

CONSERVATIVE MULTIDIMENSIONAL UPWINDING FOR THE STEADY TWO-DIMENSIONAL SHALLOW WATER EQUATIONS

M.E.Hubbard and M.J.Baines

The University of Reading

Department of Mathematics

P.O.Box 220

Whiteknights

Reading

RG6 6AX

United Kingdom

Key words:

shallow water equations, multidimensional upwinding, conservation,
fluctuation distribution, system decomposition, source terms

Subject classification:

35L65, 65M99, 76M25

Proposed running head:

Multidimensional upwinding for shallow water

Address for proofs:

Dr. M.E.Hubbard

The University of Reading

Department of Mathematics

P.O.Box 220

Whiteknights

Reading

RG6 6AX

United Kingdom

Phone: +44 (0)118 9875123 x4009

Fax: +44 (0)118 9313423

Email: M.E.Hubbard@reading.ac.uk

Abstract

In recent years upwind differencing has gained acceptance as a robust and accurate technique for the numerical approximation of the one-dimensional shallow water equations. In two dimensions the benefits have been less marked due to the reliance of the methods on standard operator splitting techniques. Two conservative genuinely multidimensional upwind schemes are presented which have been adapted from flux balance distribution methods recently proposed for the approximation of steady state solutions of the Euler equations on unstructured triangular grids. A method for dealing with source terms, such as those introduced by modelling bed slope and friction, is also suggested and results are presented for two-dimensional steady state channel flows to illustrate the accuracy and robustness of the new algorithms.

1 Introduction

In recent years, many advances have been made in the numerical solution of hyperbolic systems of conservation laws in one and more dimensions [14, 12, 1]. Of particular interest has been the prediction of discontinuous solutions to the equations, which can occur when the system is nonlinear.

In the case of the numerical solution of the shallow water equations traditional methods, such as those of Preissmann, Abbott [5] and McCormack [9] rely on central differencing and are well known to require special treatment before a realistic numerical approximation of discontinuous flows can be obtained. More recently, the concept of upwinding has been adopted from the field of gas dynamics for the modelling of shallow water flows [11, 3]. This has proved to be highly successful, particularly in one dimension, in which high order upwind schemes have been constructed which capture discontinuities sharply and smoothly. This is achieved without the addition of artificial viscosity which is normally required to stabilise central difference schemes in the vicinity of high flow gradients. Furthermore, the upwind discretisation arises naturally from the physical interpretation of hyperbolic systems of equations, also giving a framework in which boundary conditions can be applied easily. The upwinding approach is therefore ideal for the modelling of transcritical and supercritical flows.

The practical advantages of upwind schemes in higher dimensions are less clear. Historically, they have been applied to the two-dimensional shallow water equations via the use of standard operator splitting techniques, *e.g.* [3], which by implication involves the application of one-dimensional methods to a multi-dimensional system of equations, albeit in two independent directions. Recently

though, multidimensional upwind schemes have been developed for the numerical solution of the Euler equations of gas dynamics [24, 18, 21, 16] which can equally well be applied to the shallow water equations. These schemes are conservative, upwind, cell-vertex finite volume methods based on the concept of fluctuation distribution [20] and are applied on unstructured triangular grids. They differ from most standard finite volume schemes in that the underlying representation of the flow is not piecewise constant within each grid cell, as is usual, but continuous piecewise linear with the unknowns stored at the nodes of the grid, more akin to a standard finite element approximation. The resulting schemes are designed to mimic the evolution of the approximate solution within each triangular grid cell, whereas previous upwind methods concentrated on the Riemann problems arising at the discontinuities in the approximation (*e.g.* at cell edges), an inherently one-dimensional process.

Multidimensional upwind schemes are constructed from three distinct elements: a consistent, conservative linearisation of the system of equations, the decomposition of the resulting discrete system into simple (mainly scalar) components, and the subsequent evolution of the decomposed system using scalar and matrix fluctuation distribution schemes. In this paper attention is focussed on the linearisation and decomposition stages of the algorithm, both of which differ slightly from those devised for the Euler equations [18, 16]. The distribution of the components is also described but this step remains unchanged from the solution of other systems and further details can be found in [8].

Multidimensional upwind schemes for the solution of the shallow water equations have already appeared in [10, 17] but those schemes were not conservative.

In this paper a conservative formulation is presented, together with two alternative decompositions of the system of shallow water equations and a method of incorporating source terms such as those arising from the consideration of bed slope and friction. Results are presented to illustrate the quality of the numerical solutions obtained for steady state problems.

2 The Governing Equations

The shallow water equations can be used to describe the motion of ‘shallow’ free-surface flows subject to gravitational forces. The system can be obtained from the depth-averaged Navier-Stokes equations and the resulting homogeneous system represents the conservation of mass and momentum in the flow. The structure of the system is very similar to that of the Euler equations. The effects of bed slope and friction on the flow are modelled by the inclusion of source terms on the right hand side of the system which modify the momentum conservation equations.

In conservation form the unsteady shallow water equations [5], which are used for the construction of multidimensional upwind schemes even when steady state solutions are sought, are given by

$$\underline{\mathbf{U}}_t + \underline{\mathbf{F}}_x + \underline{\mathbf{G}}_y = \underline{\mathbf{q}}, \quad (2.1)$$

in which

$$\underline{\mathbf{U}} = \begin{pmatrix} h \\ hu \\ hv \end{pmatrix} \quad (2.2)$$

are the conservative variables and the corresponding flux vectors are

$$\underline{\mathbf{F}} = \begin{pmatrix} hu \\ hu^2 + \frac{gh^2}{2} \\ huv \end{pmatrix} \quad \text{and} \quad \underline{\mathbf{G}} = \begin{pmatrix} hv \\ huv \\ hv^2 + \frac{gh^2}{2} \end{pmatrix}. \quad (2.3)$$

The source terms can be written

$$\underline{\mathbf{q}} = \begin{pmatrix} 0 \\ gh(s_{\text{bx}} - s_{\text{fx}}) \\ gh(s_{\text{by}} - s_{\text{fy}}) \end{pmatrix}, \quad (2.4)$$

where the bed slope terms are defined by

$$s_{\text{bx}} = -\frac{\partial z}{\partial x} \quad \text{and} \quad s_{\text{by}} = -\frac{\partial z}{\partial y}, \quad (2.5)$$

in which z is the height of the bed above some nominal zero level, and the friction slopes are given by Manning's formula,

$$s_{\text{fx}} = \frac{n^2 u \sqrt{u^2 + v^2}}{h^{\frac{4}{3}}}, \quad s_{\text{fy}} = \frac{n^2 v \sqrt{u^2 + v^2}}{h^{\frac{4}{3}}}, \quad (2.6)$$

in which n is Manning's roughness coefficient.

The quasilinear form of the system (2.1) will also be required later. This is given by the equations

$$\underline{\mathbf{U}}_t + \mathbf{A}_U \underline{\mathbf{U}}_x + \mathbf{B}_U \underline{\mathbf{U}}_y = \underline{\mathbf{q}}, \quad (2.7)$$

where the conservative flux Jacobians are

$$\mathbf{A}_U = \frac{\partial \underline{\mathbf{F}}}{\partial \underline{\mathbf{U}}} = \begin{pmatrix} 0 & 1 & 0 \\ -u^2 + c^2 & 2u & 0 \\ -uv & v & u \end{pmatrix} \quad (2.8)$$

and

$$\mathbf{B}_U = \frac{\partial \mathbf{G}}{\partial \mathbf{U}} = \begin{pmatrix} 0 & 0 & 1 \\ -uv & v & u \\ -v^2 + c^2 & 0 & 2v \end{pmatrix}, \quad (2.9)$$

in which $c = \sqrt{gh}$ is the gravity wave speed or wave celerity. Further details about the mathematical aspects of the shallow water equations can be found in [23].

3 A Conservative Linearisation

An appropriate linearisation of the shallow water equations is required so that the decomposition and distribution stages of the algorithm give rise to a conservative scheme. Many different conservative linearisations have been constructed for the Euler equations, see for example [7, 2], but it is the most robust of these, based on Roe's one-dimensional linearisation using a set of parameter vector variables [19], which is generally used for practical calculations. This linearisation is adapted here to give an analogous discrete form of the shallow water equations.

Consider the two-dimensional homogeneous system,

$$\underline{\mathbf{U}}_t + \underline{\mathbf{F}}_x + \underline{\mathbf{G}}_y = \mathbf{0}, \quad (3.1)$$

in which the conservative variables $\underline{\mathbf{U}}$ and fluxes $\underline{\mathbf{F}}, \underline{\mathbf{G}}$ are given by (2.2) and (2.3) respectively. For a given cell in a triangular discretisation of the computational domain the flux balance is defined by

$$\begin{aligned} \underline{\Phi}_U &= - \iint_{\Delta} (\underline{\mathbf{F}}_x + \underline{\mathbf{G}}_y) \, dx \, dy \\ &= \oint_{\partial\Delta} (\underline{\mathbf{F}}, \underline{\mathbf{G}}) \cdot d\vec{n}, \end{aligned} \quad (3.2)$$

in which $d\vec{n}$ represents the inward pointing normal to the cell boundary. The numerical approximation to $\underline{\Phi}_U$ is defined to be of the form

$$\begin{aligned}\widehat{\underline{\Phi}}_U &= -S_\Delta (\widehat{\underline{F}}_x + \widehat{\underline{G}}_y) \\ &= -S_\Delta (\widehat{\underline{A}}_U \widehat{\underline{U}}_x + \widehat{\underline{B}}_U \widehat{\underline{U}}_y) ,\end{aligned}\tag{3.3}$$

where S_Δ is the cell area and $\widehat{}$ indicates a discretised quantity. The precise definition of the discrete form of the flux balance will be described below.

Multidimensional upwind schemes update flow variables stored at the nodes of the grid via the distribution of the discrete form of the flux balance, $\widehat{\underline{\Phi}}_U$ of (3.3), within each cell. Conservation requires that the overall contribution to the nodes depends only on boundary conditions, so for a linearisation such as that represented by (3.3) to be conservative the sum over the whole flow domain of the $\widehat{\underline{\Phi}}_U$ should reduce to boundary contributions alone. It follows immediately from (3.2) that a linearisation is conservative if $\widehat{\underline{\Phi}}_U = \underline{\Phi}_U$ for each grid cell, and the resulting scheme is conservative as long as the whole of each discrete flux balance is distributed to the nodes of the grid.

In keeping with the linearisation of the Euler equations, the discrete flux Jacobians in (3.3) are sought in a form which allows $\widehat{\underline{\Phi}}_U$ to be readily decomposed by the methods described in Section 4 below, *i.e.* the Jacobians are evaluated consistently from some cell-average state, $\bar{\underline{z}}$ say, so that

$$\widehat{\underline{A}}_U = \left(\frac{\partial \underline{F}}{\partial \underline{U}} \right)_{\bar{\underline{z}}} \quad \text{and} \quad \widehat{\underline{B}}_U = \left(\frac{\partial \underline{G}}{\partial \underline{U}} \right)_{\bar{\underline{z}}} .\tag{3.4}$$

The construction of a conservative linearisation of this form is aided considerably

by assuming that the components of the parameter vector

$$\underline{Z} = \sqrt{h} \begin{pmatrix} 1 \\ u \\ v \end{pmatrix} \quad (3.5)$$

vary linearly in space within each cell, *cf.* Roe's parameter vector for the Euler equations [19]. A consequence is that $\vec{\nabla} \underline{Z}$ is locally constant and so the conservative flux balance can be written as

$$\begin{aligned} \Phi_U &= - \iint_{\Delta} (\underline{E}_x + \underline{G}_y) \, dx \, dy \\ &= - \left(\iint_{\Delta} \frac{\partial \underline{F}}{\partial \underline{Z}} \, dx \, dy \right) \underline{Z}_x - \left(\iint_{\Delta} \frac{\partial \underline{G}}{\partial \underline{Z}} \, dx \, dy \right) \underline{Z}_y, \end{aligned} \quad (3.6)$$

in which

$$\frac{\partial \underline{F}}{\partial \underline{Z}} = \begin{pmatrix} \sqrt{hu} & \sqrt{h} & 0 \\ 2g(\sqrt{h})^3 & 2\sqrt{hu} & 0 \\ 0 & \sqrt{hv} & \sqrt{hu} \end{pmatrix} \quad (3.7)$$

and

$$\frac{\partial \underline{G}}{\partial \underline{Z}} = \begin{pmatrix} \sqrt{hv} & 0 & \sqrt{h} \\ 0 & \sqrt{hv} & \sqrt{hu} \\ 2g(\sqrt{h})^3 & 0 & 2\sqrt{hv} \end{pmatrix}. \quad (3.8)$$

Note that only two entries in the above Jacobian matrices - the (2,1) entry in $\frac{\partial \underline{F}}{\partial \underline{Z}}$ and the (3,1) entry in $\frac{\partial \underline{G}}{\partial \underline{Z}}$ - are not linear in the components of the parameter vector \underline{Z} .

In the case of the Euler equations it is possible to choose \underline{Z} so that each entry in the corresponding Jacobian matrices is a linear function of its components [19]. Hence, the integrals in (3.6) can be evaluated exactly in terms of a single cell-averaged value of \underline{Z} and this leads to a conservative linearisation satisfying

(3.4) [7]. The nonlinear terms in (3.7) and (3.8) mean that the linearisation of the shallow water equations cannot be constructed in precisely the same manner. In previous work [10, 17] non-conservative linearisations have been used, in which the flux balance (3.3) is evaluated consistently from an appropriate average state, but in the present work a conservative form is sought.

A conservative linearisation of the shallow water equations is achieved by evaluating the integrals in (3.6) exactly. This does not immediately give rise to linearised flux Jacobians of the form (3.4), so instead a component of (3.6) is isolated which does have this form and which therefore can be decomposed using the second stage of the algorithm, described in Section 4. Hence the numerical flux balance (3.3) is split into two parts, taking the form

$$\widehat{\Phi}_U = - \underbrace{S_\Delta \left(\frac{\partial \overline{\mathbf{F}}}{\partial \underline{\mathbf{Z}}} \overline{\underline{\mathbf{Z}}}_x + \frac{\partial \overline{\mathbf{G}}}{\partial \underline{\mathbf{Z}}} \overline{\underline{\mathbf{Z}}}_y \right)}_{(1)} - \underbrace{S_\Delta \left(\mathbf{S}_Z \overline{\underline{\mathbf{Z}}}_x + \mathbf{T}_Z \overline{\underline{\mathbf{Z}}}_y \right)}_{(2)}. \quad (3.9)$$

The overbar indicates the consistent evaluation of a quantity solely from the cell-average state given by

$$\overline{\underline{\mathbf{Z}}} = \frac{1}{3} \sum_{i=1}^3 \underline{\mathbf{Z}}_i, \quad (3.10)$$

as well as the corresponding discrete gradient (evaluated under the assumption of linearly varying $\underline{\mathbf{Z}}$)

$$\overline{\nabla \underline{\mathbf{Z}}} = \frac{1}{2S_\Delta} \sum_{i=1}^3 \underline{\mathbf{Z}}_i \vec{n}_i, \quad (3.11)$$

in which $\underline{\mathbf{Z}}_i$ are the values of the parameter vector variables at the vertices of the cell and \vec{n}_i is the inward pointing normal to the edge opposite vertex i scaled by the length of that edge.

The flux balance (3.9) can also be written simply in terms of the conservative

variables since

$$\frac{\partial \underline{\mathbf{U}}}{\partial \underline{\mathbf{Z}}} = \begin{pmatrix} 2\sqrt{h} & 0 & 0 \\ \sqrt{hu} & \sqrt{h} & 0 \\ \sqrt{hv} & 0 & \sqrt{h} \end{pmatrix} \quad (3.12)$$

is linear in the components of $\underline{\mathbf{Z}}$. It then follows that the discrete gradient of the conservative variables can be written

$$\begin{aligned} \overline{\nabla \underline{\mathbf{U}}} &= \frac{1}{S_\Delta} \iint_\Delta \frac{\partial \underline{\mathbf{U}}}{\partial \underline{\mathbf{Z}}} \overline{\nabla \underline{\mathbf{Z}}} \, dx \, dy \\ &= \frac{\overline{\partial \underline{\mathbf{U}}}}{\overline{\partial \underline{\mathbf{Z}}}} \overline{\nabla \underline{\mathbf{Z}}} \end{aligned} \quad (3.13)$$

and that, from (3.9), the discrete conservative flux balance (3.3) is given by

$$\begin{aligned} \widehat{\underline{\Phi}}_{\mathbf{U}} &= -S_\Delta \left(\frac{\overline{\partial \mathbf{F}}}{\partial \underline{\mathbf{U}}} \underline{\mathbf{U}}_x + \frac{\overline{\partial \mathbf{G}}}{\partial \underline{\mathbf{U}}} \underline{\mathbf{U}}_y \right) - S_\Delta (\mathbf{S}_{\mathbf{U}} \underline{\mathbf{U}}_x + \mathbf{T}_{\mathbf{U}} \underline{\mathbf{U}}_y) \\ &= \overline{\underline{\Phi}}_{\mathbf{U}} + \underline{\mathbf{q}}_{\mathbf{U}} \end{aligned} \quad (3.14)$$

in which

$$\underline{\mathbf{q}}_{\mathbf{U}} = -S_\Delta (\mathbf{S}_{\mathbf{U}} \underline{\mathbf{U}}_x + \mathbf{T}_{\mathbf{U}} \underline{\mathbf{U}}_y) . \quad (3.15)$$

Thus, the two components of the flux balance (3.9) reveal themselves to be (1) $\overline{\underline{\Phi}}_{\mathbf{U}}$ evaluated at the cell-average state (3.10) and (2) a small ‘source’ term $\underline{\mathbf{q}}_{\mathbf{U}}$.

The source term arising from the linearisation above can be evaluated quickly since the ‘Jacobians’ are the simple sparse matrices

$$\mathbf{S}_{\mathbf{U}} = \frac{\overline{\partial \underline{\mathbf{Z}}}}{\overline{\partial \underline{\mathbf{U}}}} \mathbf{S}_{\mathbf{Z}} = \frac{g\zeta}{\sqrt{h}} \begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad (3.16)$$

and

$$\mathbf{T}_{\mathbf{U}} = \frac{\overline{\partial \underline{\mathbf{Z}}}}{\overline{\partial \underline{\mathbf{U}}}} \mathbf{T}_{\mathbf{Z}} = \frac{g\zeta}{\sqrt{h}} \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix} , \quad (3.17)$$

where

$$\zeta = \frac{1}{S_{\Delta}} \iint_{\Delta} (\sqrt{h})^3 \, dx \, dy - (\overline{\sqrt{h}})^3 . \quad (3.18)$$

These are derived from the difference between integrating (3.6) exactly and approximating it using a one point quadrature rule. The size of the source \underline{q}_U is therefore negligible in smooth flow but may have an effect at discontinuities.

In the following section the decomposition of $\overline{\Phi}_U$ in (3.14) is described. Source terms arising from both the linearisation and the consideration of bed slope and friction are considered separately in Section 6.

4 Two Optimal Decompositions

The decomposition stage of the algorithm dictates how the first component of the flux balance, $\overline{\Phi}_U$ of (3.14), within each triangle of the grid is divided up into simple elements. In one dimension a complete decoupling of the system into scalar components is uniquely available through the transformation of the equations into characteristic variables, but in higher dimensions there are many possible ways to decompose the system.

The most successful decompositions of the Euler equations utilise a preconditioned form of the equations. Two preconditioners have been developed in the literature [18, 16], both of which are generalisations of the van Leer-Lee-Roe matrix of [25]. This was originally constructed to reduce the condition number of the system and accelerate convergence of the numerical algorithm to the steady state but here it is used to facilitate the construction of an optimal decomposition of the system in the sense that the equations of the decomposition are maximal-

ly decoupled. The corresponding analysis of the shallow water equations closely follows that of [18, 16, 25] and the resulting preconditioners are described here.

For the sake of simplifying the algebra, the homogeneous part of the system (2.1) is considered in terms of the streamwise variables, ξ and η , and the symmetrising variables \underline{Q} , defined by

$$\partial \underline{Q} = \begin{pmatrix} \frac{c}{h} \partial h \\ \partial q \\ q \partial \theta \end{pmatrix}, \quad (4.1)$$

where $q = \sqrt{u^2 + v^2}$ is the speed of the flow and $\theta = \tan^{-1} \left(\frac{v}{u} \right)$ its direction. The symmetrised form of the shallow water equations are now preconditioned by a matrix \mathbf{P} , and the resulting system written in the form

$$\underline{Q}_t + \mathbf{P} \left(\mathbf{A}_Q^S \underline{Q}_\xi + \mathbf{B}_Q^S \underline{Q}_\eta \right) = \underline{0}, \quad (4.2)$$

in which the new Jacobians are the simple symmetric matrices

$$\mathbf{A}_Q^S = \begin{pmatrix} q & c & 0 \\ c & q & 0 \\ 0 & 0 & q \end{pmatrix} \quad \text{and} \quad \mathbf{B}_Q^S = \begin{pmatrix} 0 & 0 & c \\ 0 & 0 & 0 \\ c & 0 & 0 \end{pmatrix}. \quad (4.3)$$

The superscript (and later subscript) S indicates that the streamwise coordinate system is being used.

4.1 Decomposition 1 (HELW)

Following the analysis of Mesaros and Roe for the Euler equations [16], the first preconditioning matrix suggested here is given by

$$\mathbf{P} = \frac{1}{q} \begin{pmatrix} \frac{\varepsilon F^2}{\beta \kappa} & -\frac{\varepsilon F}{\beta \kappa} & 0 \\ -\frac{\varepsilon F}{\beta \kappa} & \frac{\varepsilon}{\beta \kappa} + \varepsilon & 0 \\ 0 & 0 & \frac{\beta}{\kappa} \end{pmatrix}, \quad (4.4)$$

where $F = \frac{q}{c}$ is the local Froude number of the flow,

$$\beta = \sqrt{|F^2 - 1|}, \quad \kappa = \max(F, 1) \quad (4.5)$$

and ε is a function of the Froude number such that $\varepsilon(0) = \frac{1}{2}$ and $\varepsilon(F) = 1$ for $F \geq 1$. These restrictions on ε ensure that the decomposition is not sensitive to the flow angle in the limit as $F \rightarrow 0$ [26] and that the transition of the preconditioner through the transcritical region is smooth. Here, as in [26] ε is taken to be the C_1 function

$$\varepsilon(F) = \begin{cases} \frac{1}{2} & \text{for } F \leq \frac{1}{3} \\ -27F^3 + \frac{81}{2}F^2 - 18F + 3 & \text{for } \frac{1}{3} < F < \frac{2}{3} \\ 1 & \text{for } F \geq \frac{2}{3} \end{cases} \quad (4.6)$$

so that the first derivative also varies smoothly. The matrix \mathbf{P} in (4.4) is, in fact, precisely that of [26] with the Mach number replaced by the Froude number and without the involvement of the entropy equation. The variable κ has simply been introduced so that (4.4) is correct for both subcritical and supercritical flow.

The preconditioned system (4.2) is decomposed by transforming it into a set of characteristic equations,

$$\underline{\mathbf{W}}_t + \mathbf{A}_W^S \underline{\mathbf{W}}_\xi + \mathbf{B}_W^S \underline{\mathbf{W}}_\eta = \underline{\mathbf{0}}, \quad (4.7)$$

where the characteristic variables \underline{W} are defined by

$$\partial \underline{W}_{\text{sb}} = \begin{pmatrix} \frac{q\beta}{q} \partial h \\ q \partial \theta \\ \frac{q}{c} \partial h + F \partial q \end{pmatrix} \quad \text{and} \quad \partial \underline{W}_{\text{sp}} = \begin{pmatrix} \frac{q\beta}{c} \partial h + Fq \partial \theta \\ \frac{q\beta}{c} \partial h - Fq \partial \theta \\ \frac{q}{c} \partial h + F \partial q \end{pmatrix}, \quad (4.8)$$

for subcritical and supercritical flow respectively. The corresponding transformation matrices are given by

$$\frac{\partial \underline{W}_{\text{sb}}}{\partial \underline{Q}} = \begin{pmatrix} \frac{\beta}{F} & 0 & 0 \\ 0 & 0 & 1 \\ 1 & F & 0 \end{pmatrix} \quad \text{and} \quad \frac{\partial \underline{W}_{\text{sp}}}{\partial \underline{Q}} = \begin{pmatrix} \beta & 0 & F \\ \beta & 0 & -F \\ 1 & F & 0 \end{pmatrix} \quad (4.9)$$

and their inverses, and the resulting characteristic flux Jacobian matrices can be calculated easily using

$$\mathbf{A}_{\underline{W}}^{\text{S}} = \frac{\partial \underline{W}}{\partial \underline{Q}} \mathbf{P} \mathbf{A}_{\underline{Q}}^{\text{S}} \frac{\partial \underline{Q}}{\partial \underline{W}} \quad \text{and} \quad \mathbf{B}_{\underline{W}}^{\text{S}} = \frac{\partial \underline{W}}{\partial \underline{Q}} \mathbf{P} \mathbf{B}_{\underline{Q}}^{\text{S}} \frac{\partial \underline{Q}}{\partial \underline{W}} \quad (4.10)$$

for both subcritical and supercritical flows.

Note that the choice of variables given by (4.8) changes across the transcritical region. When the flow is supercritical the steady equations are hyperbolic and the choice of variables defined by $\partial \underline{W}_{\text{sp}}$ in (4.8) uniquely leads to the system being completely decoupled into scalar components. However, in the subcritical case only one equation can be decoupled, leaving a second component which manifests itself as a 2×2 elliptic subsystem, the form of which depends on the choice of characteristic variables, defined here by $\partial \underline{W}_{\text{sb}}$ in (4.8). The shallow water equations cannot be decoupled further in subcritical flow.

The complete decoupling of the equations in supercritical flow allows the sys-

tem (4.7) to be written in the form of three scalar advection equations, *i.e.*

$$W_t^k + \vec{\lambda}_S^k \cdot \vec{\nabla}_S W^k = 0, \quad k = 1, 2, 3, \quad (4.11)$$

in which the advection velocities in the streamwise coordinate system are

$$\vec{\lambda}_S^1 = \begin{pmatrix} \frac{\beta}{F} \\ \frac{1}{F} \end{pmatrix}_S, \quad \vec{\lambda}_S^2 = \begin{pmatrix} \frac{\beta}{F} \\ -\frac{1}{F} \end{pmatrix}_S \quad \text{and} \quad \vec{\lambda}_S^3 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}_S. \quad (4.12)$$

Hence the first component of the flux balance in (3.14) takes the form

$$\overline{\Phi_U} = -S_\Delta \sum_{k=1}^3 (\vec{\lambda}_S^k \cdot \vec{\nabla}_S W^k) \underline{\mathbf{r}}_U^k, \quad (4.13)$$

where every term on the right hand side of (4.13) is evaluated consistently from the cell-average state defined by (3.10) and (3.11), and $\underline{\mathbf{r}}_U^k$ is the k^{th} column of the matrix

$$\mathbf{R}_U = \frac{\partial \underline{\mathbf{U}}}{\partial \underline{\mathbf{Q}}} \mathbf{P}^{-1} \frac{\partial \underline{\mathbf{Q}}}{\partial \underline{\mathbf{W}}}. \quad (4.14)$$

This matrix transforms the components of the flux balance corresponding to the characteristic equations back into components of the conservative flux balance. Hence (4.13) represents a consistent decomposition of $\overline{\Phi_U}$ of (3.14), the components of which may each be distributed using a simple *scalar* scheme such as that described in Section 5.1 below.

In the case of subcritical flow the choice of characteristic variables defined by $\partial \underline{\mathbf{W}}_{\text{sb}}$ in (4.8) leads to Jacobian matrices in the system (4.7) given by

$$\mathbf{A}_W^S = \begin{pmatrix} -\varepsilon\beta & 0 & 0 \\ 0 & \beta & 0 \\ 0 & 0 & \varepsilon \end{pmatrix} \quad \text{and} \quad \mathbf{B}_W^S = \begin{pmatrix} 0 & \varepsilon & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}. \quad (4.15)$$

Hence the characteristic equations take the form of a single scalar advection equation, which is precisely the same as the $k = 3$ equation defined by (4.11) and (4.12), together with a 2×2 elliptic subsystem, so $\overline{\Phi_U}$ of (3.14) is written

$$\begin{aligned} \overline{\Phi_U} = & -S_\Delta (\underline{\mathbf{r}}_U^1, \underline{\mathbf{r}}_U^2) \left[\begin{aligned} & \begin{pmatrix} -\varepsilon\beta & 0 \\ 0 & \beta \end{pmatrix} \begin{pmatrix} W^1 \\ W^2 \end{pmatrix}_\xi + \begin{pmatrix} 0 & \varepsilon \\ 1 & 0 \end{pmatrix} \begin{pmatrix} W^1 \\ W^2 \end{pmatrix}_\eta \right] \\ & - S_\Delta (\vec{\lambda}_S^3 \cdot \vec{\nabla}_S W^3) \underline{\mathbf{r}}_U^3, \end{aligned} \quad (4.16) \end{aligned}$$

cf. the supercritical decomposition (4.13). The scalar component is treated as it was in the supercritical case but the elliptic nature of the second term suggests that an alternative to the upwind distribution scheme of Section 5.1 should be sought. Consequently the central distribution scheme proposed in [16], which involves no further decomposition of the elliptic component, is used and this is described in Section 5.2. Note though, that the choice results in a discontinuity in the distribution at the transcritical line and a loss of positivity in the subcritical region, both of which are detrimental to the robustness of the scheme. This scheme will be denoted ‘HELW’ due to the Hyperbolic/Elliptic nature of the subcritical decomposition and the Lax-Wendroff style distribution of the resulting elliptic component.

4.2 Decomposition 2 (HESUPG)

Another preconditioner which leads to a maximal decoupling of the shallow water equations corresponds to that developed by Paillère *et al.* for the Euler equations

[18] and takes the form

$$\mathbf{P} = \frac{1}{q} \begin{pmatrix} \frac{\chi F^2}{\beta_\epsilon^2} & -\frac{\chi F}{\beta_\epsilon^2} & 0 \\ -\frac{\chi F}{\beta_\epsilon^2} & \frac{\chi}{\beta_\epsilon^2} + 1 & 0 \\ 0 & 0 & \chi \end{pmatrix}, \quad (4.17)$$

where

$$\beta_\epsilon = \sqrt{\max(\epsilon^2, |F^2 - 1|)}, \quad \chi = \frac{\beta_\epsilon}{\max(F, 1)} \quad (4.18)$$

and ϵ is a nonzero constant which typically takes a value of 0.05. This matrix is again derived by following the analysis of the Euler equations [25], the result being that the Mach number is replaced by the Froude number and the entropy equation disappears.

The decoupling of the system proceeds as in the previous decomposition, leading to a set of characteristic equations (4.7) in new variables $\underline{\mathbf{W}}$, now defined by

$$\partial \underline{\mathbf{W}} = \begin{pmatrix} \frac{g\beta_\epsilon}{c} \partial h + Fq \partial \theta \\ \frac{g\beta_\epsilon}{c} \partial h - Fq \partial \theta \\ \frac{g}{c} \partial h + F \partial q \end{pmatrix}, \quad (4.19)$$

independent of the flow speed, *cf.* (4.8). The corresponding transformation matrix is given by $\frac{\partial \underline{\mathbf{W}}}{\partial \underline{\mathbf{Q}}}$ in (4.9).

The difference between the two decompositions lies in the treatment of the system for subcritical and transcritical flows ($F^2 \leq 1 + \epsilon^2$). The decision to keep the same characteristic variables in both subcritical and supercritical flow leads to Jacobian matrices in the transformed system (4.7) which are given by

$$\mathbf{A}_W^S = \begin{pmatrix} \chi\nu^+ & \chi\nu^- & 0 \\ \chi\nu^- & \chi\nu^+ & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad \text{and} \quad \mathbf{B}_W^S = \begin{pmatrix} \frac{\chi}{\beta_\epsilon} & 0 & 0 \\ 0 & -\frac{\chi}{\beta_\epsilon} & 0 \\ 0 & 0 & 0 \end{pmatrix}, \quad (4.20)$$

where

$$\nu^+ = \frac{F^2 - 1 + \beta_\epsilon^2}{2\beta_\epsilon^2} \quad \text{and} \quad \nu^- = \frac{F^2 - 1 - \beta_\epsilon^2}{2\beta_\epsilon^2}. \quad (4.21)$$

It is easy to see that in the supercritical region $\nu^- = 0$, the system is completely decoupled, and the decomposition (and subsequent distribution) reduces to precisely that given for supercritical flow in Section 4.1.

In the subcritical case the system is again decomposed into a single, independent scalar component and a pair of coupled equations, but rather than regarding the latter as a 2×2 subsystem it is instead treated as in [18], as two separate scalar equations with source terms. As a consequence, the decomposition of $\overline{\Phi_U}$ of (3.14) takes the form

$$\overline{\Phi_U} = -S_\Delta \sum_{k=1}^3 \left(\vec{\lambda}_S^k \cdot \vec{\nabla}_S W^k + q_S^k \right) \underline{\mathbf{r}}_U^k, \quad (4.22)$$

in which $\underline{\mathbf{r}}_U^k$ is the k^{th} column of the matrix \mathbf{R}_U (4.14), newly defined from the \mathbf{P} of (4.17) and the \underline{W} of (4.19),

$$\vec{\lambda}_S^1 = \begin{pmatrix} \chi \nu^+ \\ \frac{\chi}{\beta_\epsilon} \end{pmatrix}_S, \quad \vec{\lambda}_S^2 = \begin{pmatrix} \chi \nu^+ \\ -\frac{\chi}{\beta_\epsilon} \end{pmatrix}_S \quad \text{and} \quad \vec{\lambda}_S^3 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}_S, \quad (4.23)$$

and

$$q_S^1 = \chi \nu^- W_\xi^2, \quad q_S^2 = \chi \nu^- W_\xi^1 \quad \text{and} \quad q_S^3 = 0. \quad (4.24)$$

The distribution of the decoupled component of this second decomposition is once again carried out using the scalar upwind scheme of Section 5.1 below. It is possible to use the same method for the coupled components, with an appropriate modification (described in Section 5.2) to ensure that the scheme remains linearity preserving [18] in the presence of source terms. However positivity is lost as a consequence, so when the Froude number is close to unity and the advection velocities associated with each component of the nonlinear system are

very closely aligned, the distribution provides very little cross-stream diffusion and as a result lacks robustness. The actual scalar scheme used here is the SUPG scheme suggested in [18] and described at the end of Section 5.2. As in Section 4.1 the favoured distribution of the flux balance is discontinuous for transcritical flows. This approach will be denoted ‘HESUPG’ due to the Hyperbolic/Elliptic subcritical decomposition and the distribution of the elliptic component using an SUPG-type scheme.

5 Flux Balance Distribution

5.1 Distribution of Decoupled Scalar Components

Each homogeneous scalar component which results from the above decomposition can be modelled by an advection equation,

$$u_t + f_x + g_y = 0 \quad \text{or} \quad u_t + \vec{\lambda} \cdot \vec{\nabla} u = 0, \quad (5.1)$$

where $\vec{\lambda} = \left(\frac{\partial f}{\partial u}, \frac{\partial g}{\partial u} \right)^T$ defines the advection velocity.

Multidimensional upwinding for the numerical solution of the scalar equation (5.1) involves the construction of a time-stepping scheme which calculates the fluctuation,

$$\begin{aligned} \phi &= - \int \int_{\Delta} \vec{\lambda} \cdot \vec{\nabla} u \, dx \, dy \\ &= \oint_{\partial\Delta} u \vec{\lambda} \cdot d\vec{n}, \end{aligned} \quad (5.2)$$

within each cell and updates the solution at each time level by adding fractions of this quantity to the nodal values of u (see [8]). In (5.2) $\partial\Delta$ is the boundary of the cell and $d\vec{n}$ represents the inward pointing normal to the boundary. Steady

state solutions of (5.1) are calculated by repeating this update iteratively in order to approximate the solution in the limit as $t \rightarrow \infty$.

The vector $\vec{\lambda}$ in (5.2) may not be constant, in which case a conservative linearisation of the scalar advection equation (5.1) can often be constructed by treating it as a special case of the system linearisation discussed in Section 3 [8]. A conservative linearisation of this form can be constructed if a variable, z say, can be defined so that when z varies linearly in space the derivatives $\frac{\partial u}{\partial z}$, $\frac{\partial f}{\partial z}$ and $\frac{\partial g}{\partial z}$ can all be integrated exactly over any triangle by quadrature.

Assuming such a z exists and does vary linearly in space within each cell, the discrete fluctuation in a cell can be written

$$\hat{\phi} = -S_{\Delta} \hat{\vec{\lambda}} \cdot \widehat{\vec{\nabla}} u, \quad (5.3)$$

cf. the discrete flux balance in (3.3) and the decomposition of (4.13). The discrete gradient is defined to be

$$\widehat{\vec{\nabla}} u = \frac{1}{S_{\Delta}} \iint_{\Delta} \frac{\partial u}{\partial z} dx dy \vec{\nabla} z, \quad (5.4)$$

where $\vec{\nabla} z$ is evaluated exactly as in (3.11), so that the linearised advection velocity $\hat{\vec{\lambda}}$ is defined by

$$\hat{\vec{\lambda}} = \frac{\iint_{\Delta} \frac{\partial \vec{f}}{\partial z} dx dy}{\iint_{\Delta} \frac{\partial u}{\partial z} dx dy}, \quad (5.5)$$

in which $\vec{f} = (f, g)^T$. The most important consequence of choosing z to have the above properties is that $\widehat{\vec{\nabla}} u$ and $\hat{\vec{\lambda}}$ can be calculated exactly, so $\hat{\phi} = \phi$ and the linearisation is conservative.

The distribution of the discrete fluctuation of (5.3) to the grid nodes combined with an explicit forward Euler discretisation of the time derivative in (5.1) leads

to a nodal update of the form

$$u_i^{n+1} = u_i^n + \frac{\Delta t}{S_i} \sum_{\cup \Delta_i} \alpha_i^j \phi_j, \quad (5.6)$$

where S_i is the area of the median dual cell for node i (one third of the total area of the triangles with a vertex at i), α_i^j is the distribution coefficient which indicates the proportion of the fluctuation ϕ_j to be sent from cell j to node i , and $\cup \Delta_i$ represents the set of cells with vertices at node i . It can be seen from the second expression for ϕ in (5.2) that the sum of the fluctuations over the whole of the grid reduces through internal cancellation to boundary contributions alone. Therefore a nodal update of the form (5.6) leads to a conservative scheme as long as the whole of each fluctuation is distributed to the grid nodes, *i.e.*

$$\sum_i \alpha_i^j = 1 \quad \forall j. \quad (5.7)$$

For the sake of simplicity and compactness, a cell is only allowed to contribute a proportion of its fluctuation to its own vertices. The distribution coefficients α_i^j are chosen so that the resulting scheme has the following four properties:

- Upwindedness - the fluctuation within a cell is only sent to the downstream vertices of that cell, *i.e.* vertices opposite inflow edges for which $\hat{\lambda} \cdot \vec{n} > 0$, where \vec{n} is the inward pointing normal to the edge.
- Positivity - every nodal value of u at the new time level in (5.6) is a convex combination of nodal values of u at the old time level, so the scheme cannot produce new extrema in the solution at the new time-step, spurious oscillations do not appear in the solution and the scheme is stable for an appropriate time-step restriction.

- Linearity preservation - the exact steady state solution is preserved when this varies linearly in space, so no update is sent to the nodes when a cell fluctuation is zero and the scheme is second order accurate at the steady state on a regular mesh with a uniform choice of diagonals [8].
- Continuity - the contributions to the nodes, $\alpha_i^j \phi_j$ (5.6), depend continuously on the data, avoiding limit cycling as convergence is approached and improving the robustness of the scheme.

Linearity preservation should also be satisfied by the decomposition, so that no update is sent to the vertices of a cell when its flux balance is zero and the higher order accuracy possessed by the linearity preserving scalar scheme is retained by the overall algorithm. The property is obviously satisfied by the two decompositions described here because the columns of the matrix \mathbf{R}_U (4.14) are linearly independent.

A simple distribution scheme with all of the above properties is the so-called PSI scheme [8]. It is most easily described by considering a single triangular cell in isolation. If, according to the linearised advection velocity, $\hat{\lambda}$ of (5.5), the triangle has a single downstream vertex, at node i say, then that node receives the whole fluctuation, so

$$u_i^{n+1} = u_i^n + \frac{\Delta t}{S_i} \hat{\phi}, \quad (5.8)$$

while the values of u at the other two vertices remain unchanged. In the case of a triangle with two downstream vertices, at nodes i and j for example, the fluctuation is divided between these two nodes. The update due to this cell's

fluctuation can therefore be written

$$\begin{aligned} u_i^{n+1} &= u_i^n + \frac{\Delta t}{S_i} \phi_i^* , \\ u_j^{n+1} &= u_j^n + \frac{\Delta t}{S_j} \phi_j^* , \end{aligned} \quad (5.9)$$

where $\phi_i^* + \phi_j^* = \hat{\phi}$ for conservation. In the PSI scheme [22]

$$\begin{aligned} \phi_i^* &= \phi_i - L(\phi_i, -\phi_j) \\ \phi_j^* &= \phi_j - L(\phi_j, -\phi_i) , \end{aligned} \quad (5.10)$$

where

$$\phi_i = -\frac{1}{2} \hat{\lambda} \cdot \vec{n}_i (u_i^n - u_k^n) , \quad \phi_j = -\frac{1}{2} \hat{\lambda} \cdot \vec{n}_j (u_j^n - u_k^n) , \quad (5.11)$$

and L denotes the minmod limiter function,

$$L(x, y) = \frac{1}{2} (1 + \text{sgn}(xy)) \frac{1}{2} (\text{sgn}(x) + \text{sgn}(y)) \min(|x|, |y|) . \quad (5.12)$$

The PSI scheme is positive for a restriction on the time-step at a node i given by

$$\Delta t \leq \frac{S_i}{\sum_{\cup \Delta_i} \max(0, \frac{1}{2} \hat{\lambda}^j \cdot \vec{n}_i^j)} , \quad (5.13)$$

and is used in the overall algorithm for the distribution of the homogeneous scalar components which arise from the decompositions of Section 4.

5.2 Distribution of Coupled Components/Subsystems

The elliptic nature of the 2×2 subsystem which results from the decomposition of the shallow water equations in subcritical flow suggests that an upwind distribution strategy is less appropriate than for the scalar components. Two schemes are described here for the distribution of this component, one for each decomposition, following the different distributions suggested for the corresponding decompositions of the Euler equations [16, 18].

In the first decomposition (HELW) the two coupled equations are modelled by the system

$$\underline{\mathbf{u}}_t + \mathbf{A} \underline{\mathbf{u}}_x + \mathbf{B} \underline{\mathbf{u}}_y = \underline{\mathbf{0}}, \quad (5.14)$$

in which \mathbf{A} and \mathbf{B} are defined explicitly by

$$\mathbf{A} = \begin{pmatrix} -\varepsilon\beta & 0 \\ 0 & \beta \end{pmatrix} \quad \text{and} \quad \mathbf{B} = \begin{pmatrix} 0 & \varepsilon \\ 1 & 0 \end{pmatrix}, \quad (5.15)$$

and

$$\partial \underline{\mathbf{u}} = \begin{pmatrix} \frac{g\beta}{q} \partial h \\ q \partial \theta \end{pmatrix}. \quad (5.16)$$

Under the linearisation of Section 3 the component of $\overline{\Phi_U}$ (4.16) corresponding to equation (5.14) takes the form

$$\overline{\Phi} = -S_\Delta (\overline{\mathbf{A}} \underline{\mathbf{u}}_x + \overline{\mathbf{B}} \underline{\mathbf{u}}_y). \quad (5.17)$$

This quantity is distributed in a similar manner to the fluctuation in (5.2), the difference being that the distribution coefficients α_i^j now take the form of matrices which, in order to ensure conservation, must satisfy

$$\sum_i \alpha_i^j = \mathbf{I}_d \quad \forall j, \quad (5.18)$$

where d is the dimension of the subsystem (5.14).

In [16] a Lax-Wendroff type scheme (which is linearity preserving and continuous but not positive) is suggested for the distribution of $\overline{\Phi}$. The coefficients of the scheme are written

$$\alpha_i^j = \frac{1}{3} \mathbf{I}_d + \frac{\Delta t}{4S_{\Delta_j}} (\mathbf{A} n_{i_x} + \mathbf{B} n_{i_y}), \quad (5.19)$$

where $\vec{n}_i = (n_{i_x}, n_{i_y})^T$ is the scaled inward pointing normal to edge i of the triangle. This scheme is used here in the distribution of the elliptic subsystem which arises from the HELW decomposition in subcritical flow.

The second approach (HESUPG) equates the coupled subsystem with a pair of scalar advection equations with source terms of the form

$$u_t + \vec{\lambda} \cdot \vec{\nabla} u = q, \quad (5.20)$$

in which u , $\vec{\lambda}$ and q are defined by the first two entries in (4.19), (4.23) and (4.24) respectively. In [18] the quantity

$$\widehat{\phi}_q = -S_\Delta (\widehat{\vec{\lambda}} \cdot \widehat{\vec{\nabla}} u - \widehat{q}), \quad (5.21)$$

is distributed for each of the two equations using a scheme which is equivalent to a mass-lumped streamline upwind Petrov-Galerkin (SUPG) finite element scheme with additional artificial viscosity [4].

The distribution coefficients for this linearity preserving and continuous but non-positive scheme are given by

$$\alpha_i^j = \frac{1}{3} + \tau \frac{\vec{\lambda} \cdot \vec{n}_i}{2S_{\Delta_j}} + \kappa \frac{\vec{\nabla} u \cdot \vec{n}_i}{2S_{\Delta_j}}, \quad (5.22)$$

in which

$$\tau = C_1 \frac{h}{|\vec{\lambda}|}, \quad \kappa = C_2 \frac{h \operatorname{sgn}(\widehat{\phi})}{|\vec{\nabla} u| + h}. \quad (5.23)$$

The constants C_1 and C_2 are both taken to be 0.5 [13], h is a typical local length scale, *e.g.* the length of the longest edge of the cell, and $\widehat{\phi}$ is defined in (5.3). This scheme is used here for the distribution of the coupled equations which result from the subcritical HESUPG decomposition.

6 Source Terms

Source terms appear in the linearised shallow water equations both as a result of modelling bed slope and friction (2.1) and from the linearisation (3.14), and

these terms must be included in the updating of the solution.

The simplest method of treating the momentum sources, \underline{q} in (2.1), is to calculate them pointwise at each node and then add them to the conservative variables once the flux balance distribution has been completed, so

$$\underline{U}_i^{n+1} = \underline{U}_i^n + \delta \underline{U}_i + \Delta t \underline{q}_i, \quad (6.1)$$

in which $\delta \underline{U}_i^n$ is the update indicated by the distribution of the decomposed flux balance. However, it is more appropriate to the schemes presented here for all of the sources to be incorporated within the flux balance distribution itself. This is the obvious way to treat the linearisation source terms since they are inherently cell-based quantities.

One way of achieving this is to include the source terms within the decomposition, so the characteristic equations of (4.7) become

$$\underline{W}_t + \mathbf{A}_W^S \underline{W}_\xi + \mathbf{B}_W^S \underline{W}_\eta = \mathbf{R}_U^{-1} \underline{q}_{\text{tot}}, \quad (6.2)$$

where $\underline{q}_{\text{tot}}$ is the sum of the momentum and linearisation source terms consistently evaluated from the cell-average state $\bar{\underline{Z}}$. The two types of source term can be considered separately but are combined here for simplicity.

The effect of $\underline{q}_{\text{tot}}$ on the flux balance distribution can be illustrated simply by considering a scalar component of the decomposition. A characteristic equation taken from (4.11) now has the form

$$W_t^k + \vec{\lambda}_S^k \cdot \vec{\nabla}_S W^k = q_W^k, \quad (6.3)$$

where q_W^k is the k^{th} component of the vector $\mathbf{R}_U^{-1} \underline{q}_{\text{tot}}$, and the quantity which is distributed to the nodes of the grid is now

$$\widehat{\phi}_q = -S_\Delta \left(\vec{\lambda}_S^k \cdot \vec{\nabla}_S W^k - q_W^k \right). \quad (6.4)$$

A positive distribution scheme does not remain positive under this modification but the linearity preservation property is retained by calculating the distribution coefficients precisely as in the homogeneous case but then using them to distribute the quantity $\widehat{\phi}_q$. The modified updates are then transformed into increments of the conservative variables using the matrix \mathbf{R}_U (4.14) as before. The source terms which now appear in the elliptic subsystem can also be treated in this manner for both the matrix and scalar distributions.

A third method of treating the source term $\underline{q}_{\text{tot}}$ is to distribute it separately from $\underline{\Phi}_U$, and the simplest way to do this is via a symmetric distribution in which one third of $\underline{q}_{\text{tot}}$ within a cell is sent to each of its vertices. All three ways of incorporating the source terms are considered in the following section.

7 Results

Both algorithms described in the previous sections (HELW and HESUPG) have been used to solve numerically a wide variety of steady state test cases for the two-dimensional shallow water equations. In all cases the linearisation source terms are distributed separately from the rest of the flux balance by a simple central scheme since this strategy proves to be more robust than an upwind distribution and there is negligible difference between the results. The momentum sources, when they appear, are distributed in an upwind manner as part of the flux balance for the purposes of accuracy, except when robustness becomes an issue in which case they are considered separately and evaluated on a pointwise basis.

The boundary conditions are applied very simply by referring to the theory of characteristics. This determines the number and form of the physical conditions

which should be imposed at a chosen point on the boundary. One condition must be applied for each positive eigenvalue of the matrix

$$\mathbf{C}_U = \mathbf{A}_U n_x + \mathbf{B}_U n_y, \quad (7.1)$$

where $\vec{n} = (n_x, n_y)^T$ is the inward pointing normal to the boundary of the computational domain. In the case of the shallow water equations these eigenvalues are given by

$$\lambda_1 = \vec{u} \cdot \vec{n}, \quad \lambda_2 = \vec{u} \cdot \vec{n} + c \quad \text{and} \quad \lambda_3 = \vec{u} \cdot \vec{n} - c. \quad (7.2)$$

Thus, when the component of the flow normal to the boundary is supercritical either the whole solution is specified (at inflow) or none of it (outflow). For subcritical inflow two conditions are specified (total head and tangential velocity component) while for subcritical outflow a single piece of information, the depth of the flow, is set to a prespecified freestream value. At a solid wall only λ_2 is positive and this is accommodated by setting the normal velocity component to zero.

7.1 Oblique Hydraulic Jump

Few standard steady state test cases exist for the homogeneous two-dimensional shallow water equations, but there are some simple problems for which exact solutions have been calculated. One such example [3] is supercritical flow through a frictionless channel with a flat bed containing a wedge inclined at an angle θ to the direction of the flow at which an oblique hydraulic jump is induced by the interaction of the flow with the front of the wedge. The angle β which this

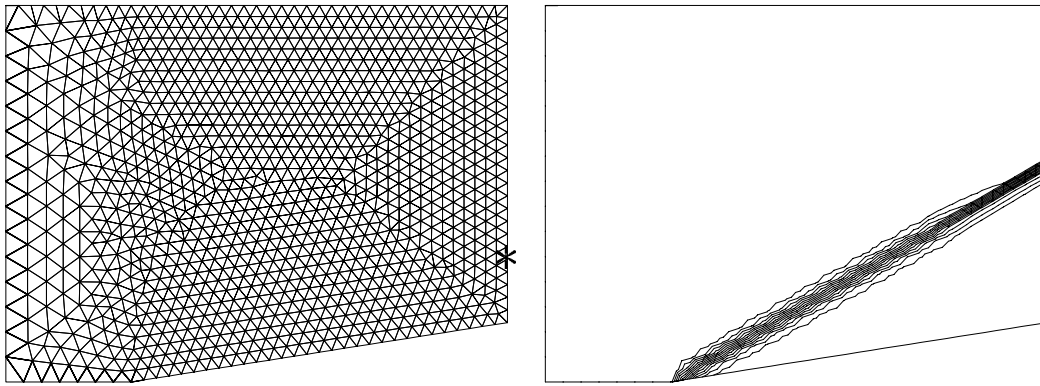


Figure 7.1: The grid (left) and local Froude number contours (right) for the oblique hydraulic jump test case. The asterisk indicates the point at which flow values have been sampled.

discontinuity makes with the direction of the freestream flow is defined by

$$\sin \beta = \frac{1}{F_u} \sqrt{\frac{h_d}{2h_u} \left(\frac{h_d}{h_u} + 1 \right)}, \quad (7.3)$$

where the subscripts u and d indicate values upstream and downstream of the jump, respectively.

In the case chosen here, the slope of the wedge is taken to be $\theta = 8.95^\circ$ and its leading edge is positioned 1m in to a $4\text{m} \times 3\text{m}$ rectangular domain. The upstream flow conditions are given as $h_u = 1\text{m}$, $u_u = 8.57\text{ms}^{-1}$ and $v_u = 0\text{ms}^{-1}$ (so $F_u = 2.74$).

The resulting flow is purely supercritical and is divided into two regions by a hydraulic jump. Upstream of the discontinuity the initial (freestream) conditions prevail at the steady state, while on the downstream side the exact solution is given by $h_d = 1.5\text{m}$ and $|\vec{u}_d| = 7.9556\text{ms}^{-1}$, so $F_d = 2.074$. The jump itself is at an angle $\beta = 30^\circ$ to the inlet flow.

The computation has been carried out on the 1175 node, 2231 cell grid shown in Figure 7.1 with the upper and lower boundaries both being treated as solid

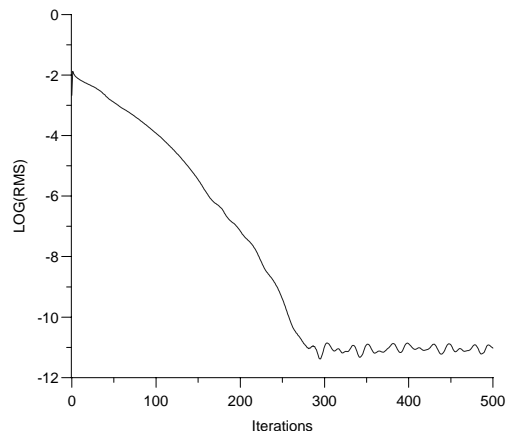


Figure 7.2: Convergence history for the oblique hydraulic jump test case.

walls. This figure also shows the local Froude number contours of the steady state solution calculated for this test case (both the HELW and the HESUPG schemes are the same for supercritical flow). The hydraulic jump can be seen to be captured sharply at the correct angle and a discontinuous water surface devoid of oscillations is obtained. The values of the flow variables downstream of the jump (sampled on the outflow boundary at the point indicated by the asterisk in Figure 7.1) are $h_d = 1.5001\text{m}$ and $|\vec{u}_d| = 7.9506\text{ms}^{-1}$ ($F_d = 2.073$), very close to the exact values.

The convergence history for this calculation pictured in Figure 7.2 shows that the numerical solution converges rapidly to machine precision. The monitor used is the root mean square of the nodal updates to the conservation of mass equation, given by

$$\text{RMS} = \sqrt{\frac{\sum_{i=1}^{N_n} \left(\frac{h_i^{n+1} - h_i^n}{h_i^n} \right)^2}{\text{CFL } N_n}}. \quad (7.4)$$

Only local time-stepping is used here to accelerate convergence. In [18] it was shown that the use of implicit and characteristic time-stepping techniques would both significantly reduce the cost of reaching the steady state in the case of the

Euler equations. Neither technique is used here but it is expected that both could be used to similar advantage.

Note that a CFL number of 0.7 has been used here but in the subsequent test cases, all of which have regions of subcritical flow, the CFL number is taken to be 0.2 which proved to be the highest value which could be taken which was stable for all of these cases. This seems to be because of the discontinuity in the distribution at the critical line and the nonorthogonality of the eigenvectors of the preconditioned system at low Froude numbers (described in more detail in [6]).

7.2 Symmetric Constricted Channel Flows

The domain for these test cases represents a channel of length 4 metres and width 1 metre with bumps of the same shape and size in the centre of either wall of the channel. The bumps are one metre in length and are defined such that the breadth of the channel is given by

$$B = B_0 - 2B_h \cos^2 \left(\frac{\pi(x - x_c)}{x_l} \right) \quad \text{for} \quad |x - x_c| \leq \frac{x_l}{2}, \quad (7.5)$$

where $B_0 = 1\text{m}$ is the breadth of the straight channel, B_h is the height of each bump, $x_c = 2\text{m}$ is the position of the centre of the constriction and $x_l = 1\text{m}$ is its length. In this particular case B_h is taken to be 0.04m. The grid on which the results have been obtained consists of 2114 nodes and 4054 cells and is shown in Figure 7.3.

For the first constricted channel test case the freestream Froude number is specified to be $F_\infty = 0.5$ while the depth is set at $h_\infty = 1\text{m}$ and the flow is perpendicular to the boundary at the inlet. The resulting steady state solution

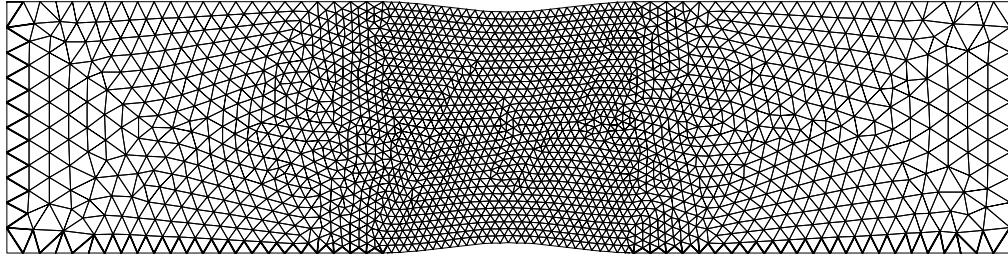


Figure 7.3: The grid for the symmetric constricted channel flow test cases.

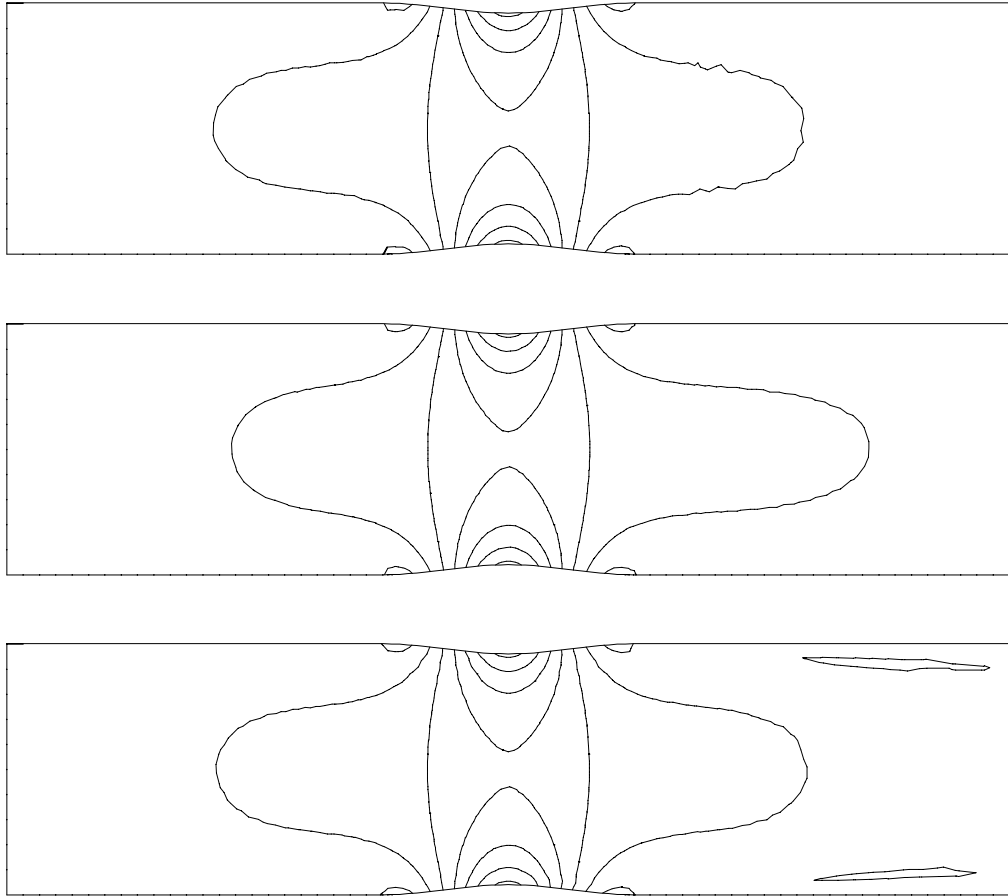


Figure 7.4: The local Froude number contours for the symmetric constricted channel flow test case with $F_\infty = 0.5$ using the HELW (top), HESUPG (middle) and Lax-Wendroff (bottom) schemes.

is completely subcritical and therefore symmetric about the centre of the constriction (the narrowest point of the channel) since the friction and bed slope are both taken to be zero.

It can be seen from Figure 7.4 that the numerical solutions obtained using the multidimensional upwind schemes differ little from the result produced by a Lax-Wendroff distribution scheme applied to the complete system without decomposition, also shown in the figure. The results from the HESUPG scheme show a very slight asymmetry, indicating that the HELW scheme is the more accurate of the two for this type of flow. This is probably because a Lax-Wendroff scheme is used to distribute the 2×2 elliptic subsystem which arises in subcritical flow rather than the more diffusive SUPG scheme. The advantages of the upwind schemes over Lax-Wendroff become clear when discontinuous flows are considered.

The freestream Froude number is now increased to $F_\infty = 0.71$ to produce a steady transcritical flow which includes a hydraulic jump downstream of the narrowest point of the channel. All other parameters remain the same. The pictures of the local Froude number contours, shown in Figure 7.5, indicate that the discontinuity is captured very sharply by both upwind schemes, within two or three cells right across the channel. The jump is sharper and slightly further downstream when the HELW scheme is used, further evidence of its less diffusive nature, but this is at the expense of small oscillations on the downstream side of the jump.

An increase in the diffusive component of the distribution coefficients (5.19) for

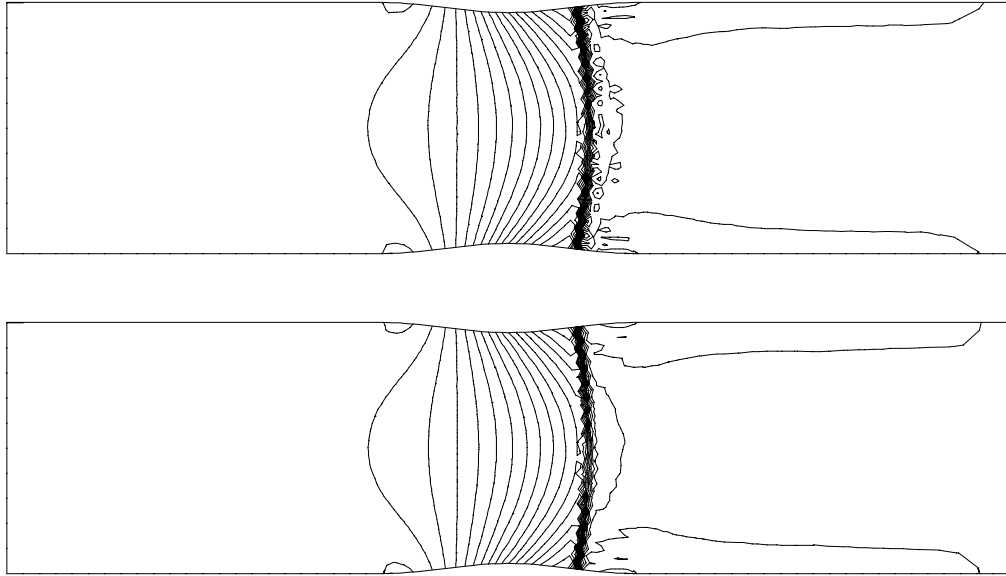


Figure 7.5: Local Froude number contours for the constricted channel flow with $F_\infty = 0.71$ for the HELW (top) and HESUPG (bottom) schemes.

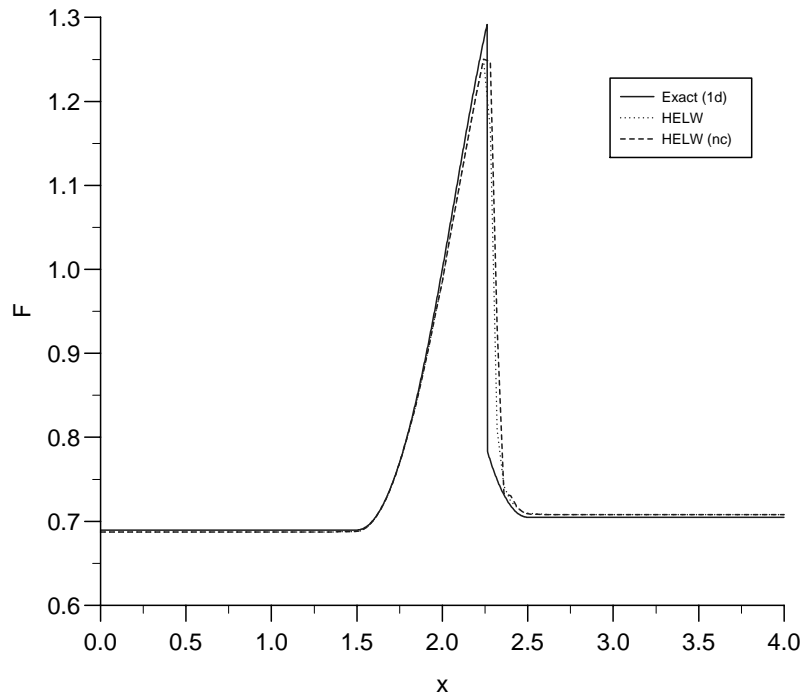


Figure 7.6: Comparison of the breadth-averaged local Froude number for the constricted channel flow test case with $F_\infty = 0.71$.

the subcritical elliptic subsystem would reduce the oscillations, but not without smearing the discontinuity as well. Although the HESUPG scheme, with its greater inherent numerical diffusion, does this automatically, neither treatment is ideal and the modelling of transcritical flows requires further consideration.

The results shown in Figure 7.6 illustrate the effect of the linearisation source terms on the solution. The values of the breadth-averaged local Froude number are plotted along the length of the channel for the HELW scheme. The small oscillations downstream of the jump are rendered almost invisible by the averaging procedure and the solutions are very close to those produced by the HESUPG scheme (not plotted here) although the latter predicts the one-dimensional discontinuity to be very slightly further upstream.

The numerical results shown are for a conservative and a non-conservative formulation in which the linearisation sources are simply ignored. Close inspection reveals that the discontinuity is predicted to be about half a cell's width further downstream by the non-conservative scheme. On a grid in which the cell edges are aligned with the discontinuity the discrepancy in jump position between the conservative and non-conservative schemes can be as much as one cell. The non-conservative formulation predicts the jump to be further away from the exact position predicted by one-dimensional theory for an open channel of varying width, the third solution shown in Figure 7.6. Thus it is important to enforce conservation for precise positioning of the discontinuity even though an adequate solution may be obtained in this case without conservation. Note also that the averaged conservative numerical approximation passes from subcritical to supercritical flow at the centre of the channel (its narrowest point) as it should

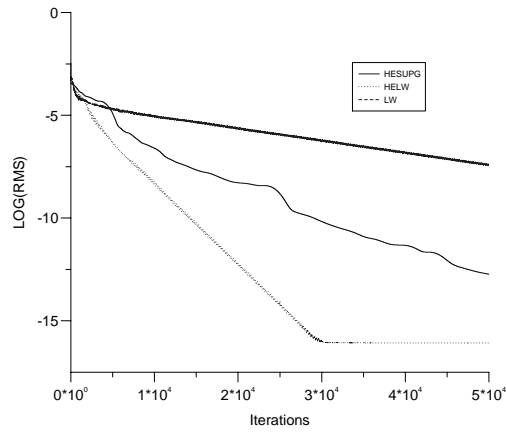


Figure 7.7: Convergence history for the symmetric constricted channel flow test case with $F_\infty = 0.5$.

do.

Figure 7.7 shows the convergence histories for the subcritical solutions presented in this section. It can be seen that in the subcritical case, although convergence is still slow without the use of acceleration techniques, both upwind schemes (particularly HELW) converge much faster than the Lax-Wendroff scheme. In the transcritical case convergence is generally very slow and in some cases machine accuracy is not achieved.

7.3 Sloping Channel Flows with Friction

Although essentially one-dimensional this test case can be used to validate the treatment of the momentum source terms on the two-dimensional triangular grid. It is one of a family of exact steady state solutions of the one-dimensional shallow water equations with bed slope and friction included which has recently been constructed for straight open channels [15].

The particular case chosen here is of flow in a rectangular channel, 1000m long and 10m wide. Manning's roughness coefficient is taken to be 0.02 and the

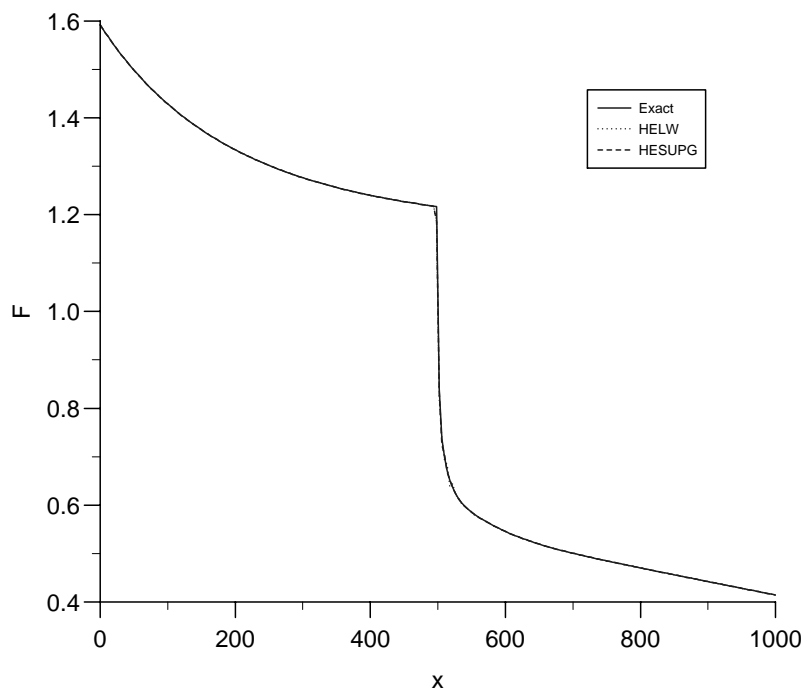


Figure 7.8: Comparison of the breadth-averaged local Froude number for the quasi-one-dimensional rectangular channel test case.

bed slope is given in terms of a hypothetical depth function \tilde{h} by

$$s_{bx} = \left(1 - \frac{4}{g\tilde{h}^3}\right) \frac{d\tilde{h}}{dx} + 0.16 \frac{(2\tilde{h} + 10)^{\frac{4}{3}}}{(10\tilde{h})^{\frac{4}{3}}} \quad \text{and} \quad s_{by} = 0, \quad (7.6)$$

where

$$\tilde{h} = \begin{cases} \left(\frac{4}{g}\right)^{\frac{1}{3}} \left(\frac{9}{10} - \frac{1}{6} \exp\left(-\frac{x}{250}\right)\right) & \text{for } 0 \leq x \leq 500 \\ \left(\frac{4}{g}\right)^{\frac{1}{3}} \left(1 + \sum_{k=1}^3 a_k \exp\left(-20k\left(\frac{x}{1000} - \frac{1}{2}\right)\right) + \frac{4}{5} \exp\left(\frac{x}{1000} - 1\right)\right) & \text{for } 500 < x \leq 1000, \end{cases} \quad (7.7)$$

in which $a_1 = -0.348427$, $a_2 = 0.552264$ and $a_3 = -0.555580$. The depth of the steady state solution is given by $h \equiv \tilde{h}$ and the discharge is taken to be $20\text{m}^3\text{s}^{-1}$. The resulting flow is supercritical at inflow and subcritical at outflow with a discontinuity half way along the channel, at $x = 500\text{m}$.

The two-dimensional solutions have been obtained on a 1004 node, 1500 cell grid in which all of the triangles have an aspect ratio of approximately one, and

the breadth-averaged local Froude number predicted by the two upwind schemes is shown in Figure 7.8, together with the exact solution. The three solutions are almost indistinguishable and even the HELW scheme exhibits no small oscillations on the subcritical side of the jump.

Both sets of results presented have been obtained using an upwind distribution of momentum sources evaluated on a cell by cell basis. Close examination of the solutions reveals that this method of treating these source terms leads to the best approximation of the exact solution, although the differences would not be visible in the figure. It should be noted though that convergence to the steady state is slightly better if the sources are incorporated at the nodal update stage, indicating greater robustness. In actual fact none of the schemes converge to machine accuracy in this transcritical case so no convergence histories are shown.

7.4 Spillway Flow

The final problem represents shallow water flow in a spillway and provides a genuinely two-dimensional test case. The flow is through a channel 10m wide with a right angled bend half way along its length. The inner and outer corners of the bend are taken to be arcs of concentric circles with radii 10m and 20m respectively, and there is 30m of straight channel both upstream and downstream of the bend. The flow is supercritical at both inflow and outflow and the inflow conditions are such that $h = 1\text{m}$, $u = 0\text{ms}^{-1}$ and $v = -5\text{ms}^{-1}$. The slope of the channel is of magnitude 1 in 5 along the straight sections and varies linearly across the channel around the bend from 1 in 5 on the outer curve to 2 in 5 on the inner curve. Manning's roughness coefficient for the flow takes the value of

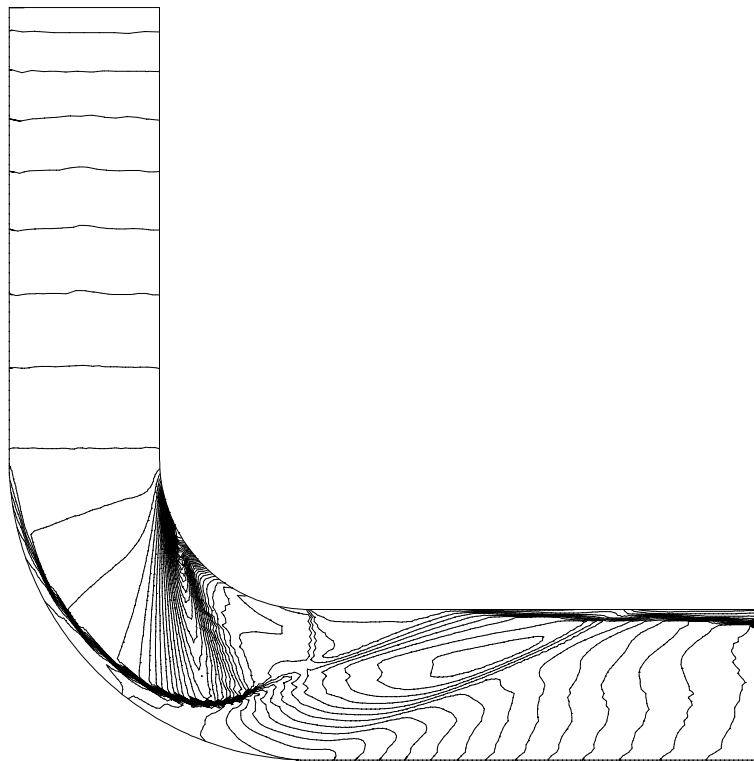
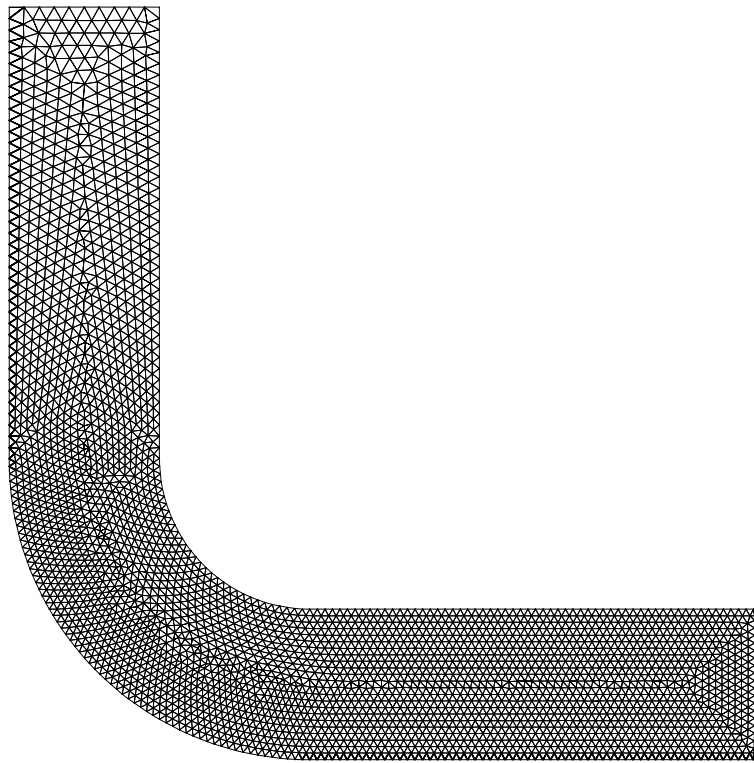


Figure 7.9: The grid (top) and local Froude number contours (bottom) for the spillway test case.

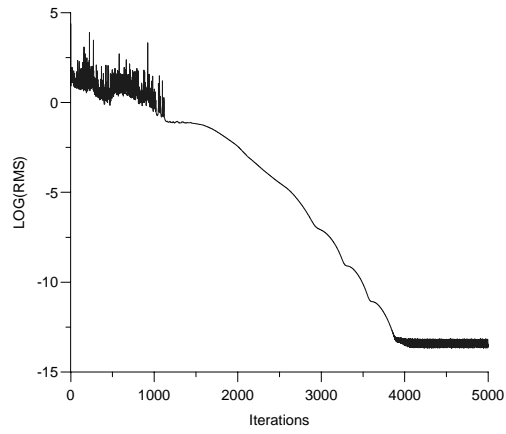


Figure 7.10: Convergence history for the spillway test case.

0.012.

The calculation has been carried out using the HESUPG scheme on the grid shown in Figure 7.9 (which is made up of 3310 nodes and 6294 cells) above the local Froude number contours of the resulting steady state solution. On the evidence of the previous test case, the momentum sources were calculated on a nodal basis and incorporated in the update after the distribution step for maximum robustness.

The hydraulic jumps in the solution are captured very sharply, across a maximum of two cells, and without spurious oscillations, despite the extreme nature of the flow in the region of the corner. Overall, the scheme appears to predict qualitatively correct flow features. The flow is completely supercritical so the HELW scheme gives precisely the same solution. Figure 7.10 shows that in this purely supercritical case the convergence to the steady state is again rapid.

8 Conclusions

Two alternative genuinely multidimensional upwind schemes have been presented for the numerical solution of the two-dimensional shallow water equations on unstructured triangular grids. Techniques which were originally developed for the solution of the Euler equations have been adapted for the approximation of the shallow water equations and a conservative formulation of the algorithm has been presented. A method of treating general source terms, appropriate for use with multidimensional upwinding, has also been suggested.

Both schemes presented here have been shown to produce high quality results for both subcritical and supercritical steady state flows and have the ability to capture discontinuities very sharply. Of the two, the HELW scheme is the less diffusive and as a result, slightly less robust. It produces sharper jumps and more accurate subcritical approximations, but at the expense of small oscillations downstream of transcritical discontinuities. These could be smoothed by adjusting the Lax-Wendroff distribution appropriately, with a consequent smearing of the discontinuity, but both schemes require further study to improve the modelling of the transition between supercritical and subcritical flows.

The treatment of the momentum sources also merits attention since it is still unclear which is the best method to use for their distribution. The most robust treatment proved to be to consider the sources on a nodal basis but distributing cell-averaged source terms in an upwind manner appears to be more accurate. It is clear though that the linearisation sources are necessary for the precise positioning of the hydraulic jumps although in many cases adequate numerical solutions can be obtained using a non-conservative formulation.

Multidimensional upwinding is still undergoing improvements for steady state flows, even though the resulting numerical solutions are now of a high quality for all flow regimes. However, it remains to construct a scheme which produces approximations of comparable accuracy to time-dependent flows.

Acknowledgements

The authors would like to thank Dr. A. Priestley, Dr. P. Garcia-Navarro and I. MacDonald for their contributions to this work and DRA Farnborough and EPSRC for providing the funding for the first author.

References

- [1] *Unstructured grid methods for advection dominated flows*, AGARD Report AGARD-R-787, 1992.
- [2] R.Abgrall, ‘Approximation of the multidimensional Riemann problem in compressible fluid mechanics by a Roe type method’, *Comptes Rendus de l’Academie des Sciences Serie 1 - Mathematique*, **319(5)**:499–504, 1994.
- [3] F.Alcrudo and P.Garcia-Navarro, ‘A high resolution Godunov-type scheme in finite volumes for the 2d shallow-water equations’, *Int. J. of Num. Methods in Fluids*, **16**:489–505, 1993.
- [4] J.-C.Carette, H.Deconinck, H.Paillère and P.L.Roe, ‘Multidimensional upwinding: it’s relationship to finite elements’, in *Proceedings of 8th International Conference on Finite Elements in Fluids*, Barcelona, 1993.

- [5] J.A.Cunge, F.M.Holly and A.Verwey, *Practical Aspects of Computational River Hydraulics*, Pitman, London, 1980.
- [6] D.L.Darmofal and P.J.Schmid, ‘The importance of eigenvectors for local preconditioners of the Euler equations’, *J. Comput. Phys.*, **127**:346–362, 1996.
- [7] H.Deconinck, P.L.Roe and R.Struijs, ‘A multi-dimensional generalization of Roe’s flux difference splitter for the Euler equations’, *J. of Computers and Fluids*, **22(2/3)**:215–222, 1993.
- [8] H.Deconinck, R.Struijs, G.Bourgois and P.L.Roe, ‘High resolution shock capturing cell vertex advection schemes for unstructured grids’, in *Computational Fluid Dynamics*, number 1994–05 in VKI Lecture Series, 1994.
- [9] R.J.Fennema and M.H.Chaudhry, ‘Explicit methods for 2-d transient free-surface flows’, *J. Hydraul. Eng.*, **116(8)**:1013–1034, 1990.
- [10] P.Garcia-Navarro, M.E.Hubbard and A.Priestley, ‘Genuinely multidimensional upwinding for the 2d shallow water equations’, *J. Comput. Phys.*, **121**:79–93, 1995.
- [11] P.Glaister, ‘Flux difference splitting for open-channel flows’, *Int. J. for Num. Methods in Fluids*, **16**:629–654, 1993.
- [12] C.Hirsch, *Computational Methods for Inviscid and Viscous Flows*, volume 2 of *Numerical Computation of Internal and External Flows*, Wiley, 1990.
- [13] C.Johnson, ‘Finite elements for flow problems’, in *Unstructured Grid Methods for Advection Dominated Flows*, AGARD Report AGARD-R-787, 1992.

- [14] R.J.Leveque, *Numerical Methods for Conservation Laws*, Birkhauser, Basel, 1992.
- [15] I.MacDonald, M.J.Baines, N.K.Nichols and P.G.Samuels, ‘Steady Open Channel Test Problems with Analytic Solutions’, Numerical Analysis Report 3/95, The University of Reading, 1995, to appear in *J. Hydraulic Eng., ASCE*.
- [16] L.M.Mesaros and P.L.Roe, ‘Multidimensional fluctuation schemes based on decomposition methods’, AIAA Paper 95–1699, 1995.
- [17] H.Paillère, ‘It is possible to solve dam-break problems using a multidimensional upwinding approach’, Thèse annexe, Université Libre de Bruxelles, 1995.
- [18] H.Paillère, E.van der Weide and H.Deconinck, ‘Multidimensional upwind methods for inviscid and viscous compressible flows’, in *Computational Fluid Dynamics*, number 1995–02 in VKI Lecture Series, 1995.
- [19] P.L.Roe, ‘Approximate Riemann solvers, parameter vectors, and difference schemes’, *J. Comput. Phys.*, **43(2)**:357–372, 1981.
- [20] P.L.Roe, ‘Fluctuations and signals - a framework for numerical evolution problems’, in *Numerical Methods for Fluid Dynamics*, edited by K.W.Morton and M.J.Baines, pages 219–257, Academic Press, 1982.
- [21] P.L.Roe, ‘Beyond the Riemann problem I’, in *Algorithmic Trends in Computational Fluid Dynamics*, Springer-Verlag, 1992.

- [22] D.Sidilkover and P.L.Roe, ‘Unification of some advection schemes in two dimensions’, ICASE Report 95–10, 1995.
- [23] J.J.Stoker, *Water Waves*, Interscience, New York, 1986.
- [24] B.van Leer, ‘Progress in multi-dimensional upwind differencing’, ICASE Report 92–43, 1992.
- [25] B.van Leer, W.-T.Lee and P.L.Roe, ‘Characteristic time-stepping or local preconditioning of the Euler equations’, AIAA Paper 91–1552–CP, 1991.
- [26] B.van Leer, E.Turkel, C.H.Tai and L.M.Mesaros, ‘Local preconditioning in a stagnation point’, AIAA Paper 95–1954, 1995.