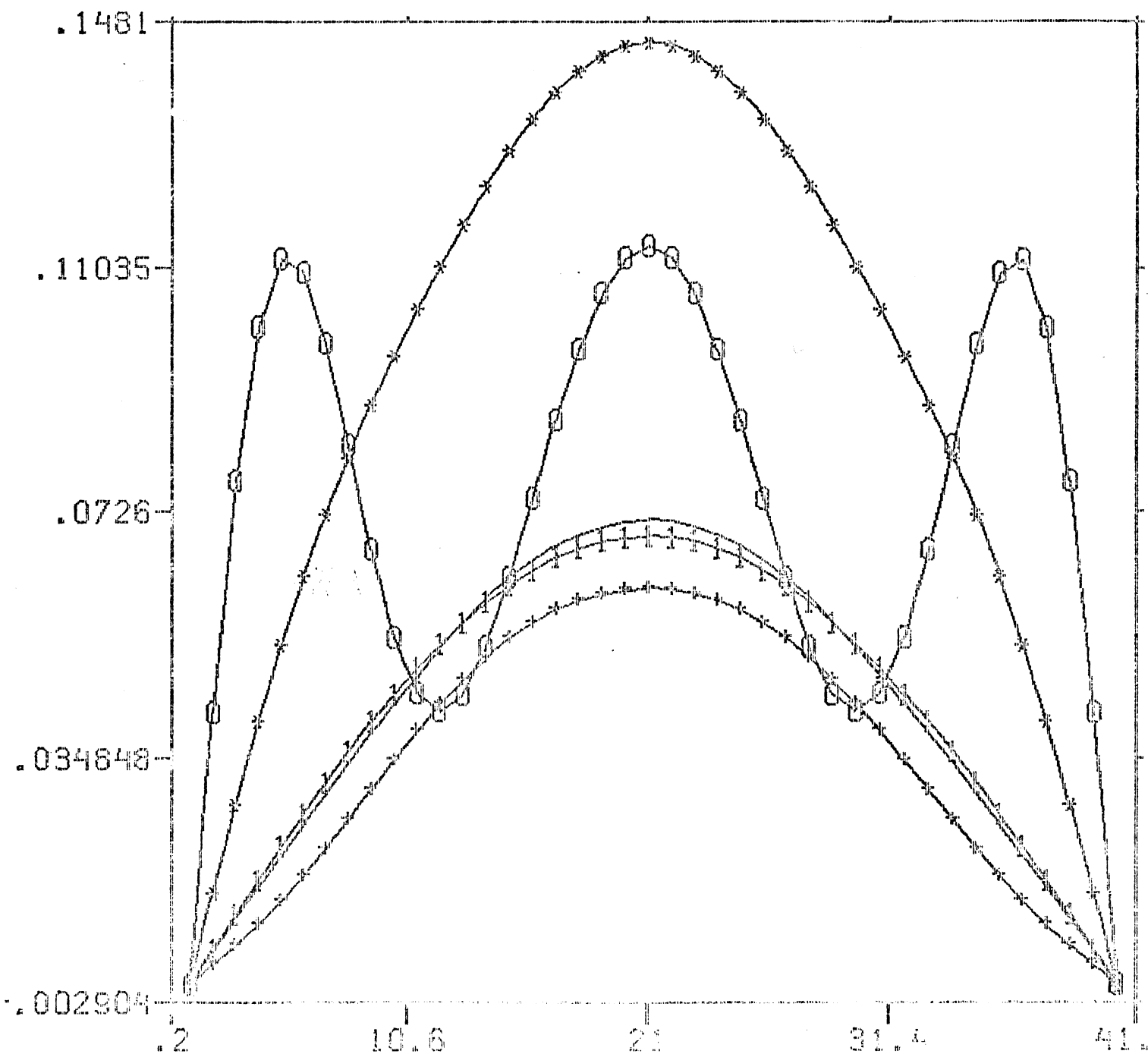


Figure 20 - Solution of problem (2.11) at  $t=1.2$ ,  $h=0.05$ ,  $\tau=0.1$  ( $r=40$ )

$\theta=1/2$ :	$\alpha=0$ ,	$\beta=1/2$ ,	$\gamma=0$ ,	$\delta=23/27$	curve (+)
	$\alpha=0$ ,	$\beta=-10/3$ ,	$\gamma=23/6$ ,	$\delta=0$	curve (*)
	$\alpha=23/12$ ,	$\beta=-17/12$ ,	$\gamma=0$ ,	$\delta=0$	curve ( $\theta$ )
	$\alpha=0$ ,	$\beta=0$ ,	$\gamma=1/2$ ,	$\delta=20/27$	curve (.)
	$\alpha=1/2$ ,	$\beta=0$ ,	$\gamma=0$ ,	$\delta=17/27$	curve (1)



## 6. Conclusions

We have extended the  $L_0$ -stable results of Lawson and Morris [5] to higher orders of accuracy. In the case of accuracies of the orders three and four, these appear to be attractive, simple methods which are  $L_0$ -stable and which produce satisfactory results for problems in which high frequency components are known to propagate when using conventional  $A_0$ -stable methods like the Crank Nicolson method.

The structure of the extrapolated algorithms we have introduced gives use, in an obvious manner, to higher order methods. For example, fifth order would be achieved by choosing the parameters suitably in

$$y(t+5\tau) = [\alpha L_T^5 + \beta L_{4T} L_T + \gamma L_{3T} L_{2T} + \delta L_{3T} L_T^2 + \epsilon L_{2T}^2 L_T + \phi L_{2T} L_T^3 + (1-\alpha-\beta-\gamma-\delta-\phi) L_{5T}] y(t)$$

However, we refrain from pursuing generalizations along these lines believing, for partial differential equations, that fourth order accuracy in time is sufficient for most (all?) applications.

A more important extension to the present work is the condition of time dependent coefficients and equations with inhomogeneous source times. This work will be reported in part III of this paper. The algorithms here apply naturally to problems in many space variables when sparse matrix algorithms are available. The question of applying the generalized methods in a splitting context is as yet an open one.

### Acknowledgement

This work was partially supported by NSERC grant A3597.

References

- [1] L. Fox. Some improvements in the use of relaxation methods for the solution of ordinary and partial differential equations, Proc. Royal Soc., London, A190 (1947), pp. 31-59.
- [2] R. Frank and C. Uberheuber. Iterated defect correction for the efficient solution of stiff systems of ordinary differential equations, BIT 17 (1977), pp. 146-159.
- [3] J.A. George and J.W.H. Liu. Computer solution of large sparse positive definite systems. Manuscript: to appear Prentice Hall, 1980.
- [4] J.D. Lambert. Computational methods in ordinary differential equations, Wiley, 1973.
- [5] J.D. Lawson and J.L. Morris. The extrapolation of first order methods for parabolic partial differential equations I, SIAM J. Numer. Anal. 15, No. 6 (1978).
- ✓✓ [6] V. Peneyra. High order finite difference solution of differential equations, Stanford Report Stan-CS-73-348, 1973.
- [7] A. Saylor. Extrapolation deferred correction, and defect correction or discrete-time Galerkin methods for linear parabolic problems. Ph.D. Thesis, University of Kentucky, 1979.
- [8] L.F. Richardson and J.A. Gaunt. The deferred approach to the limit II interpenetrating lattices, Philo. Trans., Roy. Soc. London, Ser. A, 226 (1979), pp. 350-361.
- [9] H. Stetter. The defect correction principle and discretization methods, Numer. Math. 29 (1978), pp. 425-443.
- [10] P. Zadunaisky. On the estimation of errors propagated in the numerical integration of ordinary differential equations. In proceedings of Conference on the Numerical Solutions of Differential Equations, Ed. G.A. Watson, Springer Verlag Lecture Notes in Mathematics, 362, 1974.