# NUMERICAL ANALYSIS OF TIME-DEPENDENT BOUSSINESQ MODELS

P. J. VAN DER HOUWEN

*Centre for Mathematics and Computer Science, Post Box 4079, 1009 AB Amsterdam, The Netherlands*

J. MOOIMAN

*Delft Hydraulics, Post Box 177, 2600 MH Delft, The Netherlands*

AND

F. W. WUBS

*University of Groningen, Post Box 800, 9700AV Groningen, The Netherlands*

## SUMMARY

In this paper we analyse numerical models for time-dependent Boussinesq equations. These equations arise when so-called Boussinesq terms are introduced into the shallow water equations. We use the Boussinesq terms proposed by Katapodes and Dingemans. These terms generalize the constant depth terms given by Broer. The shallow water equations are discretized by using fourth-order finite difference formulae for the space derivatives and a fourth-order explicit time integrator. The effect on the stability and accuracy of various discrete Boussinesq terms is investigated. Numerical experiments are presented in the case of a fourth-order Runge–Kutta time integrator.

KEY WORDS  Numerical analysis  Stability  Boussinesq equations

## 1. INTRODUCTION

Boussinesq equations describe the behaviour of fairly long, low-amplitude waves in flow models. The starting point is the shallow water model where terms are added which take into account the effects of wave dispersion. If we define the characteristic parameters

$$\mu := \left( \frac{hk}{2\pi} \right)^2, \qquad \varepsilon := \frac{a}{h},$$

where $h$ is the depth function and $k$ and $a$ are respectively the spatial frequency and amplitude of the waves, then these terms are $O(\mu + \varepsilon)$ with $\mu$ and $\varepsilon$ of the same order of magnitude. The order-$\mu$ and order-$\varepsilon$ terms are respectively related to the frequency and amplitude dispersion.

For very low frequencies $O(\mu)$-terms are negligible so that the Boussinesq model reduces to the shallow water model, while for very-low-amplitude waves the $O(\varepsilon)$-terms are negligible, leading to the linearized wave equations. If both types of terms are neglected, then the wave equation with constant celerity for all waves is recovered.

In practice, Boussinesq models represent a significant improvement over shallow water models because they allow (moderate) curvature of the free surface, non-depth-averaged velocities, non-hydrostatic pressure and wave dispersion.

In this paper we are concerned with the case where $\mu$ and $\varepsilon$ are of magnitude at most $\frac{1}{50}$, so that frequencies and amplitudes satisfying the conditions

$$k = \sqrt{\left(\frac{4\mu\pi^2}{h^2}\right)} \leqslant \frac{\pi\sqrt{2}}{5h}, \qquad a = \varepsilon h \leqslant \frac{h}{50} \tag{1}$$

should be described accurately by the Boussinesq models to be used. Frequencies less than $1/h$ (say) will be considered as the *relevant* frequencies.

In our analysis we concentrate on one-dimensional Boussinesq models, but the analysis can straightforwardly be extended to two-dimensional models. Consider the equations

$$\frac{\partial}{\partial t}\begin{pmatrix} u \\ z \end{pmatrix} = -\begin{pmatrix} 0 & g\partial/\partial x \\ z\partial/\partial x + L^{-1}(\partial/\partial x)h & 0 \end{pmatrix}\begin{pmatrix} u \\ z \end{pmatrix} - u\frac{\partial}{\partial x}\begin{pmatrix} u \\ z \end{pmatrix}, \tag{2a}$$

where $z$ is the free surface elevation, $u$ is the horizontal velocity at the free surface, $h$ is assumed to be independent of $t$, $g$ is the acceleration due to gravity and $L$ is a linear spatial differential operator characterizing the particular form of the Boussinesq approximation. If $L$ equals the identity operator, then (2a) reduces to the one-dimensional shallow water equations. One of the forms of the operator $L$ proposed by Katopodes and Dingemans[1] for describing Boussinesq models reads

$$L := 1 - \frac{1}{2}h\frac{\partial^2}{\partial x^2}h + \frac{1}{6}h^2\frac{\partial^2}{\partial x^2}. \tag{2b}$$

This operator generalizes the operator used by Broer[2] for the constant depth case; i.e. if $h$ does not depend on $x$, then the Boussinesq model defined by (2) reduces to the model derived by Broer.

Since we are mainly interested in the low-frequency range of the solution space of (1), the operator $L$ defined by (2b) may be considered as a perturbation of the identity operator; i.e. the norm of the operator $1 - L$ is small on the space of low frequencies. This property will become important in designing numerical approximations to $L$.

Following the method-of-lines approach, we replace the spatial domain by a discrete set of grid points and approximate the continuous functions $u$, $v$, $z$ and $h$ on these grid points by grid functions $U$, $V$, $Z$ and $H$. Furthermore, the differential operator $\partial/\partial x$ is approximated by difference operators $D_x$ which are defined on the space of grid functions. In this paper we shall assume that the boundary conditions are given by periodicity conditions and that the spatial grid consists of uniformly spaced grid points $j\Delta x$.

Let $L^*$ denote a discretization of $L$; then we are led to a semidiscretization of (2) given by the system of ordinary differential equations (ODEs)

$$\frac{d}{dt}\begin{pmatrix} U \\ Z \end{pmatrix} = -\begin{pmatrix} 0 & gD_x \\ ZD_x + (L^*)^{-1}D_xH & 0 \end{pmatrix}\begin{pmatrix} U \\ Z \end{pmatrix} - UD_x\begin{pmatrix} U \\ Z \end{pmatrix}. \tag{2*}$$

Since one usually wants high-order discretizations in Boussinesq models, the discretization stencils defining $D_x$ are rather large. In fact, in this paper it is assumed that a fourth-order discretization both in time and space is desired. As a consequence, the blocks in the Jacobian matrix associated with the semidiscretization (2*) contain a considerable number of non-zero diagonals. If implicit time integrators are used, then such Jacobian matrices imply a rather computationally intensive linear algebra problem for solving the implicit relations. In one spatial dimension this linear algebra problem does not prevent us from using implicit integrators; however, in two spatial dimensions explicit time integrators seem to be more attractive. Since it is

our aim to extend the results of this study to the two-dimensional case, we shall use *explicit* time integration methods.

In the actual integration of the system (2*) by explicit ODE solvers, one or more right-hand-side evaluations of (2*) are needed in each integration step. Hence in each step the quantity

$$B := (L^*)^{-1} D_x H U$$

is to be computed, so that in each step we have to solve the equation

$$L^* B = D_x H U. \tag{3}$$

The most obvious definition of $L^*$ is the difference operator

$$L_0^* := I - \tfrac{1}{2} H D H + \tfrac{1}{6} H^2 D, \tag{4}$$

where $D$ denotes a discretization of the operator $\partial^2 / \partial x^2$. (However, we will see that the corresponding semidiscretization (2*) becomes easily unstable for negative values of the elevation, so that alternative discretizations are desirable; see Section 5.) Putting aside the particular discretization we use for $L$, we will always be faced with the problem of solving equation (3), in spite of our restriction to explicit time integrators. As already observed, in one spatial dimension this is not a severe problem. However, in two spatial dimensions it requires special attention.

In the remainder of this paper the following aspects will be discussed. Section 2 deals with the stability of the continuous problem (2) and of the semidiscretization (2*). In Section 3 the stability condition associated with explicit integration methods is derived. Sections 4 and 5 respectively treat the difference operators $D_x$ and the discretization of the operator $L$. Finally, in Section 6 we present numerical results for the case of the standard fourth-order Runge–Kutta time integrator.

## 2. STABILITY

Before selecting an ODE solver for integrating the system (2*), we investigate the stability properties of both the continuous problem (2) and the semidiscretization (2*). This will be done in the case where the depth function $h$ is constant and with respect to the function space spanned by complex exponentials.

### 2.1. Stability of the continuous problem

We shall investigate the *local* stability of problem (2) by substituting continuous harmonic data

$$\begin{pmatrix} u \\ z \end{pmatrix} = \mathbf{a}(t) \exp(\mathrm{i} k x) \tag{5}$$

into (2) at some fixed point $x$ in the domain of definition (this is often called the method of 'frozen coefficients'). Here $k$ is the real-valued spatial frequency and $\mathbf{a}$ does not depend on $x$. We readily find that (5) satisfies (2) if $\mathbf{a}(t)$ is a solution of the ODE

$$\frac{\mathrm{d}}{\mathrm{d}t} \mathbf{a} = -A \mathbf{a} - \alpha \mathbf{a}, \qquad A := \mathrm{i} \begin{pmatrix} 0 & gk \\ [z + h/\lambda(L)]k & 0 \end{pmatrix}, \qquad \alpha := \mathrm{i} u k, \tag{6}$$

where $u$ and $h$ are defined at the point $x$ and where for any linear operator $M$, $\lambda(M)$ denotes an eigenvalue of $M$. The eigenvalues of $L$ can be expressed in terms of $k$, i.e.

$$\lambda(L) = 1 - \frac{1}{3} \lambda \left( \frac{\partial^2}{\partial x^2} \right) h^2 = 1 + \tfrac{1}{3} k^2 h^2.$$

The condition for (local) stability of the system (6) requires that the eigenvalues of the matrix $-(A + \alpha I)$ are located in the non-positive half-plane. The non-trivial eigenvalues of this matrix are given by

$$\lambda(-A - \alpha I) = -i \left[ uk \pm \sqrt{\left( gk^2 \frac{h + z\lambda(L)}{\lambda(L)} \right)} \right]. \tag{7}$$

The values of $1/\lambda(-A - \alpha I)$ are called *time constants* and depend on the frequency $k$. Thus we have stability if the time constants are located in the non-positive half-plane, i.e. $\mathrm{Re}\lambda(-A - \alpha I) \leqslant 0$. We shall say that the problem (2) is *dissipative* if $\mathrm{Re}\lambda(-A - \alpha I) < 0$ and *zero-dissipative* or *conservative* if $\mathrm{Re}\lambda(-A - \alpha I)$ vanishes. The following theorem is now immediate.

*Theorem 1*

If and only if $h \geqslant -z\lambda(L)$, then (6) is stable and at the same time zero-dissipative.     □

From this theorem it follows that for negative $z$ we only have stability with respect to spatial frequencies satisfying the inequality $-zk^2 \leqslant 3(h + z)/h^2$. For positive $z$ this condition is always satisfied, but for negative $z$ it prescribes an upper bound for the spatial frequencies. Recalling that the range of relevant frequencies is bounded by $1/h$, we conclude that the relevant frequencies always satisfy the above stability condition.

It may be of interest to express the stability condition of Theorem 1 in terms of the characteristic parameters $\varepsilon$ and $\mu$ introduced in Section 1. Introducing the wave amplitude $a := |z|_{\max}$, we find

$$k^2 \leqslant 3 \frac{h - a}{ah^2} = \frac{3(1 - \varepsilon)}{\varepsilon h^2}, \qquad \varepsilon := \frac{a}{h},$$

and substitution of

$$k^2 = \frac{4\mu\pi^2}{h^2}$$

yields

$$\frac{8\mu\pi^2}{h^2} \leqslant \frac{3(1 - \varepsilon)}{\varepsilon h^2},$$

so that the Boussinesq model (2) is stable if the solution space is restricted to frequencies for which $\varepsilon\mu$ is bounded by $3(1 - \varepsilon)/8\pi^2 \approx 3/80$.

## 2.2. Stability of the semidiscrete problem

Following the above approach, we insert into the semidiscretization (2*), at a fixed grid point $x$, the harmonic data

$$\begin{pmatrix} U \\ Z \end{pmatrix} = \mathbf{a}(t)\exp(\mathrm{i}kx_j), \quad 0 \leqslant k \leqslant \frac{\pi}{\Delta x}. \tag{5*}$$

Here $x_j$ runs through the grid points and $\mathbf{a}$ again only depends on $t$. The frequency $k$ is restricted to the interval $[0, \pi]/\Delta x$ because the grid cannot 'resolve' higher frequencies. In practice, accurate solutions can only be obtained for frequencies contained in an interval which is an order

of magnitude smaller, say $[0, 0.2]/\Delta x$. As stated in Section 1, the relevant frequencies for the Boussinesq model range from zero to $1/H$, so that $1/H$ should be less than $0.2/\Delta x$, i.e. $\Delta x \leqslant H/5$.

In the following analysis it is assumed that the grid functions $\exp(ikx_j)$ are eigenfunctions of the operators $D_x$ and $L^*$ with eigenvalues $\delta_x$ and $\lambda(L^*)$ respectively. Then the analogue of (6) becomes

$$\frac{d}{dt}\mathbf{a} = -A^*\mathbf{a} - \alpha^*\mathbf{a}, \qquad A^* := \begin{pmatrix} 0 & g\delta_x \\ [Z + H/\lambda(L^*)]\delta_x & 0 \end{pmatrix}, \qquad \alpha^* := U\delta_x. \quad (6^*)$$

The condition for stability of the system (6*) requires that $\lambda(-A^* - \alpha^* I)$ lies in the non-positive half-plane. By 'freezing' the coefficients in (6*), we find that the non-trivial eigenvalues of $-(A^* + \alpha^* I)$ are given by

$$\lambda(-A^* - \alpha^* I) = -U\delta_x \pm \sqrt{\left( g\delta_x^2 \frac{H + Z\lambda(L^*)}{\lambda(L^*)} \right)}. \quad (7^*)$$

As for the continuous problem, we say that the semidiscretization (2*) is *dissipative* if $\mathrm{Re}\lambda(-A^* - \alpha^* I) < 0$ and *zero-dissipative* or *conservative* if $\mathrm{Re}\lambda(-A^* - \alpha^* I)$ vanishes. Let $\rho(L^*)$ denote the spectral radius of $L^*$; then the analogue of Theorem 1 becomes as follows.

*Theorem 1\**

If the eigenvalues of the discretization $D_x$ are purely imaginary and if the eigenvalues $\lambda(L^*)$ are positive and satisfy the condition $H \geqslant -Z\rho(L^*)$, then and only then is the semidiscretization (6*) stable and at the same time zero-dissipative. □

The condition $H \geqslant -Z\rho(L^*)$ shows that in the case of negative elevation waves the quantity $1/\rho(L^*)$ may be interpreted as an upper bound for $|Z/H|$. Since $|Z/H|$ is bounded by the characteristic parameter $\varepsilon$, we conclude that $1/\rho(L^*)$ should not be less than $\varepsilon$. In order to see its implications, we consider the case where $L^*$ is defined by (4). Then the eigenvalues of $L^*$ are

$$\lambda(L_0^*) = 1 - \tfrac{1}{3}\lambda(D)H^2.$$

Assuming that $D$ has negative eigenvalues, we find

$$\rho(L_0^*) = 1 + \tfrac{1}{3}\rho(D)H^2. \quad (8)$$

Since $\rho(D)$ is usually extremely large, we see that in the case of large negative elevation waves the magnitude of $1/\rho(L^*)$ is easily less than $\varepsilon$. (We recall that the order of magnitude of $\varepsilon$ and $\mu$ is at most $\tfrac{1}{50}$, so that $\rho(L^*)$ should not exceed 50.) In Section 5 we return to the problem of discretizing the operator $L$ by better-conditioned difference operators than the operator $L_0^*$ defined in (4).

*2.3. Artificial stabilizing terms*

In Section 2.2 we have seen that the system (6*) is stable if and only if the eigenvalues of $D_x$ and $L^*$ satisfy the conditions of Theorem 1* and that the corresponding time constants lie on the imaginary axis which separates the regions of stability and instability. This 'marginal' stability property of the semidiscretization causes a numerical integration process to become easily unstable. Therefore it may be necessary to stabilize the system (2*) by adding artificial stabilizing terms.

*2.3.1. Artificial diffusion.* The most simple way to achieve additional stabilization is the introduction of an artificial diffusion term into (2*):

$$\frac{\mathrm{d}}{\mathrm{d}t}\begin{pmatrix} U \\ Z \end{pmatrix} = -\begin{pmatrix} 0 & gD_x \\ ZD_x + (L^*)^{-1}D_x H & 0 \end{pmatrix}\begin{pmatrix} U \\ Z \end{pmatrix} - [UD_x - d(\Delta x)^p D]\begin{pmatrix} U \\ Z \end{pmatrix},$$

where $d$ and $p$ are positive and $D$ denotes the discretization of the operator $\partial^2/\partial x^2$. As a result of this term the system is changed by a $p$th-order perturbation. The time constants are now given by

$$\lambda(-A^* - \alpha^* I) + \Delta\lambda, \qquad \Delta\lambda := d(\Delta x)^p \lambda(D). \tag{7'}$$

Assuming that the conditions of Theorem 1* are satisfied, these values are located on a curve in the left half-plane and no longer on the imaginary axis. Since $\Delta\lambda$ is negative, the semidiscretization has become dissipative.

*2.3.2. Fischer-type semidiscretization.* An alternative way to introduce dissipation is the following: let **S** denote the state vector $(U, Z)^{\mathrm{T}}$ and write (2*) in the compact form

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbf{S} = -Q\mathbf{S}.$$

Furthermore, let the matrix operator $Q$ be split according to $Q = T + (Q - T)$, where $T$ is the strictly lower triangular part of $Q$, and define the operator $P := I + q(\Delta x)^p T$, where $q$ and $p$ are positive. Instead of the semidiscretization (2*) we now consider the *preconditioned* semidiscretization

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbf{S} = -P^{-1}Q\mathbf{S} = -[I + qT(\Delta x)^p]^{-1}Q\mathbf{S}.$$

Since $P^{-1}$ is triangular, the evaluation of the right-hand-side function $-P^{-1}Q\mathbf{S}$ does not require more computational effort than that of $-Q\mathbf{S}$. Evidently, this preconditioned system is an $O[(\Delta x)^p]$-perturbation of the original system.

The method used by Fischer[3] for solving the shallow water equations can be interpreted as the explicit Euler method applied to the above preconditioned semidiscretization with $p = q = 1$ and $\Delta x = \Delta t$. Fischer showed that the resulting method is conditionally stable whereas application of explicit Euler to the original semidiscretization would lead to an unconditionally unstable method. The reason is of course that explicit Euler possesses an empty imaginary stability interval. As we shall see below, the preconditioning trick forces the time constants of the semidiscretization into the left half-plane, where explicit Euler has a non-empty stability region. Instead of using explicit Euler, one may use any ODE solver for integrating the preconditioned semidiscretization. We shall call this particular type of semidiscretization *Fischer-type semidiscretization*.

It is easily verified that the time constants $\lambda(-P^{-1}Q)$ associated with the Fischer-type semidiscretization are the roots of the equation

$$(\lambda + \alpha^*)^2 + \gamma q(\Delta x)^p \lambda + \gamma = 0, \qquad \alpha^* := U\delta_x, \qquad \gamma := -g\delta_x^2 \frac{H + Z\lambda(L^*)}{\lambda(L^*)}.$$

Writing $\lambda(-P^{-1}Q) = \lambda(-A^* - \alpha^* I) + \Delta\lambda$, we obtain

$$\lambda(-P^{-1}Q) = \lambda + \Delta\lambda, \qquad \Delta\lambda := \frac{\gamma q(\Delta x)^p \lambda(A^* + \alpha^* I)}{2[-\lambda(A^*) + \gamma q(\Delta x)^p]}. \tag{7''}$$

It can be shown that $\text{Re}(\Delta\lambda)$ and hence $\text{Re}\lambda(-P^{-1}Q)$ is negative, so that the Fischer-type semidiscretization is dissipative. The advantage of the above preconditioning over adding artificial diffusion lies in its lower computational costs.

## 3. STABILITY OF EXPLICIT TIME INTEGRATORS

It has already been observed that we shall use *explicit* time integration methods for the integration of the semidiscrete system (2*) in order to avoid the rather computationally intensive linear algebra involved in integrating (2*) by implicit methods. We also recall that we cannot completely avoid the solution of implicit equations because we always have to solve the system (3) defining the quantity $B$.

### 3.1. Stability condition of the zero-dissipative semidiscretization

If the conditions of Theorem 1* are satisfied, then the time constants of the system of ODEs (2*) are purely imaginary and (2*) is therefore stable. Hence the integration process used for integrating (2*) is (linearly) stable if its time step $\Delta t$ satisfies the stability condition

$$\Delta t \leqslant \frac{\beta}{\rho(A^* + \alpha^* I)},$$

where $\beta$ is the imaginary stability boundary of the time integrator used. This leads us to the following theorem.

### Theorem 2

If the time integrator chosen for integrating (2*) has imaginary stability boundary $\beta$ and if the discretization $D_x$ has imaginary eigenvalues, then a sufficient condition for linear stability is

$$\Delta t \leqslant \frac{\beta}{\rho(D_x)|U| + \sqrt{(g\rho(D_x^2[(L^*)^{-1}H + Z]))}}. \qquad \Box$$

We remark that this condition on $\Delta t$ reduces to the familiar stability condition for shallow water models if $L^* = I$. In order to get some insight into the actual step size limitations of this condition, we consider the case where $L^*$ has eigenvalues greater than or equal to unity, so that $Z \leqslant \lambda((L^*)^{-1}(H + ZL^*)) \leqslant H + Z$. Hence

$$\rho(D_x^2(L^*)^{-1}(H + ZL^*)) \leqslant \rho(D_x^2)\rho((L^*)^{-1}(H + ZL^*)) = \rho(D_x^2)\,\text{Max}\,\{H + Z, |Z|\}.$$

Thus we have the following corollary of Theorem 2.

### Corollary

If the eigenvalues of $L^*$ are greater than or equal to unity, then a sufficient condition for linear stability is

$$\Delta t \leqslant \frac{\beta}{\rho(D_x)[|U| + \sqrt{(g\,\text{Max}\{H + Z, |Z|\})}]}. \qquad \Box$$

Notice that the operator $L^*$ does not appear in this condition. If the eigenvalues of $L^*$ are less than unity, then $\rho((L^*)^{-1}(H + ZL^*)) > H + Z$, leading to a smaller maximum time step. Furthermore, we should bear in mind that the condition of this corollary may be more restrictive than that of Theorem 2.

Let us present the time step condition of the corollary in the form

$$\Delta t \leqslant \frac{\beta \Delta x}{C_x[|U| + \sqrt{(g \operatorname{Max}\{H + Z, |Z|\})}]}, \qquad C_x := \Delta x \rho(D_x), \qquad (9)$$

where $C_x$ is a constant depending on the particular discretization formula used. By a judicious choice of the discretization $D_x$ the constant $C_x$ can be minimized, thereby relaxing the stability condition. In a typical case we have

$$|U| = 1 \text{ m s}^{-1}, \qquad Z = 4 \cdot 4 \text{ m}, \qquad H = 10 \text{ m}, \qquad g = 10 \text{ m s}^{-2}, \qquad \Delta x = 2 \text{ m},$$

so that the stability condition becomes

$$\Delta t \leqslant \frac{2\beta}{13 C_x}.$$

### 3.2. Stability condition for the dissipative semidiscretization

In the case where artificial diffusion is added to the dissipative semidiscretization (2*) the stability condition (9) is changed to (see Section 2.3.1)

$$\Delta t \leqslant \frac{\beta \Delta x}{\sqrt{([d C_D (\Delta x)^{p-1}]^2 + \{C_x[|U| + \sqrt{(g \operatorname{Max}\{H + Z, |Z|\})}]\}^2)}}, \qquad C_D := \rho(D)(\Delta x)^2, \qquad (9')$$

where $\beta$ is the radius of the half-circle that can be inscribed in the stability region of the ODE solver used. On substitution of the numerical values given above we obtain

$$\Delta t \leqslant \frac{\beta \Delta x}{\sqrt{\{[d C_D (\Delta x)^{p-1}]^2 + (13 C_x)^2\}}},$$

showing that for $p = 4$ and $\Delta x \leqslant 2$ the value $d = \frac{1}{4}$ (say) seems to be suitable in the sense that the denominator in (9') is only slightly larger than that of (9) (here we assume that the constants $C_x$ and $C_D$ do not differ much in magnitude). We remark that the introduction of artificial diffusion does not relax the stability condition, but it improves the stability behaviour of the integration process because of dissipation of the higher frequencies.

## 4. THE OPERATOR $D_x$

We shall use discretizations of the symmetric form

$$D_x := \frac{1}{2\Delta x} \sum_{j=0}^{m} d_j (E_x^j - E_x^{-j}). \qquad (10a)$$

Here $d_j$ are scalar weights and $E_x$ is the shift operator in the $x$-direction, i.e. for any function $g(x)$ we define

$$E_x g(x, y) := g(x + \Delta x). \qquad (10b)$$

It is sometimes convenient to present the difference operator $D_x$ by so-called stencils (or molecules). For example, for $m = 3$ such a stencil is given by

$$D_x = \frac{1}{\Delta x} |-d_3 \quad -d_2 \quad -d_1 \quad 0 \quad d_1 \quad d_2 \quad d_3|.$$

We shall call $m$ the dimension of the discretization stencil. For the general difference operator of the type (10) we have the following theorem.

*Theorem 3*

The following assertions hold.

(a) The eigenvalues of $D_x$ are purely imaginary.
(b) The spectral radius of $D_x$ is given by $\rho(D_x) = C_x/\Delta x$, where

$$C_x \leqslant 2\|c(p)\|, \qquad c(p) := \sum_{j=0}^{m} d_j \sin(jp), \tag{11}$$

with $\|\cdot\|$ denoting the maximal norm with respect to all values of $p$ and $q$.

(c) The discretization (10) is fourth-order-accurate if

$$\sum_{j=0}^{m} 2j d_j = 1, \qquad \sum_{j=0}^{m} j^3 d_j = 0. \tag{12}$$

*Proof.* (a) Since

$$E_x \exp(ikx_j) = e^{ip} \exp(ikx_j), \qquad p := k_x \Delta x,$$

we find that

$$D_x \exp(ikx_j) = \delta_x \exp(ikx_j), \qquad \delta_x = \delta_x(p) = \frac{2i}{\Delta x} \sum_{j=0}^{m} d_j \sin(jp), \tag{13}$$

showing that $\exp(ikx_j)$ is an eigenfunction of $D_x$ with purely imaginary eigenvalues $\delta_x$.

(b) This estimate is immediate from the expression for $\delta_x$.

(c) Let $g(x)$ be a sufficiently differentiable function; then we can write

$$D_x g(x) = \frac{\partial}{\partial x} X\left(\Delta x \frac{\partial}{\partial x}\right) g(x), \qquad X(x, y) := 2 \sum_{j=0}^{m} d_j \frac{\sinh(jx)}{x}.$$

It is straightforwardly verified that

$$X(x) = \sum_{j=0}^{m} (2j + \tfrac{1}{3}j^3 x^2) d_j + O(x^4),$$

from which the theorem is immediate. □

By means of this theorem, fourth-order difference operators $D_x$ can straightforwardly be constructed. However, because the actual implementation of these operators will be based on staggered grids (i.e. the components of $U$ and $Z$ will be computed at distinct grid points), we shall distinguish two special cases:

(i) $m = 2$, no restrictions on the weights $d_j$
(ii) $m = 3$, $d_j = 0$ if $j$ is odd.

In the first case we deduce from Theorem 2 that fourth-order accuracy is possible for $m = 2$. This leads to the conventional four-point 'line' discretization

$$D_x := \frac{1}{12\Delta x} |1 \quad -8 \quad 0 \quad 8 \quad -1|, \quad \text{with } C_x = 1\cdot37149 \ldots. \tag{14}$$

On staggered grids we need discretizations of $D_x$ with $m = 3$ and $d_j = 0$ if $j$ is odd. This leads to the conventional four-point 'line' discretization

$$D_x := \frac{1}{48\Delta x} |1 \quad 0 \quad -27 \quad 0 \quad 27 \quad 0 \quad -1|, \quad \text{with } C_x = \tfrac{7}{6}. \tag{14'}$$

## 5. THE OPERATOR $L^*$

In Section 2.2 (Theorem 1*) it was shown that the discretization $L^*$ of the operator $L$ should satisfy the condition $H \geqslant -Z\rho(L^*)$ in order to achieve stability, i.e. in the case of negative elevation waves $1/\rho(L^*)$ should not be less than the parameter $\varepsilon$ characterizing the Boussinesq model. We recall that the 'natural' discretization $L_0^*$ defined by (4) may lead to severe restrictions on $|Z/H|$. In this section we therefore consider alternative discretizations which are in fact approximations to $L_0^*$ with reduced spectral radius. We shall discuss 'low-frequency' approximations to $L_0^*$ and preconditioning (or smoothing) of the operator $L_0^*$. In both cases the spectral radius of the resulting operator $L^*$ is reduced considerably while the whole spectrum is bounded below by unity. Of course, the defect $L^* - L_0^*$ should be small for accuracy reasons. In order to measure this defect, we consider the quantity

$$\delta(k) := (L^* - L_0^*)\exp(ikx_j) = \lambda(L^*) - \lambda(L_0^*) \tag{15}$$

as a function of $k$. Since we are only interested in the lower frequencies, i.e. $k$ in the interval $[0, \pi\sqrt{2/5}H]$ (see (1)), it is justified to restrict our considerations to this range of low frequencies. Let $\| \cdot \|$ denote the maximum norm with respect to all frequencies less than $k_0$; then we define

$$\Delta(k_0) := \| \delta(k) \| \tag{16}$$

as a measure for the-low frequency defect.

### 5.1. Low-frequency approximations to $L_0^*$

Consider the operator

$$L^* = L^*(\omega, \theta) := (I + \omega HDH - \theta H^2 D)^{-1}[I - \tfrac{1}{2}(1 - 2\omega)HDH + \tfrac{1}{6}(1 - 6\theta)H^2 D], \tag{17}$$

where $\omega$ and $\theta$ are free parameters (notice that $L^*(\omega, \theta) = L_0^*$ for $\omega = \theta = 0$).

First we derive an expression for the spectral radius of $L^*$. It is readily verified that in the constant coefficient case the eigenvalues of $L^*$ corresponding to the eigenfunctions $\exp(ikx)$ are given by

$$\lambda(L^*) = \lambda(L^*(\omega, \theta)) = \frac{1 - (\theta - \omega + \tfrac{1}{3})\lambda(D)H^2}{1 - (\theta - \omega)\lambda(D)H^2},$$

so that

$$\delta(k) = \lambda(L^*) - \lambda(L_0^*) = \frac{-(\theta - \omega)[\lambda(D)H^2]^2}{3[1 - (\theta - \omega)\lambda(D)H^2]},$$

where $\lambda(D)$ denotes the eigenvalues of the discretization of the operator $\partial^2/\partial x^2$. Using staggered grids, we define the operator $D$ by the fourth-order-accurate difference formula

$$D = \frac{1}{48(\Delta x)^2} |-1 \quad 0 \quad 16 \quad 0 \quad -30 \quad 0 \quad 16 \quad 0 \quad -1|. \tag{18}$$

We find that the eigenvalues of $D$ are given by

$$\lambda(D) = (2\Delta x)^{-2} d(\xi), \qquad d(\xi):= -\xi - \tfrac{1}{12}\xi^2, \qquad \xi:= 4\sin^2(k_x \Delta x), \qquad (19a)$$

so that $0 \leqslant \xi \leqslant 4$ and

$$\rho(D) = \frac{4}{3(\Delta x)^2}. \qquad (19b)$$

Because $D$ has negative eigenvalues, the expression for $\lambda(L^*)$ shows that $\lambda(L^*(\omega, \theta))$ is bounded below by unity if $\theta - \omega \geqslant 0$. For $\theta - \omega \geqslant 0$ we find

$$\rho(L^*) = \frac{1 + (\theta - \omega + \tfrac{1}{3})\rho(D)H^2}{1 + (\theta - \omega)\rho(D)H^2} = \frac{1 + 16(\theta - \omega + \tfrac{1}{3})Q^2/3}{1 + 16(\theta - \omega)Q^2/3}, \qquad Q:= \frac{H}{2\Delta x}. \qquad (20)$$

Next we compute the defect

$$\Delta(k_0) = \|\delta(k)\|, \qquad \delta(k) = \frac{(\theta - \omega)[\lambda(D)H^2]^2}{3[1 - (\theta - \omega)\lambda(D)H^2]} = \frac{(\theta - \omega)[d(\xi)Q^2]^2}{3[1 - (\theta - \omega)d(\xi)Q^2]}.$$

Let $k$ be less than $k_0$; then the variable $\xi$ is bounded by $2[1 - \cos(k_0 \Delta x)]$. Hence, by taking the maximum norm with respect to this range of $\xi$-values, we obtain

$$\Delta(k_0) = \left| \frac{(\theta - \omega)[d(\xi_0)Q^2]^2}{3[1 - (\theta - \omega)d(\xi_0)Q^2]} \right|, \qquad \xi_0:= 2[1 - \cos(2k_0 \Delta x)]. \qquad (21)$$

Hence

$$\Delta(k_0) \approx \left| \frac{4(\theta - \omega)(Hk_0)^4}{3[1 + 2(\theta - \omega)(Hk_0)^2]} \right| \quad \text{as } k_0 \Delta x \to 0, \qquad (21')$$

showing that for the relevant frequencies $k \leqslant k_0 \approx 1/H$ the defect is bounded by

$$4(\theta - \omega)/[3 + 6(\theta - \omega)].$$

Thirdly we consider the system (3) for computing the quantity $B$ in the case of (17):

$$[I - \tfrac{1}{2}(1 - 2\omega)HDH + \tfrac{1}{6}(1 - 6\theta)H^2 D]B = (I + \omega HDH - \theta H^2 D)D_x HU.$$

In general, solving this system has the same computational complexity as that of the system arising for $L^* = L_0^*$ ($\omega = \theta = 0$). The values $\omega = \tfrac{1}{2}$ and $\theta = \tfrac{1}{6}$ seem to be of interest because $B$ is then explicitly defined. However, since $\theta - \omega$ should be non-negative, this choice is excluded. Another attractive choice which does preserve stability is $\omega = 0$ and $\theta = \tfrac{1}{6}$, leading to

$$(I - \tfrac{1}{2}HDH)B = (I - \tfrac{1}{6}H^2 D)D_x HU.$$

Next we compute in the case $\theta - \omega = \tfrac{1}{6}$ the spectral radius $\rho(L^*)$ and the defect $\Delta(k_0)$ for a few values of $Q$ and $k_0$. Furthermore, as a reference, we also list the values of $\rho(L_0^*)$. Table I clearly demonstrates the considerable reduction of the spectral radius by using low-frequency approximations to the operator $L_0^*$. Since stability requires that $\rho(L^*)$ should be less than $1/\varepsilon$ (see Section 2.2), we conclude that in the range $3 \leqslant Q \leqslant 13$ the discretization $L^*(\theta, \theta - \tfrac{1}{6})$ allows waves with $\varepsilon$-values as large as 0·33, whereas the discretization $L_0^* = L^*(0, 0)$ allows waves with $\varepsilon$-values varying from 0·06 to 0·0033. However, Table I also shows that the defect in the range of relevant frequencies is rather large and cannot be improved by decreasing $\Delta x$.

Table I. Spectral radius and defect of the operator $L^*$ defined by (17) for $\theta - \omega = \frac{1}{6}$

| $Q = H/2\Delta x$: | 3 | 5 | 7 | 9 | 11 | 13 | $\infty$ |
|---|---|---|---|---|---|---|---|
| $\rho(L_0^*)$ | 17·0 | 45·4 | 88·1 | 145·0 | 216·1 | 301·4 | $\infty$ |
| $\rho(L^*)$ | 2·78 | 2·91 | 2·96 | 2·97 | 2·98 | 2·99 | 3·0 |
| $\Delta(\pi\sqrt{2}/5H)$ | 0·03 | 0·03 | 0·03 | 0·03 | 0·03 | 0·03 | 0·03 |
| $\Delta(1/H)$ | 0·05 | 0·05 | 0·05 | 0·05 | 0·05 | 0·05 | 0·05 |
| $\Delta(2/H)$ | 0·53 | 0·53 | 0·53 | 0·53 | 0·53 | 0·53 | 0·53 |

Finally we derive the time step condition according to Theorem 2 in the case $\theta - \omega = \frac{1}{6}$ with $D_x$ defined by (14′). We deduce from (13) that

$$\lambda(D_x^2) = (\delta_x)^2 = -(2\Delta x)^2 \xi(1 + \tfrac{1}{24}\xi)^2,$$

where $\xi$ is defined as before. A comparison with $\lambda(D)$ defined in (19a) reveals that

$$\frac{\lambda(D_x^2) - \lambda(D)}{\lambda(D)} = \frac{\xi^2}{48(12 + \xi)},$$

showing that we may replace $\lambda(D_x^2)$ by $\lambda(D)$ without introducing large errors. Hence

$$\lambda(D_x^2(L^*)^{-1}(H + ZL^*)) \approx \lambda(D(L^*)^{-1}(H + ZL^*)) = \lambda(D)\left(\frac{1 - \lambda(D)H^2/6}{1 - \lambda(D)H^2/2}H + Z\right).$$

It can be shown that this expression is monotone in $\lambda(D)$, so that it follows from (19b) that

$$\rho(D_x^2(L^*)^{-1}(H + ZL^*)) \approx \frac{4}{3(\Delta x)^2}\left(\frac{2H^2 + 9(\Delta x)^2}{6H^2 + 9(\Delta x)^2}H + Z\right) \approx \frac{4}{9(\Delta x)^2}(H + 3Z).$$

Upon substitution into the condition of Theorem 2 we obtain

$$\Delta t \leqslant \frac{6\beta\Delta x}{7|U| + 4\sqrt{[g(H + 3Z)]}}. \tag{22}$$

We remark that the Corollary of Theorem 2 would result in the more restrictive condition

$$\Delta t \leqslant \frac{6\beta\Delta x}{7|U| + 7\sqrt{(g\,\mathrm{Max}\{H + Z, |Z|\})}}. \tag{22′}$$

### 5.2. Preconditioning of $L_0^*$

The preconditioned discretizations of this subsection possess a smaller defect than those of the preceding subsection but at the cost of larger $\rho(L^*)$-values.

Consider the discretization

$$L^* = SL_0^*, \qquad S := (I + qQ^2 D_S)^{-1}, \qquad Q := \frac{H}{2\Delta x}, \tag{23}$$

where $q$ is a free parameter and $D_S$ is a difference operator. The system (3) that has to be solved in each call of the right-hand side of the semidiscretization (2*) now assumes the form

$$L_0^* B = (1 + qQ^2 D_S)D_x HU, \tag{24}$$

showing that the computational complexity is hardly increased by introducing the preconditioner $S$.

Let the operator $D$ be defined by the fourth-order line molecule given by (18) and let $D_S$ be defined by

$$D_S = |a_2 \ 0 \ a_1 \ 0 \ 1 \ 0 \ a_1 \ 0 \ a_2 \ |. \tag{25}$$

The eigenvalues of $D$ (with respect to the eigenfunctions $\exp(ikx)$) are given by

$$\lambda(D) = (\Delta x)^{-2} d(\xi), \qquad d(\xi) := -(\xi + \tfrac{1}{12}\xi^2), \qquad \xi := 2[1 - \cos(2k\Delta x)]. \tag{26}$$

Hence the eigenvalues of $L^*$ can be expressed as

$$\lambda(L^*) = \frac{1 - \tfrac{1}{3}d(\xi)Q^2}{1 + q\lambda(D_S)Q^2} = \frac{1 + \tfrac{1}{3}(\xi + \tfrac{1}{12}\xi^2)Q^2}{1 + q\lambda(D_S)Q^2}, \quad 0 \leqslant \xi \leqslant 4, \tag{27}$$

and the defect function becomes

$$\delta(k) = \frac{-q\lambda(D_S)Q^2}{1 + q\lambda(D_S)Q^2} [1 + \tfrac{1}{3}(\xi + \tfrac{1}{12}\xi^2)Q^2].$$

where

$$\lambda(D_S) = (2a_1 + 2a_2 + 1) + (-a_1 - 4a_2)\xi + a_2\xi^2.$$

Suppose that we choose the parameters in (25) such that the first two terms in $\lambda(D_S)$ vanish, i.e. $a_1 = -\tfrac{2}{3}$ and $a_2 = \tfrac{1}{6}$; then

$$D_S = \tfrac{1}{6}|1 \ 0 \ -4 \ 0 \ 6 \ 0 \ -4 \ 0 \ 1|,$$

$$\lambda(L^*) = \frac{36 + (12\xi + \xi^2)Q^2}{36 + 6q\xi^2Q^2}, \qquad \delta(k) = \frac{-q\xi^2Q^2}{6 + q\xi^2Q^2} [1 + \tfrac{1}{3}(\xi + \tfrac{1}{12}\xi^2)Q^2], \qquad 0 \leqslant \xi \leqslant 4.$$

From these expressions it can be derived that $\lambda(L^*)$ is never less than unity if $q$ assumes values in the range $[\tfrac{1}{6}, \tfrac{2}{3}]$ and that in this range of $q$-values the magnitude of $\rho(L^*)$ is minimized for $q = \tfrac{2}{3}$. Introducing the new variable $x = \xi Q$, we may write

$$\lambda(L^*) = \frac{36 + 12xQ + x^2}{36 + 4x^2}, \qquad \delta(k) = \frac{-2x^2}{18 + 2x^2}(1 + \tfrac{1}{3}xQ + \tfrac{1}{36}x^2), \qquad 0 \leqslant x \leqslant 4Q. \tag{28}$$

For larger values of $Q$ the spectrum function $\lambda(L^*)$ behaves as $3xQ/(9 + x^2)$, which assumes its maximum at the point $x = 3$, so that $\rho(L^*) \approx Q/2$. From this result the analogue of Table I becomes Table II.

Table II. Spectral radius and defect of the operator $L^*$ defined by (23) and (25) with $a_1 = -\tfrac{2}{3}$, $a_2 = \tfrac{1}{6}$ and $q = \tfrac{2}{3}$

| $Q = H/2x$: | 3 | 5 | 7 | 9 | 11 | 13 | $\infty$ |
|---|---|---|---|---|---|---|---|
| $\rho(L_0^*)$ | 17·0 | 45·4 | 88·1 | 145·0 | 216·1 | 301·4 | $\infty$ |
| $\rho(L^*)$ | 2·17 | 3·15 | 4·14 | 5·14 | 6·06 | 7·11 | $\infty$ |
| $\Delta(\pi\sqrt{2}/5H)$ | 0·01 | 0·003 | 0·002 | 0·001 | 0·001 | 0·001 | 0 |
| $\Delta(1/H)$ | 0·16 | 0·006 | 0·003 | 0·002 | 0·001 | 0·001 | 0 |
| $\Delta(2/H)$ | 0·36 | 0·15 | 0·08 | 0·05 | 0·034 | 0·024 | 0 |

Finally we derive the stability condition according to Theorem 2. From (28) we obtain in terms of $x$

$$\lambda(D_x^2[(L^*)^{-1}H+Z]) = -(\Delta x)^{-2}\left(\frac{x}{Q}+\frac{x^2}{12Q^2}\right)\left(\frac{36+4x^2}{36+12xQ+x^2}H+Z\right),$$

which assumes its maximum value at $x = 4Q$. Hence $\rho[(L^*)^{-1}] \approx 1$ and Theorem 2 yields the same condition as stated in the Corollary, i.e. condition (22').

## 6. NUMERICAL EXPERIMENTS

In order to test the theory developed in this paper, we add to the right-hand side of equation (2a) some source function such that a prescribed function is identical to the exact solution. This enables us to determine the accuracy of the numerical solutions. Let us first rewrite (2a) in the form

$$\frac{\partial}{\partial t}\begin{pmatrix} u \\ Lz \end{pmatrix} = -\begin{pmatrix} 0 & g\partial/\partial x \\ Lz\partial/\partial x +(\partial/\partial x)h & 0 \end{pmatrix}\begin{pmatrix} u \\ z \end{pmatrix} - \begin{pmatrix} u(\partial/\partial x)u \\ Lu(\partial/\partial x)z \end{pmatrix}.$$

Then, by introducing the source function $s(x, t)$, we obtain

$$\frac{\partial}{\partial t}\begin{pmatrix} u \\ Lz \end{pmatrix} = -\begin{pmatrix} 0 & g\partial/\partial x \\ Lz\partial/\partial x +(\partial/\partial x)h & 0 \end{pmatrix}\begin{pmatrix} u \\ z \end{pmatrix} - \begin{pmatrix} u(\partial/\partial x)u \\ Lu(\partial/\partial x)z \end{pmatrix} + \begin{pmatrix} s_1(x, t) \\ s_2(x, t) \end{pmatrix}. \quad (29)$$

By prescribing the exact solution $u(x, t)$ and $z(x, t)$, we deduce from (29) that the corresponding source function $s(x, t)$ is defined by

$$s_1(x, t) = \frac{\partial}{\partial t}u + g\frac{\partial}{\partial x}z + u\frac{\partial}{\partial x}u, \qquad s_2(x, t) = L\left(\frac{\partial}{\partial t}z + \frac{\partial}{\partial x}uz\right) + \frac{\partial}{\partial x}hu. \quad (30)$$

In our numerical experiments we always prescribed the exact solution

$$u(x, t) = -\sin(cx^2)\sin(dt), \qquad z(x, t) = \cos(cx^2)\cos(dt), \qquad c = \frac{4\pi}{b^2}, \qquad d = \frac{2\pi}{T}, \quad (31)$$

where $[0, b]$ is the spatial domain and $[0, T]$ is the integration interval.

As a consequence of the introduction of the source function $s(x, t)$, we now have to solve in each integration step the equation (cf. (3))

$$L^*B = D_xHU - S_2, \quad (32)$$

where $S_2$ is the discretization of $s_2(x, t)$.

Equation (29) was discretized using the staggered grid difference approximation (14') on a uniform grid with mesh size $\Delta x$. The operator $L$ defined in (2b) was discretized according to formula (17) with $\omega = 0$ and $\theta = \frac{1}{6}$ and according to (23) and (25) with $a_1 = \frac{2}{3}$, $a_2 = \frac{1}{6}$ and $q = \frac{2}{3}$. The time integration was performed using the standard fourth-order Runge–Kutta method with constant step size $\Delta t$.

Since the imaginary stability boundary of the standard Runge–Kutta method is given by $2\sqrt{2}$, the stability conditions (22) and (22') take respectively the form

$$\Delta t \leqslant \frac{12\Delta x\sqrt{2}}{7|U|+4\sqrt{[g(H+3Z)]}}, \qquad \Delta t \leqslant \frac{12\Delta x\sqrt{2}}{7|U|+7\sqrt{[g(H+Z)]}}. \quad (33)$$

In view of Theorem 1*, there is an additional stability condition reading

$$H \geqslant -Z\rho(L^*), \tag{34}$$

where $\rho(L^*)$ is given in Tables I and II depending on the discretization used for $L$.

In the tables of results we have listed the accuracy obtained at the end point $t = T$. The accuracy is measured by the number of correct decimal digits, i.e. by writing the elevation error in the form

$$|Z(x_j, T) - z(x_j, T)| = 10^{-\Delta}. \tag{35}$$

### 6.1. Constant depth function

In our first test the following input data were used:

| | |
|---|---|
| domain of definition | $0 \leqslant x \leqslant b := 1000, \qquad 0 \leqslant t \leqslant T := 60;$ |
| initial values | $u(x, 0) = 0, \qquad z(x, 0) = \cos(cx^2);$ |
| boundary values | $u(0, t) = u(b, t) = 0;$ |
| coefficient functions | $g = 9.81, \qquad h(x) = 10;$ |
| source function | $s_1(x, t) = -\sin(cx^2)[(d + 2gcx)\cos(dt) - 2cx\cos(cx^2)\sin^2(dt)],$ |
| | $s_2(x, t) = -\sin(dt)\{(d + 2hcx)\cos(cx^2) + 2cx[\cos^2(cx^2)$ |
| | $\qquad - \sin^2(cx^2)]\cos(dt) - \tfrac{2}{3}cdh^2[\sin(cx^2) + 2cx^2\cos(cx^2)]$ |
| | $\qquad - \tfrac{16}{3}c^2xh^2[3\cos(cx^2)\sin(cx^2) + 2cx^2[\cos(cx^2)$ |
| | $\qquad - \sin(cx^2)]\cos(dt)\};$ |
| mesh size | $\Delta x = 1.$ |

For these data we found that the operator $L_0^*$ leads to instabilities irrespective the value of $\Delta t$. The reason is that the stability condition (34) is violated. However, when using the operators defined by (17) and (23), (25), this condition is always satisfied; hence the step size condition (33)

Table III. Values of the number of correct digits of $Z$ at $T = 60$ for discretization (17) with $\omega = 0$ and $\theta = \tfrac{1}{6}$

| $\Delta t$ | $x = b/8$ | $x = 2b/8$ | $x = 3b/8$ | $x = 4b/8$ | $x = 5b/8$ | $x = 6b/8$ | $x = 7b/8$ |
|---|---|---|---|---|---|---|---|
| 0·1 | 2·1 | 3·2 | 1·9 | 1·8 | 1·9 | 2·5 | 1·8 |
| 0·2 | 2·1 | 2·8 | 1·9 | 1·8 | 1·9 | 2·8 | 1·8 |
| 0·3 | 2·1 | 2·6 | 1·9 | 1·8 | 1·9 | 3·6 | 1·8 |
| 0·4 | 2·1 | 2·4 | 1·9 | 1·8 | 1·9 | 3·0 | 1·8 |
| 0·5 | * | * | * | * | * | * | * |

Table IV. Values of the number of correct digits of $Z$ at $T = 60$ for discretization (23) and (25) with $a_1 = \tfrac{2}{3}$, $a_2 = \tfrac{1}{6}$ and $q = \tfrac{2}{3}$

| $\Delta t$ | $x = b/8$ | $x = 2b/8$ | $x = 3b/8$ | $x = 4b/8$ | $x = 5b/8$ | $x = 6b/8$ | $x = 7b/8$ |
|---|---|---|---|---|---|---|---|
| 0·1 | 2·1 | 3·2 | 1·9 | 1·8 | 1·9 | 2·5 | 1·8 |
| 0·2 | 2·1 | 2·8 | 1·9 | 1·8 | 1·9 | 2·8 | 1·8 |
| 0·3 | * | * | * | * | * | * | * |

completely determines the stability. Since the maximal numerical values assumed by $|Z|$ and $|U|$ are approximately unity, we may expect the results to be stable if respectively $\Delta t \leqslant 0\cdot 32$ and $\Delta t \leqslant 0\cdot 21$. We obtained the results listed in the Tables III and IV ($*$ indicates instability).

## REFERENCES

1. N. D. Katopodes and M. W. Dingemans, 'Stable equations for surface wave propagation on an uneven bottom', *Progress Report Draft II*, Delft Hydraulics, 1988.
2. L. J. F. Broer, 'Approximate equations for long water waves', *Appl. Sci. Res.*, **31**, 377–395 (1975).
3. G. Fischer, 'A numerical method for tidal computations in inner seas' (in German), *Tellus*, **11**, 60–76 (1959).
4. L. J. F. Broer, 'On the hamiltonian theory of surface waves', *Appl. Sci. Res.*, **30**, 430–446 (1974).
5. L. J. F. Broer, E. W. C. van Groessen and J. M. W. Timmers, 'Stable model equations for long water waves', *Appl. Sci. Res.*, **32**, 619–639 (1976).
6. L. J. F. Broer and J. A. Kobussen, 'Canonical transformations and generating functionals', *Physica*, **61**, 275–288 (1972).
7. M. W. Dingemans, 'Water waves over an uneven bottom', *Report R729-II*, Delft Hydraulics, 1973.