# An Analysis of Two-Step Time Discretizations in the Solution of the Linearized Shallow Water Equations

M. G. G. FOREMAN

*Institute of Ocean Sciences, P. O. Box 6000,
Sidney, British Columbia V8L 4B2, Canada*

Received June 24, 1982; revised February 4, 1983

A Fourier analysis for evaluating the accuracy of finite element methods which solve the linearized shallow water equations is extended to include group velocity. The technique is first applied to the selection of a spatial discretization and then, assuming a specific finite element discretization, to finding the most accurate two-step time stepping method.

## INTRODUCTION

Within the last decade, finite element methods have become popular for solving the shallow water equations. Many such methods are available since typically, each combines a spatial discretization with a time stepping or spectral method. The spatial discretization is determined by the particular finite element approach (e.g., Galerkin), the approximating basis functions, and the size, shape, and configuration of the spatial elements. It has the effect of reducing the governing partial differential equations (PDEs) to a system of ordinary differential equations (ODEs) in time. These can then be solved by one of many methods discussed in texts such as Gear [4] or Lambert [7]. In light of the many method combinations, it is natural to ask which is the most accurate.

Several techniques are available for evaluating the accuracy of a finite element (or finite difference) method. Assuming a specific spatial discretization, a traditional analysis of ODE methods involves investigating the absolute stability region and calculating the truncation errors. Stability of a prospective time stepping method is determined by ensuring that the spectrum of the Jacobian of the ODE system, when scaled by $\Delta t$, lies in the absolute stability region of the method. Accuracy is determined by evaluating the local truncation error, or by comparing the principal root of the method's characteristic polynomial to the exponential function (which is the analytic result). This approach is discussed in Gear [4] and used by Praagman [17] for analysing finite element solutions of the shallow water equations.

A second popular technique for evaluating numerical methods which solve hyperbolic PDEs was developed by Leendertse [9]. It is based on propagation factors. These are ratios of the computed wave to the analytic wave after the time it takes for

454

the analytic wave to propagate one wavelength. Gray and Lynch [5] apply this analysis to finite element solutions of the shallow water equations. In one dimension, they assess various time stepping schemes in combination with a Galerkin finite element method and piecewise linear basis functions.

The technique employed in this paper is an extension of a third common type (e.g., [12]). It is a Fourier (phase/space) analysis whereby the amplitude and phase of the discrete problem are studied. The method extension is through examination of group velocity. Group velocity is important in all wave problems since it describes the speed and direction of energy propagation. For tsunami problems, it is vital since the wave packet speed rather than that of an individual wave determines the arrival time [14]. Although shallow water waves have virtually the same phase and group velocity, their numerical model representations may not. It is therefore important to study the properties of both in assessing the merits of a numerical scheme. As it will be seen, a method which most accurately represents phase velocity may not be best for group velocity.

Previous studies of group velocity in numerical methods are not common. Vichnevetsky [22] shows that zero group velocity characterizes a cutoff frequency beyond which wave solutions exhibit a spurious amplitude decay. Vichnevetsky and Peiffer [23] demonstrate that spurious $2\Delta x$ waves generated by mesh refinement or near-discontinuities in the exact solution, travel at the group speed. Schoenstadt [19] and Williams [27] include group velocity in their evaluation of several numerical methods for solving the atmospheric shallow water equations. In all these studies, only the effects of the spatial discretization are considered.

A recent paper by Trefethen [21] surveys and illustrates the relevance of group velocity in numerical schemes. Among the important points that he discusses are the following:

   (i)   although wave crests travel at the phase velocity, wave packets travel at the group velocity,

   (ii)   energy travels at the group velocity,

   (iii)   group speed is the only meaningful speed for studying parasitic numerical solutions,

   (iv)   instability of an initial boundary value problem is related to the possibility that waves radiating out from a boundary with positive group velocity may not be stimulated by incoming waves with negative group velocity,

   (v)   zero group velocity defines a cutoff frequency for transmission through an interface.

In brief, Trefethen demonstrates that there is more to the inaccuracy of a numerical scheme than its truncation error.

In addition to the inclusion of group velocity, this analysis approach has other advantages. Calculations to determine accuracy and stability are closely correlated and expressed in terms of amplitude, phase velocity, and group velocity. Unlike stability regions, truncation errors, and propagation factors, these are familiar

concepts for the physcal oceanographer. Furthermore, the same analysis technique can be used to evaluate a method both before and after the ODE is solved. That is, the analysis can asses the merits of the spatial discretization as well as the time stepping method.

However, the analysis does have limited application. It requires that the PDE be linear, and have constant coefficients and periodic boundary conditions. (In some cases, analyses are possible for nonconstant coefficients. See Section 7.) A constant time step and a regular mesh configuration are usually assumed as well. Although few problems are this simple, it is important to understand numerical behaviour in such a setting before introducing the additional complexities of boundary conditions and varying coefficients. In order to simplify the analysis, only the one-dimensional shallow water equations are examined in this paper. However, extensions to two dimensions are straightforward and have also been done.

This paper first examines eight finite element and finite difference spatial discretizations. It then studies the effects of combining an ODE solver from the class of linear two-step methods with the particular spatial discretization, a Galerkin finite element method with piecewise linear basis functions. Although these choices may not produce the most accurate numerical method, they do effectively illustrate the analysis technique. In particular, two-step methods permit an accuracy optimization through their parameterization and illustrate the problems that can arise from a spurious dispersion relationship. A Galerkin finite element method with piecewise linear basis functions illustrates the additional problems associated with cutoff frequencies and $2\Delta x$ waves. Similar analysis have been applied to other spatial discretizations.

In this study, little attention is given to program storage requirements and economy of the numerical calculations. In two dimensions, these are probably the most important criteria for selecting a numerical method. Therefore a complete evaluation of a numerical method should include not only the accuracy considerations studied here but also cost estimates of its implementation. In terms of computational cost, Weare [26], for example, shows that finite element methods cannot compare favourably to finite difference methods as long as band algorithms are used in their solution. Of course there are many alternatives to band algorithms. A popular one involves "mass lumping" the matrix (e.g., [20]). When used in combination with an explicit time stepping method, this procedure diagonalizes the band matrix, thereby greatly reducing the storage and computational cost. In fact, the resultant technique should be economically competitive with explicit finite difference methods. However, a corresponding accuracy loss [2, 13] can be expected and warrants further investigation with the present analysis technique.

This paper is divided into seven sections. Section 1 specifies the shallow water equations and their analytic solution. It also defines dispersion relationship, phase velocity, and group velocity. Section 2 calculates the phase and group velocities arising from eight finite element and finite difference spatial discretizations and discusses their relative merits. Section 3 introduces the class of two-step methods for solving an ODE. Section 4 uses these methods in combination with a specific spatial

descretization to solve the governing equations. Dominant phase and group velocities, and dominant wavenumbers are also defined and illustrated. Section 5 defines three accuracy measure or error functions and uses them to find the most accurate two-step method. Section 6 validates these accuracy measure functions with numerical tests and a truncation error analysis. Finally, Section 7 summarizes and briefly discusses the results.

## 1. MATHEMATICAL BACKGROUND

The one-dimensional linearized shallow water equations are

$$\frac{\partial z}{\partial t} + \frac{\partial (hu)}{\partial x} = 0, \tag{1a}$$

$$\frac{\partial u}{\partial t} + g\frac{\partial z}{\partial x} + \tau u = 0, \tag{1b}$$

where $z(x, t) = $ elevation above mean sea level, $u(x, t) = $ velocity, $h(x) = $ mean sea depth, $g = $ gravity, $\tau = $ linear bottom friction coefficient.

In the present analysis, (1a) and (1b) are solved on an infinite channel (i.e., with periodic boundary conditions) subject to initial conditions.

Assuming a constant depth and travelling wave solutions of the form

$$\begin{pmatrix} z(x, t) \\ u(x, t) \end{pmatrix} = \mathrm{Re}\left[ \begin{pmatrix} z_0 \\ u_0 \end{pmatrix} \exp(ikx + i\omega t) \right], \tag{2}$$

where $\omega$ is frequency and $k$ is wavenumber, the dispersion relationship is

$$\omega = i\tfrac{1}{2}\tau \pm (ghk^2 - (\tfrac{1}{2}\tau)^2)^{1/2} = i\tfrac{1}{2}\tau \pm \omega_r. \tag{3}$$

The solutions are then

$$z(x, t) = z_0 \exp(-\tfrac{1}{2}\tau t)\cos(kx \pm \omega_r t),$$

$$u(x, t) = -z_0\left(\frac{g}{h}\right)^{1/2}\exp(-\tfrac{1}{2}\tau t)\cos(kx + \omega_r t + \theta)$$

$$= z_0\left(\frac{g}{h}\right)^{1/2}\exp(-\tfrac{1}{2}\tau t)\cos(kx - \omega_r t - \theta), \tag{4}$$

where

$$\theta = \arctan(\tfrac{1}{2}\tau, \omega_r).$$

Using (3), phase and group velocities are found from their respective definitions,

$$C = \frac{\omega_r}{k} \quad \text{and} \quad G = \frac{\partial \omega_r}{\partial k}. \tag{5}$$

Waves whose propagation speed $C$ is not independent of $k$ are said to be dispersive. This implies $|C| \neq |G|$. When travelling in a group, dispersive waves seem to be created at the leading or trailing edge of the packet, and disappear at the other edge.

Changes in the solution over the interval $\Delta t$ are expressed through the amplitude and phase of the function

$$\lambda = \exp(i\omega\Delta t) \tag{6}$$

with $\omega$ substituted from (3).

Dispersion relationships, phase velocities, and group velocities can also be found for numerical methods by requiring nontrivial travelling wave solutions to the discretized versions of (1a) and (1b) with constant $\Delta t$ and $\Delta x$. In this case, $\lambda$ is the root of the associated characteristic polynomial. It is also an eigenvalue of the amplification matrix resulting from a linear stability analysis [18]. Due to this relationship, $\lambda$ will be referred to as an eigenvalue throughout this paper. Richtmeyer and Morton [18] call it an amplification factor.

Although analytic waves may be nondispersive, all discrete models of them are dispersive [21].

## 2. AN ANALYSIS OF SPATIAL DISCRETIZATIONS

Dispersion relationships may be obtained for the system of ODEs which arise from a particular spatial discretization of (1a) and (1b). The resulting phase and group velocities may be interpreted as arising from a numerical scheme where the time dependency can be solved exactly. They thus provide a measurement of inaccuracy solely due to the spatial discretization. However, this does not mean that a subsequent time discretization will contribute further errors. It is possible that some cancellation may occur and the fully discretized equations may be more accurate.

The analysis approach is similar to the spatial Fourier transform method used by Schoenstadt [19] and Williams [27]. Their analyses are for simplified two-dimensional versions of (1a) and (1b) and include some of the following discretizations studied here:

(D1)   a centred finite difference method with an unstaggered grid,

(D2)   a centred finite difference method with a staggered grid,

(D3)   a Galerkin finite element method with piecewise linear basis functions for both variables and unstaggered elements,

(D4)   a Galerkin finite element method with piecewise linear basis functions for both variables and staggered elements,

(D5)   a residual least squares finite element with piecewise linear basis functions for both variables and unstaggered elements,

(D6)   a Galerkin finite element method with unstaggered elements, piecewise constant basis functions for one variable and piecewise linear for the other,

(D7)   a Galerkin finite element method with unstaggered elements, piecewise linear basis functions for one variable, piecewise quadratic for the other and

   (a)   $\Delta x =$ distance between adjacent "linear variables" or

   (b)   $\Delta x =$ distance between adjacent "quadratic variables,"

(D8)   a Galerkin finite element method with unstaggered elements, piecewise quadratic basis functions for both variables, and

   (a)   $\Delta x =$ distance between nodes of the same type, (i.e., between mid-element nodes or end-element nodes),

   (b)   $\Delta x =$ distance between adjacent nodes.

For $\tau = 0$, the spatially discretized equations and their corresponding dispersion relationships are listed in Tables I and II, respectively. Figure 1 plots the nondimensional phase and group velocities versus $k\Delta x/\pi$. Both analytic velocities are identically equal to 1.0.

The $(0, \pi]$ range for $k\Delta x$ reflects grid sampling per wavelength. The upper value corresponds to the shortest resolvable wavelength, namely, $2\Delta x$, while the lower value represents infinite sampling. Numerical models are usually designed so that desired wavelengths are at least $20\Delta x$ (i.e., $k\Delta x/\pi \leqslant 0.1$). Figure 1 shows that most of the selected discretizations are quite accurate in this range.

The interpretation of $\Delta x$ necessitates two representations for (D7) and (D8). In both cases, representation (a) is simply the first half of representation (b) stretched by a factor of 2. Piecewise quadratic approximation requires two types of basis function
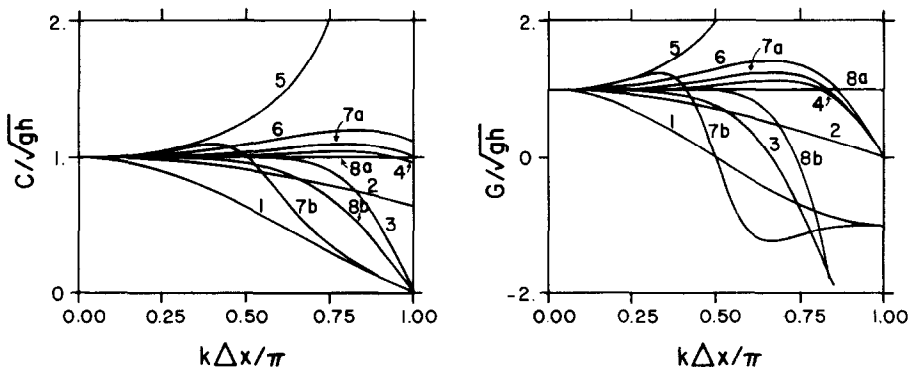


FIG. 1.   Nondimensional phase and group velocities for several spatial discretizations. Analytic values are identically equal to 1.0 and scarcely distinguishable from those of (D8a).

## TABLE I
### Spatially Discretized Shallow Water Equations

Continuity Equation(s):
$$\frac{\partial}{\partial t}\sum_{i=-1}^{1} a_i z_{j+i} + \frac{h}{\Delta x}\sum_{i=-3/2}^{3/2} b_i u_{j+i} = 0$$

Momentum Equation(s):
$$\frac{\partial}{\partial t}\sum_{i=-1}^{1} c_i u_{j+1} + \frac{g}{\Delta x}\sum_{i=-3/2}^{3/2} d_i z_{j+i} = 0$$

| Spatial Discretization | References | $a_{-1}$ | $a_{-1/2}$ | $a_0$ | $a_{1/2}$ | $a_1$ | $b_{-3/2}$ | $b_{-1}$ | $b_{-1/2}$ | $b_0$ | $b_{1/2}$ | $b_1$ | $b_{3/2}$ | $c_{-1}$ | $c_{-1/2}$ | $c_0$ | $c_{1/2}$ | $c_1$ | $d_{-3/2}$ | $d_{-1}$ | $d_{-1/2}$ | $d_0$ | $d_{1/2}$ | $d_1$ | $d_{3/2}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| D1 | [19, 23] | | | 1 | | | | $-\frac{1}{2}$ | | | | $\frac{1}{2}$ | | | | 1 | | | | $-\frac{1}{2}$ | | | | $\frac{1}{2}$ | |
| D2 | [19] | | | 1 | | | | | $-1$ | | $1$ | | | | | 1 | | | | | $-1$ | | $1$ | | |
| D3 | [19, 23] | $\frac{1}{6}$ | | $\frac{2}{3}$ | | $\frac{1}{6}$ | | $-\frac{1}{2}$ | | | | $\frac{1}{2}$ | | $\frac{1}{6}$ | | $\frac{2}{3}$ | | $\frac{1}{6}$ | | $-\frac{1}{2}$ | | | | $\frac{1}{2}$ | |
| D4 | [19] | $\frac{1}{6}$ | | $\frac{2}{3}$ | | $\frac{1}{6}$ | $-\frac{1}{8}$ | | $-\frac{5}{8}$ | | $\frac{5}{8}$ | | $\frac{1}{8}$ | $\frac{1}{6}$ | | $\frac{2}{3}$ | | $\frac{1}{6}$ | $-\frac{1}{8}$ | | $-\frac{5}{8}$ | | $\frac{5}{8}$ | | $\frac{1}{8}$ |
| D5 | | $\frac{-1}{2\Delta x}$ | | | | $\frac{1}{2\Delta x}$ | | $\frac{1}{\Delta x}$ | | $\frac{-2}{\Delta x}$ | | $\frac{1}{\Delta x}$ | | $\frac{-1}{2\Delta x}$ | | | | $\frac{1}{2\Delta x}$ | | $\frac{1}{\Delta x}$ | | $\frac{-2}{\Delta x}$ | | $\frac{1}{\Delta x}$ | |
| D6 | [28] | | | 1 | | | | | $-1$ | | $1$ | | | $\frac{1}{6}$ | | $\frac{2}{3}$ | | $\frac{1}{6}$ | | | $-1$ | | $1$ | | |
| D7a | [16] | $-\frac{1}{10}$ | $\frac{1}{5}$ | $\frac{4}{5}$ | $\frac{1}{5}$ | $-\frac{1}{10}$ | $\frac{1}{6}$ | | $-\frac{2}{3}$ | | $\frac{2}{3}$ | | $-\frac{1}{6}$ | $-\frac{1}{10}$ | $\frac{1}{5}$ | $\frac{4}{5}$ | $\frac{1}{5}$ | $-\frac{1}{10}$ | $\frac{1}{6}$ | | $-\frac{2}{3}$ | | $\frac{2}{3}$ | | $-\frac{1}{6}$ |
| D8a | [2] | $\frac{1}{10}$ | $\frac{1}{5}$ | $\frac{4}{5}$ | $\frac{1}{5}$ | $\frac{1}{10}$ | $\frac{1}{2}$ | | $-2$ | $-1$ | $2$ | | $-\frac{1}{2}$ | $\frac{1}{10}$ | $\frac{1}{5}$ | $\frac{4}{5}$ | $\frac{1}{5}$ | $\frac{1}{10}$ | $\frac{1}{2}$ | | $-2$ | $-1$ | $2$ | | $-\frac{1}{2}$ |

*Note.* $\Delta x = x_{j+1} - x_j = x_{j+1/2} - x_{j-1/2}$ is the distance between end-element nodes and mid-element nodes.

TABLE II

Dispersion Relationships for the Spatial Discretizations of Fig. 1

| Spatial Discretization | $\omega \Delta x/(gh)^{1/2}$ |
|---|---|
| D1 | $\pm \sin(k\Delta x)$ |
| D2 | $\pm 2 \sin\left(k\,\dfrac{\Delta x}{2}\right)$ |
| D3 | $\pm \dfrac{3 \sin(k\Delta x)}{(2 + \cos(k\Delta x))}$ |
| D4 | $\pm \dfrac{3}{4}\left(\dfrac{\sin(\frac{3}{2}k\Delta x) + 5 \sin(k\Delta x/2)}{2 + \cos(k\Delta x)}\right)$ |
| D5 | $\pm 2\,\dfrac{(1 - \cos(k\Delta x))}{\sin(k\Delta x)}$ |
| D6 | $\pm \left[\dfrac{6(1 - \cos(k\Delta x))}{2 + \cos(k\Delta x)}\right]^{1/2}$ |
| D7a | $0,\ \pm 2 \sin\left(\dfrac{k\Delta x}{2}\right)\left[\dfrac{2(4 - \cos(k\Delta x))}{(2 + \cos(k\Delta x))(3 - \cos(k\Delta x))}\right]^{1/2}$ |
| D7b | $0,\ \pm\sin(k\Delta x)\left[\dfrac{2(4 - \cos(2k\Delta x))}{(2 + \cos(2k\Delta x))(3 - \cos(2k\Delta x))}\right]^{1/2}$ |
| D8a | $\pm 2 \sin\left(\dfrac{k\Delta x}{2}\right)\left[\dfrac{(10 - \cos^2(k\Delta x/2))^{1/2} \pm 2 \cos(k\Delta x/2)}{2 - \cos^2(k\Delta x/2)}\right]$ |
| D8b | $\pm\sin(k\Delta x)\left(\dfrac{(10 - \cos^2(k\Delta x))^{1/2} \pm 2 \cos(k\Delta x)}{2 - \cos^2(k\Delta x)}\right)$ |

and the introduction of mid-element nodes (e.g., [20]). Consequently, waves of length $\Delta x$ may exist in the approximated variables. In order to represent these waves in Fig. 1, either the upper limit for $k\Delta x$ should be extended to $2\pi$, or $\Delta x$ should be halved. The latter approach was adopted.

Ideally, the phase and group velocities of a spatial discretization should be close to their analytic values. Few of the discretizations shown in Fig. 1 are close, particularly for large wavenumbers. Since $2\Delta x$ waves are frequently troublesome in shallow water models [24], their behaviour is important. Figure 1 shows that $2\Delta x$ waves for (D1), (D3), (D7b), and (D8b) have zero phase velocity and thus do not propagate. Their corresponding group velocities are negative. Hence the energy associated with these waves is moving, but in the wrong direction. One might therefore see the same generation of spurious waves at an interface with these discretizations as was demonstrated by Trefethen [21]. Furthermore, an inappropriate choice of boundary

conditions could also cause the instability that he mentions. Zero group velocities for discretizations (D2), (D4), (D6), and (D7a) indicate that although $2\Delta x$ waves are propagating, the associated energy is not.

Cutoff frequencies exist for (D1), (D3), (D7b), and (D8b) since they all attain zero group velocity for waves longer than $2\Delta x$. Precise values for these frequencies can be calculated from the dispersion relationships in Table II. As demonstrated by Trefethen [21], when using one of these discretizations in a problem containing an interface (e.g., due to a mesh refinement or a change of coefficient), it may happen that a wave incident from one side has a frequency which is not sustainable on the other. As shown by Vichnevetsky [22], difficulties may also arise with time-varying boundary conditions which oscillate at frequencies higher than the cutoff.

Graphically it would seem that (D8a) is the best spatial discretization. In actual computations, it may not be. Representation (D8a) ignores waves shorter than twice the distance between end-element nodes. This is valid provided measures such as artificial viscosity can effectively eliminate these waves. Otherwise, intra-element oscillations can exist and may contaminate the highly accurate longer waves. An additional complication for (D8) is that only two of its four dispersion relationships (as shown in Table II) have nondimensional phase (and group) velocities whose magnitudes tend to 1.0 as $k\Delta x$ tends to zero. The other two tend to $\pm 5$, and thus are not consistent with the analytic solution. If waves represented by these spurious curves are generated and sustained in a numerical model, further inaccuracies can be expected. Cullen [2] investigates (D8) in more detail.

Provided intra-element oscillations can be avoided, (D7) is another promising spatial discretization. Walters and Carey [24] recommend the linear basis functions for $z(x, t)$, and the quadratic functions for $u(x, t)$, since this choice generates fewer spurious modes than vice versa. Consistent with this analysis, they remark that a small amount of dissipation may be necessary in the nonlinear equations to remove $2\Delta x$ waves in the velocity field.

Figure 1 also indicates good accuracy with (D4) and (D6). In fact (D4) is superior to the three vorticity–divergence formulations investigated by Williams [27]. Unfortunately, a convenient triangular element analog in two dimensions is not apparent, especially for the case of irregular geometry. Discretization (D6) is also difficult to extend to two dimensions since the piecewise constant variable is discontinuous at the inter-element nodes [24]. However, the "wave equation" finite element method of Gray and Lynch [5, 10] has the same dispersion relationship (thus phase and group velocity) as (D6) in one dimension [3], and has been extended to two dimensions. In some sense it may therefore be viewed as a two-dimensional version of (D6).

A realistic positive value for $\tau$ would have little effect on the plots of Fig. 1. The analytic nondimensional phase velocity would become slightly less than 1. for all wavenumbers, and the associated group velocity would become slightly greater than 1. All velocities for the eight spatial discretizations would also exhibit small shifts in varying degrees. All velocities would equal zero for small $k$, since a wave solution to (1a) and (1b) cannot be supported there. As seen from (3), this occurs when $\omega_r$ is imaginary. However, a more significant change would occur with (D7). Its secondary

or spurious dispersion relationship would no longer be zero, thereby permitting the existence of associated spurious waves in the numerical solution.

The preceding analysis illustrates how phase and group velocity accuracy can aid in the selection of a spatial discretization. Because of its restrictive nature (i.e., one-dimensional linearized equations, constant $\Delta x$ and $h$, $\tau = 0$.) and the fact that economy of the calculations has been ignored, an analysis of this type should be only one part of the selection process. It must also be stressed that implementation of a time stepping technique can change the relative accuracy of two spatial discretizations. Analyses of the fully discretized equations should therefore always accompany analyses of spatial discretizations.

## 3. A Class of ODE Methods: Linear Two-Step Methods

Pinder and Gray [16] and Walters and Cheng [25] recommend a centered time stepping scheme for solving the system of ODEs that arise from the spatially discretized versions of (1a) and (1b). With the one-step method

$$y^{n+1} - y^n = \Delta t [\theta f^{n+1} + (1 - \theta) f^n] \tag{7}$$

for solving the ODE

$$\frac{\partial y}{\partial t} = f(y), \tag{8}$$

a centered scheme is characterized by identical central times for the left and right sides of (7). In this case, the left and right sides are respectively centered over $t_{n+1/2}$ and $t_{n+\theta}$. Therefore $\theta = 0.5$ represents a centered scheme. With uncentered versions of (7), these authors found either excessive damping of the solution or substantially incorrect phases. This result is not surprising since (7) is second order accurate (i.e., the truncation error is $O(\Delta t^3)$) when $\theta = 0.5$, and first order otherwise. In order to have the desirable behavior associated with centered schemes it is therefore wise to insist on at least second order accuracy.

A broad class of numerical techniques for solving ODEs are multistep methods. As the trapezoid scheme is the only second order one-step method, within this class (outside this class, the second order Runge–Kutta method is also one-step), the larger class of two-step methods will be considered. For solving (7), all two-step methods are characterized by the formula [4, 7]

$$a_2 y^{n+2} + a_1 y^{n+1} + a_0 y^n = \Delta t (b_2 f^{n+2} + b_1 f^{n+1} + b_0 f^n), \tag{9}$$

where $a_2$, $a_1$, $a_0$, $b_2$, $b_1$, and $b_0$ are real numbers. In the sense that both sides of (9) can be multiplied by any constant and not alter the relationship, this equation

requires a normalization. Lambert [7] suggests $a_2 = 1$, while Gear [4] recommends

$$b_0 + b_1 + b_2 = 1. \tag{10}$$

The latter convention will be adopted here.

For at least second order accuracy, only two parameters remain free. Choosing them to be $a_2$ and $b_2$, the others are specified as

$$
\begin{aligned}
a_0 &= a_2 - 1, & b_0 &= \tfrac{1}{2} - a_2 + b_2, \\
a_1 &= 1 - 2a_2, & b_1 &= \tfrac{1}{2} + a_2 - 2b_2.
\end{aligned}
\tag{11}
$$

These relationships are derived in [4, 7, 8].

All second order two-step methods are centered. Some familiar ones with their $(a_2, b_2)$ values are: trapezoid or Crank–Nicolson $(1, \tfrac{1}{2})$; Gear stiffly stable [4], $(\tfrac{3}{2}, 1)$; Adams–Bashforth, $(1, 0)$; Adams–Moulton, $(1, \tfrac{5}{12})$; Milne, $(\tfrac{1}{2}, \tfrac{1}{6})$; and leapfrog, $(\tfrac{1}{2}, 0)$. Explicit methods are characterized by $b_2 = 0$. Third order methods include Adams–Moulton and have the additional constraint

$$b_2 = \tfrac{1}{2}a_2 - \tfrac{1}{12}. \tag{12}$$

Milne's method is fourth order.

A necessary and sufficient condition for stability of a second order two-step method is $a_2 \geqslant 0.5$. This follows from the root condition [4]. Multistep methods which satisfy this condition are called zero-stable [7].

## 4. THE FULLY DISCRETIZED EQUATIONS

This section studies the effects of combining an ODE from the class of second order two-step methods with the particular spatial discretization (D3). Although Section 2 shows that (D3) is not the most accurate discretization, it is commonly used and effectively illustrates the analysis technique. Similar analyses have also been performed for (D7) and the "wave equation" approach of Gray and Lynch [5, 10].

The four numerical eigenvalues arising from a two-step method solution of the ODEs resulting from (D3) are

(i)   for $k\Delta x = \pi$,

$$\lambda_1 = \lambda_2 = 1, \qquad \lambda_3 = \lambda_4 = \frac{(a_2 - 1)}{a_2}; \tag{13a}$$

(ii)   otherwise,

$$\lambda_{1,2} = \frac{-T_1 \pm (T_1^{\,2} - 4T_0 T_2)^{1/2}}{2T_2}$$

$$\lambda_{3,4} = \frac{-R_1 \pm (R_1^{\,2} - 4R_0 R_2)^{1/2}}{2R_2}, \tag{13b}$$

where

$$T_0 = a_0 + b_0 S_+, \qquad T_1 = a_1 + b_1 S_+, \qquad T_2 = a_2 + b_2 S_+,$$

$$R_0 = a_0 + b_0 S_-, \qquad R_1 = a_1 + b_1 S_-, \qquad R_2 = a_2 + b_2 S_-,$$

$$S_\pm = \frac{1}{2}\tau\Delta t \pm i\left[gh\left(\frac{\Delta t}{\Delta x}\right)^2\left(\frac{3\sin(k\Delta x)}{2+\cos(k\Delta x)}\right)^2 - \left(\frac{1}{2}\tau\Delta t\right)^2\right]^{1/2}. \qquad (13c)$$

Complex eigenvalues occur in conjugate pairs corresponding to progressive and retrogressive travelling waves. Two progressive waves arise when all four roots have nonzero imaginary parts. In this case, only the principal root represents the desired solution, the other is called spurious or parasitic. Real valued eigenvalues signify a nonpropagating wave and frequently arise for $2\Delta x$ waves (when $k\Delta x = \pi$).

Assuming a travelling wave solution and no multiple eigenvalues, the component of $z(x, t)$ (or $u(x, t)$) with wavenumber sampling $k\Delta x$ has the following complex valued amplitude at time step $n$:

$$z_n(k\Delta x) = \sum_{j=1}^{4} P_j(k\Delta x)(\lambda_j(k\Delta x))^n \qquad (14)$$

for some functions $P_j(k\Delta x)$. As $n$ increases this amplitude is dominated by the eigenvalue with the largest modulus. For stability, it is necessary that the dominant eigenvalue have modulus less than or equal to 1.0 for all $k\Delta x$. This is a special case of the von Neumann stability condition

$$|\lambda| \leqslant 1 + O(\Delta t). \qquad (15)$$

The $O(\Delta t)$ term is usually omitted (e.g., [12, 18]) when the exact solution does not grow exponentially. Since $h(x)$ is constant and $\tau \geqslant 0$, this is the case here.

Each numerical eigenvalue has its own dispersion relationship and thus phase and group velocity. Dominant eigenvalues imply dominant dispersion relationships and dominant velocities. Since the same eigenvalue may not be dominant for all $k\Delta x$, switch points may exist. At these points the dominant dispersion relationship is usually multivalued. Numerical difficulties can be expected at wavenumbers where the parasitic dispersion relationship dominates. If through boundary conditions, initial conditions, or an interface, parasitic waves or wave packets are generated at such wavenumbers, they will eventually overshadow principal waves of the same length.

Associated with each dominant dispersion relationship is a dominant or favoured wavenumber. At this $k\Delta x$ value, the amplitude of the dominant eigenvalue is maximum. A favoured wavenumber therefore denotes the wave which grows most rapidly, or decays most slowly, as time advances. Dissipative schemes such as Lax–Wendroff have amplitudes curves which decrease with increasing wavenumber [12]. Small wavenumbers therefore dominate and shorter waves are increasingly damped. Schemes where $k\Delta x = \pi$ is favoured can expect problems with $2\Delta x$ waves.

Figure 2 illustrates the numerical eigenvalues, dispersion curves, and phase and group velocities for three two-step methods. Values for the analytic and spatially discretized solutions, and the discrete numerical solution (i.e., from the eigenvalues of the matrix equation solved at each time step) arising from a ring domain test model with 10 grid points, are also included. Results are parameterized in terms of

$$f_1 = \frac{\tau \Delta x}{(gh)^{1/2}}, \qquad \text{and} \qquad f_2 = (gh)^{1/2} \frac{\Delta t}{\Delta x}. \tag{16}$$

The latter parameter is commonly referred to as the Courant number.

The eigenvalue spectra diagrams show ranges of the four numerical eigenvalues for $k\Delta x$ in the interval $(0, \pi]$. The unit circle is included as a reference for stability. All eigenvalue paths lie entirely in either the nonnegative imaginary half plane and correspond to progressive wave solutions, or in the nonpositive imaginary half plane and correspond to retrogressive waves. For these examples, paths of the principal numerical solutions lie almost entirely in either the first or fourth quadrants, while the spurious numerical solutions are in the second and third. As $k\Delta x$ increases from zero, the principal progressive numerical eigenvalue moves in a counterclockwise direction from the positive real axis. When $k\Delta x$ is approximately $2\pi/3$, this excursion reverses and returns to the real axis along exactly the same path. Platzman $|15|$ refers to the $k$ value at this turning point as the folding wavenumber $k_f$ and discusses the aliasing problems that result from its existence. At the folding wavenumber, the real part of the principal progressive dispersion curve is maximum and the corresponding group velocity is zero. The associated frequency is therefore a cutoff frequency. Although the analytic and principal numerical eigenvalue paths are close when $k < k_f$, it cannot be determined from this diagram if adjacent points in these paths arise for the same $k\Delta x$ value.

The second series of diagrams in Fig. 2 permits such a comparison by plotting angular displacement (real part of the dispersion relationship) as a function of $k\Delta x$. Only the progressive wave solutions have been shown. Notice that curves arising from the fully discretized numerical solution are determined to a large extent by those solely due to the spatial discretization. However, the principal numerical dispersion curve in the second example (leapfrog method) does illustrate that a subsequent time discretization can improve accuracy for some range of $k\Delta x$. For larger $k\Delta x$ values, the first and third examples demonstrate that the spurious numerical solution can provide a better approximation to the analytic dispersion curve than the principal numerical solution.

The third series of diagrams permits determination of instability and the dominant numerical eigenvalue. The first example shows a switch of dominance between the principal and spurious numerical eigenvalues. Specifically, the principal eigenvalue is only dominant for $0.325 < k\Delta x/\pi < 0.880$. In the second example, the spurious eigenvalue is both dominant and unstable, while in the third, the principal eigenvalues is dominant and unstable in the neighbourhood of the folding wavenumber.

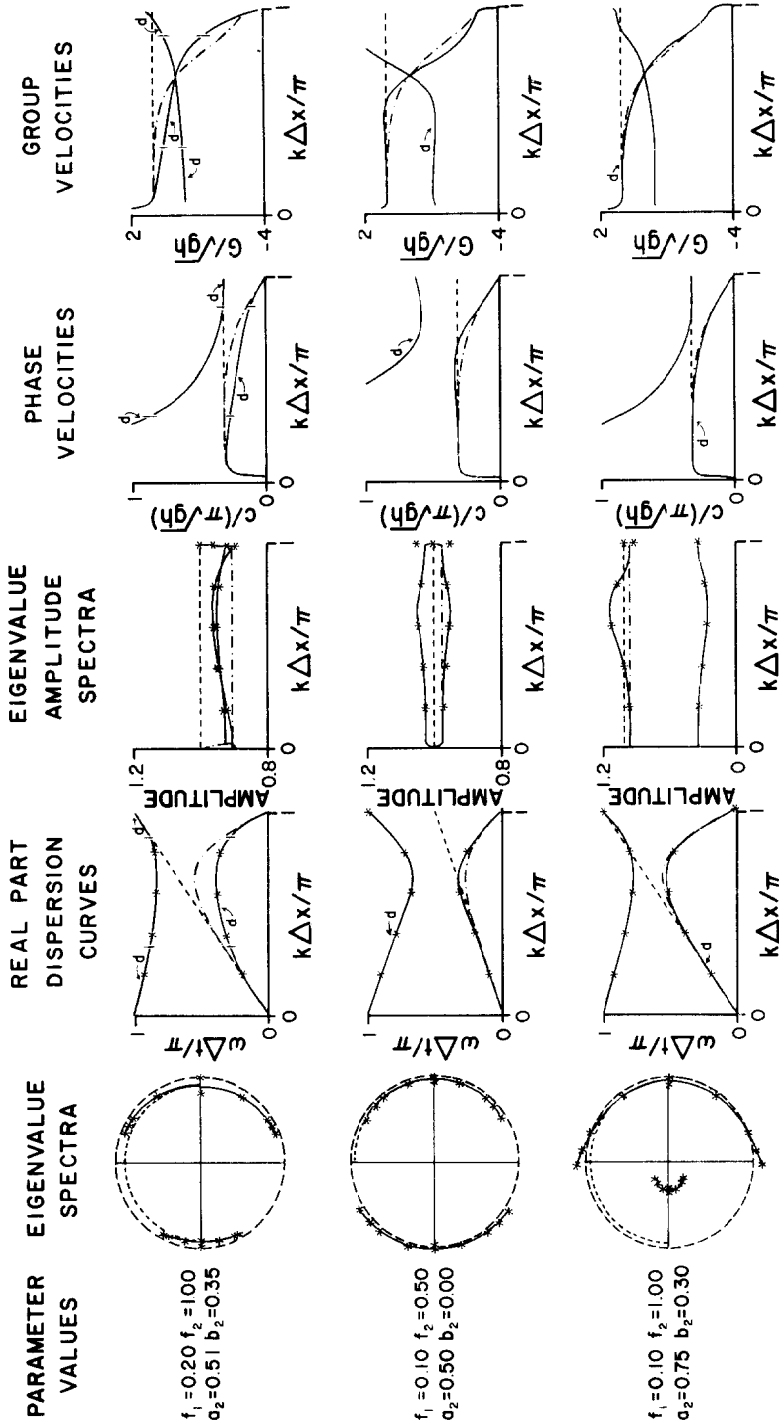In each example the dominant and folding wavenumbers are identical and approx-

FIG. 2. Eigenvalue spectra, dispersion curves, and phase and group velocities for three finite element methods.(———) stability limit, (- - -) analytic values, (——) values from spatial discretization, (—) fully discretized continuous spectrum, $d$ = dominant curve, $*$ = ringed domain solution.

imately equal to $2\pi/3$. This is apparent from the eigenvalue spectra plots since the reversal point in each eigenvalue path corresponds to the maximum amplitude.

The fourth and fifth diagrams plot the corresponding nondimensional phase and group velocities. (They do not have a limiting value of 1.0 at $k\Delta x = 0$ because nonzero friction does not permit a wave solution there.) Phase velocities for the first example are less than the analytic values thereby indicating that the numerical wave solutions travel too slowly. In fact, $2\Delta x$ waves do not travel at all. Group velocities for these cases are also too small and are seen to become negative beyond the folding wavenumber. This indicates energy propagating in the wrong direction. Dominance of the spurious numerical eigenvalue in the second example results in substantially inaccurate phase and group velocities. Switching of the dominant eigenvalue in the first example produces double-valued phase and group velocities at the switch points.

## 5. AN ACCURACY ANALYSIS

The preceding discussion suggests three functions to measure accuracy of the numerical solution; one for each the amplitude, phase velocity, and group velocity. Their respective definitions are

$$M_A = \left| \frac{\lambda_n}{\lambda_a} \right|, \tag{17a}$$

$$M_C = \frac{C_n - C_a}{C_a}, \tag{17b}$$

$$M_G = \frac{G_n - G_a}{G_a}, \tag{17c}$$

where $\lambda_n$ is the dominant progressive numerical eigenvalue, $\lambda_a$ is the analytic progressive eigenvalue, and $C_n$, $C_a$, $G_n$, $G_a$ are the corresponding phase and group velocities.

The velocity accuracy measures are simply relative errors. Negative values denote waves travelling too slowly while zero values are optimal. For example, $-0.01$ denotes a numerical velocity which is 1% too slow. The amplitude measure is a ratio denoting the growth (or decay) factor per time step relative to the analytic solution. Values greater than the optimum of 1. signify a solution which will decay too slowly or grow too rapidly. After $n$ time steps, the ratio of the numerical amplitude to the analytic will be $(M_A)^n$.

Figures 3 and 4 show the accuracy measure contours as functions of the second order two-step parameters $a_2$ and $b_2$ (designated there by $A_2$ and $B_2$). In all plots, a dotted line represents third order methods while asterisks locate the six familiar methods listed in Section 3. The stability region is bounded to the left by $a_2 = 0.5$ and from below by the heavy solid line. All methods corresponding to $(a_2, b_2)$ values

FIG. 3.   Accuracy measure values for $f_1 = 0.1$ and $f_2 = 1.0$.

outside this region have a dominant eigenvalue modulus greater than 1.0 for some $k\Delta x$. They will therefore be unstable.

Figure 3 shows accuracy measure changes as $k\Delta x$ increases and $f_1$ and $f_2$ remain fixed at 0.1 and 1.0, respectively. Notice that the most accurate methods may not coincide for all three measures or even lie within the stability region. Thus a method which is most accurate for one $k\Delta x$ value may be unstable for others. Also notice that a method which has more accurate phase velocity may not have more accurate group velocity, and vice versa. For example, with $k\Delta x/\pi = 0.4$ Adams–Bashforth

FIG. 4.   Accuracy measure values for $f_2 = 1.0$ and $k\Delta x \geqslant 0.1$.

$((a_2, b_2) = (1, 0))$ has a better phase velocity than Adams–Moulton $((a_2, b_2) = (1, \frac{5}{12}))$; but the latter has a better group velocity. (Since both methods are unstable, this is admittedly a poor example.) High accuracy measure values along the lower $b_2$ axis arise because the parasitic eigenvalue is dominant.

In most numerical models, desired waves have $k\Delta x < \pi/10$ (wavelengths longer than $20\Delta x$). Figure 4 illustrates the accuracy measure changes as $f_1$ increases with $k\Delta x$ and $f_2$ fixed at $\pi/10$ and $1.0$, respectively. The stability region has the lower boundary $b_2 = 0.5a_2$ for $f_1 = 0.0$ and becomes less restrictive as $f_1$ increases. In all cases, the most accurate and stable methods lie on, or very close to this line. With

$f_2 = 0.5$, the same series of plots reveals similar patterns but a less restrictive lower boundary for stabilty. Optimal accuracy now occurs along the dotted line or as close to it as stability permits.

Choosing the most accurate two-step method will depend on $f_1, f_2, k\Delta x$, and the relative importance of amplitude and velocity. In most cases, accuracy measure and stability results indicate that all methods along the line $b_2 = 0.5a_2$ are very good choices. In fact, the four numerical eigenvalues for all methods in this subset are

$$\lambda_1 = \frac{1 - S_+/2}{1 + S_+/2}, \tag{18a}$$

$$\lambda_2 = \frac{1 - S_-/2}{1 + S_-/2}, \tag{18b}$$

and

$$\lambda_3 = \lambda_4 = \frac{a_2 - 1}{a_2}, \tag{18c}$$

were $S_+$ and $S_-$ are defined in (13b). The first two are principal numerical eigenvalues and are independent of $a_2$. They are identical for all methods. The other two spurious eigenvalues vary with the two-step method and are constant for all $k\Delta x$. Hence the associated numerical solution will not propagate. Provided $\lambda_1$ and $\lambda_2$ dominate, all accuracy measures (and numerical solutions after many time steps) for this subset of methods will be identical. However, for some wavenumbers, $\lambda_3$ may dominate and the accuracy measures and numerical solution will vary with $a_2$. From this perspective, the Crank–Nicolson method ($a_2 = 1$) is optimal within the subset because both its spurious eigenvalues are zero and thus can never dominate. Furthermore, being a one-step method it should also have the most economical storage requirements. However, it is implicit and may be expensive with regard to computing time. No second order explicit method exists within the subset.

## 6. VALIDATION OF THE ACCURACY MEASURES

In order to validate the previous accuracy measure analysis, truncation errors were calculated and several numerical tests were performed. Depth, $\Delta x$ and $\Delta t$ were constant throughout each test and the additional complication of boundary conditions was avoided by choosing a ring as the test domain. All tests were initial value problems where the propagation characteristics of one or two progressive waves were studied as they travelled around the ring. Numerical solutions were obtained with a Galerkin finite element method which combined (D3) with a second order two-step method for solving the ODEs in time.

Two series of tests were made. The first was designed for checking only amplitude and phase velocity and was characterized by initial conditions which were spatially

sampled values of a travelling wave (as given by (3)) with wavelength equal to the ring circumference.

Five test problems were selected, each with $f_1, f_2$, and $k\Delta x/\pi$ values corresponding to one of the plots in Figs. 3, 4, or the counterpart to Fig. 4 with $f_2 = 0.5$. Wavelength and depth were chosen so that the resultant problem would be realistic for semi-diurnal tides along a one-dimensional continental shelf.

Each test problem was run for approximately ten periods and solved with ten different second order two-step methods. If the numerical solution remained stable, a spectral analysis of the $z(x, t)$ and $u(x, t)$ values over the ring was first used to determine if the original travelling wave had dispersed into other wavelengths. As expected for linear equations, this never occurred. The amplitude and phase lag for the wave were then calculated and compared to the analytic result. The amplitude change per time step and the nondimensional phase velocity were also calculated and compared to the values predicted by a dispersion analysis of the numerical method. From these model values, ratios were formed as in (17) and compared to the accuracy measure values.

The ten second order two-step methods which were used to solve the five test problems were loosely selected upon the following criteria:

(i)   representation of most regions in the domain $0.5 \leqslant a_2 \leqslant 2.5$, $0.0 \leqslant b_2 \leqslant 1.5$;

(ii)   inclusion of some well-known methods;

(iii)   inclusion of some expected unstable methods (i.e., those for which $b_2 < \frac{1}{2}a_2$),

(iv)   inclusion of some methods with small truncation errors.

The chosen methods are listed in Table III and shown in Fig. 5.

The truncation error which arises from combining a second order two-step method with (D3) is

$$
\Delta x \, \Delta t^3 \frac{\partial^2}{\partial t^2} \left( 1 + \frac{\Delta x^2}{6} \frac{\partial^2}{\partial x^2} \right) \left[ \frac{1}{6} \frac{\partial z}{\partial t} + \frac{1}{2} \left( \frac{1}{2} - a_2 + 2b_2 \right) h \frac{\partial u}{\partial x} \right]
$$

$$
+ \Delta x \, \Delta t^4 \left( \frac{2a_2 - 1}{24} \right) \frac{\partial^3}{\partial t^3} \left( 1 + \frac{\Delta x^2}{6} \frac{\partial^2}{\partial x^2} \right) \left[ \frac{\partial z}{\partial t} + 2h \frac{\partial u}{\partial x} \right]
$$

$$
+ O(\Delta x^5) \, O(\Delta t) + O(\Delta t^5) \, O(\Delta x), \tag{19}
$$

where $z = z(x, t)$ and $u = u(x, t)$ are the true solutions to (1a) and (1b). Methods whose parameterization satisfies (12) therefore have a smaller truncation error, and Milne's method is minimal. The test methods denoted by $(a_2, b_2) = (1.0, .417)$ and $(2.0, .917)$ have smaller error constants than the others. In the subsequent discussion they will be referred to as M3 and M4, respectively.

Results for the five test problems are given in Table IV. Initial conditions at times 0 and $\Delta t$ were specified exactly. A run was judged unstable when the absolute value

TABLE III

Second Order Two-Step Methods Used in the Numerical Tests

| Parameter Value | Crank–Nicolson | | Adams–Moulton | | | | Stable | | | Leapfrog |
|---|---|---|---|---|---|---|---|---|---|---|
| $a_2$ | 1.0 | 0.75 | 1.0 | 2.0 | 0.6 | 2.5 | 1.5 | 1.0 | 0.5 | 0.5 |
| $b_2$ | 0.5 | 0.75 | 0.417 | 0.917 | 0.3 | 1.25 | 1.0 | 1.25 | 1.0 | 0.0 |

of the first elevation point became greater than ten times the initial amplitude $z_0$ of 1.0. Only the Adams–Moulton and leapfrog methods became unstable and only the latter was unstable for all tests. Instability can occur even though $|\lambda| \leqslant 1.$ for the wavelength of the initial travelling wave. During the numerical computations round-off errors produce signals at all wavelengths. So if $|\lambda| > 1.$ for any $k\Delta x$, this signal will grow withut bound and eventually dominate the initial wave.

Table IV shows that the dispersion analysis and test model results are very close. In most cases, differences in the amplitude changes and nondimensional phase velocities occurred in the fifth digit. Consequently, accuracy measures calculated from the test models were virtually the same as those from the dispersion analysis. In fact, only for the second problem and the method $(a_2, b) = (0.5, 1.0)$ are the discrepancies as large as 1%.

The relative performance of M3 and M4 varied with each test. With tests 4 and 5 they most accurately represented phase velocity and amplitude decay. With test 2 and
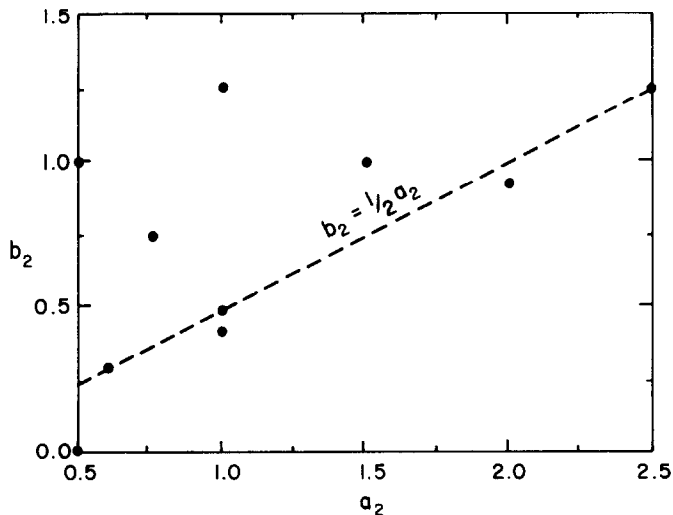


FIG. 5. The $(a_2, b_2)$ coordinates of the second order two-step methods used in numerical tests.

TABLE IV

Results for the first series of numerical tests.

| | | Problem Number and Parameter Values | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | | 2 | | 3 | | 4 | | 5 | |
| | | $t_1$ / $t_2$ / $\frac{k\Delta x}{\pi}$ | | $t_1$ / $t_2$ / $\frac{k\Delta x}{}$ | | $t_1$ / $t_2$ / $\frac{k\Delta x}{}$ | | $t_1$ / $t_2$ / $\frac{k\Delta x}{\pi}$ | | $t_1$ / $t_2$ / $k\Delta x$ | |
| | | .10 / 1.0 / .1 | | .10 / 1.0 / .4 | | .00 / 1.0 / .1 | | .05 / .5 / .1 | | .20 / .5 / .1 | |
| Two-step method parameters $(a_2,b_2)$ | Source of results | $|\lambda|$ | $\frac{c}{(gh)^{\frac12}}$ | $|\lambda|$ | $\frac{c}{(gh)^{\frac12}}$ | $|\lambda|$ | $\frac{c}{(gh)^{\frac12}}$ | $|\lambda|$ | $\frac{c}{(gh)^{\frac12}}$ | $|\lambda|$ | $\frac{c}{(gh)^{\frac12}}$ |
| (1.0,0.5) | analysis | .95234 | .97999 | .96444 | .88055 | 1.0000 | .99184 | .98765 | .99478 | .95148 | .94677 |
| | model | .95234 | .98000 | .96445 | .88055 | 1.0000 | .99184 | .98765 | .99478 | .95148 | .94677 |
| (.75,.75) | analysis | .95613 | .94836 | .94865 | .65918 | .99925 | .95831 | .98793 | .98591 | .95257 | .94122 |
| | model | .95613 | .94836 | .94906 | .65931 | .99925 | .95831 | .98793 | .98591 | .95257 | .94122 |
| (1.0,.417) | analysis | .95157 | .98782 | 1.02658 | .95197 | 1.0004 | .99976 | .98760 | .99681 | .95124 | .94808 |
| | model | unstable | | 1.02658$^a$ | .95197 | unstable | | .98760 | .99681 | .95124 | .94808 |
| (2.0,.917) | analysis | .95234 | .98782 | 1.01216 | .90152 | 1.0010 | .99843 | .98765 | .99681 | .95127 | .94837 |
| | model | .95234 | .98782 | 1.01210$^a$ | .90158 | 1.0010$^a$ | .99843 | .98765 | .99681 | .95127 | .94837 |
| (0.6,0.3) | analysis | .95234 | .97999 | .96444 | .88055 | 1.0000 | .99184 | .98765 | .99478 | .95148 | .94677 |
| | model | .95234 | .98000 | .96480 | .88046 | 1.0000 | .99184 | .98765 | .99478 | .95148 | .94677 |
| (2.5,1.25) | analysis | .95234 | .97999 | .96444 | .88055 | 1.0000 | .99184 | .98765 | .99478 | .95148 | .94677 |
| | model | .95234 | .98000 | .96440 | .88080 | 1.0000 | .99184 | .98765 | .99478 | .95148 | .94677 |
| (1.5,1.0) | analysis | .95346 | .95726 | .87190 | .75200 | .99805 | .97066 | .98773 | .98875 | .95216 | .94241 |
| | model | .95346 | .95726 | .87191 | .75201 | .99805 | .97066 | .98773 | .98875 | .95216 | .94241 |
| (1.0,1.25) | analysis | .95810 | .91886 | .92298 | .56281 | .99747 | .92954 | .98809 | .97721 | .95327 | .93506 |
| | model | .95810 | .91886 | .92257 | .56259 | .99747 | .92954 | .98809 | .97721 | .95327 | .93506 |
| (0.5,1.0) | analysis | .96039 | .92224 | .98897 | .56060 | 1.0000 | .92857 | .98829 | .97751 | .95365 | .93632 |
| | model | .96039 | .92221 | .98127 | .56505 | 1.0009 | .92830 | .98829 | .97750 | .95365 | .93632 |
| (0.5,0.0) | analysis | 1.05397 | 8.9977 | 1.96666 | 1.30461 | 1.0000 | 1.01717 | 1.01274 | 18.999 | 1.05184 | 19.0497 |
| | model | unstable | | unstable | | unstable | | .99379 | .98044 | unstable | |
| Analytic solution | analysis | .95123 | .98725 | .95123 | .99921 | 1.0000 | 1.0000 | .98758 | .99683 | .95123 | .94799 |
| | model | .95123 | .98725 | .95123 | .99921 | 1.0000 | 1.0000 | .98758 | .99683 | .95123 | .94799 |

$^a$ Going unstable but does not satisfy instability criterion.

3, they had accurate velocities but were unstable. With test 1, M4 was most accurate while M3 was unstable. The truncation error analysis therefore predicted the high accuracy, but did not foresee the potential stability problems. This is to be expected.

Figure 6 shows the $z(x,t)$ and $u(x,t)$ profiles around the ring domain for test problem 1 when solved analytically and with the Gear method. The wave is moving leftward and the numerical solution is seen to be too slow (by 3% as calculated from values in Table IV). After 42 time steps, this translates to a phase discrepancy of 21.9° between the numerical and analytic solutions. It is also evident that the numerical amplitude is not decaying as quickly as it should. An error of 0.234% in $|\lambda|$ (from Table IV) in this case compounds to an amplitude error of 10% after 42 time steps.

The second series of tests is similar to the first but permits checking of the group velocity calculations. Two travelling waves of equal amplitude but different

FIG. 6. Elevation and velocity profiles for problem 1 in the first series of numerical tests. (—) numerical solution, (— — —) analytic solution.

wavelengths were now initially specified on the ring domain. As time progresses their combined effect is a short wavelength carrier wave moving inside and at a different speed than a long wavelength envelope (e.g., see Fig. 8). Algebraically, this is seen |1| by considering two close frequency/wavenumber coordinates, $(\omega_1, k_1)$ and $(\omega_2, k_2)$, on a dispersion curve as shown in Fig. 7. Defining

$$\omega_0 = \tfrac{1}{2}(\omega_1 + \omega_2), \qquad k_0 = \tfrac{1}{2}(k_1 + k_2),$$
$$\Delta\omega = \tfrac{1}{2}|\omega_2 - \omega_1|, \qquad \Delta k = \tfrac{1}{2}|k_2 - k_1|, \tag{20}$$



FIG. 7. A sample dispersion curve.

the combined effect of two equal amplitude progressive waves at these frequencies is then

$$A \cos(\omega_1 t - k_1 x) + A \cos(\omega_2 t - k_2 x) = 2A \cos(\omega_0 t - k_0 x) \cos(\Delta k x - \Delta \omega t). \quad (21)$$
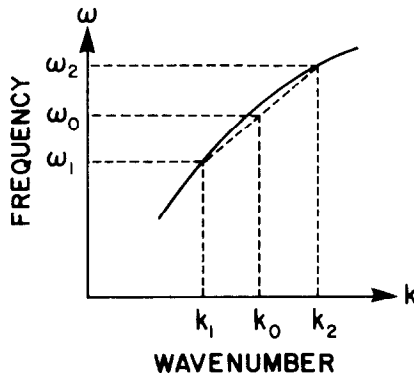
This represents an envelope with wavenumber $\Delta k$ and a carrier wave with wavenumber $k_0$. As $k_1$ and $k_2$ approach $k_0$, the speeds of the envelope and carrier waves approximate the group and phase velocity, respectively, since

$$C(k_0) = \frac{\omega(k_0)}{k_0} = \lim_{k_1, k_2 \to k_0} \left( \frac{\omega_0}{k_0} \right) \quad \text{and} \quad G(k_0) = \frac{\partial \omega(k_0)}{\partial k} = \lim_{k_1, k_2 \to k_0} \left( \frac{\Delta \omega}{\Delta k} \right). \quad (22)$$

So if $k_1$ and $k_2$ are sufficiently close, speeds of the envelope and carrier waves are approximately $G(k_0)$ and $C(k_0)$, respectively.

Numerical tests to measure the group velocity using the preceding approach have an additional complication. Since the eigenvalue amplitudes for wavenumbers $k_1$ and $k_2$ are generally not the same, the two waves do not decay (or grow) at the same rate. So even though they may have equal amplitudes initially, after one step, there is a slight difference. In order that (21) be a reasonable representation of the two waves, all numerical tests were run for only a few time steps.

In all numerical experiments wavelengths $L_1$ and $L_2$ of the two travelling waves were chosen so that the ring circumference $L$ was one lobe of the envelope (as in Fig. 8) and $L/L_1$ and $L/L_2$ were both integer valued. Consequently, for this series of
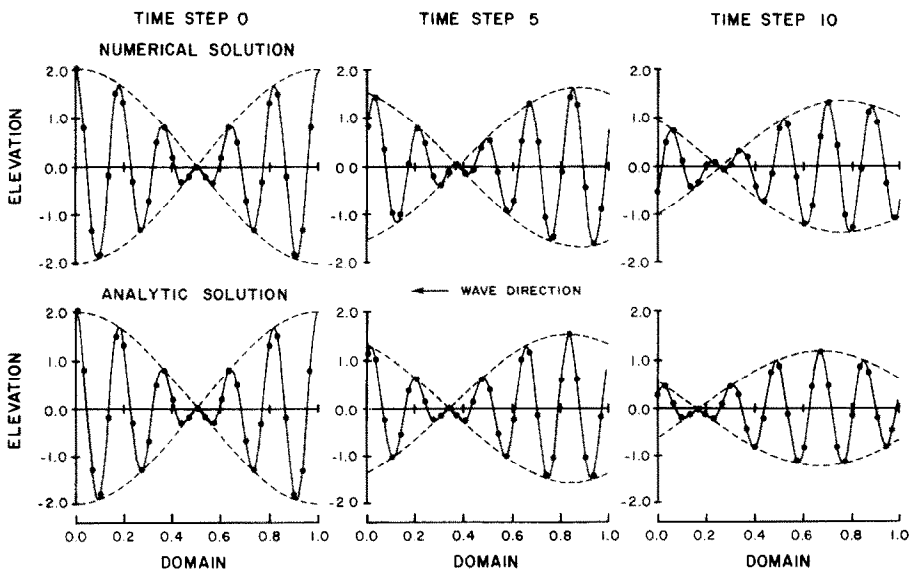


FIG. 8. Numerical and analytic elevation profiles for problem 6 in the second series of numerical tests. The Crank–Nicolson method is used for time stepping. (——) carrier wave, (– – –) envelope, (●) $Z$ value.

tests the parameter values $f_1, f_2$, and $k\Delta x/\pi$ for the six selected test problems could only approximate those for one of the accuracy measure plots in Figs. 3 and 4, or the counterpart to Fig. 4 with $f_2 = 0.5$.

Assuming the $z(x, t)$ and $u(x, t)$ profiles around the domain can be approximated at any time step by

$$A \cos(k_0 x - \phi_1) \cdot B \cos(\Delta kx - \phi_2), \qquad (23)$$

where $k_0$ and $\Delta k$ are specified, the parameters $A$, $B$, $\phi_1$, and $\phi_2$ then characterize the wave packet. Specifically, if $k_1$ and $k_2$ are sufficiently close, then $AB$ is the amplitude of the envelope and $(\phi_1/k_0 t)$ and $(\phi_2/\Delta kt)$, respectively, approximate the phase and group velocity. Values for these parameters were calculated from nonlinear least squares fits to the $z(x, t)$ and $u(x, t)$ profiles. In all cases, velocities and amplitude changes were the same for both variables.

All tests were for only ten time steps with the initial conditions at times 0 and $\Delta t$ specified exactly. The same ten second-order two-step methods were tested in this series as before.

Results for these numerical tests are presented in Table V. Dispersion analysis and test values are not as close as before but due to the several approximations involved, this was expected. Comparisons between tabulated results with the same $f_1, f_2$, and $k\Delta x/\pi$ values (e.g., test 1 in the first series and test 2 in the second) provide an estimate of the error associated with these approximations. In all cases, increasing the number of grid points in the domain would decrease $\Delta k\Delta x$ and reduce this error.

For most tests, differences between the dispersion analysis and test model estimates of the phase and group velocity, and the eigenvalue amplitude were less than 1 %. The two-step method with the poorest correspondence was $(a_2, b_2) = (0.5, 1.0)$. This was also poorest for the first series of tests and is because the parasitic eigenvalue is only slightly smaller than the principal eigenvalue. Many time steps are therefore required before the energy assigned to the parasitic solution by the initial conditions becomes insignificant. In fact, were it not for round-off errors and initial conditions which are, in varying degrees, inconsistent with each numerical method, the results from the first series of tests would be exactly the same as those predicted by the principal numerical eigenvalue.

Throughout the second series of tests, M3 and M4 most accurately represented the phase and group velocity. In fact, only for tests 1 and 6 were their amplitude decay factors not the most accurate. These two tests have the highest $k\Delta x$ values thereby

similar to those of Fig. 2 confirm this. The $|\lambda|$ values for test 6 indicate future instability. Although those for test 1 suggest stability, Fig. 3 indicates magnitudes greater than 1. at other wavelengths. Hence eventual instability can be expected here also.

Due to the shortness of the tests, instability (judged as before) occurred only once. Had the runs been longer, dominant eigenvalues with magnitudes greater than 1.0 would have caused other numerical solutions to become unstable.

TABLE V

Results for the second series of numerical tests

Problem Number and Parameter Values

| Two-step method parameters $(a_2,b_2)$ | Source of results | P1 $\lvert\lambda\rvert$ | P1 $\frac{C}{(gh)^{1/2}}$ | P1 $\frac{G}{(gh)^{1/2}}$ | P2 $\lvert\lambda\rvert$ | P2 $\frac{C}{(gh)^{1/2}}$ | P2 $\frac{G}{(gh)^{1/2}}$ | P3 $\lvert\lambda\rvert$ | P3 $\frac{C}{(gh)^{1/2}}$ | P3 $\frac{G}{(gh)^{1/2}}$ | P4 $\lvert\lambda\rvert$ | P4 $\frac{C}{(gh)^{1/2}}$ | P4 $\frac{G}{(gh)^{1/2}}$ | P5 $\lvert\lambda\rvert$ | P5 $\frac{C}{(gh)^{1/2}}$ | P5 $\frac{G}{(gh)^{1/2}}$ | P6 $\lvert\lambda\rvert$ | P6 $\frac{C}{(gh)^{1/2}}$ | P6 $\frac{G}{(gh)^{1/2}}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Parameters | | $t_1=.10$ | $f_2=1.0$ | $\frac{k_0\Delta x}{\pi}=.268$, $\frac{\Delta k\Delta x}{\pi}=.042$ | $t_1=.10$ | $f_2=1.0$ | $\frac{k_0\Delta x}{\pi}=.104$, $\frac{\Delta k\Delta x}{\pi}=.021$ | $t_1=.00$ | $t_2=1.0$ | $\frac{k_0\Delta x}{\pi}=.104$, $\frac{\Delta k\Delta x}{\pi}=.021$ | $f_1=.05$ | $t_2=.5$ | $\frac{k_0\Delta x}{\pi}=.104$, $\frac{\Delta k\Delta x}{\pi}=.021$ | $f_1=.20$ | $f_2=.5$ | $\frac{k_0\Delta x}{\pi}=.104$, $\frac{\Delta k\Delta x}{\pi}=.021$ | $f_1=.10$ | $f_2=1.0$ | $\frac{k_0\Delta x}{\pi}=.367$, $\frac{\Delta k\Delta x}{\pi}=.033$ |
| (1.0,0.0) | analysis | .956 | .963 | .902 | .952 | .980 | .986 | 1.000 | .991 | .974 | .988 | .995 | .996 | .952 | .951 | 1.044 | .963 | .898 | .710 |
| | model | .956 | .959 | .901 | .952 | .974 | .987 | 1.000 | .990 | .973 | .988 | .994 | .996 | .952 | .946 | 1.047 | .963 | .896 | .710 |
| (.75,.75) | analysis | .961 | .851 | .672 | .956 | .946 | .883 | .999 | .955 | .873 | .988 | .985 | .967 | .953 | .944 | 1.018 | .952 | .690 | .347 |
| | model | .959 | .844 | .640 | .956 | .941 | .884 | .999 | .951 | .874 | .988 | .984 | .967 | .953 | .941 | 1.020 | .940 | .694 | .365 |
| (1.0,.417) | analysis | .955 | .996 | .991 | .952 | .989 | 1.013 | 1.000 | 1.000 | .999 | .988 | .997 | 1.003 | .951 | .952 | 1.051 | 1.009 | .966 | .825 |
| | model | .960 | .995 | .989 | .952 | .988 | 1.014 | 1.000 | 1.000 | .999 | .988 | .997 | 1.003 | .951 | .950 | 1.053 | 1.011 | .963 | .823 |
| (2.0,.917) | analysis | .965 | .983 | .935 | .953 | .989 | 1.009 | 1.001 | .998 | .992 | .986 | .997 | 1.002 | .951 | .953 | 1.051 | 1.003 | .920 | .725 |
| | model | .967 | .976 | .934 | .953 | .988 | 1.010 | 1.001 | .998 | .992 | .988 | .997 | 1.003 | .951 | .950 | 1.054 | 1.007 | .916 | .714 |
| (0.6,0.3) | analysis | .956 | .963 | .902 | .952 | .980 | .986 | 1.000 | .991 | .974 | .988 | .995 | .996 | .952 | .951 | 1.044 | .963 | .898 | .710 |
| | model | .957 | .960 | .902 | .952 | .979 | .988 | 1.000 | .990 | .974 | .988 | .995 | .996 | .952 | .948 | 1.047 | .988 | .896 | .703 |
| (2.5,1.25) | analysis | .956 | .963 | .902 | .952 | .980 | .986 | 1.000 | .991 | .974 | .988 | .995 | .996 | .952 | .951 | 1.044 | .963 | .898 | .710 |
| | model | .957 | .959 | .896 | .953 | .980 | .988 | 1.000 | .990 | .974 | .988 | .995 | .997 | .952 | .948 | 1.048 | .965 | .891 | .695 |
| (1.5,1.0) | analysis | .942 | .892 | .743 | .953 | .956 | .916 | .998 | .968 | .912 | .988 | .988 | .977 | .952 | .946 | 1.024 | .887 | .776 | .510 |
| | model | .942 | .884 | .739 | .953 | .953 | .917 | .998 | .966 | .911 | .988 | .987 | .977 | .952 | .943 | 1.027 | .887 | .771 | .503 |
| (1.0,1.25) | analysis | .955 | .777 | .499 | .958 | .914 | .800 | .997 | .924 | .795 | .988 | .976 | .940 | .954 | .937 | .990 | .929 | .594 | .239 |
| | model | .952 | .765 | .497 | .958 | .907 | .798 | .997 | .918 | .794 | .988 | .973 | .939 | .954 | .933 | .992 | .918 | .590 | .240 |
| (0.5,1.0) | analysis | .976 | .778 | .493 | .961 | .918 | .805 | 1.000 | .923 | .791 | .988 | .976 | .941 | .954 | .939 | .995 | .988 | .592 | .235 |
| | model | .973 | .773 | .529 | .959 | .909 | .797 | .998 | .917 | .776 | .988 | .976 | .946 | .954 | .935 | .998 | .972 | .613 | .288 |
| (0.5,0.0) | analysis | .930¹ | 1.082 | 1.307 | .949¹ | 1.005 | 1.064 | 1.000¹ | 1.019 | 1.058 | .987¹ | 1.001 | 1.013 | .955¹ | .955 | 1.059 | .694 | 1.442 | -.311 |
| | model | .910 | 1.073 | 1.226 | .949 | 1.007 | 1.077 | 1.000 | 1.021 | 1.059 | .987 | 1.001 | 1.014 | .951 | .952 | 1.059 | | unstable | |
| Analytic solution | analysis | .951 | .997 | 1.003 | .951 | .988 | 1.012 | 1.000 | 1.000 | 1.000 | .988 | .997 | 1.003 | .951 | .952 | 1.050 | .951 | .999 | 1.001 |
| | model | .951 | .997 | 1.003 | .951 | .988 | 1.012 | 1.000 | 1.000 | 1.000 | .988 | .997 | 1.003 | .951 | .950 | 1.053 | .951 | .999 | 1.001 |

¹calculated from the sub-dominant eigenvalue.

Figure 8 illustrates the results of solving test problem 6 with the Crank–Nicolson method. Both travelling waves are again moving leftward and decaying at the rate of 4% per time step. The phase velocity is larger than the group velocity causing the carrier wave to move leftward inside the envelope. Analytic values are also shown and after 10 time steps have the following features relative to the numerical solution:

(i)  an envelope amplitude which is about 11% smaller,

(ii)  a carrier wave which is about 61° further advanced,

(iii)  an envelope which is about 16° further advanced.

Even though the numerical group velocity error is larger than the numerical phase velocity (29% vs. 10% from Table V), the envelope has less phase error after 10 time steps because its frequency is smaller by a factor of 11.

The results of both sets of numerical tests validate the accuracy measure calculations of Section 5. Only for the method $(a_2, b_2) = (0.5, 1.0)$ were there notable discrepancies between the test results and the accuracy measure calculations. These can be attributed to the fact that the spurious and principal numerical eigenvalues had virtually the same magnitude. Hence, over the test period, neither one dominated the other.

The performance of methods M3 and M4 confirms the high accuracy predicted by their truncation errors. An investigation of their absolute stability regions could be expected to predict the instability. The relatively good performance of methods in the subset $b_2 = \frac{1}{2}a_2$ is also substantiated by (19). They all have the same error constant. In fact, for each constant $c$, all methods related by

$$b_2 = \tfrac{1}{2}a_2 + c \tag{24}$$

have the same truncation error. This explains the general tendency toward contour lines of this slope in Figs. 3 and 4, and further validates the analysis of Section 5.

## 7. Summary and Conclusions

Let us summarize some highlights of the preceding sections.

In Section 2, several spatial discretizations were examined for the accuracy of their phase and group velocities. Each of the four most accurate were shown to have drawbacks which could affect their performance or implementation in two dimensions.

In Section 4, the class of second order two-step methods was combined with the particular spatial discretization, a Galerkin finite element method with piecewise linear basis functions. The concepts of dominant dispersion relationship, dominant phase and group velocity, and dominant or favoured wavenumber were defined and illustrated. It was shown that the same dispersion relationship may not be dominant for all wavenumbers, and the dominant dispersion relationship may be multivalued at some points.

In Section 5, three accuracy measure functions were defined to facilitate the search for an optimally accurate two-step method. It was shown that the most accurate methods for wave amplitude, phase velocity, and group velocity may not coincide. In particular, it was demonstrated that the best method for phase velocity may not be best for group velocity, and vice versa. Furthermore, a method which most accurately represents either velocity may be unstable. In general, the choice of an optimally accurate method depends on $f_1$, $f_2$, $k\Delta x$, and the relative importance of amplitude, phase velocity, and group velocity.

In Section 6, numerical tests validated the phase velocity, group velocity, and amplitude decay factors which were calculated in Section 5. Only in cases where the spurious and principal eigenvalues had approximately the same magnitude were there significant discrepancies between the analysis and test results. Truncation errors were also calculated and correctly predicted the most accurate methods, when they remained stable.

For a Galerkin finite element method with piecewise linear basis functions, the most accurate and stable two-step methods are characterized by $b_2 = 0.5a_2$. Crank–Nicolson ($a_2 = 1$) is the best among these since it has no spurious eigenvalues. However, it is implicit and may be expensive with respect to computing time. Although a similar analysis has shown that Crank–Nicolson is also the most accurate with (D7), it is not best for all spatial discretizations. Due to second derivatives in their continuity equation, a variation of the linear two-step methods introduced in Section 3 is required for the Gray and Lynch "wave equation" method. An accuracy measure analysis of this approach shows that the Crank–Nicolson analog is not the most accurate [3]. However, in this case, the most accurate methods are virtually independent of wavenumber.

Again it must be emphasized that in the preceding analysis, accuracy was the only consideration in determining a good method. In two-dimensional problems this is no longer a sufficient criterion. Storage requirements and computational costs are now at least as important and may necessitate the use of a method which is less accurate but more economical.

Travelling wave solutions of the form (2) do not exist when the depth in (1a) is assumed nonconstant. With a forcing frequency $\omega$, solutions can now be expected to have the form

$$\begin{pmatrix} z(x, t) \\ u(x, t) \end{pmatrix} = \mathrm{Re} \left[ \begin{pmatrix} z_0(x) \\ u_0(x) \end{pmatrix} \exp(i\omega t) \right], \tag{25}$$

where $z_0(x)$ and $u_0(x)$ are complex functions representing the spatial amplitude and

absence of friction and with a linear depth and specific boundary conditions, Lamb [6] shows that

$$|z_0(x)| \propto J_0(2(\kappa x)^{1/2}), \qquad \text{where} \quad h(x) = h_0 x \quad \text{and} \quad \kappa = \omega^2/gh_0. \tag{26}$$

Lynch and Gray [11] extend this result to the case $h(x) = h_0 x^n$ for integer $n$, and include linear friction as in (1b).

In general, the depth dependency of the solution will be such that waves of constant frequency will have their wavelength decrease and their amplitude increase as they enter shallow water. Phase and group velocity will also become spatially dependent. This same behaviour can be expected in a numerical model, although it may not be accurately represented. Unfortunately, the model will not differentiate between spurious and principal waves; all will become shorter and grow. Although it may not be the case analytically, it is possible that with particular numerical schemes and depth variations, shorter waves will grow more quickly. This could be disastrous, for if the short waves are spurious, they may eventually contaminate the numerical solution.

For some depth variations, it is possible to forecast the rapid growth of short waves with an analysis similar to that of Section 4. Since amplitude is now a function of both space and time, spatial growth curves (with $k\Delta x$ along the abscissa) are required in addition to the temporal growth curves of Fig. 2. In fact, it may be necessary to produce these curves for several depth characteristics (e.g., ratios of depth gradient to depth). Numerical schemes which favour high wavenumbers could then be expected to exhibit rapid growth of short waves and should be avoided.

In the absence of nonlinear terms, short waves may be generated numerically by boundary conditions, an interface, round-off errors, or arise naturally such as through a transition from deep to shallow water. Intuitively, this last source can be controlled by maintaining the same sampling rate per wavelength everywhere in the model. This requires a constant $k\Delta x$ for each wave as it moves throughout the model domain. Therefore any transition from deep to shallow water would not correspond to a rightward shift on a spatial amplitude growth curve which has $k\Delta x$ as the abscissa and which may favour large wavenumbers. Using the dispersion relationship for constant depth, a first approximation to uniform sampling is attained by choosing $\Delta x$ proportional to $(h(x))^{1/2}$. This choice has further appeal. Stability conditions when they arise are frequently in the form

$$\Delta t \leqslant c\, \Delta x/(h(x))^{1/2} \qquad (27)$$

for some constant $c$. Therefore a constant value for $\Delta x/(h(x))^{1/2}$ implies that deep regions of the model where there may be little variation in the numerical solution, are not dictating the largest possible time step. This would be the case with constant $\Delta x$.

However, choosing $\Delta x$ proportional to $(h(x))^{1/2}$ will not affect the generation of short waves due to round-off errors, boundary conditions, or an interface. It may only control their subsequent wavenumber transitions. If an amplitude growth curve shows that these waves will grow faster than the desired longer waves, numerical difficulties can be expected.

We may conclude that the preceding extended Fourier analysis is a valid technique for evaluating the accuracy of both a spatial discretization and a time stepping method. Although the analysis was illustrated with finite element solutions of the

shallow water equations, the concepts are sufficiently general that they could be applied to other numerical techniques (e.g., finite differences) and wave equations. Furthermore, the analysis can be extended to two-dimensional equations, as in |13|. Dispersion curves are then replaced by dispersion surfaces and the phase and group velocities become vectors. Consequently, both magnitude and direction errors may be introduced by a numerical method. Limited experience with such analyses indicates that the numerical phase and group velocity are generally not co-directional, as they should be for shallow water waves.

## ACKNOWLEDGMENTS

## REFERENCES

1. L. BRILLOUIN, "Wave Propagation and Group Velocity," Academic Press, New York, 1960.
2. M. J. P. CULLEN, *J. Comput. Phys.* **45** (1982), 221.
3. M. G. G. FOREMAN, *J. Comput. Phys.*, in press.
4. C. W. GEAR, "Numerical Initial Value Problems in Ordinary Differential Equations," Prentice–Hall, Englewood Cliffs, N.J., 1971.
5. W. G. GRAY AND D. R. LYNCH, *Advan. Water Resour.* **1** (1977), 83.
6. H. LAMB, "Hydrodynamics," 6th ed., Dover, New York, 1932.
7. J. D. LAMBERT, "Computational Methods in Ordinary Differential Equations," Wiley, London, 1973.
8. L. LAPIDUS AND J. H. SEINFELD, "Numerical Solution of Ordinary Differential Equations," Academic Press, New York, 1971.
9. J. J. LEENDERTSE, "Aspects of a Computational Model for Long-Period Water-Wave Propagation," Rand Memorandum, RM–5294–PR, 1967.
10. D. R. LYNCH AND W. G. GRAY, *Comput. Fluids* **7** (1979), 207.
11. D. R. LYNCH AND W. G. GRAY, *J. Hydraul. Div. Amer. Soc. Civ. Eng.* **104** (H10) (1978), 1409.
12. F. MESINGER AND A. ARAKAWA, "Numerical Methods Used in Atmospheric Models," Vol. 1, WMO–ICSU Joint Organizing Committee, GARP Publication Series, No. 17, 1976.
13. R. MULLEN AND T. BELYTSCHKO, *Int. J. Numer. Methods Eng.* **18** (1982), 11.
14. T. S. MURTY, "Seismic Sea Waves—Tsunamis," Department of Fisheries and the Environment, Ottawa, 1977.
15. G. W. PLATZMAN, *J. Comput. Phys.* **40** (1981), 36.
16. G. F. PINDER AND W. G. GRAY, "Finite Element Simulation in Surface and Subsurface Hydrology," Academic Press, New York, 1977.
17. N. PRAAGMAN, "Numerical Solution of the Shallow Water Equations by a Finite Element Method," Ph. D. thesis, Delft University of Technology, 1979.
18. R. D. RICHTMEYER AND K. W. MORTON, "Difference Methods for Initial-Value Problems," Wiley–Interscience, New York, 1967.
19. A. L. SCHOENSTADT, *Mon. Weather Rev.* **108** (1980), 1248.

20. G. STRANG AND G. J. FIX, "An Analysis of the Finite Element Method," Prentice–Hall, Englewood Cliffs, N. J., 1973.
21. L. N. TREFETHEN, *SIAM Rev.* **24** (1982), 113.
22. R. VICHNEVETSKY, *Math. Comput. Simulation* **22** (1980), 98.
23. R. VICHNEVETSKY AND B. PEIFFER, Error waves in finite element and finite difference methods for hyperbolic equations, *in* "Advances in Computer Methods for Partial Differential Equations" (R. Vichnevetsky, Ed.), Assoc. Int. Calcul. Analogique, Ghent, 1975.
24. R. A. WALTERS AND G. F. CAREY, "Analysis of Spurious Oscillation Modes for the Shallow Water and Navier–Stokes Equations," TICOM Report, No. 81–3, 1981.
25. R. A. WALTERS AND R. T. CHENG, *Advan. Water Resour.* **2** (1979), 177.
26. T. J. WEARE, *Comput. Methods Appl. Mech. Eng.* **7** (1976), 351.
27. R. T. WILLIAMS, *Mon. Weather Rev.* **109** (1981), 463.
28. R. T. WILLIAMS AND O. C. ZIENKIEWICZ, *Int. J. Numer. Methods Fluids* **1** (1981), 81.