



**HAL**  
open science

# On the numerical solution of sparse linear systems emerging in finite volume discretizations of 2D Boussinesq-type models

Maria Kazolea, M Gaitani, Argiris I. Delis

► **To cite this version:**

Maria Kazolea, M Gaitani, Argiris I. Delis. On the numerical solution of sparse linear systems emerging in finite volume discretizations of 2D Boussinesq-type models. [Research Report] RR-8778, INRIA Bordeaux, équipe CARDAMOM. 2015, pp.28. hal-01202983

**HAL Id: hal-01202983**

**<https://inria.hal.science/hal-01202983>**

Submitted on 22 Sep 2015

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



# On the numerical solution of sparse linear systems emerging in finite volume discretizations of 2D Boussinesq-type models

M. Kazolea, M. Gaitani, A.I. Delis

**RESEARCH  
REPORT**

**N° 8778**

September 2015

Project-Teams CARDAMOM





# On the numerical solution of sparse linear systems emerging in finite volume discretizations of 2D Boussinesq-type models

M. Kazolea\*, M. Gaitani†, A.I. Delis‡

Project-Teams CARDAMOM

Research Report n° 8778 — September 2015 — 25 pages

**Abstract:** This work supplements the realization and validation of a higher-order well balanced finite volume (FV) scheme developed for numerically simulating, on triangular meshes, weakly non-linear weakly dispersive water waves over varying bathymetries. The scheme has been recently presented by Kazolea et al. (*Coastal Eng.* 69:42-66, 2012 and *J. Comp. Phys.* 271:281-305, 2014). More precisely, we investigate and develop solution strategies for the sparse linear system that occurs during this FV discretisation of a set of Boussinesq-type equations on unstructured meshes. The resultant system of equations must be solved at each time step as to recover the actual velocity field of the flow. The system's coefficient matrix is sparse, un-symmetric and often ill-conditioned. Its characteristics are affected by physical quantities of the problem to be solved, such as the undisturbed water depth and the mesh topology. This work investigates the application of different iterative techniques, with and without the usage of preconditioners and reordering, for the solution of this sparse linear system. Two different iterative methods, three preconditioning techniques, including different ILU factorizations and two different reordering techniques are implemented and discussed. An optimal strategy, in terms of computational efficiency and robustness, is proposed.

**Key-words:** Boussinesq-type equations, Finite Volumes, Sparse matrices, preconditioning, Reordering

---

\* maria.kazolea@inria.fr

† School of Production Engineering & Management, Technical University of Crete, mgaitani@isc.tuc.gr

‡ School of Production Engineering & Management, Technical University of Crete, adelis@science.tuc.gr

**RESEARCH CENTRE  
BORDEAUX – SUD-OUEST**

351, Cours de la Libération  
Bâtiment A 29  
33405 Talence Cedex

## Sur la solution numerique des systems creuses et linéaires émergents de la discretization volume finis des modeles 2D de type Boussinesq

**Résumé :** Ce travail concerne la réalisation et la validation d'un schéma Volumes Finis d'ordre élevé pour la simulation des vagues en régime faiblement non-linéaire et faiblement dispersif sur bathymétries variables. Le schéma implémenté est celui proposé récemment par Kazolea et al. (*Coastal Eng* 69: 42-66, 2012 et *J. Phys Comp* 271:.. 281-305, 2014). Plus précisément, nous étudions et développons des stratégies de solution pour le système linéaire creux qui se produit au cours de la discrétisation des équations de Boussinesq sur maillages non structurés. Le système d'équations résultant doit être résolu à chaque pas de temps pour récupérer la vitesse. La matrice du système est creuse, non symétrique et souvent mal conditionnée. Ses caractéristiques sont affectées par des quantités physiques tels que la profondeur de l'eau au repos et la topologie du maillage. Ce travail étudie l'application de différentes techniques itératives, avec et sans l'utilisation de préconditionneurs et de re-numérotation, pour la solution de ce système linéaire creux. Deux méthodes itératives différentes, trois techniques de préconditionnement, y compris les différentes factorisations ILU et deux techniques de réordonnement différentes sont mises en oeuvre et évaluées. Une stratégie optimale, en termes d'efficacité de calcul et de robustesse, est proposée.

**Mots-clés :** équations de Boussinesq, Volume Finis, système linéaire creux, préconditionneurs, re-numérotation

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Governing equations</b>	<b>4</b>
<b>3</b>	<b>The numerical model</b>	<b>6</b>
3.1	The FV approximation . . . . .	6
3.2	The Resulting Linear System . . . . .	7
3.3	Matrix properties . . . . .	9
<b>4</b>	<b>Iterative methods, preconditioning and reordering</b>	<b>12</b>
4.1	Application of Iterative Methods . . . . .	12
4.2	Application of Preconditioning methods . . . . .	13
4.2.1	The ILU(0) preconditioner . . . . .	14
4.2.2	The ILU( $k$ ) preconditioner . . . . .	15
4.2.3	ILUT( $p, \tau$ ) preconditioner: . . . . .	16
<b>5</b>	<b>Reorderings</b>	<b>17</b>
5.1	Spatial accuracy and efficiency . . . . .	19
<b>6</b>	<b>Conclusions</b>	<b>20</b>

## 1 Introduction

Boussinesq-type (BT) equations have been widely used in the past few decades for the description of water wave propagation and transformation along the cost line. Starting from the classical Boussinesq equations [24], which are limited to relatively shallow water, and up to now-days, numerous researchers have contributed to the development of analytical theories and their numerical approximation as to simulate the various wave phenomena such as, wave shoaling, diffraction, refraction and breaking. As such, significant attempts have been made to extend the applicability of the Boussinesq-type models to deeper water, for example in [22, 23, 21]. These extended models give a more accurate representation of the waves's phase and group velocities in intermediate waters, with water depth to wavelength ratio up to  $1/2$ , and sometimes are referred to as low-order enhanced BT equations. Moreover, significant effort has been made in recent years into advancing the nonlinear and dispersive properties of BT models by including high-order nonlinear and dispersion terms, we refer for example to [17, 20, 7, 36]. An extensive review, which describes the state of the art in water wave modelling by means of BT models, can be found in [8].

For the numerical solution of BT equations different methods have been proposed such as, the Finite-Difference (FD) method, [23, 36, 15, 21] and the Finite-Element method, [13, 35, 25]. In the last ten years the use of the Finite-Volume (FV) method has become the most widely used due to its effectiveness in the approximation of hyperbolic conservation laws. Applications and advances along this line of research can be found, for example, in [12, 10, 30, 27, 33, 18, 19]. Usually in the numerical solution of BT equations the inversion of one or more matrices, i.e. solution(s) of linear system(s), is required as to recover the actual velocity field from the solution variables. This is an essential

procedure (especially in two-dimensional simulations) if the original dispersive BT equations are re-written in a conservative-like form and the vector of the unknown variables includes part of the dispersive terms, [8]. In the most simple cases e.g. when structured computational meshes are used the matrices can be of tridiagonal shape, but when using e.g. unstructured meshes the matrices that occur are more complicated, unsymmetric and have variable sparsity patterns [18, 14, 35]. For this reason further investigations for the efficient and accurate solution of the resulting linear systems in two-dimensional computations on unstructured meshes are in needed.

This work is complementary to [18, 19] where, for the first time, a high-order well-balanced unstructured finite volume (FV) scheme on triangular meshes was presented for modeling weakly nonlinear and weakly dispersive water waves over slowly varying bathymetries, as described by the 2D depth-integrated extended Boussinesq equations of Nwogu [23] rewritten in a conservation law form. The FV scheme numerically solves the conservative form of the equations following the median dual node-centered approach, for both the advective and dispersive part of the equations. For the advective fluxes, the scheme utilizes an approximate Riemann solver along with a well-balanced topography source term upwinding. Higher order accuracy in space and time is achieved through a MUSCL-type reconstruction technique and through a strong stability preserving explicit Runge-Kutta time stepping. After each step in the time marching scheme a sparse linear system must be solved in order to recover the velocity field. This work examines, in depth, the solution process that must be followed for this sparse and large linear system. The complexity of the system is such that iterative methods are necessary to obtain efficient solutions for even moderately sized problems. The applicability of different iterative methods, preconditioning and reordering possibilities are examined as to find an efficient solution procedure. The effect of the mesh's topology and of certain physical conditions, e.g. the flow field's reference water depth, to the solution procedure is also investigated.

The outline of this paper is as follows: The BT model equations are described in Section 2 while the numerical model used is briefly presented in Section 3 along with the derivation of the sparse matrix and its properties. The iterative methods used for the solution of the resulting linear system, the preconditioning methods tested and a detailed comparison of their performance are presented in Section 4, while Section 5 describes the reordering methods. Conclusions are drawn in Section 6.

## 2 Governing equations

A number of BT models have been developed to describe the transformation and propagation of waves in coastal regions. In this work, the model equations solved are the extended BT equations of Nwogu [23] which are based on the assumption that the wave height ( $A$ ) is much smaller than the water depth  $h$ . The equations derived by Nwogu [23] using the velocity vector  $[u, v]^T = \mathbf{u} \equiv \mathbf{u}_a$  at an arbitrary distance,  $z_a$ , from a still water level,  $h$ , as the velocity variable, instead of the commonly used depth-averaged velocity. The elevation of the velocity variable  $z_a$  became a free parameter used to optimize the equations and making them applicable to a wider range of water depths, compared to the classical Boussinesq

equations. The equations of Nwogu describe weakly non-linear and weakly dispersive water waves in variable water depth.  $\epsilon := A/h$ , which measures the weight of nonlinear effects, and the square water depth to wave length ( $L$ ) ratio  $\mu^2 := h^2/L^2$ , which represents the dimension of the dispersive effects, is of the same order with, i.e. the Stokes number  $S := \epsilon/\mu^2 = O(1)$ . The equations provide accurate linear dispersion and shoaling characteristics for values of  $kh \approx 3$  (intermediate water depths), where  $k$  is the wave number and  $kh$  is essentially a scale of the value of  $\mu$ , providing a correction of  $O(\mu^2)$  to the shallow water theory. The equations are presented here in a conservative-like form and as such are numerically approximated by an unstructured FV scheme. Following [18] the vector conservative form of the equations reads as:

$$\partial_t \mathbf{U} + \nabla \cdot \mathcal{H}(\mathbf{U}^*) = \mathbf{S} \quad \text{on } \Omega \times [0, t] \subset \mathbb{R}^2 \times \mathbb{R}^+, \quad (1)$$

where  $\Omega \times [0, t]$  is the space-time Cartesian domain over which solutions are sought,  $\mathbf{U}^* = [H, Hu, Hv]^T$  are the physically conservative variables,  $\mathbf{U}$  is the vector of the actual solution variables.  $H = \eta + h$  is the total water depth, with  $\eta$  the free surface elevation and  $h$  the still water depth.  $\mathcal{H} = [\mathbf{F}, \mathbf{G}]$  are the nonlinear flux vectors given as

$$\mathbf{U} = \begin{bmatrix} H \\ P_1 \\ P_2 \end{bmatrix}, \quad \mathbf{F} = \begin{bmatrix} Hu \\ Hu^2 + \frac{1}{2}gH^2 \\ Huv \end{bmatrix}, \quad \mathbf{G} = \begin{bmatrix} Hv \\ Huv \\ Hv^2 + \frac{1}{2}gH^2 \end{bmatrix},$$

where

$$\mathbf{P} = \begin{bmatrix} P_1 \\ P_2 \end{bmatrix} = \begin{bmatrix} Hu + HA \\ Hv + HB \end{bmatrix} = H \left[ \frac{z_a^2}{2} \nabla(\nabla \cdot \mathbf{u}) + z_a \nabla(\nabla \cdot h\mathbf{u}) + \mathbf{u} \right]. \quad (2)$$

The source term vector,  $\mathbf{S} = \mathbf{S}_b + \mathbf{S}_f + \mathbf{S}_d$ , includes the bed topography's ( $b$ ) slope  $\mathbf{S}_b$ , the bed friction effects  $\mathbf{S}_f$ , given in this work in terms of the Manning coefficient  $n_m$ , and the dispersive terms  $\mathbf{S}_d$ . These terms read as

$$\mathbf{S}_b = \begin{bmatrix} 0 \\ -gH\partial_x b \\ -gH\partial_y b \end{bmatrix}, \quad \mathbf{S}_d = \begin{bmatrix} -\psi_c \\ -u\psi_c + \psi_{M_x} \\ -v\psi_c + \psi_{M_y} \end{bmatrix} \quad \text{and} \quad \mathbf{S}_f = \begin{bmatrix} 0 \\ -gn_m^2 u \|\mathbf{u}\| h^{-1/3} \\ -gn_m^2 v \|\mathbf{u}\| h^{-1/3} \end{bmatrix},$$

with

$$\psi_c = \mathcal{C} + \mathcal{D} = \nabla \cdot \left[ \left( \frac{z_a^2}{2} - \frac{h^2}{6} \right) h \nabla(\nabla \cdot \mathbf{u}) + \left( z_a + \frac{h}{2} \right) h \nabla(\nabla \cdot h\mathbf{u}) \right], \quad (3)$$

$$\psi_M = \begin{bmatrix} \psi_{M_x} \\ \psi_{M_y} \end{bmatrix} = \begin{bmatrix} H_t \mathcal{A} \\ H_t \mathcal{B} \end{bmatrix} = \partial_t H \frac{z_a^2}{2} \nabla(\nabla \cdot \mathbf{u}) + \partial_t H z_a \nabla(\nabla \cdot h\mathbf{u}) \quad (4)$$

where

$$\mathcal{A} = \left( \frac{z_a^2}{2} (u_{xx} + v_{yx}) + z_a ((hu)_{xx} + (hv)_{yx}) \right), \quad (5)$$

$$\mathcal{B} = \left( \frac{z_a^2}{2} (u_{xy} + v_{yy}) + z_a ((hu)_{xy} + (hv)_{yy}) \right), \quad (6)$$

$$\mathcal{C} = \left[ \left( \frac{z_a^2}{2} - \frac{h^2}{6} \right) h (u_x + v_y)_x + \left( z_a + \frac{h}{2} \right) h ((hu)_x + (hv)_y)_x \right]_x \quad \text{and} \quad (7)$$

$$\mathcal{D} = \left[ \left( \frac{z_a^2}{2} - \frac{h^2}{6} \right) h (u_x + v_y)_y + \left( z_a + \frac{h}{2} \right) h ((hu)_x + (hv)_y)_y \right]_y. \quad (8)$$

Equations (1) have flux terms identical as those in the non-linear shallow water (or Saint-Venant) equations and variables  $\mathbf{P}$  contain all time derivatives in the momentum equations, including part of the dispersion terms. The dispersion vector  $\mathbf{S}_d$  contains only spatial derivatives since  $\partial_t H$  is explicitly defined by the mass equation.

### 3 The numerical model

To numerically solve equations (1) we use the FV scheme proposed in [18, 19]. The FV approach is of the node-centered median-dual type where the control volumes (see Fig. 1) are elements dual to the primal triangular mesh. The locations of the discrete solutions are called data points  $N$  which essentially correspond to the number of vertices of the mesh. Referring to Fig. 1, the boundary  $\partial C_P$  of a control volume (cell),  $C_P$ , around an internal node  $P$ , is defined by connecting the barycenters of the surrounding triangles (having  $P$  as a common vertex) with the mid-points of the corresponding edges that meet at node  $P$ .

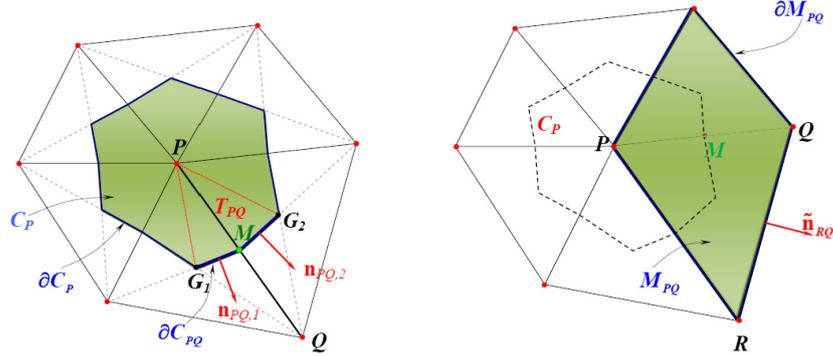


Figure 1: Median-dual computational cell (left) and computational cell for the gradient of the divergence (right)

#### 3.1 The FV approximation

After integration of (1) over each computational cell and application of the Gauss divergence theorem the semi-discrete form of the FV scheme reads as:

$$\frac{\partial \mathbf{U}_P}{\partial t} = -\frac{1}{|C_P|} \sum_{Q \in K_P} \Phi_{PQ} - \frac{1}{|C_P|} \Phi_{P,\Gamma} + \frac{1}{|C_P|} \iint_{C_P} \mathbf{S}_d \Omega, \quad (9)$$

where  $\mathbf{U}_P$  is the volume-averaged value of the conserved-like quantities at a given time,  $K_P$  is the set of the neighboring nodes to  $P$ , i.e.  $K_P := \{Q \in \mathbb{N} \mid \partial C_P \cap \partial C_Q \neq \emptyset\}$ ,  $\Gamma$  is the boundary of the computational domain  $\Omega$  and  $\Phi_{PQ}$ ,  $\Phi_{P,\Gamma}$  are the numerical flux vectors across each internal face,  $\partial C_P$ , and boundary face respectively.

The numerical fluxes are evaluated solving a Riemann problem at cell interfaces using the approximate Riemann solver of Roe [26]. To reach higher-order spatial accuracy an extension of the MUSCL methodology of Van Leer [34] is used. This extension relies on the evaluation of the fluxes with extrapolated physical variables  $[h, u, v]$  at the midpoint  $M$  of an edge  $PQ$ . Each component of the physical variables and bed topography  $b$  are extrapolated using extrapolation gradients which are obtained using a combination of centered and upwind gradients [1, 9, 18, 31] as to increase accuracy of the basic MUSCL reconstruction [9]. In this way a third-order spatial accuracy is obtained [18]. For the reduction of oscillations in cases where non-linearity prevails (e.g. when the dispersive terms become negligible) the use of a slope limiting procedure is necessary and the edge-based non-linear slope limiter of Van Albada-Van-Leer is used. Details for the numerical model used such as wet/dry front treatment, boundary conditions and discretization of the dispersive terms can be found in [18, 19].

### 3.2 The Resulting Linear System

Concerning the time discretization an optimal third order explicit Strong Stability Preserving Runge-Kutta (SSP-RK) method was adopted [32, 18] under the usual CFL stability restriction. In each time step on the RK scheme the values of the velocities  $u, v$  must be extracted from the new solution variable  $\mathbf{P} = [P_1, P_2]^T$  following from (2). The FV discretization of  $\mathbf{P}$  results to a sparse matrix. The linear system  $\mathbf{A}\mathbf{V} = \mathbf{C}$  with  $\mathbf{A} \in \mathbb{R}^{2N \times 2N}$ ,  $\mathbf{V} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_N]^T$  and  $\mathbf{C} = [\mathbf{P}_1 \ \mathbf{P}_2 \ \dots \ \mathbf{P}_N]^T$ , has to be solved in each step of the RK time marching scheme.

Keeping in mind that  $\mathbf{u}$  is our unknown velocity vector at each mesh node, each two rows of the matrix correspond to a node  $P \in \{1, 2, \dots, N\}$  on the grid and for each such node we have,

$$H_P \left[ \frac{z_a^2}{2} \nabla(\nabla \cdot \mathbf{u}) + z_a \nabla(\nabla \cdot h\mathbf{u}) + \mathbf{u} \right]_P = \mathbf{P}_P. \quad (10)$$

Now the gradients  $\nabla(\nabla \cdot \mathbf{u})$  and  $\nabla(\nabla \cdot h\mathbf{u})$  must be computed. For that reason, we use the average of the gradient  $(\nabla w)_P$  in a cell [2, 3, 4]:

$$(\nabla w)_P = \frac{1}{|C_P|} \sum_{Q \in K_P} w_M \mathbf{n}_{PQ}. \quad (11)$$

which is computed in each mesh node  $P$  by applying the Green-Gauss theorem in the region  $\Omega_P$ , i.e. the union of all triangles which share the vertex  $P$  and  $w_M$  is the value of  $w$  at the midpoint of the edge  $PQ$ . The outward normal vector to  $\partial C_{PQ}$  as  $\mathbf{n}_{PQ} = [n_{PQx}, n_{PQy}]^T$ , while  $\tilde{\mathbf{n}}_{PQ}$  is the corresponding unit vector. Using (11) equation (10) reads as:

$$H_P \left[ \frac{(z_a^2)_P}{2} \frac{1}{|C_P|} \sum_{Q \in K_P} (\nabla \cdot \mathbf{u})_M \mathbf{n}_{PQ} + \frac{(z_a)_P}{|C_P|} \sum_{Q \in K_P} (\nabla \cdot h\mathbf{u})_M \mathbf{n}_{PQ} + \mathbf{u}_P \right] = \mathbf{P}_P, \quad (12)$$

and is now obvious that we have to compute  $\nabla \cdot \mathbf{u}$  and  $\nabla \cdot h\mathbf{u}$  at  $M$ . Referring again to figure 1, a new computational cell  $M_{PQ}$  constructed by the union of

two triangles which share the edge  $PQ$  is defined. Hence, the discrete averages of the divergence can be computed as follows for  $(\nabla \cdot \mathbf{u})$ ,

$$(\nabla \cdot \mathbf{u})_M \approx \frac{1}{|M_{PQ}|} \sum_{\substack{R, Q \in K_{PQ} \\ R \neq Q}} \frac{1}{2} (\mathbf{u}_R + \mathbf{u}_Q) \cdot \mathbf{n}_{RQ} \quad (13)$$

where  $K_{PQ} = \{R \in \mathbb{N} \mid R \text{ is a vertex of } M_{PQ} \text{ and } RQ \in \partial M_{PQ}\}$ . A similar computation is performed for the approximation of  $(\nabla \cdot h\mathbf{u})$ . We refer to [18] for more details on the discretization. By performing the above approximations we restrict the unknown information used in (10), i.e. values of  $\mathbf{u}$ , only to that coming from the nodes that are neighbors of node  $P$ , i.e. nodes  $Q \in K_P$ .

Substituting the above equation to (12) and for each node, gives:

$$\begin{aligned} & \frac{(z_a^2)_P}{2} \frac{H_P}{|C_P|} \sum_{Q \in K_P} \left[ \frac{1}{|M_{PQ}|} \sum_{\substack{R, Q \in K_{PQ} \\ R \neq Q}} \frac{1}{2} (\mathbf{u}_R + \mathbf{u}_Q) \cdot \mathbf{n}_{RQ} \right] \mathbf{n}_{PQ} + \\ & \frac{(z_a)_P H_P}{|C_P|} \sum_{Q \in K_P} \left[ \frac{1}{|M_{PQ}|} \sum_{\substack{R, Q \in K_{PQ} \\ R \neq Q}} \frac{1}{2} (h_R \mathbf{u}_R + h_Q \mathbf{u}_Q) \cdot \mathbf{n}_{RQ} \right] \mathbf{n}_{PQ} = \mathbf{P}_P, \quad P = 1, \dots, (14) \end{aligned}$$

which can be further rewritten as:

$$\begin{aligned} & \frac{(z_a^2)_P}{2|C_P|} \sum_{Q \in K_P} \left[ \frac{1}{2|M_{PQ}|} \left( \sum_{\substack{R, Q \in K_P \cap K_{PQ} \\ R \neq Q}} \mathbf{u}_R \cdot (\mathbf{n}_{PR} + \mathbf{n}_{RQ}) + \mathbf{u}_P \cdot (\mathbf{n}_{SP} + \mathbf{n}_{PR}) \right) \right] \mathbf{n}_{PQ} + \\ & \frac{(z_a)_P}{|C_P|} \sum_{Q \in K_P} \left[ \frac{1}{2|M_{PQ}|} \left( \sum_{\substack{R, Q \in K_P \cap K_{PQ} \\ R \neq Q}} h_R \mathbf{u}_R \cdot (\mathbf{n}_{PR} + \mathbf{n}_{RQ}) + h_P \mathbf{u}_P \cdot (\mathbf{n}_{SP} + \mathbf{n}_{PR}) \right) \right] \mathbf{n}_{PQ} \\ & = \frac{1}{H_P} \mathbf{P}_P, \quad P = 1, \dots, N \quad (15) \end{aligned}$$

After some calculations the sparse  $2N \times 2N$  linear system to be solved can be presented in a compact form, as:

$$\frac{(z_a^2)_P}{2|C_P|} \sum_{Q \in K_P} (\mathbf{A}_Q \mathbf{u}_Q + \mathbf{A}_P \mathbf{u}_P) + \frac{(z_a)_P}{|C_P|} \sum_{Q \in K_P} (\mathbf{B}_Q \mathbf{u}_Q + \mathbf{B}_P \mathbf{u}_P) + \mathbf{I}_2 \mathbf{u}_P = \frac{1}{H_P} \mathbf{P}_P, \quad P = 1 \dots N, \quad (16)$$

where the sub-matrices  $\mathbf{A}_Q, \mathbf{A}_P, \mathbf{B}_Q, \mathbf{B}_P$ , depend only on geometric quantities and area  $|M_{PQ}|$ , and are given as

$$\begin{aligned} \mathbf{A}_Q &= \frac{1}{2|M_{PQ}|} \begin{bmatrix} \sum_{R, Q \in K_P \cap K_{PQ}} (n_{PRx} + n_{RQx}) n_{PQx} & \sum_{R, Q \in K_P \cap K_{PQ}} (n_{PRy} + n_{RQy}) n_{PQx} \\ \sum_{R, Q \in K_P \cap K_{PQ}} (n_{PRx} + n_{RQx}) n_{PQy} & \sum_{R, Q \in K_P \cap K_{PQ}} (n_{PRy} + n_{RQy}) n_{PQy} \end{bmatrix}, \\ \mathbf{B}_Q &= \frac{1}{2|M_{PQ}|} \begin{bmatrix} \sum_{R, Q \in K_P \cap K_{PQ}} h_R (n_{PRx} + n_{RQx}) n_{PQx} & \sum_{R, Q \in K_P \cap K_{PQ}} h_R (n_{PRy} + n_{RQy}) n_{PQx} \\ \sum_{R, Q \in K_P \cap K_{PQ}} h_R (n_{PRx} + n_{RQx}) n_{PQy} & \sum_{R, Q \in K_P \cap K_{PQ}} h_R (n_{PRy} + n_{RQy}) n_{PQy} \end{bmatrix}, \end{aligned}$$

Inria

$$\mathbf{A}_P = \frac{1}{2|M_{PQ}|} \begin{bmatrix} (n_{SP_x} + n_{PR_x}) n_{PQ_x} & (n_{SP_y} + n_{PR_y}) n_{PQ_x} \\ (n_{SP_x} + n_{PR_x}) n_{PQ_y} & (n_{SP_y} + n_{PR_y}) n_{PQ_y} \end{bmatrix} \quad \text{and} \quad \mathbf{B}_P = h_P \mathbf{A}_P.$$

The number of geometrical entries in each summation is always two, while the number of entries in the summation  $\sum_{Q \in K_P}$  in (16) is equal to the number of the neighbors of  $P$ . This means that the maximum non-zero elements of the matrix  $\mathbf{M}$  in each row  $P$  in (16) are two times the number of the neighbors of  $P$  plus one.

It is important to state here that, for the linear system (16) its coefficient matrix is constant in time. So it is constructed and stored, in a compressed sparse row CSR format [29], at a pre-processing stage at the beginning of each simulation. The present work concerns the solution process after the sparse matrix has been stored.

### 3.3 Matrix properties

System's coefficient matrix, named  $\mathbf{A}$ , is an un-symmetric but structurally symmetric matrix, in terms of its non-zero entries. It is often ill-conditioned and also mesh dependent. This means that the sparsity pattern of the matrix depends on the ordering of the nodes on each grid. Different grids lead to different matrix structures. In this work four type of grids are used, see figure 2. For a

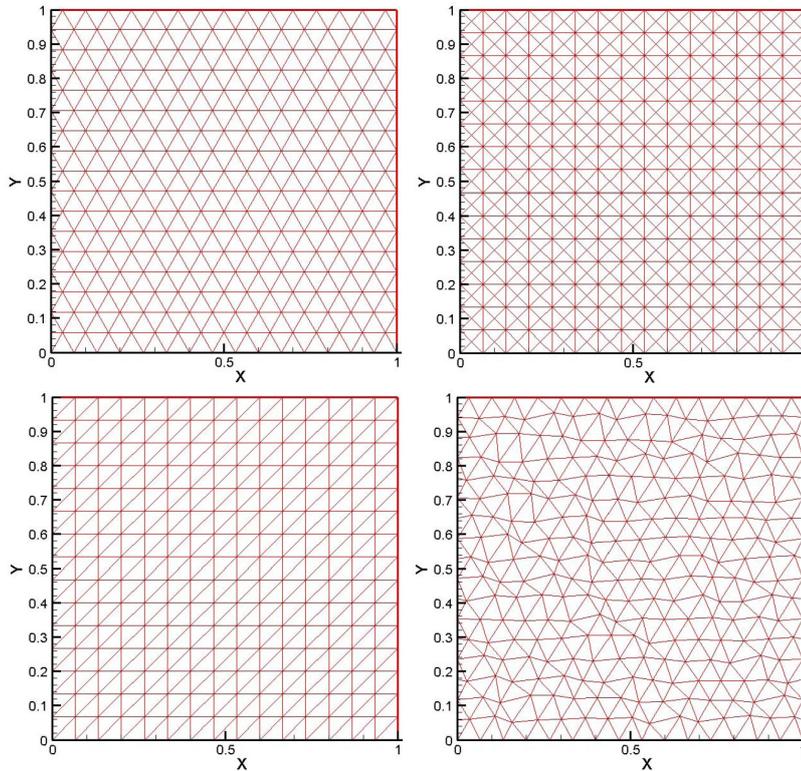


Figure 2: Representative grid types: Equilateral, Orthogonal I, Orthogonal II, Distorted (left to right)

given computational domain, and with out loss of generality, with dimensions  $L_x \times L_y$  in the  $x$ - and  $y$ - direction respectively, we define a subdivision of  $L_x$  by  $N_x$  line segments, namely  $D_x = L_x/N_x$  and depending on the grid type the corresponding subdivision  $D_y$  of  $L_y$  can be easily determined. As such, we define the characteristic length (effective mesh size) for each grid as  $h_N = \sqrt{\frac{L_x \times L_y}{N}}$ .

Figure 3 shows the structure of each matrix using the corresponding grids of figure 2. Each grid used has 15 nodes in the  $x$ -axis ( $N_x = 15$ ) and the resulting

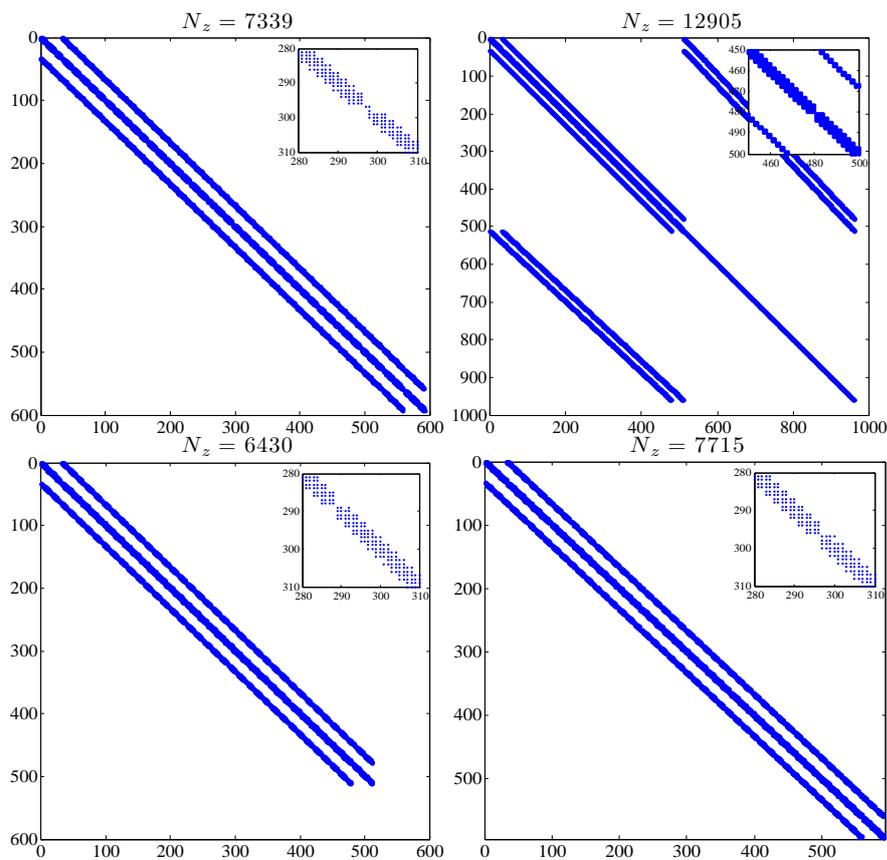


Figure 3: Sparsity patterns for the different mesh types for  $N_x = 15$

matrix has  $N_z$  non-zero elements shown. The resulting structure using the Orthogonal I type of grid is quite different from that of the other types which have a much smaller bandwidth. Further, for the Orthogonal I grid type the number of the non-zero elements,  $N_z$ , in the system's matrix is almost double of that obtained from the other grids. Matrices structure remains the same while refining the meshes. Table 1 shows the characteristics of indicative matrices produced using for  $L_x = L_y = 1$  and different nodes in  $x$ -axis ( $N_x$ ) while Table 2 show the  $h_N$  values for each matrix reported in Table 1.

The properties of the matrix are also affected by the physical situation of the problem examined. The most important parameters are the still water level,

$N_x$	Equilateral	Orthogonal I	Orthogonal II	Distorted
15	352836 (7339)	925444 (12905)	262144 (6430)	352836 (7715)
30	4857616 (28436)	13853284 (50806)	3694084 (25163)	4857616 (29746)
60	74132100 (118293)	214388164 (201931)	55383364 (99258)	74132100 (118293)

Table 1: Total ( $2N \times 2N$ ) and non-zero elements ( $N_z$ ) for the matrices produced for  $L_x = L_y = 1m$  with different  $N_x$ .

$N_x$	Equilateral	Orthogonal I	Orthogonal II	Distorted
15	0.058	0.0456	0.0625	0.058
30	0.0301	0.0232	0.0323	0.0301
60	0.0152	0.017	0.0164	0.0152

Table 2: The  $h_N$  values for the matrices produced for  $L_x = L_y = 1m$  with different  $N_x$ .

see (15) with relation to the nodes used on the grid, i.e.  $h_N$ . To illustrate the dependence on the ratio  $\frac{h}{h_N}$  the spectrum of eigenvalues for six matrices are shown in Figures 4 and 5. The six matrices examined have been produced using two type of grids, equilateral and Orthogonal I. Figure 4 shows the eigenvalues of three matrices produced using the equilateral type of Grid. The first matrix (on the left) has  $h = 1m$  and  $h_N = 0.058$  the second (center) uses  $h = 1m$ ,  $h_N = 0.0301$  and the third one (right) uses  $h = 100m$  and  $h_N = 0.058$ . All matrices have minimum eigenvalues near zero and it is noted that as the ratio  $\frac{h}{h_N}$  grows the spectrum of the matrix has a much larger spread of eigenvalues through the right half of the complex plane. The same conclusions can be obtained from figure 5 which depicts the spectrum of eigenvalues for three matrices, obtained using the Orthogonal I type of grid. The parameters used are the same as in figure 4.

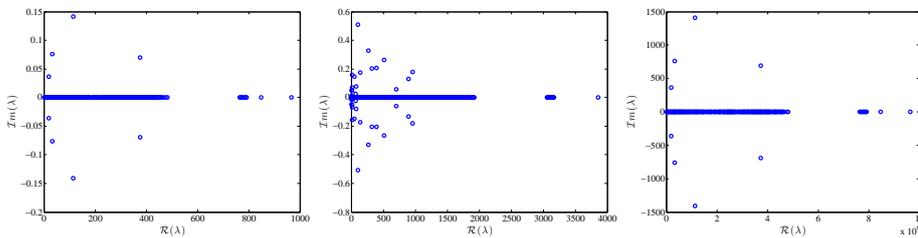


Figure 4: Eigenvalues of three matrices using the equilateral type of grid.  $\frac{h}{h_N} = \frac{1}{0.058}$  (left),  $\frac{1}{0.031}$  (middle),  $\frac{100}{0.058}$  (right)

The above behavior gives evidence that preconditioning is crucial in this type of problems. These conclusions are further reflected in the respective condition numbers of the three matrices which are between  $O(10^3)$ - $O(10^5)$ . For example the matrix produced by the distorted grid and  $h = 1m$ ,  $N_x = 15$  (see figure 2)

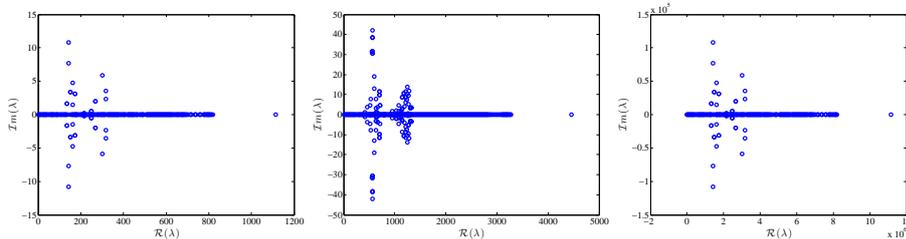


Figure 5: Eigenvalues of three matrices using the Orthogonal I type of grid.  $\frac{h}{h_N} = \frac{1}{0.0456}$  (left),  $\frac{1}{0.0232}$  (middle),  $\frac{100}{0.0456}$  (right)

has a condition number  $7.7732E+05$  while the one produced from the equilateral grid has condition number  $5.5348E+03$ . The matrices become even more ill-conditioned as the depth is further increased, or the grid is refined. Similar behavior has been mentioned in [14] where the finite difference method is used to solve a Boussinesq model in two dimensions.

## 4 Iterative methods, preconditioning and reordering

### 4.1 Application of Iterative Methods

In this presentation, we want to develop an optimal strategy for the solution of (16) abbreviated to  $\mathbf{Ax}=\mathbf{b}$ , and we use a "toy" problem in which the right hand side vector  $\mathbf{b}$  is computed adding the columns of the matrix  $\mathbf{A}$ . The initial guess used for the iterative methods was the zero vector. From now on all test cases presented will solve this "toy" problem unless otherwise said. We must mention again that the matrix depends on the numbering and the geometrical quantities of the mesh nodes. So, and since it is not depending on the unknown quantities of the BT model we can construct it and store it only once in the beginning of our simulation.

One common solution method for sparse systems  $\mathbf{Ax} = \mathbf{b}$  is to use a sparse direct solution algorithm via a complete factorization of the matrix. However, for large scale problems concluding to large sparse linear systems, the computational time required for the factorization along with the storage requirements could be an insurmountable problem. An alternative to direct solution methods is the use of iterative ones which have significantly lower storage demands. Two common Krylov Subspace iterative methods are used in this work. The Generalized Minimal Residual (GMRES) and the Biconjugate Gradient Stabilized (BiCGStab) algorithms. Both methods are implemented using the SPARSKIT package [28]. The optimal solution strategy depends on factors such as problem size, sparsity pattern and the matrix's eigenspectrum.

Figure 6 depicts the performance (in terms of CPU time) of each method for two different type of grids, of the Equilateral type (left) and the Orthogonal I type. Each sub-figure depict the computational time needed as to solve the sparse system while the nodes are refined and the depth is increased. For both cases we can see that the computational time grows dramatically not only as the

h

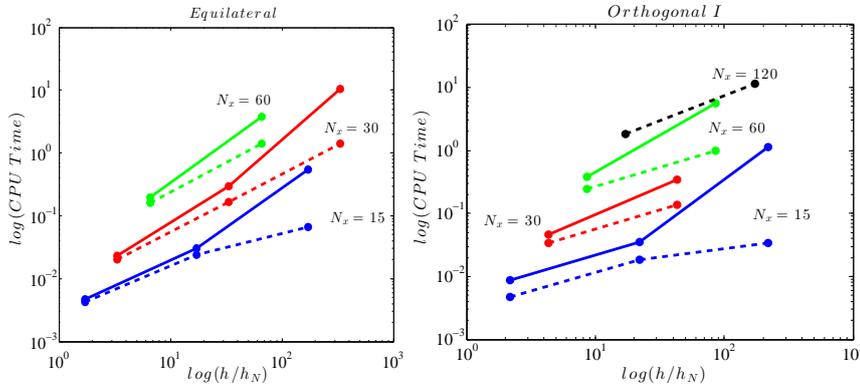


Figure 6: CPU time versus variable still water level to  $h_N$  ratio for GMRES (solid line) and BiCGStab (dashed line)

mesh is refined (which is expected) but also as the still water level  $h$  (with respect to  $h_N$ ) is increased. Furthermore, comparing the CPU times needed for each iterative method we can see that they start to vary as the water depth to  $h_N$  ratio is increased. The BiCGStab method solves the linear systems produced (for both types of grid) in less time than the GMRES. In some cases both methods could not reach convergence for higher  $h/h_N$  values. This was more prominent for GMRES. For both iterative methods a relative residual error tolerance  $r = \|b - Ax\|_2 / \|b\|_2 \leq 10^{-6}$  was used as the convergence criterion. Following from the above, we can clearly understand that the usage of a preconditioning method is mandatory as to reduce the number of iterations needed and consequently the total computational time.

## 4.2 Application of Preconditioning methods

Although Krylov Subspace iterative methods are well suited to solve relatively sparse system, they can exhibit slow convergence. Thus, it is essential to use a good preconditioning strategy as to enhance their convergence properties and reduce the computational cost. A survey of preconditioning techniques for large linear systems can be found in [5]. So, this section introduces the different techniques of preconditioning tested.

A good preconditioner must approximate matrix  $\mathbf{A}$  well, while at the same time being easy to solve. Let  $\mathbf{M}$  be a new non-singular matrix which is a good approximation of  $\mathbf{A}$ . The preconditioned system  $\mathbf{M}^{-1}\mathbf{A}\mathbf{x} = \mathbf{M}^{-1}\mathbf{b}$  will have the same solution as system  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , and it can be solved easier. The non-singular matrix  $\mathbf{M}$  is called preconditioner. In the present work all preconditioning is done from the left even though preconditioning from the right has been found to be equally effective. Three preconditioning techniques, which are freely available as part of the SPARSKIT [28] package, are presented and tested. The ILU(0), ILU( $k$ ) and ILUT preconditioners.

Since convergence of an iterative method such as GMRES and BiCGStab depends on the eigenvalues of the matrix, generally speaking, preconditioning

attempts to improve the spectral properties of  $\mathbf{A}$  [5]. A "good" preconditioner transforms the matrix so that the original sparse linear system can be solved easily, with low storage demands, in low computational expense. The preconditioned matrix should have eigenvalues away from zero. Even for non-symmetric matrices with a cluster of the eigenvalue spectrum away from zero can still lead to a rapid convergence.

One way to approach this is to use a direct method such as  $\mathbf{LU}$ -decomposition, i.e.  $\mathbf{M} = \mathbf{LU}$ . The system can then be solved in two steps using the factor  $\mathbf{L}$  and  $\mathbf{U}$ . The drawback of this method is that during the factorization process, matrices  $\mathbf{L}$  and  $\mathbf{U}$  are dense so the computer storage demands and the dimension of the problem may become huge. One of the simplest ways of defining a preconditioner is to perform an incomplete factorization of the original matrix  $\mathbf{A}$ . This entails a decomposition of the form  $\mathbf{A} = \mathbf{LU} - \mathbf{R}$  where  $\mathbf{L}$  and  $\mathbf{U}$  have the same nonzero structure as the lower and the upper part of  $\mathbf{A}$ , respectively and  $\mathbf{R}$  is the residual matrix of the factorization. This incomplete factorization is known as  $\text{ILU}(0)$  and often leads to a crude approximation which may result to many iterations to converge. To remedy this, several alternative incomplete factorizations have been developed by allowing more fill in  $\mathbf{L}$  and  $\mathbf{U}$ . In practice and as to find the  $\mathbf{L}$  and  $\mathbf{U}$  the Gaussian elimination process is used and a level of fill is attributed to each element which is dropped or not according to this value. [29]. Generally the more accurate  $\text{ILU}$  factorizations require fewer iterations to converge but of course the preprocessing cost to compute the factors is higher. Incomplete factorizations that rely on the level of fill are blind to numerical values because elements that are dropped depend only on the structure of the Matrix. Some methods are available based on dropping elements in the Gaussian elimination process according to their magnitude rather than their locations.

Different factorization algorithms have different rules that govern the dropping or fill-in in the incomplete factors. In this work two rules will be used; the level of fill approach and the drop tolerance parameter. From now on applying a dropping rule to an element will only mean replacing the element with zero if it satisfies a set of criteria. The drop tolerance is a positive number  $\tau$  which is used in a dropping criterion. The dropping strategy depends on the matrix, the non-zero numbers and the conditional number. The philosophy of that is to accept only the values greater (in absolute value) than the tolerance in the new fill-ins. Generally, it is difficult to choose the right value for the drop tolerance, because one can't predict the amount of storage that will be needed to store [5]. Three preconditioning techniques, which are freely available as part of the SPARSKIT [28] package, are presented and tested in this work. The  $\text{ILU}(0)$ ,  $\text{ILU}(k)$  and  $\text{ILUT}$  preconditioners.

#### 4.2.1 The $\text{ILU}(0)$ preconditioner

When the level of fill is equal to zero then the  $\text{ILU}(0)$  preconditioner, in the  $\text{ILU}$  factorization, is recovered. So,  $\text{ILU}(0)$  refers to a full factorization, with no reduction or fill-in, and it is also called the no-fill  $\text{ILU}$  preconditioner. We tested this preconditioner for systems produced for two types of grids, namely the Equilateral and Orthogonal of type I, using the two iterative methods (GMRES and BiCGStab). The computational values in the toy problem were  $h = 1m$  and  $N_x = 15, 30, 60, 120$ . For the Equilateral type of grid none of the two iterative

methods was able to converge while for the Orthogonal I type of grid only the BiCGStab method for  $N_x = 30$  converged after 34957 iterations. Even though the ILU(0) preconditioner is easy to implement and its computation is inexpensive, it is quite effective mainly for simpler problems. By simpler problems we mean, for example, low-order discretizations of scalar elliptic PDEs leading to non-singular matrices and diagonally dominated matrices. These are the type of problems for which these preconditioners were originally proposed [5]. For this case study the non-fill factorization resulted in too crude approximation of the matrix  $\mathbf{A}$ , so more sophisticated preconditioners must be used.

#### 4.2.2 The ILU( $k$ ) preconditioner

The computation of the ILU( $k$ ) preconditioner requires only the fill-in criterion. So, if  $k$  is a non-negative integer, all the fill-ins whose level is greater than  $k$  are dropped. The limitation of this process is that for matrices that are not diagonal dominated ILU( $k$ ) requires to store many fill-ins that are small in absolute value leading to expensive computations and less quality for the preconditioner. In this work we used the value  $k = 300$  elements per row (in the factors) which has been found to be a good universal parameter for the problems presented here.

Figure 7 demonstrates how the performance of the iterative methods is affected by the application of the ILU( $k$ ) strategy in terms of the increasing still water level to  $h_N$  ratio for system matrices produced by Equilateral type of grids. The iterations needed for convergence vary between 2 and 4 for both methods. Comparing with figure 6 (on the left) a substantial improvement in computational time can be observed. The results show that the total time needed for convergences is not affected by the ratio  $h/h_N$  but, as expected, is highly increased while increasing the number of unknowns of the problem (for larger  $N_x$  values i.e. smaller  $h_N$  values). The two iterative methods needed the same computational time.

Table 3 presents the actual times and iterations needed for the solution of the linear systems produced using the Orthogonal I type of mesh. We must specify that whenever the results are omitted or a dash is placed then the iterative method used failed to converge. It can be concluded from Table 3 that, as the ratio  $h/h_N$  increases from a grid refinement, i.e. lower values of  $h_N$  are produced, the computational time dramatically increases and in some cases convergence can not be obtained for system resulting from this type of grid. The main reasons for this are, the large amount of nodes in the mesh, compared to the other types of meshes (see Table 2) and the different coefficient matrix structure, which has a bigger bandwidth (please refer to figure 3). For this grid type, and although it is a structured one, the ILU( $k$ ) preconditioner fails to improve the condition of the matrix for finer meshes even though the iterations needed for the converged are significantly decreased at coarser ones.

It should be noted here that, the total time needed for the computation of the preconditioner is independent of the still water depth  $h$ . The actual CPU times needed are for matrices resulting from refined meshes for the two grid types are presented in Table 4. These times are very small compared to the total time needed for the iterative methods to reach convergence.

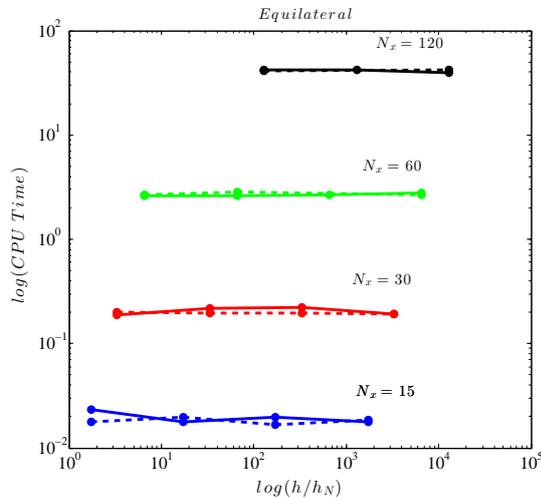


Figure 7: CPU time versus  $h/h_N$  for GMRES (solid line) and BiCGStab (dashed line) using the  $ILU(k)$  preconditioner with  $k = 300$ .

$h/h_N$	GMRES (sec) / iterations	BiCGStab (sec) / iterations
0.1/0.0456	1.588515759 / 2	1.588515759 / 3
1.0/0.0456	1.525697947 / 2	1.581536055 / 3
10/0.0456	1.580311775 / 3	1.591470003 / 3
0.1/0.0232	87.64499593 / 2	87.37235212 / 3
10/0.0232	-	87.47821903 / 3

Table 3: CPU time for the solution of linear systems resulting from Orthogonal I type of grids using the  $ILU(k)$  preconditioner, with  $k = 300$ , preconditioner

#### 4.2.3 $ILUT(p, \tau)$ preconditioner:

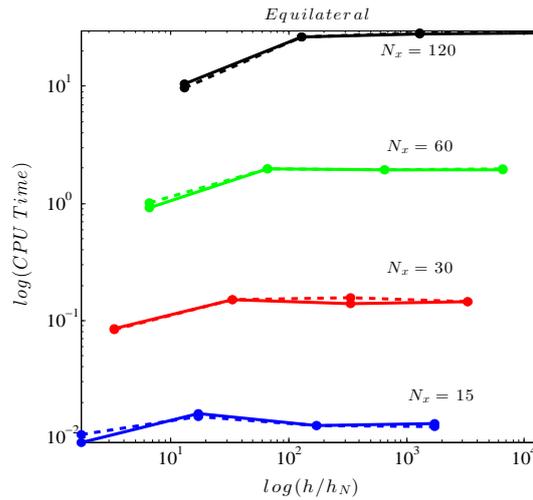
The general algorithm of preconditioner ILU with threshold,  $ILUT(p, \tau)$ , includes a set of rules for dropping small elements. Thus, the used of  $ILUT$  produces a lower storage factorization. The main idea here is to replace an element with zero if its value is smaller than a threshold value  $\tau$ , set by the user. The dropping rule can be applied to a row, by checking all the elements and keeping/ignoring them depending on their arithmetic value [29]. By increasing the drop tolerance (threshold) the sparsity of the preconditioner increases and the amount of work that will be needed for applying the preconditioner decreases. The limitation here is that the application of the  $ILUT$  can affect the solution's accuracy thus the iterative methods may require more iterations. An additional dropping rule which can be applied is to keep only the  $p$ -th largest elements in the  $\mathbf{L}$  part of the row and the  $p$ -th largest elements in the  $\mathbf{U}$  part of the row in addition to the diagonal element, which is always kept. In this work we use always  $p = 300$ . The goal of this dropping step is to control the number of elements per row.

Figure 8 presents, like before, how the performance of the  $ILUT$  strategy is affected by the variable water depth to  $h_N$  ratio, for the matrices produced by

$N_x$	Equilateral (sec)	Orthogonal I (sec)
15	0.014	0.012
30	0.19	0.15
60	2.62	2.24
120	-	-

Table 4: Times needed for the ILU( $k$ ) preconditioners's computation.

the equilateral type of grid. The threshold value used here was set to  $\tau = 10^{-5}$ . The iterations needed for convergence varied for 3-13 but the CPU time needed is slightly less compared to the one obtained using the ILU( $k$ ) preconditioner (as presented in figure 7) for higher  $h/h_N$  values. For smaller ratios some time improvement can be noticed. We must acknowledge here that for the systems produced by the Orthogonal I type of grid both iterative methods failed to converge. Lowering the threshold value  $\tau$  lead to no substantial improvement in computational efficiency or convergence of the iterative methods.

Figure 8: CPU time versus  $h/h_N$  for GMRES (solid line) and BiCGStab (dashed line) applying the ILUT preconditioner with  $\tau = 10^{-5}$ .

## 5 Reorderings

Since all preconditioners applied in the previous Section failed to substantially improve the condition of the produced matrices as to decrease the iterations for convergence and/or the computational time needed and in some cases no convergence could be obtained, we must use a reordering technique as to improve the stability of the incomplete factorizations. It is well known that, the incomplete factorization preconditioners are sensitive to the ordering of unknowns and equations [5]. Optimal reordering strategies can be used with the dual purpose of limiting the bandwidth of the discrete operator  $\mathbf{A}$  and to reduce excessive fill-in in the factorization of the involved operators [11]. In most cases,

reordering techniques tend to affect the rate of convergence of preconditioned Krylov subspace methods [5]. These algorithms aim to minimize the bandwidth of the matrix  $\mathbf{A}$ . Two different approaches are used in this work namely, the Cuthill-McKee (CMK) and the Reverse Cuthill-McKee (RCM) permutations.

Figure 9 shows the CPU time versus the relative water depth for both type of grids using again the GMRES and BiCGStab algorithms along with the  $ILU(k)$  preconditioner with  $k = 300$  and CMK reordering. As can be seen in figure 9 the usage of CMK reordering in addition to the  $ILU(k)$  preconditioner improves the convergence process, compared to the single usage of the preconditioner, for all the matrices produced with Orthogonal I type of grids. It can be also observed that on grid with the same degrees for freedom, increasing the water depth does not affect the computational time needed.

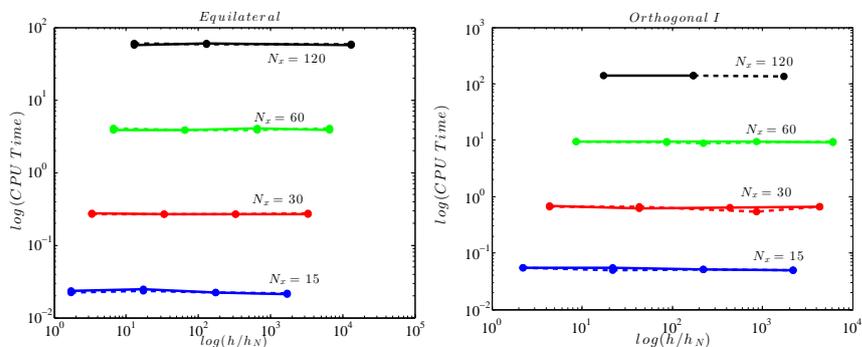


Figure 9: CPU time versus  $h/h_N$  for GMRES (solid line) and BiCGStab (dashed line) using the  $ILU(k)$  preconditioner and CMK reordering.

Using the ILUT preconditioner and the CMK reordering (see figure 10) the computational results vary. We tested different values of drop tolerance  $\tau$  for this preconditioner. For values greater than  $10^{-5}$  none of the iterative methods was able to converge. Figure 10 shows the CPU time versus the relative water depth using CMK and ILUT with drop tolerance  $10^{-5}$  (left column). Even though we observe convergence to the solution for most of the test cases, the computational time increased as the relative water depth was increased. This makes the choice of the drop tolerance unsuitable for physical problems with relative large water depth. Using a smaller value for  $\tau = 10^{-10}$  (see figure 10 right column) the behavior of the linear systems, especially when using GMRES, closely approximates the behavior shown in figure 9 since the CPU time is not affected by the variation of the ratio  $(h/h_N)$  (using the same degrees of freedom). The results of this work confirm also those found in [14], where non-linear and highly dispersive Boussinesq equations are solve using a finite difference scheme. The ILUT preconditioner works quite well in shallow to intermediate water depths, however may rapidly lose effectiveness as the depth is further increased. However, a big improvement was that, for systems produced from Orthogonal I type grid in this case both iterative methods converged with satisfactory performance.

One of the most common reordering techniques, used mostly with the finite element method, is the Reverse Cuthill-McKee (RCM) ordering [16]. The work

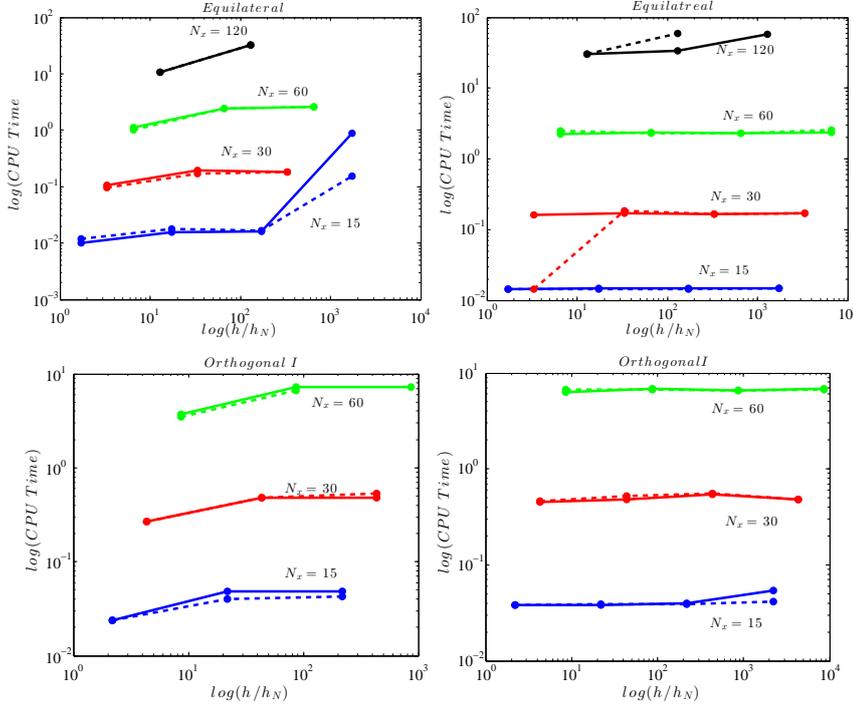


Figure 10: CPU time versus variable water depth for GMRES (solid line) and BiCGStab (dashed line) using ILUT preconditioner with threshold  $10^{-5}$  (left) and  $10^{-10}$  (right) and CMK reordering. For equilateral type of grid (up) and for Orthogonal I (down).

of Benzi et. al [6], on ordering for incomplete factorization, revealed that the use of RCM method was advantageous for non-symmetric problems.

Figure 11 depicts the cpu time versus the ratio  $h/h_N$  using the  $ILU(k)$  preconditioner and the RCM reordering while figure 12 shows the same but with the usage of the ILUT preconditioner and the RCM reordering. Again two different tolerances have been used. The same behavior as the one presented in figures 9 and 10 is observed.

## 5.1 Spatial accuracy and efficiency

We performed studies for the accuracy and efficiency of the FV numerical from [18] considering the propagation of a solitary wave over an undisturbed depth. The numerical test case consists of a solitary wave of amplitude  $A = 0.1m$  which propagates over a flat topography of depth  $h = 1m$ . The computational domain was  $(x, y) \in [0, 300m] \times [0, 5m]$ . To prove the validity of the current study we examined the total and per time-step CPU time need to advance the model in one time-step and additionally the total and per time-step CPU times for the solutions of the  $2N \times 2N$  sparse linear system, using the BiCGStab method (and  $ILUT(p, \tau)$  preconditioner with  $p = 300$  and  $\tau = 10^{-10}$  and RCMK reordering).

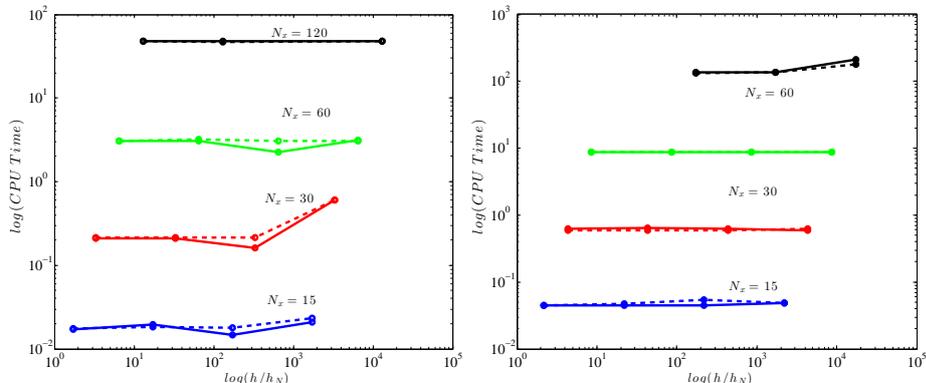


Figure 11: CPU time versus variable water depth for GMRES (solid line) and BiCGStab (dashed line) using  $ILU(k)$  preconditioner and RCM reordering. for equilateral type of grids (left) and for Orthogonal I types (right).

As shown in figure 13 the CPU time grows like  $O(\|E(\eta)\|^{-1})$  (linearly) while the time needed by the BiCGStab (as to solve the linear system) like  $O(\|E(\eta)\|^{-0.85})$ , for the finer grids where  $\|E(\eta)\|$  is the error of the free surface elevation measured in  $L_2$  norm. However, and due to the increase of the number of time steps needed on finer grids, the total CPU time grows approximately like  $O(\|E(\eta)\|^{-1.25})$  while the total time needed by the BiCGStab like  $O(\|E(\eta)\|^{-1.5})$  and starts to dominate the overall time, as grids get refined.

To assess the effect the increase of the number of grid points  $N$  has to the storage requirements of the non-zero elements ( $N_z$ ) of the  $2N \times 2N$  sparse linear system and to the computational efficiency, we present relevant comparison in figure 14. As expected, the  $N_z$  entries grow almost linearly with respect to  $N$ . The BiCGStab CPU time per time step scales like  $O(N^{3/2})$  however, the total CPU time per time step is growing like  $O(N^{5/4})$ , close to linear.

All computations were performed on a shared memory machine HP DL180G6 consists of two 6-core Xeon X5660@2.8GHz type processor with 12 MB Level 3 cache memory. The total memory is 64 GB and the operating system is Oracle Linux version 6.1. The application is developed in double precision Fortran code using Oracle Solaris Studio compilers version 12.2.

## 6 Conclusions

A numerical study is presented for numerically solving the sparse linear system that results from a finite volume discretization of the extended Boussinesq-type equations of Nwogu on unstructured triangular meshes. The linear system arises due to the reformulation of the model's equations in conservative-like form and its solution is essential to recover the actual velocity field in the flow. The systems's coefficient matrices are structurally symmetric and sparse but in most cases ill-conditioned. Its numerical solution consumes much of the actual computation time needed for the numerical scheme to advance one discrete time step, especially as the meshes get finer. The effect of the grids resolution along

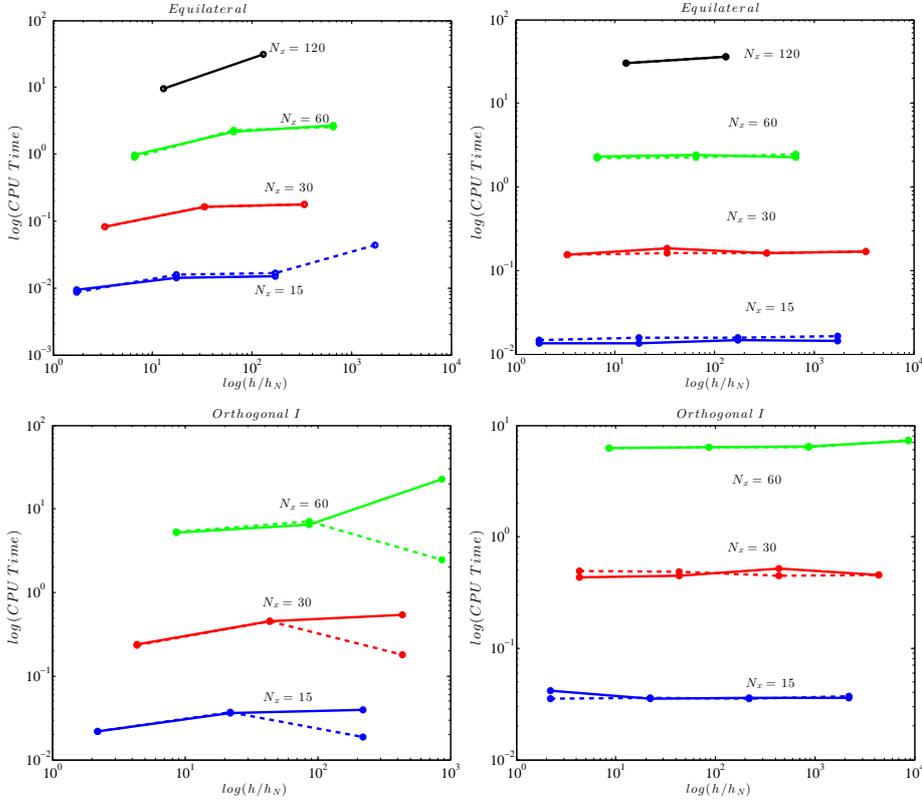


Figure 12: CPU time versus variable water depth for GMRES (solid line) and BiCGStab (dashed line) using ILUT preconditioner with threshold  $10^{-5}$  (left) and  $10^{-10}$  (right) and RCM reordering for Equilateral type of grids (up) and Orthogonal I types (down).

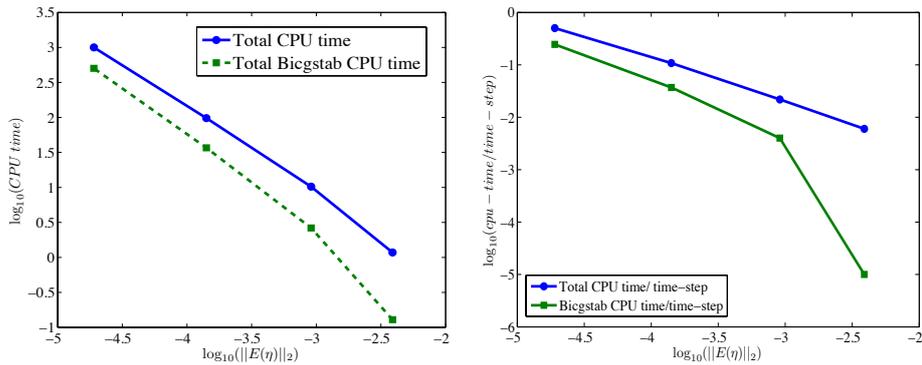


Figure 13: Non-zero elements as a function of  $N$  (left) and CPU times as a function of  $N$  (right)

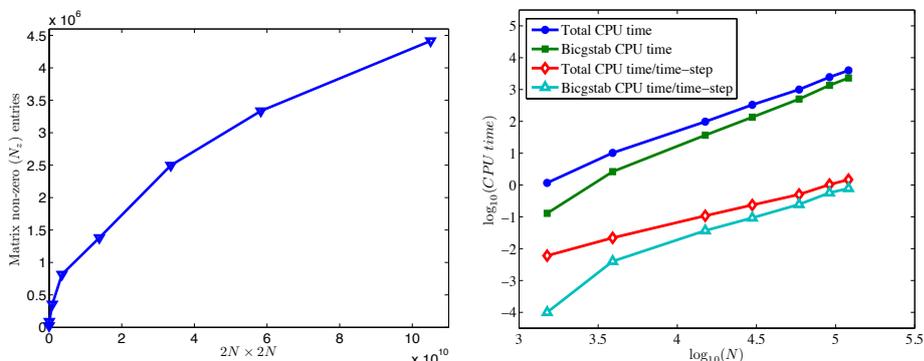


Figure 14: Non-zero elements as a function of  $N$  (left) and CPU times as a function of  $N$  (right)

with the physical value of the still water reference in the problem was investigated. Different iterative methods, precondition and reordering techniques were investigated as to to conclude to a an optimal and robust strategy for the system's solution. The following major conclusions can be drawn from this work:

- BiCGSTAB and GMRES iterative methods give almost similar results for the resulting systems, with the BiCGSTAB to have been proven more robust in some cases and is the method of choice following from this work.
- The usage of preconditioning and/or reordering is mandatory as to achieve convergence for the different mesh types used.
- Using preconditioning and reordering we gained convergence for (all) systems in every water depth. Using only preconditioning we were able to solve efficiently systems that have a small condition number (usually derived from equilateral grids).
- Using a drop tolerance  $\tau = 10^{-5}$  (for  $\text{ILU}(k)$  and  $\text{ILUT}$  preconditioners): CPU time using  $\text{ILUT}$  is less than that of using  $\text{ILU}(k)$  in average water depths. The usage of  $\text{ILU}(k)$  maybe more expensive in time but results on an overall the same CPU time in any water depth for the same grid resolution for convergence.
- As to correct the limitation of  $\text{ILUT}$  we decrease the drop tolerance and we observed that systems produced for deeper water depths both iterative methods converge, but of course with an additional time cost. Like before the CPU time is independent on the relative water depth on each matrix.

## References

## References

- [1] T. J. Barth. A 3-D upwind Euler solver for unstructured meshes. *AIAA paper 91-1548CP*, 1991.

- 
- [2] T. J. Barth. Aspects of unstructured grids and finite volume solvers for the Euler and Navier-Stokes equations. In Special Course on Unstructured Grid Methods for advection Dominated Flows, AGARD report 787, 1992.
- [3] T. J. Barth. *Numerical Methods and Error Estimation for Conservation laws on Structured and Unstructured Meshes*. VKI Computational Fluid Dynamics Lecture Series, 2003.
- [4] T. J. Barth and M. Ohlberger. Finite volume methods: foundation and analysis. In E. Stein, R. de Borst, and T.R. Huges, editors, *Encyclopedia of Computational Mechanics*. John Wiley and Sons Ltd., 2004.
- [5] M. Benzi. Preconditioning Techniques for Large Linear Systems: A Survey. *J. Comp. Phys.*, 182:418–477, 2002.
- [6] M. Benzi, D. B. Szyld, and A. V. Duin. Orderings for incomplete factorization preconditioning of nonsymmetric problems. *SIAM Journal of Scientific Computing*, 20(5):1652–1670, 1999.
- [7] P. Bonneton, E. Barthelemy, J. D. Carter, F. Chazel, and S. T. Cien. Fully nonlinear weakly dispersive modelling of wave transformation, breaking and runup. *Preprint submitted to European Journal of Mechanics - B/Fluids*, 2011.
- [8] M. Brocchini. A reasoned overview on Boussinesq-type models: the interplay between physics, mathematics and numerics. *Proc. R. Soc. A*, 469, 2013.
- [9] A. I. Delis, I. K. Nikolos, and M. Kazolea. Performance and comparison of cell-centered and node-centered unstructured finite volume discretizations for shallow water free surface flows. *Archives of Computational Methods in Engineering*, 18:57–118, 2011.
- [10] W. El Asmar and O. Nwogu. Finite volume solution of Boussinesq-type equations on an unstructured grid. In J. McKee Smith, editor, *Proceedings of the 30th International Conference on Coastal Engineering*, pages 73–85, Sam Diego, CA, USA, 2006.
- [11] A.-P. Engsig-Karup. *Unstructured Nodal DG-FEM Solution of High-order Boussinesq-type equations*. PhD thesis, Technical University of Denmark, 2006.
- [12] K. S. Erduran. Further application of hybrid solution to another form of Boussinesq equations and comparisons. *Int. J. Numer. Methods Fluids*, 53:827–849, 2007.
- [13] C. Eskilsson and S. J. Sherwin. Spectral/*hp* discontinuous Galerkin methods for modelling 2D Boussinesq equations. *J. Comp. Phys.*, 210:566, 2006.
- [14] D. R. Fuhrman and B. A. Bingham. Numerical solutions of fully non-linear and highly dispersive Boussinesq equations in two horizontal dimensions. *Int. J. Numer. Meth. Fluids*, 44:231–255, 2004.
- [15] D. R. Fuhrman and P. A. Madsen. Simulation of nonlinear wave run-up with a high-order Boussinesq model. *Coast. Eng.*, 55:139–154, 2008.

- [16] A. George and J. W. H Liu. *Computer solution of Large Sparce Positive Definite Systems*. Prentice Hall: Englewood Cliffs, N.J., 1981.
- [17] M. F. Gobbi, J. T. Kirby, and G. Wei. A fully non-linear Boussinesq model for surface waves. Part 2. Extension to  $o(kh^4)$ . *Journal of Fluid Mechanics*, 405:181, 2000.
- [18] M. Kazolea, A. I. Delis, I. A Nikolos, and C. E. Synolakis. An unstructured finite volume numerical scheme for extended 2D Boussinesq-type equations. *Coastal Eng.*, 69:42–66, 2012.
- [19] M. Kazolea, A. I. Delis, and C. E. Synolakis. Numerical treatment of wave breaking on unstructured finite volume approximations for extended Boussinesq-type equations. *J. Comp. Phys.*, 271:281–305, 2014.
- [20] P. Lynett and P. L. F Liu. A two-layer approach to wave modeling. In *Proceedings of the Royal Society of London Series A* 460, 2637, 2004.
- [21] P. A. Madsen and O. R. Sørensen. A new form of the Boussinesq equations with improved linear dispersion characteristics. Part 2: A slowing varying bathymetry. *Coast. Eng.*, 18:183–204, 1992.
- [22] P.A. Madsen, O. R. Sørensen, and H. A. Schäffer. Surf zone dynamics simulated by a Boussinesq-type model: Part I. Model description and cross-shore motion of regular waves. *Coast. Eng.*, 32:255–288, 1997a.
- [23] O. Nwogu. An alternative form of the Boussinesq equations for nearshore wave propagation. *Journal of Waterway, Port, Coastal, and Ocean Engineering*, 119:618–638, 1994.
- [24] D. H. Peregrine. Long waves on a beach. *J. Fluid Mech.*, 27:815–882, 1967.
- [25] M. Ricchiuto and Filippini A. G. Upwind residual discretization of enhanced Boussinesq equations for wave propagation over complex bathymetries. *J. Comp. Phys.*, 271:306–341, 2014.
- [26] P. L. Roe. Approximate Riemann solvers, parameter vectors, and difference schemes. *J. Comp. Phys.*, 43:357–372, 1981.
- [27] V. Roeber, K. F. Cheung, and M. H. Kobayashi. Shock-capturing Boussinesq-type model for nearshore wave processes. *Coast. Eng.*, 57:407–423, 2010.
- [28] Y. Saad. SPARSIT: A basic tool kit for sparse matrix computations, version 2. Available: [www.cs.umn.edu/research/arpa/SPARSKIT/sparsekit.html](http://www.cs.umn.edu/research/arpa/SPARSKIT/sparsekit.html), 1994.
- [29] Y. Saad. *Iterative Methods for Sparse Linear Systems*. PWS, 1996.
- [30] F. Shi, J. T. Kirby, J. C. Harris, J. D. Geiman, and S. T. Grilli. A high-order adaptive time-stepping tvd solver for boussinesq modeling of breaking waves and coastal inundation. *Ocean Modelling*, 43-44:36–51, 2012.
- [31] T. M. Smith, M. F. Barone, and R. B. Bond. Comparison of reconstruction techniques for unstructured mesh vertex centered finite volume schemes. *18th AIAA Computational Fluid Dynamics Conference*, pages 1–22, 2007.

- 
- [32] R. J. Spiteri and S. J. Ruuth. A new class of optimal high-order strong-stability-preserving time discretization methods. *SIAM J. Numer. Anal.*, 40:469, 2002.
- [33] M. Tonelli and M. Petti. Hybrid finite-volume finite-difference scheme for 2DH improved Boussinesq equations. *Coast. Eng.*, 56:609–620, 2009.
- [34] B. van Leer. Towards the ultimate conservative difference scheme V. A second order sequel to Godunov’s method. *J. Comp. Phys.*, 32:101, 1979.
- [35] M. Walkey and M. Berzins. A finite element method for the two-dimensional extended Boussinesq equations. *Int. J. Numer. Meth. Fluids*, 39:865, 2002.
- [36] G. Wei and J. T. Kirby. A time-dependent numerical code for extended Boussinesq equations. *Journal of Waterway, Port, Coastal, and Ocean Engineering*, 120:251–261, 1995.



**RESEARCH CENTRE  
BORDEAUX – SUD-OUEST**

351, Cours de la Libération  
Bâtiment A 29  
33405 Talence Cedex

Publisher  
Inria  
Domaine de Voluceau - Rocquencourt  
BP 105 - 78153 Le Chesnay Cedex  
[inria.fr](http://inria.fr)

ISSN 0249-6399