# Flux Difference Splitting and the Balancing of Source Terms and Flux Gradients. *

M.E.Hubbard

The University of Reading, Department of Mathematics,

P.O.Box 220, Whiteknights, Reading, Berkshire, RG6 6AX, U.K.

and

P.Garcia-Navarro

Area de Mecanica de Fluidos, Centro Politecnico Superior,

C/Maria de Luna 3, 50015 Zaragoza, España.

21st May 1999

## Abstract

Flux difference splitting methods are widely used for the numerical approximation of homogeneous conservation laws where the flux depends only on the conservative variables. However, in many practical situations this is not the case. Not only are source terms commonly part of the mathematical model, but the flux can vary spatially even when the conservative variables do not. It is the discretisation of the additional terms arising from these two situations which is addressed in this work, given that a specific flux difference splitting method has been used to approximate the underlying conservation law. The discretisation is constructed in a manner which retains an exact balance between the flux gradients and the source terms when this is appropriate.

The effectiveness of these new techniques, in both one and two dimensions, is illustrated using the shallow water equations, in which the additional terms arise from the modelling of bed slope and, in one dimension, breadth variation. Roe's scheme is chosen for the approximation of the conservation laws and appropriate discrete forms are constructed for the additional terms, not only in the first order case (which has been done before) but also in the presence of flux limited and slope limited high resolution corrections. The method is then extended to two-dimensional flow where it can be applied on both quadrilateral and triangular grids.

# 1  Introduction

There has been much research in CFD into the accurate and efficient solution of homogeneous systems of conservation laws. More recently, as numerical models become more complicated and the areas of application of these methods widens, it has become important that other aspects of the discretisation be given due attention. This is certainly true in the field of computational hydraulics where the modelling can be dominated by the effects not only of source terms, but also of quantities which vary spatially but independently of the flow variables.

It can be argued that the presence of these effects warrants the construction of new numerical schemes which are appropriate to the nature of the equations, not one of the many which have been constructed for the simple, homogeneous case. However, this work is concerned with how the additional terms should be discretised, given that a specific scheme has been used to approximate the flux terms. This approach has been taken previously by a number of authors, and applied in a variety of different situations. For example, Smolarkiewicz has adapted his own MPDATA scheme to solve inhomogeneous equations arising from geophysical flows [14], LeVeque has incorporated the modelling of source terms for shallow water flows within his wave-propagation algorithm [9], and Roe's scheme [11] has been modified by a number of authors to include source terms, the research of Glaister [6], Vázquez-Cendón [16], Bermúdez and Vázquez [1] and Bermúdez *et al.* [2] being of particular relevance to this work.

In each of the aforementioned papers discussing Roe's scheme, the discrete form of the source terms has been deliberately constructed along similar lines to the numerical fluxes. This is done to ensure that equilibria which occur in the mathematical model are retained by the numerical model, and that in the absence of additional terms, the conservative fluxes are retrieved for accurate modelling of discontinuous solutions. However, all of the previous work deals only with the first order scheme. The intention of this paper is to provide an extension of these ideas to higher order Total Variation Diminishing (TVD) versions of Roe's scheme (using both flux limiting and slope limiting techniques) and to describe a source term approximation which has each of the above properties on all types of

regular and irregular grids in any number of dimensions. Furthermore, following on from [4], a new formulation is presented for the discretisation of the flux in the case where it depends on a spatially varying quantity which is independent of the solution.

The shallow water equations have been chosen to demonstrate the effectiveness of these new techniques in one and two dimensions, by modelling the effects of a sloping bed and, in one dimension only, the inclusion of breadth variation in an open channel. The one-dimensional discretisation is described first, in Section 2, for a general system of conservation laws, followed by its application to the shallow water equations and a wide selection of results to show its accuracy. In Section 3 the generalisation to two dimensions, illustrated using unstructured triangular grids, is presented and again applied to the shallow water equations. The final section contains some brief conclusions obtained from the work.

## 2　One dimension

The one-dimensional equations representing a system of conservation laws with source terms may be written

$$\underline{U}_t + \underline{F}_x \;=\; \underline{S}\,, \tag{2.1}$$

where $\underline{U}$ is the vector of conservative variables, $\underline{F}$ is the conservative flux vector and $\underline{S}$ includes all of the source terms. In this section it is assumed that $\underline{F} = \underline{F}(\underline{U})$; in Section 2.2 the flux will be assumed to depend not only on the conservative variables but also another independent, spatially varying quantity, i.e. $\underline{F} = \underline{F}(\underline{U}, B(x))$.

Using the standard finite volume approximation of the flux terms in (2.1), combined with a simple, forward Euler discretisation of the time derivative leads to a difference scheme which can be written

$$\underline{U}_i^{n+1} \;=\; \underline{U}_i^n - \frac{\Delta t}{\Delta x_i}\left(\underline{F}^*_{i+\frac{1}{2}} - \underline{F}^*_{i-\frac{1}{2}}\right) + \frac{\Delta t}{\Delta x_i}\,\mathbf{S}^*_i\,, \tag{2.2}$$

in which $\underline{F}^*$ represents a numerical flux evaluated at an interface between control volumes and $\mathbf{S}^* \approx \int \underline{S}\,\mathrm{d}x$ is a numerical source integral over the control volume,

4

which has yet to be approximated. For convenience, a cell centre scheme in which the control volumes coincide with the mesh cells has been considered throughout this work, although the ideas may be applied to other types of scheme in a similar manner.

At first sight, the second term on the right hand side of (2.2) looks like a discrete flux derivative for cell $i$. However, for the purposes of this work it is more convenient to consider it from the point of view of the numerical fluxes being constructed from an approximation to the integral of the flux derivatives over dual cells and providing contributions to the cell updates ($\Delta x_i$ comes from the integration of the original equations over the control volume).

Commonly, $\underline{\mathbf{S}}_i^*$ is evaluated pointwise, taking the value $\Delta x_i \, \underline{S}(\underline{U}_i)$, or split symmetrically, giving an expression of the form

$$\underline{\mathbf{S}}_i^* \;=\; \frac{\Delta x_i}{2} \left( \underline{S}(\underline{U}_{i-\frac{1}{2}}) + \underline{S}(\underline{U}_{i+\frac{1}{2}}) \right) , \tag{2.3}$$

but a more sophisticated approach is sought here, based on the approach of Glaister [6], which accounts for the form of the numerical fluxes.
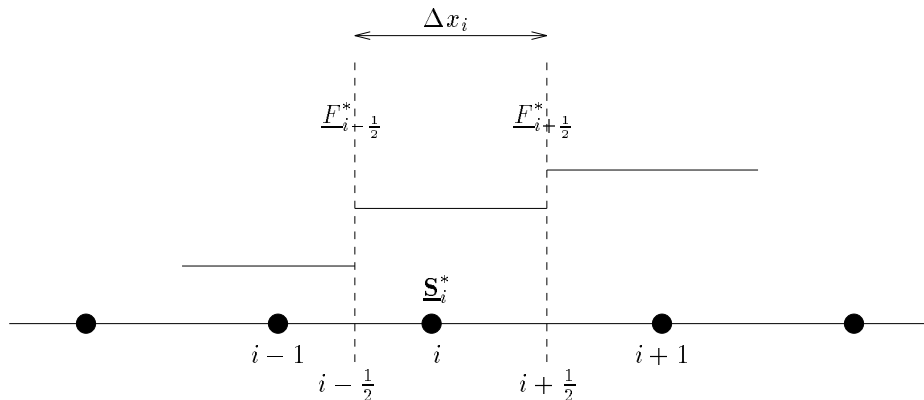


Figure 2.1: Numerical fluxes and sources for the cell centre scheme.

Note that in the absence of source terms the scheme given by (2.2) reduces to a conservative discretisation of the homogeneous system. Also, (2.2) has been written with irregular grids in mind, and as a consequence the mesh spacing $\Delta x_i = x_{i+\frac{1}{2}} - x_{i-\frac{1}{2}}$ relates to the cells, not the nodes (see Figure 2.1).

Roe's scheme [11] is one of the most commonly used examples of the conservative finite volume method mentioned above. This is an upwind scheme which uses

an approximate Riemann solver to decompose the flux terms into characteristic components by diagonalisation of the homogeneous part of a linearised form of the system (2.1), which is

$$\underline{U}_t + \tilde{\mathbf{A}}\underline{U}_x = \underline{0},\tag{2.4}$$

where $\tilde{\mathbf{A}} \approx \frac{\partial \underline{F}}{\partial \underline{U}}$ is the linearised flux Jacobian of the system. The Riemann problems arise at the interfaces between the control volumes (the mesh nodes in this case) where discontinuities occur in the discrete representation of the solution.

Application of Roe's Riemann solver results in a decoupling of the linearised equations that splits the flux difference so that it can be written in a number of equivalent forms, *i.e.* at an interface

$$\Delta \underline{F}_{i+\frac{1}{2}} = (\tilde{\mathbf{A}} \, \Delta \underline{U})_{i+\frac{1}{2}} = (\tilde{\mathbf{R}} \, \tilde{\mathbf{\Lambda}} \tilde{\mathbf{R}}^{-1} \Delta \underline{U})_{i+\frac{1}{2}} = \left( \sum_{k=1}^{N_w} \tilde{\alpha}_k \tilde{\lambda}_k \tilde{\underline{r}}_k \right)_{i+\frac{1}{2}},\tag{2.5}$$

in which $\Delta \underline{F}$ represents the jump in $\underline{F}$ across the edge of a control volume, $\tilde{\mathbf{R}}$ is the matrix whose columns are the right eigenvectors $\tilde{\underline{r}}_k$ of $\tilde{\mathbf{A}}$, $\tilde{\mathbf{\Lambda}}$ is the diagonal matrix of eigenvalues $\tilde{\lambda}_k$ of $\tilde{\mathbf{A}}$, and the components of $\tilde{\mathbf{R}}^{-1} \Delta \underline{U}(= \Delta \underline{W})$ are the 'strengths' $\tilde{\alpha}_k$ associated with each component of the decomposition ($\underline{W}$ being the vector of characteristic variables of the system). The final expression indicates how the flux difference is decomposed into $N_w$ characteristic components (or waves of the Riemann problem), where $N_w$ is the number of equations of the system (2.4). In both (2.4) and (2.5), $\tilde{\phantom{i}}$ denotes the evaluation of a quantity at its Roe-average state [11, 5]. This is a special average state of the flow variables which is constructed so that (2.5) is always satisfied for the given system.

Having obtained the decomposition (2.5), Roe's scheme for a homogeneous system of equations is constructed from (2.2) by taking the numerical fluxes to be

$$\underline{F}^*_{i+\frac{1}{2}} = \frac{1}{2}(\underline{F}_{i+1} + \underline{F}_i) - \frac{1}{2}\left( \tilde{\mathbf{R}}|\tilde{\mathbf{\Lambda}}|\tilde{\mathbf{R}}^{-1}\Delta \underline{U} \right)_{i+\frac{1}{2}},\tag{2.6}$$

where $|\tilde{\mathbf{\Lambda}}| = \mathrm{diag}(|\tilde{\lambda}_k|)$: the source terms have been temporarily ignored. A similar expression can be written down for $\underline{F}^*_{i-\frac{1}{2}}$.

Generally, $\tilde{\phantom{i}}$ in (2.6) can represent any consistent approximation to the specified variables, and the resulting scheme will be conservative. However, forcing it
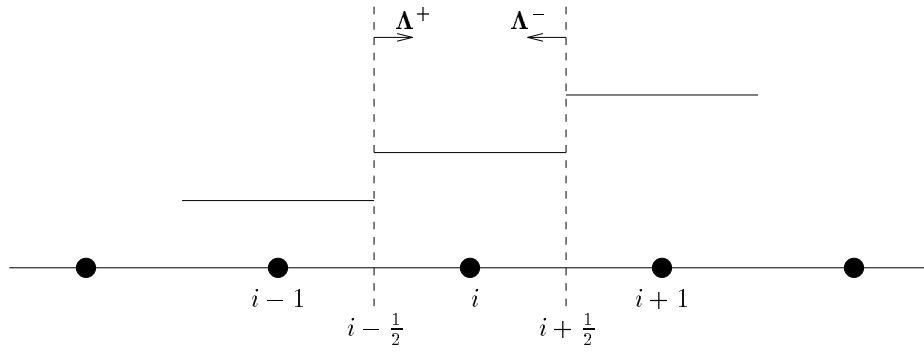
Figure 2.2: Wave propagation directions in a control volume.

Choosing the Roe-average state (represented by $\tilde{\cdot}$) to satisfy (2.5) means that the resulting approximate Riemann solver is an exact solver for this local linearisation of the Riemann problem. More importantly, in the context of this work, when (2.5) is combined with (2.6) the nodal update scheme given by (2.2) is equivalent to the fluctuation-signal scheme [12] given by

$$\underline{U}_i^{n+1} = \underline{U}_i^n - \frac{\Delta t}{\Delta x_i}\left[\left(\tilde{\mathbf{R}}\tilde{\mathbf{\Lambda}}^-\tilde{\mathbf{R}}^{-1}\Delta\underline{U}\right)_{i+\frac{1}{2}} + \left(\tilde{\mathbf{R}}\tilde{\mathbf{\Lambda}}^+\tilde{\mathbf{R}}^{-1}\Delta\underline{U}\right)_{i-\frac{1}{2}}\right] + \frac{\Delta t}{\Delta x_i}\,\underline{\mathbf{S}}_i^*\,, \quad (2.7)$$

in which

$$\tilde{\mathbf{\Lambda}}^\pm = \frac{1}{2}(\tilde{\mathbf{\Lambda}} \pm |\tilde{\mathbf{\Lambda}}|)\,. \quad (2.8)$$

This splits the update into contributions related to right-going ($+$) and left-going ($-$) characteristics in the decomposition. It follows that the solution is updated using only contributions from the wave perturbations of the Riemann problems at the nodes which enter the cell under consideration, as illustrated in Figure 2.2. It remains to choose an appropriate form for the numerical source term integral $\underline{\mathbf{S}}^*$.

## 2.1   Source terms

This work follows much recent research into source term discretisation, see for example [5, 6, 1, 4], which has concentrated on the use of a characteristic decomposition of the type shown in (2.5). This similarly projects the source term integral onto the eigenvectors of the flux Jacobian $\tilde{\mathbf{A}}$, so that in its linearised

form it can be expressed as

$$\int_{x_i}^{x_{i+1}} \underline{S} \, \mathrm{d}x \;\approx\; \tilde{\underline{\mathbf{S}}}_{i+\frac{1}{2}} \;=\; \left( \tilde{\mathbf{R}} \, \tilde{\mathbf{R}}^{-1} \tilde{\underline{\mathbf{S}}} \right)_{i+\frac{1}{2}} \;=\; \left( \sum_{k=1}^{N_w} \tilde{\beta}_k \tilde{\underline{r}}_k \right)_{i+\frac{1}{2}} , \qquad (2.9)$$

where $\tilde{\beta}_k$, the coefficients of the decomposition, are the components of the vector $\tilde{\mathbf{R}}^{-1}\tilde{\underline{\mathbf{S}}}$. Note that the integral approximated in (2.9) is over a dual cell of the mesh (associated with the interface $i + \frac{1}{2}$), and can be easily incorporated within the fluctuation-signal form of the finite volume scheme given by (2.7). $\underline{\mathbf{S}}_i^*$ will be constructed out of contributions from both ends of the cell, with consistency assured as long as the whole of each dual cell integral (2.9) is distributed.

It is useful (though less so than in higher dimensions) to note here that the analytical form of the source term can be split up into components which can be discretised separately, *i.e.*

$$\underline{S} \;=\; \underline{S}^0 + \sum_j \underline{S}_j^1 \frac{\partial S_j^2}{\partial x} , \qquad (2.10)$$

so that its integral can be approximated consistently by

$$\int_{x_i}^{x_{i+1}} \underline{S} \, \mathrm{d}x \;\approx\; \tilde{\underline{\mathbf{S}}}_{i+\frac{1}{2}} \;=\; \left( \Delta x \underline{\tilde{S}}^0 + \sum_j \underline{\tilde{S}}_j^1 \, \Delta S_j^2 \right)_{i+\frac{1}{2}} , \qquad (2.11)$$

and comparison with (2.9) leads directly to the coefficients $\tilde{\beta}_k$ of the characteristic decomposition of $\tilde{\underline{\mathbf{S}}}_{i+\frac{1}{2}}$.

The terms within the sum on the right hand side of (2.11) may be called upon to balance components of the flux difference $\Delta \underline{F}$ (2.5) so they must be linearised in the same way to ensure that, for the chosen equilibrium state,

$$\underline{F}_x - \underline{S} \equiv \underline{0} \quad \Rightarrow \quad \Delta \underline{F}_{i+\frac{1}{2}} - \tilde{\underline{\mathbf{S}}}_{i+\frac{1}{2}} \;=\; \underline{0} \qquad (2.12)$$

throughout the domain. This follows because at this equilibrium the decompositions (2.5) and (2.9) have been constructed to give $\tilde{\mathbf{\Lambda}}\tilde{\mathbf{R}}^{-1} \, \Delta \underline{U} = \tilde{\mathbf{R}}^{-1}\tilde{\underline{\mathbf{S}}}$ (or alternatively $\tilde{\alpha}_k \tilde{\lambda}_k = \tilde{\beta}_k$). Hence $\tilde{\cdot}$ still represents the evaluation of a quantity at the Roe-average state.

The first term on the right hand side of (2.11) contains only contributions which provide no exact balance with the flux derivatives (*e.g.* bed friction terms in the shallow water equations), so the precise form of their linearisation is not

prescribed by the above arguments. However, it seems sensible that they should also be evaluated at the same state, given by the Roe-average.

As a result of the characteristic decomposition (2.9), the source terms may be discretised in an 'upwind' manner (although, since none of the components has an inherent upwind direction, this must be taken from the corresponding flux component). This leads straightforwardly to an appropriate upwind fluctuation-signal formulation for the first order scheme (2.7) with source terms, given by

$$
\underline{U}_i^{n+1} = \underline{U}_i^n - \frac{\Delta t}{\Delta x_i} \left[ \left( \tilde{\mathbf{R}}(\tilde{\mathbf{\Lambda}}^- \tilde{\mathbf{R}}^{-1}\Delta\underline{U} - \mathbf{I}^- \tilde{\mathbf{R}}^{-1}\underline{\tilde{\mathbf{S}}}) \right)_{i+\frac{1}{2}} \right.
$$
$$
\left. + \left( \tilde{\mathbf{R}}(\tilde{\mathbf{\Lambda}}^+ \tilde{\mathbf{R}}^{-1}\Delta\underline{U} - \mathbf{I}^+ \tilde{\mathbf{R}}^{-1}\underline{\tilde{\mathbf{S}}}) \right)_{i-\frac{1}{2}} \right], \quad (2.13)
$$

in which $\mathbf{I}^\pm = \tilde{\mathbf{\Lambda}}^{-1}\tilde{\mathbf{\Lambda}}^\pm$. The correct balance follows immediately from (2.12).

It is not immediately clear though, how the discretisation of the source term implied by (2.13) can be converted into a numerical source integral $\underline{\mathbf{S}}_i^*$ so that the same balance can be achieved within the flux-based form of the scheme (2.2), particularly when it is extended to higher order. Previous attempts have only proved successful for first order schemes.
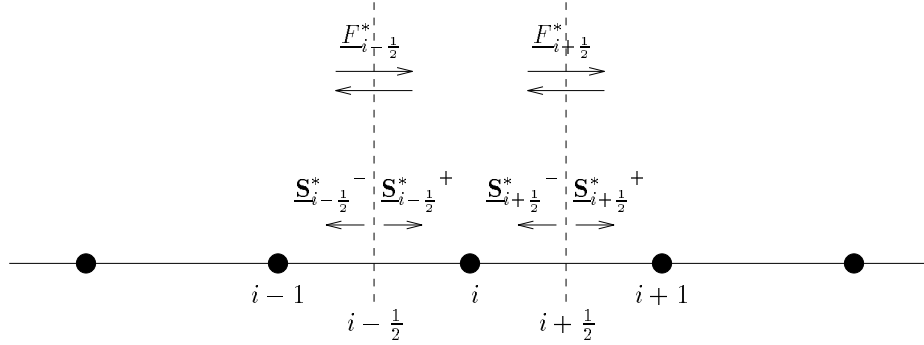


Figure 2.3: Flux and source distribution within a control volume.

The difficulties which arise (and the solution to the problem) can be highlighted by following the transformation of (2.13) into an equation corresponding to (2.2). With a small amount of algebraic manipulation (2.13) becomes

$$
\underline{U}_i^{n+1} = \underline{U}_i^n - \frac{\Delta t}{2\Delta x_i} \left[ \left( \tilde{\mathbf{R}}(\tilde{\mathbf{\Lambda}}\tilde{\mathbf{R}}^{-1}\Delta\underline{U} - \tilde{\mathbf{R}}^{-1}\underline{\tilde{\mathbf{S}}}) \right)_{i+\frac{1}{2}} \right.
$$
$$
\left. + \left( \tilde{\mathbf{R}}(\tilde{\mathbf{\Lambda}}\tilde{\mathbf{R}}^{-1}\Delta\underline{U} - \tilde{\mathbf{R}}^{-1}\underline{\tilde{\mathbf{S}}}) \right)_{i-\frac{1}{2}} \right]
$$
$$
+ \frac{\Delta t}{2\Delta x_i} \left[ \left( \tilde{\mathbf{R}}(|\tilde{\mathbf{\Lambda}}|\tilde{\mathbf{R}}^{-1}\Delta\underline{U} - \text{sgn}(\mathbf{I})\tilde{\mathbf{R}}^{-1}\underline{\tilde{\mathbf{S}}}) \right)_{i+\frac{1}{2}}
$$

9

$$- \left( \tilde{\mathbf{R}} (|\tilde{\mathbf{\Lambda}}| \tilde{\mathbf{R}}^{-1} \Delta \underline{U} - \text{sgn}(\mathbf{I}) \tilde{\mathbf{R}}^{-1} \tilde{\underline{S}}) \right)_{i-\frac{1}{2}} \bigg] , \qquad (2.14)$$

in which $\text{sgn}(\mathbf{I}) = \tilde{\mathbf{\Lambda}}^{-1} |\tilde{\mathbf{\Lambda}}|$. Since (2.5) and (2.9) hold, and

$$\Delta \underline{F}_{i+\frac{1}{2}} + \Delta \underline{F}_{i-\frac{1}{2}} = (\underline{F}_{i+1} + \underline{F}_i) - (\underline{F}_i + \underline{F}_{i-1}) , \qquad (2.15)$$

it follows that the scheme (2.14) can be simplified to

$$\underline{U}_i^{n+1} = \underline{U}_i^n - \frac{\Delta t}{\Delta x_i} \left( \underline{F}_{i+\frac{1}{2}}^* - \underline{F}_{i-\frac{1}{2}}^* \right) + \frac{\Delta t}{\Delta x_i} \left( \underline{\mathbf{S}}_{i+\frac{1}{2}}^{*\ -} + \underline{\mathbf{S}}_{i-\frac{1}{2}}^{*\ +} \right) . \qquad (2.16)$$

The numerical fluxes, $\underline{F}_{i+\frac{1}{2}}^*$ and $\underline{F}_{i-\frac{1}{2}}^*$, are precisely those defined by (2.6) and the numerical source term integral of (2.2) is given by

$$\underline{\mathbf{S}}_i^* = \underline{\mathbf{S}}_{i+\frac{1}{2}}^{*\ -} + \underline{\mathbf{S}}_{i-\frac{1}{2}}^{*\ +} , \qquad (2.17)$$

where

$$\underline{\mathbf{S}}_{i+\frac{1}{2}}^{*\ -} = \frac{1}{2} \left( \tilde{\mathbf{R}} (\mathbf{I} - \text{sgn}(\mathbf{I})) \tilde{\mathbf{R}}^{-1} \tilde{\underline{S}} \right)_{i+\frac{1}{2}} = \left( \tilde{\mathbf{R}} \mathbf{I}^- \tilde{\mathbf{R}}^{-1} \tilde{\underline{S}} \right)_{i+\frac{1}{2}} \qquad (2.18)$$

and

$$\underline{\mathbf{S}}_{i-\frac{1}{2}}^{*\ +} = \frac{1}{2} \left( \tilde{\mathbf{R}} (\mathbf{I} + \text{sgn}(\mathbf{I})) \tilde{\mathbf{R}}^{-1} \tilde{\underline{S}} \right)_{i-\frac{1}{2}} = \left( \tilde{\mathbf{R}} \mathbf{I}^+ \tilde{\mathbf{R}}^{-1} \tilde{\underline{S}} \right)_{i-\frac{1}{2}} . \qquad (2.19)$$

Note that because the numerical source integral cannot, in general, be written as a difference, nothing similar to (2.15) can be applied to it to allow it to be included within the numerical flux (2.6). This means that the balance which is sought between flux derivatives and sources in the flux-based scheme can only be obtained locally by balancing non-zero fluxes through the edges of a control volume, and not by setting each edge flux to zero. One important consequence of this is that the most sensible method of applying the boundary conditions to the numerical scheme is through the addition of ghost cells, since this requires no further correction to maintain the balance which is sought. The distribution of the numerical fluxes and source term components is shown in Figure 2.3.

It is of course possible to overcome the above problem when the source term takes the form of a derivative. If this is the case the source simply augments the conservative flux in the scheme (2.2), *i.e.* given $\underline{S} = \underline{G}_x$ then

$$\underline{F}^* \rightarrow \underline{F}^* - \underline{G}^* , \qquad (2.20)$$

and $\underline{\mathbf{S}}^*$ becomes obsolete. In some cases it may also be possible to incorporate some part of the source term which can be expressed as a derivative within the numerical flux, and then apply an appropriate discretisation to the remaining component of the source.

### 2.1.1  Flux limited schemes

The approach presented in the previous section is no different to the standard upwind technique for approximating source terms when a first order upwind flux discretisation is being used [5]. The only new aspect is the way it has been written, splitting the dual cell source integral into two parts. Usually though, accuracy of higher than first order is required for practical calculations.

The accuracy of Roe's scheme is improved, without introducing spurious oscillations into the solution, by the application of flux limiting techniques [15, 8]. These ensure second order accuracy in smooth regions of the flow, whilst enforcing a Total Variation Diminishing (TVD) property. It is achieved by including a high order correction term in the numerical flux, which becomes [13]

$$\underline{F}^*_{i+\frac{1}{2}} = \frac{1}{2}\left(\underline{F}_{i+1} + \underline{F}_i\right) - \frac{1}{2}\left(\tilde{\mathbf{R}}|\tilde{\mathbf{\Lambda}}|\mathbf{L}\tilde{\mathbf{R}}^{-1}\Delta\underline{U}\right)_{i+\frac{1}{2}}, \tag{2.21}$$

in which $\mathbf{L} = \mathrm{diag}(1 - L(r_k)(1 - |\nu_k|))$, where $\nu_k = \tilde{\lambda}_k\Delta t/\Delta x$ is the Courant number associated with the $k^{\mathrm{th}}$ component of the decomposition, $L$ is a nonlinear flux limiter function, as described in [8, 15], and

$$r_k = \frac{\tilde{\alpha}_k^{\mathrm{upwind}}}{\tilde{\alpha}_k^{\mathrm{local}}}. \tag{2.22}$$

It is clear that a corresponding high order correction must also be made to the source term approximation, and its form can be derived simply by comparing the numerical sources of (2.18,2.19) with the numerical fluxes in (2.6), all of which have been split into two parts which are balanced separately. The flux limiter is only applied to the second part of the numerical flux, so the flux limited numerical source which maintains the balance achieved by the first order discretisation takes the form

$$\underline{\mathbf{S}}^{*}_{i+\frac{1}{2}}{}^{-} = \frac{1}{2}\left(\tilde{\mathbf{R}}(\mathbf{I} - \mathrm{sgn}(\mathbf{I})\mathbf{L})\tilde{\mathbf{R}}^{-1}\underline{\tilde{\mathbf{S}}}\right)_{i+\frac{1}{2}}, \tag{2.23}$$

with a similar expression for $\underline{\mathbf{S}}^*_{i-\frac{1}{2}}{}^+$ in (2.17). Note that since (2.23) is an edge-based quantity, it is simple to evaluate with the fluxes and include within the numerical model.

At this point it should be emphasised that the TVD condition which the flux limiter has been constructed to satisfy applies to the homogeneous system of conservation laws, and the inclusion of source terms means that spurious oscillations may appear in the final solution. The same is true of the slope limited schemes of the next section. This problem has not been addressed in this work.

### 2.1.2 Slope limited schemes

The same balance is slightly more difficult to achieve when the high resolution scheme is constructed using a MUSCL-type slope limiting approach [17]. This is because the underlying representation of the solution is now taken to be linear within each cell so that (2.15) is no longer true. It can though, be replaced by the more general expression,

$$\Delta \underline{F}_{i+\frac{1}{2}} + \Delta \underline{F}_{i-\frac{1}{2}} \;=\; (\underline{F}^{\mathrm{R}}_{i+\frac{1}{2}} + \underline{F}^{\mathrm{L}}_{i+\frac{1}{2}}) - (\underline{F}^{\mathrm{R}}_{i-\frac{1}{2}} + \underline{F}^{\mathrm{L}}_{i-\frac{1}{2}}) - 2(\underline{F}^{\mathrm{L}}_{i+\frac{1}{2}} - \underline{F}^{\mathrm{R}}_{i-\frac{1}{2}}), \quad (2.24)$$

where the superscripts $\cdot^{\mathrm{R}}$ and $\cdot^{\mathrm{L}}$ represent evaluation on, respectively, the right and left hand sides of the interface indicated by the associated subscript (as shown in Figure 2.4). The corresponding numerical flux is

$$\underline{F}^*_{i+\frac{1}{2}} \;=\; \frac{1}{2}\left(\underline{F}^{\mathrm{R}}_{i+\frac{1}{2}} + \underline{F}^{\mathrm{L}}_{i+\frac{1}{2}}\right) - \frac{1}{2}\left(\tilde{\mathbf{R}}|\tilde{\mathbf{\Lambda}}|\tilde{\mathbf{R}}^{-1}\Delta\underline{U}\right)_{i+\frac{1}{2}}, \quad (2.25)$$

in which the Roe-averages are now evaluated from the reconstructed piecewise linear solution. An appropriate correction must therefore be made to the numerical source within each cell, and this leads to

$$\underline{\mathbf{S}}^*_i \;=\; \left(\underline{\mathbf{S}}^*_{i+\frac{1}{2}}{}^- + \underline{\mathbf{S}}^*_{i-\frac{1}{2}}{}^+\right) - \tilde{\underline{\mathbf{S}}}\left(\underline{U}^{\mathrm{L}}_{i+\frac{1}{2}}, \underline{U}^{\mathrm{R}}_{i-\frac{1}{2}}\right). \quad (2.26)$$

The first term on the right hand side is evaluated precisely as before, in (2.17), except that the interface values are now those of the MUSCL reconstruction of the solution within each cell. $\tilde{\underline{\mathbf{S}}}$ is simply the source term integral approximated over the mesh cell (*cf.* (2.11)), and hence evaluated at the Roe-average of the left and right states of the linear reconstruction of the solution within the cell. In

12

terms of the approximations (2.12) and (2.11) the extra term can be thought of as a correction to the integral of the source term over the dual cell arising from the linear variation of the approximation.
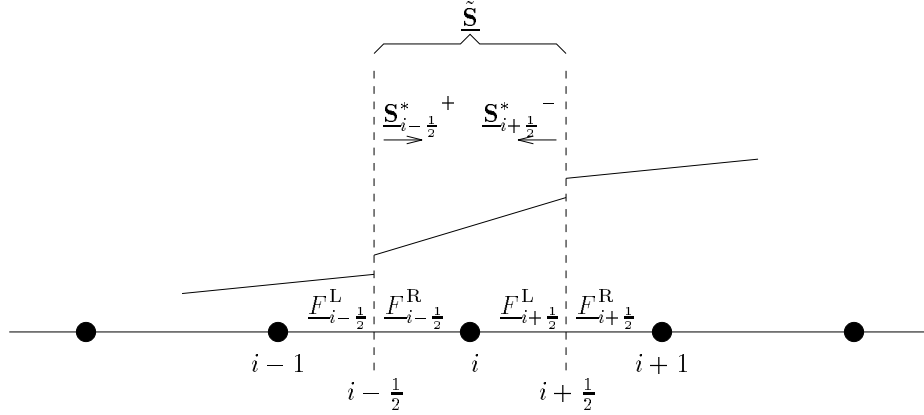


Figure 2.4: Flux and source evaluation for the MUSCL scheme.

## 2.2  Spatially dependent fluxes

In some situations the flux may depend on quantities other than the flow variables and the numerical scheme needs to be modified appropriately. Only one extra spatially varying quantity will be considered here but the approach is easily extended to any number. Returning to Equation (2.1) and taking $\underline{F} = \underline{F}(\underline{U}, B(x))$, where $B$ varies independently of $\underline{U}$, requires a modification to the characteristic decomposition, so that (2.5) becomes

$$
\begin{aligned}
\Delta \underline{F}_{i+\frac{1}{2}} &= \left( \tilde{\mathbf{A}} \, \Delta \underline{U} + \underline{\tilde{V}} \right)_{i+\frac{1}{2}} \\
&= \left( \tilde{\mathbf{R}} \tilde{\mathbf{\Lambda}} \tilde{\mathbf{R}}^{-1} \, \Delta \underline{U} + \tilde{\mathbf{R}} \, \tilde{\mathbf{R}}^{-1} \underline{\tilde{V}} \right)_{i+\frac{1}{2}} \\
&= \left( \sum_{k=1}^{N_w} \tilde{\alpha}_k \tilde{\lambda}_k \underline{\tilde{r}}_k + \sum_{k=1}^{N_w} \tilde{\gamma}_k \underline{\tilde{r}}_k \right)_{i+\frac{1}{2}} ,
\end{aligned}
\tag{2.27}
$$

where $\underline{\tilde{V}} \approx \frac{\partial F}{\partial B} \Delta B$ and $\tilde{\gamma}_k$, the coefficients of the decomposition of this extra term, are the components of $\tilde{\mathbf{R}}^{-1} \underline{\tilde{V}}$. (2.27) gives a set of $N_w$ equations in $N_w + 1$ unknowns, which are taken to be a set of consistent Roe-averaged independent variables from which $\underline{\tilde{U}}$ and $\tilde{B}$ (and all other variables) can be evaluated. This leaves one degree of freedom which can be used by enforcing $\Delta \underline{F} - \underline{\tilde{S}} = \underline{0}$ at an

13

appropriate state of equilibrium (*cf.* (2.11) in which the equilibrium was achieved automatically using the original averages).

Following the same steps as in Section 2 to transform the fluctuation-signal scheme to the flux-based scheme, but including this extra term in the flux difference, leads to precisely the same form for the scheme when approximating the homogeneous system as shown in (2.2), but with new expressions for the numerical fluxes, given by

$$\underline{F}^*_{i+\frac{1}{2}} = \frac{1}{2}(\underline{F}_{i+1} + \underline{F}_i) - \frac{1}{2}\left(\tilde{\mathbf{R}}|\tilde{\mathbf{\Lambda}}|\tilde{\mathbf{R}}^{-1}\,\Delta\underline{U} + \tilde{\mathbf{R}}\,\mathrm{sgn}(\mathbf{I})\tilde{\mathbf{R}}^{-1}\tilde{\underline{V}}\right)_{i+\frac{1}{2}}, \qquad (2.28)$$

in the first order case,

$$\underline{F}^*_{i+\frac{1}{2}} = \frac{1}{2}(\underline{F}_{i+1} + \underline{F}_i) - \frac{1}{2}\left(\tilde{\mathbf{R}}|\tilde{\mathbf{\Lambda}}|\mathbf{L}\tilde{\mathbf{R}}^{-1}\,\Delta\underline{U} + \tilde{\mathbf{R}}\,\mathrm{sgn}(\mathbf{I})\mathbf{L}\tilde{\mathbf{R}}^{-1}\tilde{\underline{V}}\right)_{i+\frac{1}{2}}, \qquad (2.29)$$

when the flux limited high resolution scheme is being used, or

$$\underline{F}^*_{i+\frac{1}{2}} = \frac{1}{2}\left(\underline{F}^{\mathrm{R}}_{i+\frac{1}{2}} + \underline{F}^{\mathrm{L}}_{i+\frac{1}{2}}\right) - \frac{1}{2}\left(\tilde{\mathbf{R}}|\tilde{\mathbf{\Lambda}}|\tilde{\mathbf{R}}^{-1}\,\Delta\underline{U} + \tilde{\mathbf{R}}\,\mathrm{sgn}(\mathbf{I})\tilde{\mathbf{R}}^{-1}\tilde{\underline{V}}\right)_{i+\frac{1}{2}}, \qquad (2.30)$$

for the MUSCL scheme, where the averages are now calculated from the linearly reconstructed solution. By including this extra term in the numerical flux it is possible to avoid altering the form of the source term, as was suggested in [10].

## 2.3  Shallow water flows

The shallow water equations have been chosen as the system of equations to illustrate the use of these new techniques. In one dimension, shallow water flow through a rectangular open channel of varying breadth and bed slope. The effects of bed friction may also be included and, as described in Section 2.1, are simple to treat within the new framework without disturbing balance between the other source terms and the flux derivatives. However, since it is this balance which the new discretisation has been constructed to maintain, friction is not included in the following discussion. The remaining system can be modelled by the equations

$$\begin{pmatrix} bd \\ bdu \end{pmatrix}_t + \begin{pmatrix} bdu \\ bdu^2 + \frac{1}{2}gbd^2 \end{pmatrix}_x = \begin{pmatrix} 0 \\ \frac{1}{2}gd^2 b_x + gbdh_x \end{pmatrix}, \qquad (2.31)$$

which, when compared with (2.1) to find $\underline{U}$, $\underline{F}$ and $\underline{S}$, ultimately leads to

$$\frac{\partial \underline{F}}{\partial \underline{U}} = \begin{pmatrix} 0 & 1 \\ gd - u^2 & 2u \end{pmatrix}, \quad \frac{\partial \underline{F}}{\partial B} = \begin{pmatrix} 0 \\ -\frac{1}{2}gd^2 \end{pmatrix}. \qquad (2.32)$$

14

In these equations $d$ is the depth of the flow, $h$ is the depth of the bed below a nominal still water level, $b = b(x)$ is the channel breadth, $u$ is the flow velocity, and $g$ is the acceleration due to gravity. These quantities are depicted in Figure 2.5.
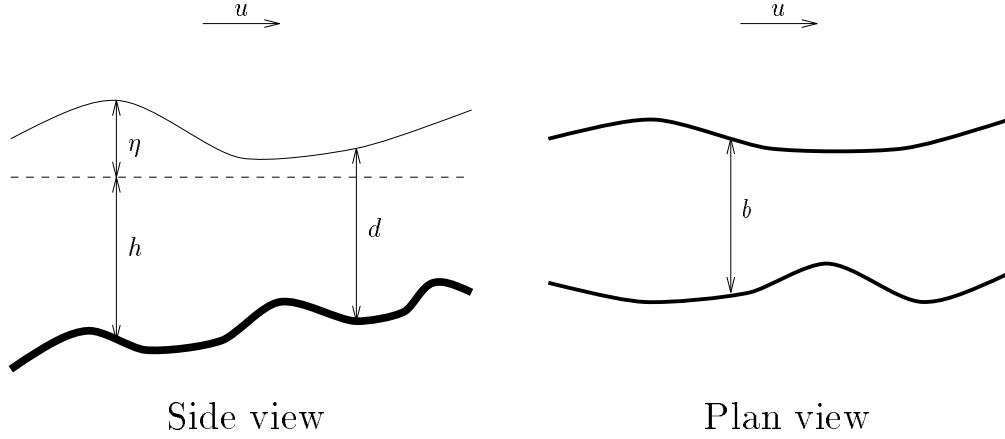


Figure 2.5: The shallow water flow variables.

Equation (2.31) provides an example which includes source terms and a spatial dependence on channel breadth which is independent of the flow. Furthermore, the balance which has been sought in previous sections is illustrated by the steady state represented by still water ($d \equiv h$ and $u \equiv 0$), in which case the system (2.31) reduces to

$$\left(\frac{1}{2}gbd^2\right)_x = \frac{1}{2}gd^2 b_x + gbdh_x \,. \tag{2.33}$$

Previously it has only been possible to maintain this steady state numerically when first order schemes have been used.

The characteristic decomposition (2.27) for the one-dimensional shallow water equations (2.31) and (2.32) is completely defined by

$$\tilde{\alpha}_1 = \frac{\Delta(bd)}{2} + \frac{1}{2\tilde{c}}\left(\Delta(bdu) - \tilde{u}\,\Delta(bd)\right), \quad \tilde{\alpha}_2 = \frac{\Delta(bd)}{2} - \frac{1}{2\tilde{c}}\left(\Delta(bdu) - \tilde{u}\,\Delta(bd)\right)$$

$$\tilde{\lambda}_1 = \tilde{u} + \tilde{c}, \quad \tilde{\lambda}_2 = \tilde{u} - \tilde{c}$$

$$\underline{\tilde{r}}_1 = \begin{pmatrix} 1 \\ \tilde{u} + \tilde{c} \end{pmatrix}, \quad \underline{\tilde{r}}_2 = \begin{pmatrix} 1 \\ \tilde{u} - \tilde{c} \end{pmatrix}$$

$$\tilde{\gamma}_1 = -\frac{1}{4g}\tilde{c}^3\Delta b, \quad \tilde{\gamma}_2 = \frac{1}{4g}\tilde{c}^3\Delta b, \tag{2.34}$$

and it is easily shown that (2.27) is satisfied exactly when

$$\tilde{u} = \frac{\sqrt{b^{\mathrm{R}}d^{\mathrm{R}}}u^{\mathrm{R}} + \sqrt{b^{\mathrm{L}}d^{\mathrm{L}}}u^{\mathrm{L}}}{\sqrt{b^{\mathrm{R}}d^{\mathrm{R}}} + \sqrt{b^{\mathrm{L}}d^{\mathrm{L}}}} \, , \quad \tilde{c}^2 = g\left(\frac{\sqrt{b^{\mathrm{R}}}d^{\mathrm{R}} + \sqrt{b^{\mathrm{L}}}d^{\mathrm{L}}}{\sqrt{b^{\mathrm{R}}} + \sqrt{b^{\mathrm{L}}}}\right) , \tag{2.35}$$

which reduce to the Roe-averages for one-dimensional shallow water flow described in [5] in the absence of breadth variation (*i.e.* when $b^{\mathrm{R}} = b^{\mathrm{L}}$). The corresponding decomposition of the source terms (2.9) then leads to

$$\tilde{\beta}_1 = \frac{1}{4g}\tilde{c}^3\Delta b + \frac{1}{2}\tilde{b}\tilde{c}\Delta h = -\tilde{\beta}_2 \tag{2.36}$$

In order for (2.18) and (2.19) to maintain the correct balance, *i.e.*

$$\tilde{\alpha}_k\tilde{\lambda}_k + \tilde{\gamma}_k - \tilde{\beta}_k = 0 \quad \forall k \tag{2.37}$$

or equivalently,

$$\tilde{\mathbf{R}}\left(\tilde{\boldsymbol{\Lambda}}\tilde{\mathbf{R}}^{-1}\,\Delta\underline{U} + \tilde{\mathbf{R}}^{-1}\underline{\tilde{V}} - \tilde{\mathbf{R}}^{-1}\underline{\tilde{\mathbf{S}}}\right) = \underline{0} \tag{2.38}$$

when the flow is quiescent, $\tilde{b}$ is constructed so that it satisfies

$$\tilde{b}\Delta h = \Delta(bh) - \tilde{h}\Delta b \, , \tag{2.39}$$

where $\tilde{h}$ is evaluated in a similar manner to $\tilde{d}$,

$$\tilde{h} = \frac{\sqrt{b^{\mathrm{R}}}d^{\mathrm{R}} + \sqrt{b^{\mathrm{L}}}d^{\mathrm{L}}}{\sqrt{b^{\mathrm{R}}} + \sqrt{b^{\mathrm{L}}}} \, , \tag{2.40}$$

so that $d \equiv h \Rightarrow \tilde{d} = \tilde{h}$ throughout the domain. Note that this also requires that $d$ and $h$ are reconstructed in the same manner if the MUSCL high resolution scheme is used.

### 2.3.1 Numerical results

The results presented in this section have been chosen to illustrate the improvement in the approximation using the new techniques by focusing on the following:

- the ability to maintain quiescent flow,

- the accuracy of approximations to both continuous and discontinuous steady state solutions,

- the accuracy of simple time-dependent approximations.

16

These have been studied using a variety of channel geometries.

The geometry for the first test case was proposed by the Working Group On Dam-Break Modelling [3], and the bed and breadth variation of the channel (of length 1500) are depicted in Figure 2.6. The upwind source term treatment described in this paper is compared with a much simpler pointwise discretisations in Figure 2.7 (using a uniform 600 cell grid, so that $\Delta x = 2.5$), which show graphs of water surface level and unit discharge for the numerical steady states which result from quiescent initial conditions ($\eta = d - h = 12.0$ and $u = 0.0$), and applying simple non-reflecting boundary conditions. In this case the initial (still water) conditions should be maintained indefinitely by the numerical scheme.
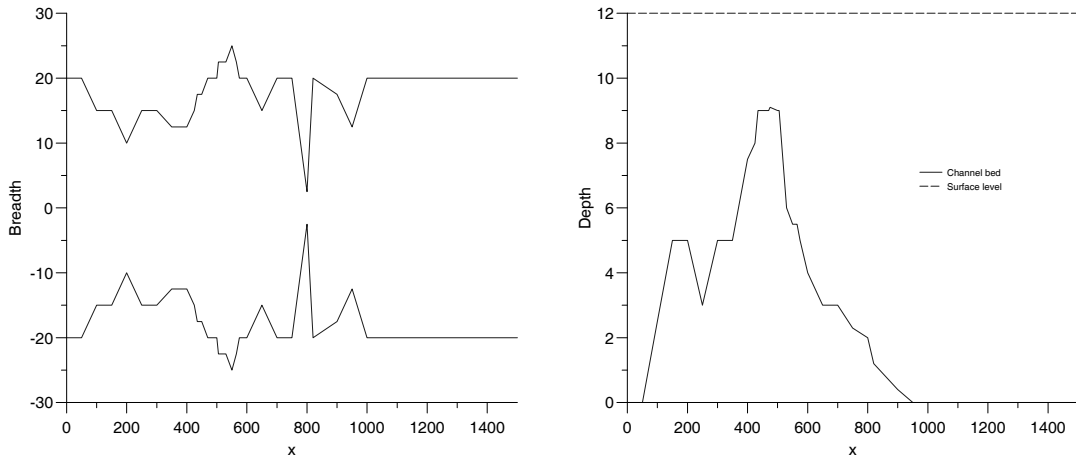


Figure 2.6: Breadth (left) and bed (right) variation for the 'tidal flow' test case.

The comparison is made between first order, slope limited and flux limited schemes combined with pointwise and upwind source term discretisations: in all high resolution cases the Minmod limiter [15] has been applied. The upwind source term discretisations always produce the correct steady state solution, exact to machine accuracy and indistinguishable from the exact solution in the graphs. This is not only true for the first order scheme (which has been achieved previously) but also for the high resolution TVD schemes using any flux or slope limiter on any grid in the presence of bed slope and breadth variations. The pointwise discretisations show small discrepancies (a central discretisation of the source term was also tried but produced even worse results than the pointwise approximation and isn't presented here), most notably in the unit discharge, a

17

quantity which depends on the flow velocity. In each case the method described in Section 2.2 is used to discretise the fluxes where the channel breadth varies.

The second channel geometry which will be used in this work is defined over the interval $[0.0, 3.0]$ and has a smoothly varying depth and breadth, given by

$$
b(x) = \begin{cases} 1.0 - (1.0 - b_{\min}) \cos^2(\pi(x - 1.5)) & \text{for} \quad |x - 1.5| \leq 0.5 \\ 1.0 & \text{otherwise}, \end{cases}
$$

(2.41)

where $b_{\min}$ is the minimum channel breadth, and

$$
h(x) = \begin{cases} 1.0 - z_{\max} \cos^2(\pi(x - 1.5)) & \text{for} \quad |x - 1.5| \leq 0.5 \\ 1.0 & \text{otherwise}, \end{cases}
$$

(2.42)

in which $z_{\max}$ is the maximum height of the bed above the level $\eta = 0.0$. This has been chosen as a simple channel geometry for which exact steady state solutions to the one-dimensional shallow water equations are available for comparison [7]. The parameters chosen to define the channel here are $z_{\max} = 0.1$ and $b_{\min} = 0.9$. A uniform 150 cell grid has been used for each computation.

Three flows are compared:

- $F_\infty = 0.5$, $d_\infty = 1.0$, giving purely subcritical flow which is symmetric about the throat of the constriction (the most narrow point, $x = 1.5$),

- $F_\infty = 0.6$, $d_\infty = 1.0$, giving transcritical flow with a stationary hydraulic jump downstream of the throat and a critical point at the throat,

- $F_\infty = 1.7$, $d_\infty = 1.0$, giving purely supercritical flow which is symmetric about the throat.

The subscript $\cdot_\infty$ represents the freestream flow values 'at infinity' which are used in the application of simple characteristic boundary conditions at inflow and outflow. The results of the comparisons for each of the schemes are shown in Figures 2.8–2.10. The graphs show the variation of total discharge $Q$ through the channel, a quantity which should remain constant at the steady state.

In each case the upwinded source terms can be seen to model the solution better than the pointwise evaluation. This is particularly noticeable in the subcritical
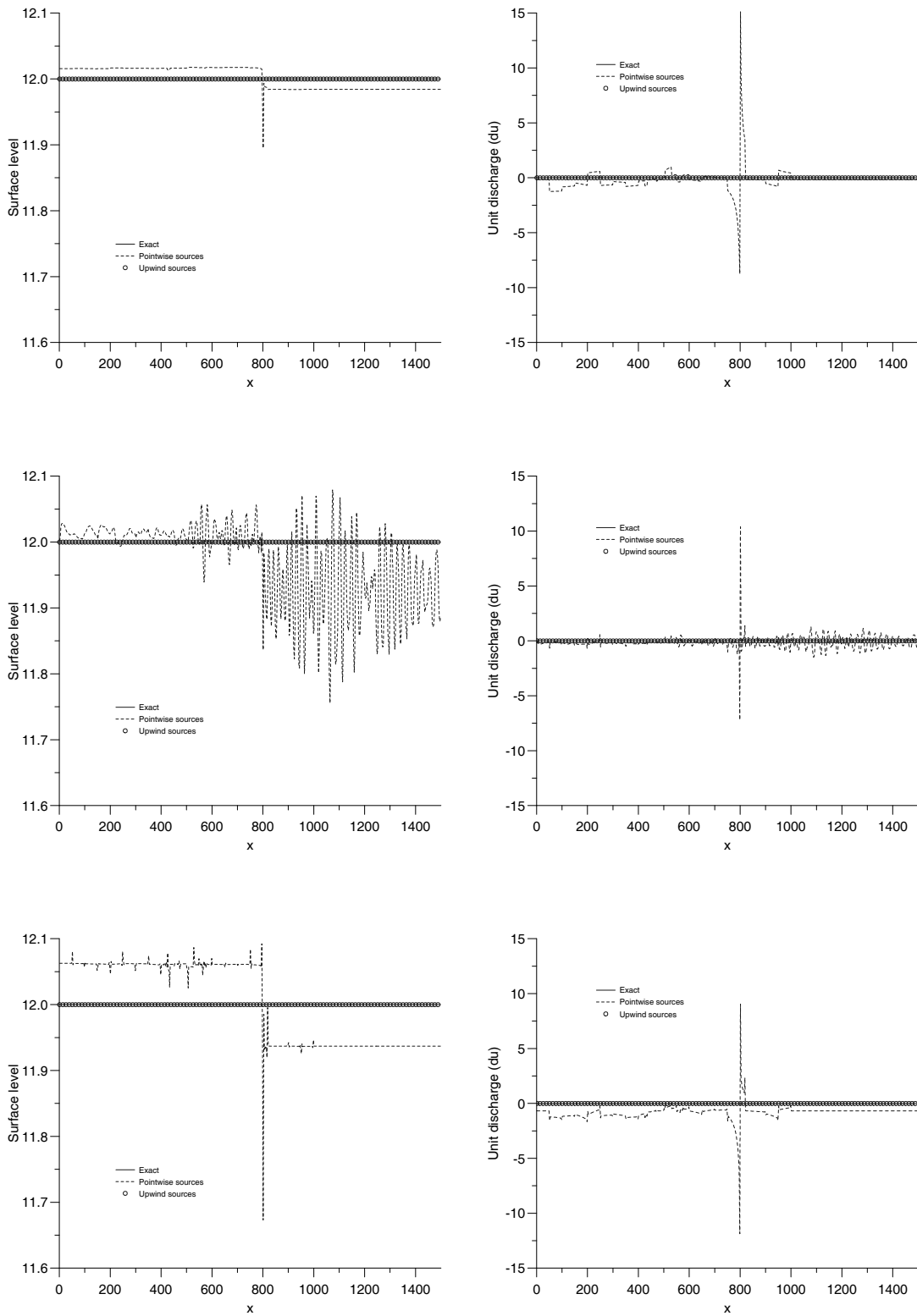
18

Figure 2.7: Water surface level and unit discharge for quiescent flow in a channel with variable bed and breadth, see Figure 2.6, for first order (top) and high resolution slope limited (centre) and flux limited (bottom) schemes ($t = 1000$).
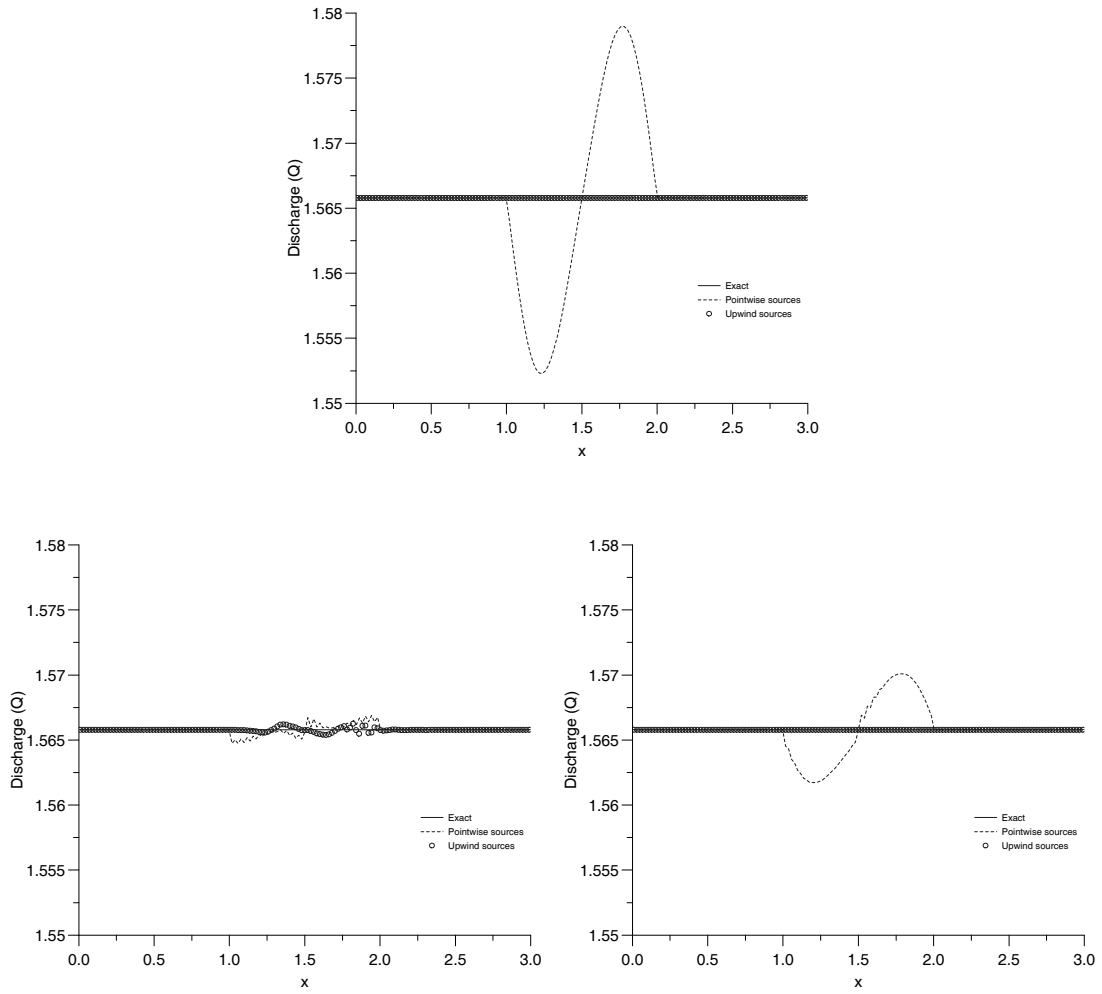
Figure 2.8: Discharge for the steady, subcritical, symmetric constricted channel test case for first order (top) and high resolution slope limited (bottom left) and flux limited (bottom right) schemes.
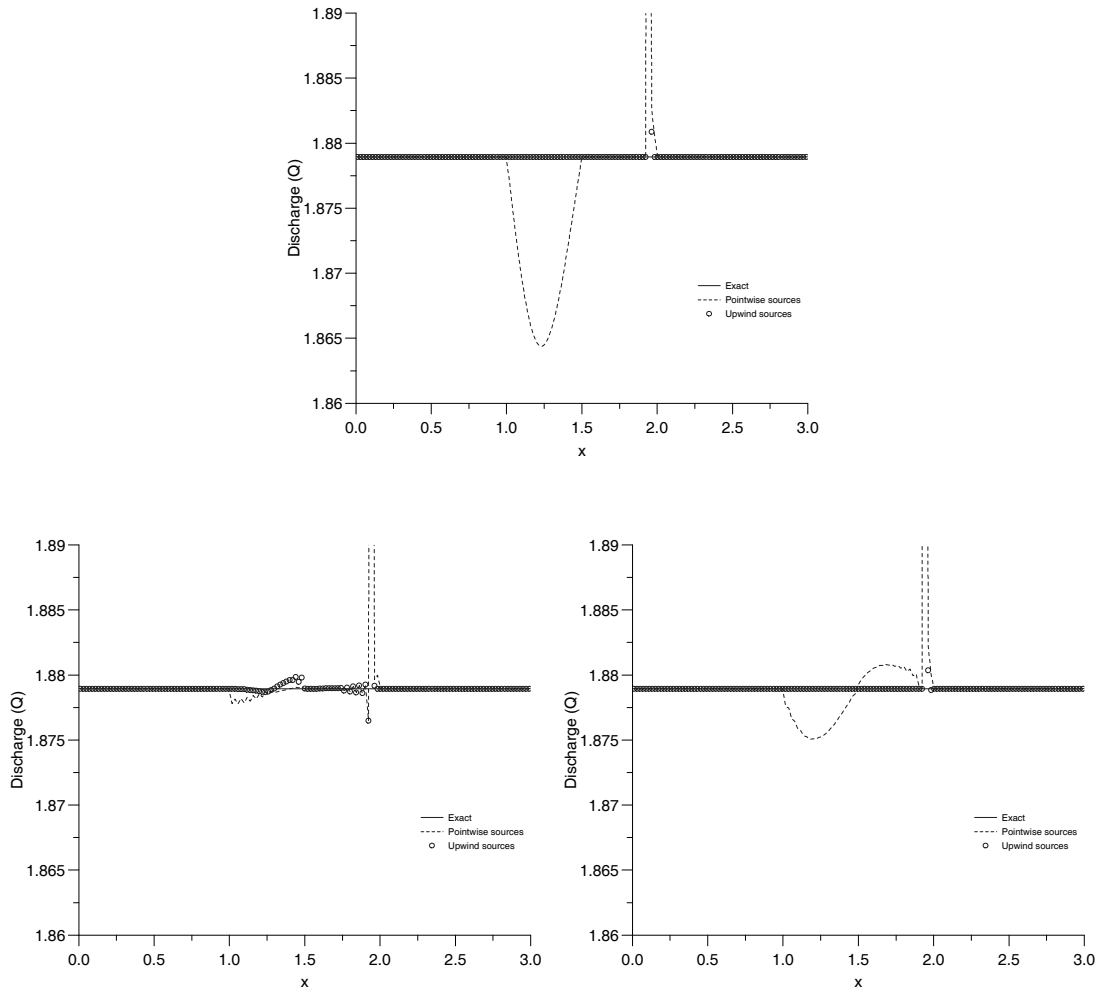
Figure 2.9: Discharge for the steady, transcritical, symmetric constricted channel test case for first order (top) and high resolution slope limited (bottom left) and flux limited (bottom right) schemes.
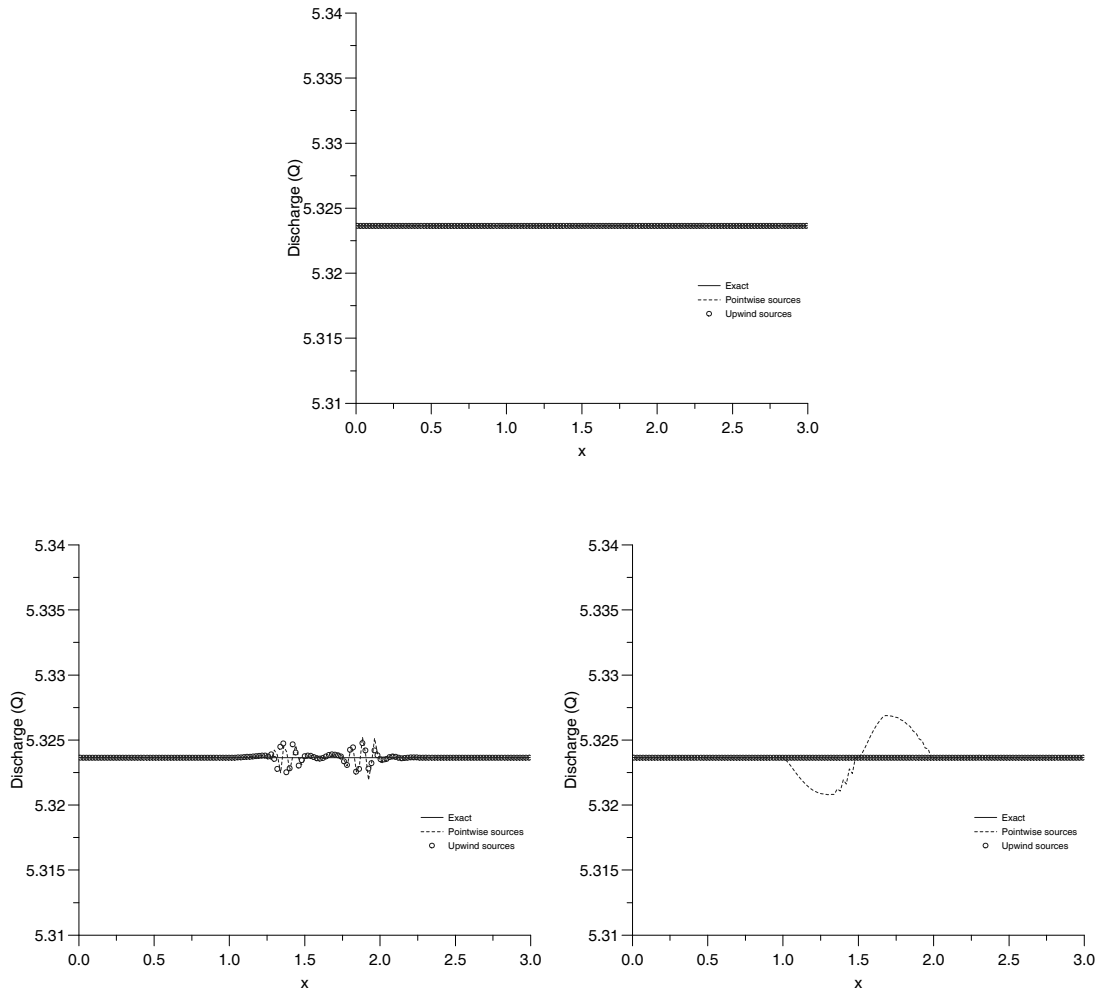
Figure 2.10: Discharge for the steady, supercritical, symmetric constricted channel test case for first order (top) and high resolution slope limited (bottom left) and flux limited (bottom right) schemes.

case where the latter is unable to attain a symmetric solution. The position and strength of the hydraulic jump is predicted accurately by all of the schemes, although there is a small discrepancy in the discharge at the discontinuity in every case. Note that in the second order case small oscillations appear in the 'up-winded', slope limited solution. These are not prohibited by enforcing the TVD condition because this only applies to the homogeneous equations, although they appear in neither the first order nor the flux limited results. This indicates that the correction term of (2.26) may require modification away from the still water steady state.

In [16] it is shown that, for a 'short' channel (of length $L$, taken here to be 1500) and 'low-speed' flow, given the initial conditions

$$d(x,0) \;=\; h(x)\,, \quad q(x,0) \;=\; 0\,, \tag{2.43}$$

where $q = du$ and $h(x)$ is indicated in Figure 2.6, and the boundary conditions

$$d(0,t) \;=\; h(0) + \phi(t)\,, \quad q(L,t) \;=\; \psi(t)\,, \tag{2.44}$$

then a first order approximate solution to the equations (2.31) can be expressed as

$$
\begin{aligned}
d(x,t) \;&=\; h(x) + \phi(t) \\
q(x,t) \;&=\; \psi(t) + \frac{\phi'(t)}{b(x)} \int_x^L b(s)\,\mathrm{d}s\,.
\end{aligned}
\tag{2.45}
$$

The quiescent flow case considered earlier corresponds to taking $\phi(t) \equiv \psi(t) \equiv 0$.

A time-dependent 'tidal' flow test case was suggested in [16] for which

$$\phi(t) \;=\; 4 + 4\sin\left(\frac{(t - 10800)\pi}{21600}\right) \tag{2.46}$$

and $\psi(t) \equiv 0$, the asymptotically exact solution being given by (2.45). The 'exact' and numerical solutions (all computed on the same regular 600 cell grid) to this problem when $t = 10800$ are compared in Figure 2.11. The agreement is very close for the first order and both of the higher order schemes when the upwind source discretisation is used. The main disadvantage of the higher order method is that there is a much stricter practical bound on the CFL number for
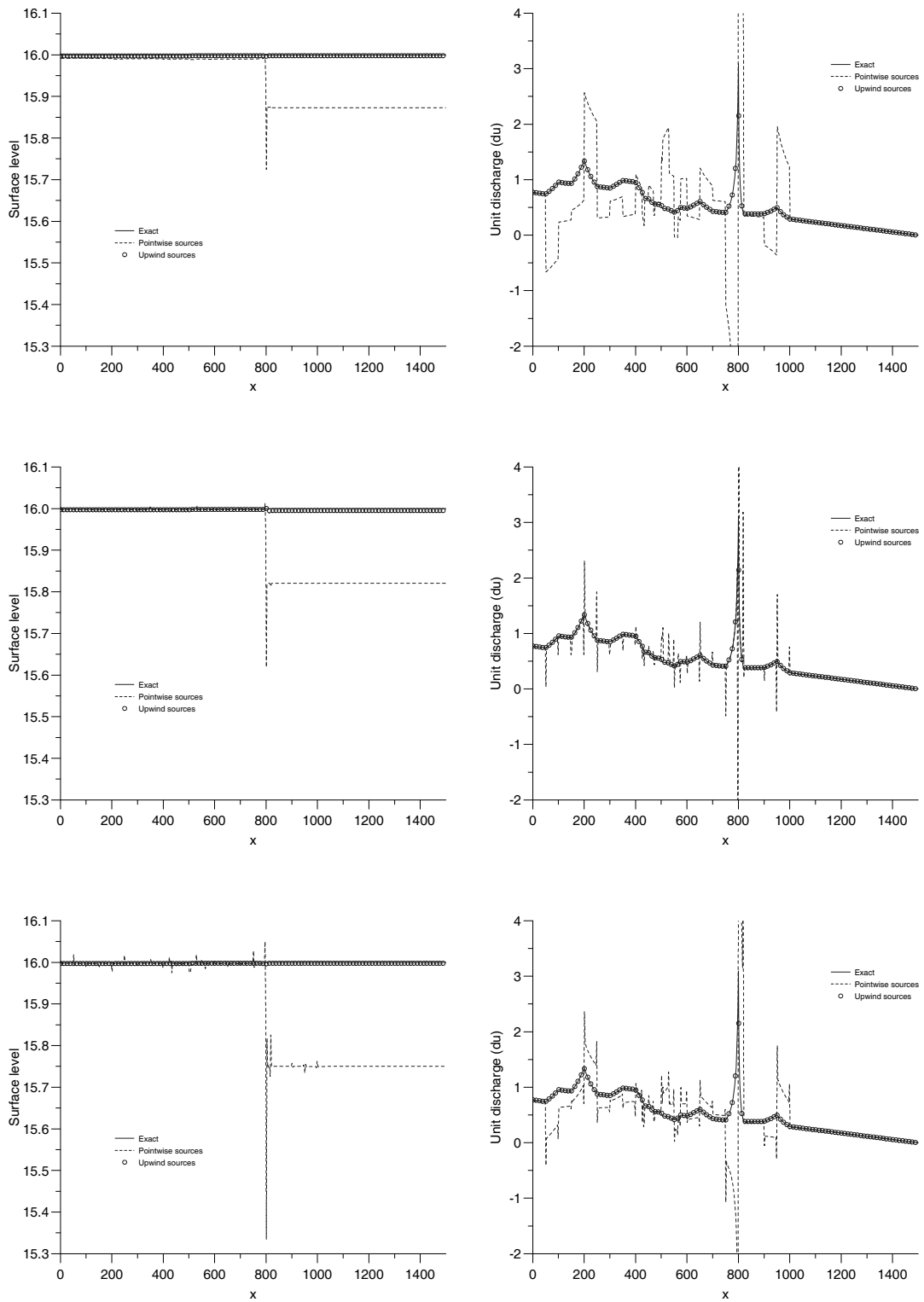
Figure 2.11: Water surface level and unit discharge for the tidal flow test case for first order (top) and high resolution slope limited (centre) and flux limited (bottom) schemes.

the solution to remain free of unwanted oscillations (a value of 0.1 was used in the second order case compared with 0.8 for the first order scheme). At higher CFL numbers the accuracy of the solutions is comparable to the accuracy of those obtained with the simpler source term discretisations. This is due to the fact that the TVD condition only applies in the absence of source terms. Note though that, as in the still water test, even though the pointwise source discretisation gives a reasonable approximation to the depth, it is very poor at predicting the flow velocity.

# 3　Higher dimensions

The following analysis is presented for the two-dimensional case but can be applied simply in three dimensions as well. The conservative form of a system of conservation laws with additional source terms is expressed as

$$\underline{U}_t + \underline{F}_x + \underline{G}_y = \underline{S} \, , \tag{3.1}$$

in which there are now two flux vectors, denoted by $\underline{F} = \underline{F}(\underline{U})$ and $\underline{G} = \underline{G}(\underline{U})$. The case where the fluxes depend on a quantity other than the flow variables is not presented here, having no obvious application to two-dimensional shallow water flows, but can be dealt with in a similar manner to the one-dimensional case presented in Section 2.2.

A combination of a standard finite volume approximation of the flux terms on an arbitrary polygonal mesh (although only triangular and quadrilateral meshes will be considered in the results) and a forward Euler discretisation of the time derivative leads to the conservative difference scheme,

$$\underline{U}_i^{n+1} = \underline{U}_i^n - \frac{\Delta t}{V_i} \sum_{l=1}^{N_e} L_{il} \left( \underline{F}_{il}^*, \underline{G}_{il}^* \right) \cdot \hat{\vec{n}}_{il} + \frac{\Delta t}{V_i} \mathbf{S}_i^* \tag{3.2}$$

where $V_i$ is the area of the chosen control volume, $N_e$ is the number of edges it has, $\hat{\vec{n}}_{il}$ is the outward pointing unit normal to the edge common to cells $i$ and $l$ (where $l$ represents a generic neighbouring cell) and $L_{il}$ is the length of that edge (as shown for a triangular mesh cell in Figure 3.1). $\mathbf{S}^* \approx \iint_{\text{cell}} \underline{S} \, \mathrm{d}x \, \mathrm{d}y$ is once more a numerical approximation to the source integral over the control volume.

For simplicity the scheme will again be assumed to be a cell centre discretisation in which the control volumes coincide with the mesh cells, although the techniques may also be applied to other types of scheme. The following analysis runs along similar lines to that presented in previous sections for the one-dimensional case.
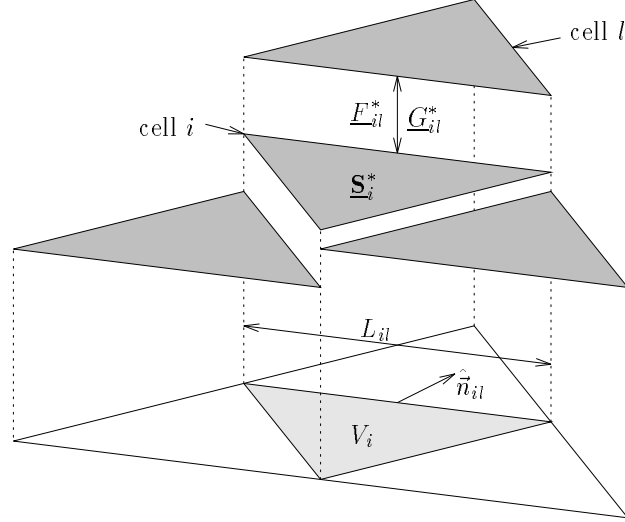


Figure 3.1: Numerical fluxes and sources for the cell centre scheme.

## 3.1 The first order scheme

The numerical fluxes which lead to the first order Roe's scheme in two dimensions are given by

$$(\underline{E}_{il}^*, \underline{G}_{il}^*) \cdot \hat{\tilde{n}}_{il} = \frac{1}{2}(\underline{E}_i + \underline{E}_l, \underline{G}_i + \underline{G}_l) \cdot \hat{\tilde{n}}_{il} - \frac{1}{2}\left(\tilde{\mathbf{R}}|\tilde{\mathbf{\Lambda}}|\tilde{\mathbf{R}}^{-1}\Delta\underline{U}\right)_{il}, \qquad (3.3)$$

in which the eigenvectors and eigenvalues which are needed to construct $\tilde{\mathbf{R}}$ and $\tilde{\mathbf{\Lambda}}$ are now those of the matrix $\tilde{\mathbf{C}}_n = (\tilde{\mathbf{A}}, \tilde{\mathbf{B}}) \cdot \hat{\tilde{n}}$, where

$$\tilde{\mathbf{A}} \approx \frac{\partial \underline{F}}{\partial \underline{U}} \quad \text{and} \quad \tilde{\mathbf{B}} \approx \frac{\partial \underline{G}}{\partial \underline{U}} \qquad (3.4)$$

are the linearised flux Jacobians. It can be seen that the numerical flux is similar in form to that used in one dimension (2.6). In particular, $\tilde{\ }$ again denotes the evaluation of a quantity at its Roe-average state.

Since the two-dimensional scheme is based on Riemann solvers oriented perpendicular to the edges of the grid cells the decomposition also bears a strong resemblance to the one-dimensional case. Once more, as long as the quantities

denoted $\tilde{\cdot}$ are evaluated at the appropriate Roe-average state [11] then the flux differences can be written in the decomposed form

$$\Delta(\underline{F}, \underline{G}) \cdot \hat{\vec{n}} \;=\; \tilde{\mathbf{C}}_n \Delta\underline{U} \;=\; \tilde{\mathbf{R}}\tilde{\mathbf{\Lambda}}\tilde{\mathbf{R}}^{-1}\Delta\underline{U} \;=\; \sum_{k=1}^{N_w} \tilde{\alpha}_k \tilde{\lambda}_k \tilde{\underline{r}}_k \tag{3.5}$$

from which it follows in much the same way as in one dimension that the scheme (3.2) is equivalent to

$$\underline{U}_i^{n+1} \;=\; \underline{U}_i^n - \frac{\Delta t}{V_i}\sum_{l=1}^{N_e} L_{il}\left(\tilde{\mathbf{R}}\tilde{\mathbf{\Lambda}}^-\tilde{\mathbf{R}}^{-1}\Delta\underline{U}\right)_{il} + \frac{\Delta t}{V_i}\,\underline{\mathbf{S}}_i^* \,, \tag{3.6}$$

where the superscript $\cdot^-$ now indicates the incoming characteristics at the appropriate edge of the control volume (see Figure 3.2). It is easily seen that this reduces to (2.7) when restricted to one dimension.
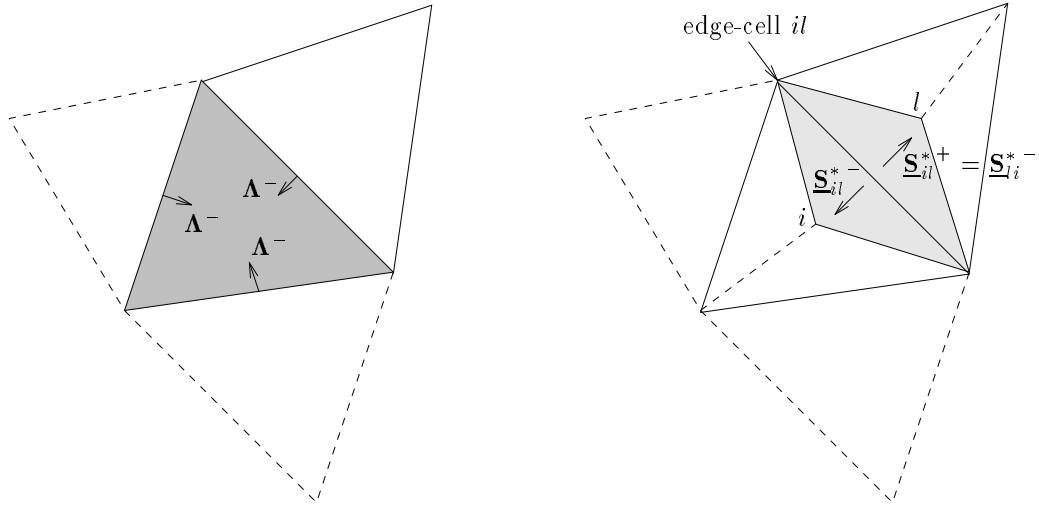


Figure 3.2: Wave propagation directions (left) and source distribution (right) within a triangular cell in two dimensions.

As in one dimension the analytical form of the source term can be split into components to be integrated separately (*cf.* (3.30)) so that

$$\underline{S} \;=\; \underline{S}^0 + \sum_j \underline{S}_j^1 \, \vec{\nabla} \cdot (S_j^x, S_j^y) \,. \tag{3.7}$$

Hence, integrating over an edge-cell and applying the divergence theorem to the terms within the sum leads naturally to the approximation

$$\iint_{\diamond_{il}} \underline{S}\,\mathrm{d}x\,\mathrm{d}y \;\approx\; \tilde{\underline{\mathbf{S}}}_{il} \;=\; \left(V_\diamond \tilde{\underline{S}}^0 + \sum_j \tilde{\underline{S}}_j^1 \oint_{\partial\diamond}(S_j^x, S_j^y)\cdot\mathrm{d}\vec{n}\right)_{il}, \tag{3.8}$$

in which $V_\diamond$ is the area of the edge-cell. Now, given that the solution has already been assumed to be constant in each part of the edge-cell for the purposes of the Riemann solver, and hence the flux evaluation, the approximation reduces to

$$\underline{\tilde{\mathbf{S}}}_{il} \;=\; \left( V_\diamond \underline{\tilde{S}}^0 + \sum_j \underline{\tilde{S}}^1_j \, \Delta(S^x_j, S^y_j) \cdot \vec{n} \right)_{il} , \qquad (3.9)$$

where $\vec{n}$ is the normal to the edge, scaled by its length, see also (3.31). The terms within the sum may again be required to balance the flux difference, so the same Roe linearisation is used in their evaluation, and it follows that

$$\underline{F}_x + \underline{G}_y - \underline{S} \;\equiv\; \underline{0} \quad \Rightarrow \quad \left( \Delta(\underline{F},\underline{G}) \cdot \vec{n} - \underline{\tilde{\mathbf{S}}} \right)_{il} \;=\; \underline{0} \qquad (3.10)$$

throughout the domain: $\diamond_{il}$ is the edge-cell corresponding to the edge between cells $i$ and $l$, as shown in Figure 3.2. The three-dimensional case is similar, with all the approximations being carried out over a face-cell with the solution being assumed constant on either side.

The two-dimensional source term can now be written as a characteristic decomposition similar to that of the flux difference (3.5), *i.e.* its linearisation can take the form

$$\underline{\tilde{\mathbf{S}}}_{il} \;=\; \left( \tilde{\mathbf{R}} \, \tilde{\mathbf{R}}^{-1} \underline{\tilde{\mathbf{S}}} \right)_{il} \;=\; L_{il} \left( \sum_{k=1}^{N_w} \tilde{\beta}_k \underline{\tilde{r}}_k \right)_{il} . \qquad (3.11)$$

Evaluating this at the same Roe-average state as the flux difference means that the correct balance is attained because, at equilibrium, the decompositions give $L\boldsymbol{\Lambda} \mathbf{R}^{-1} \, \Delta\underline{U} = \mathbf{R}^{-1}\underline{S}$. $\underline{S}^*_i$ will be constructed out of contributions from each edge of the cell, with consistency assured as long as the whole of each edge-cell integral (3.11) is distributed.

The decomposition has been carried out so that, when (3.6) is combined with (3.11) to give

$$\underline{U}^{n+1}_i \;=\; \underline{U}^n_i - \frac{\Delta t}{V_i} \sum_{l=1}^{N_e} \left( \tilde{\mathbf{R}}(L\tilde{\boldsymbol{\Lambda}}^- \tilde{\mathbf{R}}^{-1}\Delta\underline{U} - \mathbf{I}^- \tilde{\mathbf{R}}^{-1}\underline{\tilde{\mathbf{S}}}) \right)_{il} , \qquad (3.12)$$

a precise balance can be achieved when one is sought between the sources and the flux gradients.

The relationship between the two forms of the finite volume scheme, (3.2) and (3.6), can now be exploited. Substituting for $\mathbf{I}^-$ in (3.12) gives

$$\underline{U}^{n+1}_i \;=\; \underline{U}^n_i - \frac{\Delta t}{2V_i} \sum_{l=1}^{N_e} \left( \tilde{\mathbf{R}}(L\tilde{\boldsymbol{\Lambda}}\tilde{\mathbf{R}}^{-1}\Delta\underline{U} - \tilde{\mathbf{R}}^{-1}\underline{\tilde{\mathbf{S}}}) \right)_{il}$$

$$-\frac{\Delta t}{2V_i} \sum_{l=1}^{N_e} \left( \tilde{\mathbf{R}}(L|\tilde{\mathbf{\Lambda}}|\tilde{\mathbf{R}}^{-1} \Delta \underline{U} - \mathrm{sgn}(\mathbf{I})\tilde{\mathbf{R}}^{-1}\underline{\tilde{\mathbf{S}}}) \right)_{il} . \quad (3.13)$$

In addition, it is easily shown that

$$\sum_{l=1}^{N_e} \Delta(\underline{F}_{il}, \underline{G}_{il}) \cdot \vec{n}_{il} = \sum_{l=1}^{N_e} (\underline{F}_i + \underline{F}_l, \underline{G}_i + \underline{G}_l) \cdot \vec{n}_{il} , \quad (3.14)$$

in which $\Delta \underline{F}_{il} = \underline{F}_l - \underline{F}_i$ is the jump in $\underline{F}$ across the $l^{\text{th}}$ edge of cell $i$ (and similarly for $\underline{G}$). Therefore, since $\tilde{\cdot}$ indicates evaluation at the Roe-average state, (3.5) holds and (3.13) can be rewritten as

$$\underline{U}_i^{n+1} = \underline{U}_i^n - \frac{\Delta t}{V_i} \sum_{l=1}^{N_e} (\underline{F}_{il}^*, \underline{G}_{il}^*) \cdot \vec{n}_{il} + \frac{\Delta t}{V_i} \underline{\mathbf{S}}_i^* , \quad (3.15)$$

in which the numerical fluxes are given by (3.3) and the numerical source is

$$\underline{\mathbf{S}}_i^* = \sum_{l=1}^{N_e} \underline{\mathbf{S}}_{il}^{*\,-} , \quad (3.16)$$

where

$$\underline{\mathbf{S}}_{il}^{*\,-} = \frac{1}{2} \left( \tilde{\mathbf{R}}(\mathbf{I} - \mathrm{sgn}(\mathbf{I}))\tilde{\mathbf{R}}^{-1}\underline{\tilde{\mathbf{S}}} \right)_{il} = \left( \tilde{\mathbf{R}}\mathbf{I}^-\tilde{\mathbf{R}}^{-1}\underline{\tilde{\mathbf{S}}} \right)_{il} . \quad (3.17)$$

These expressions bear a close resemblance to the numerical fluxes and can be incorporated into the flux-based scheme in a similar manner. As in one dimension it is not possible to combine the source term completely with the numerical fluxes.

## 3.2 High resolution schemes

When the accuracy of the scheme is increased by the use of a flux limiting technique the numerical flux takes the form

$$(\underline{F}_{il}^*, \underline{G}_{il}^*) \cdot \hat{\vec{n}}_{il} = \frac{1}{2}(\underline{F}_i + \underline{F}_l, \underline{G}_i + \underline{G}_l) \cdot \hat{\vec{n}}_{il} - \frac{1}{2} \left( \tilde{\mathbf{R}}\tilde{\mathbf{\Lambda}}\mathbf{L}\tilde{\mathbf{R}}^{-1} \Delta \underline{U} \right)_{il} , \quad (3.18)$$

and the appropriate discretisation of the source term can be shown to be

$$\underline{\mathbf{S}}_{il}^{*\,-} = \frac{1}{2} \left( \tilde{\mathbf{R}}(\mathbf{I} - \mathrm{sgn}(\mathbf{I})\mathbf{L})\tilde{\mathbf{R}}^{-1}\underline{\tilde{\mathbf{S}}} \right)_{il} \quad (3.19)$$

by similar arguments to those used in one dimension.

For a MUSCL-type slope limited higher order numerical scheme, the numerical fluxes take the form

$$(\underline{F}_{il}^*, \underline{G}_{il}^*) \cdot \hat{\vec{n}}_{il} = \frac{1}{2}(\underline{F}_{Il} + \underline{F}_{iL}, \underline{G}_{Il} + \underline{G}_{iL}) \cdot \hat{\vec{n}}_{il} - \frac{1}{2} \left( \tilde{\mathbf{R}}\tilde{\mathbf{\Lambda}}\tilde{\mathbf{R}}^{-1} \Delta \underline{U} \right)_{il} , \quad (3.20)$$

in which the subscripts $\cdot_{Il}$ and $\cdot_{iL}$ represent evaluation of the piecewise linear reconstruction of the solution on, respectively, the inside and the outside of the edge between cells $i$ and $l$, relative to cell $i$ (indicated in Figure 3.3), giving new values from which the Roe-averages at the interface are calculated. Now, instead of (3.14) the flux differences satisfy the more general expression

$$\sum_{l=1}^{N_e} \Delta(\underline{F}_{il}, \underline{G}_{il}) \cdot \vec{n}_{il} = \sum_{l=1}^{N_e} (\underline{F}_{Il} + \underline{F}_{iL}, \underline{G}_{Il} + \underline{G}_{iL}) \cdot \vec{n}_{il}$$
$$-2 \sum_{l=1}^{N_e} (\underline{F}_{Il} - \underline{F}_i, \underline{G}_{Il} - \underline{G}_i) \cdot \vec{n}_{il} \,. \qquad (3.21)$$

Consequently, the numerical source term appropriate to this type of scheme is given by

$$\underline{S}_i^* = \sum_{l=1}^{N_e} \left( \underline{S}_{il}^{*\,-} - \tilde{\underline{S}}(\underline{U}_{Il}, \underline{U}_i) \right) \,, \qquad (3.22)$$

where $\tilde{\phantom{.}}$ indicates the evaluation of the source term integral ($cf.$ (3.9)) at the Roe-average of the specified conservative variables (taken from the linear reconstruction at the midpoints of the cell edges) and $\underline{S}_{il}^{*\,-}$ is taken directly from (3.17). This can again be considered as applying a higher order correction to the integral of the source term over the edge-cell.
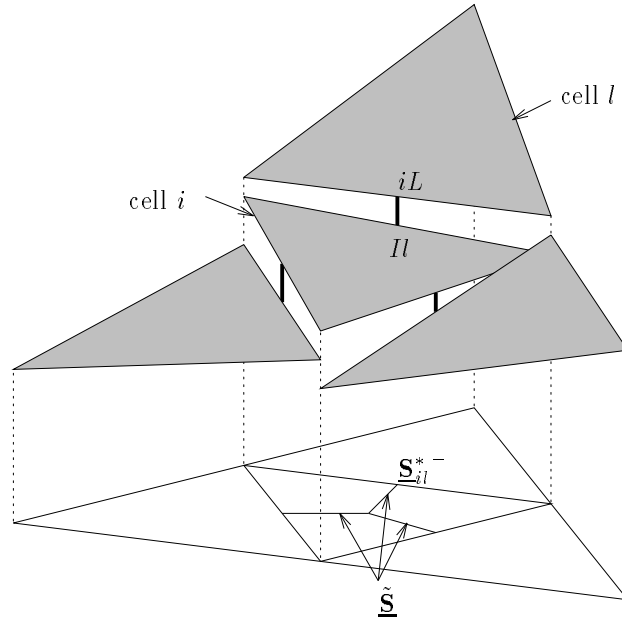


Figure 3.3: Flux and source evaluations for a two-dimensional MUSCL-type scheme on triangles.

## 3.3 Shallow water flows

In two dimensions the shallow water equations including the effects of varying bed slope are obtained by substituting

$$
\underline{U} = \begin{pmatrix} d \\ du \\ dv \end{pmatrix} , \quad \underline{F} = \begin{pmatrix} du \\ du^2 + \frac{gd^2}{2} \\ duv \end{pmatrix} , \quad \underline{G} = \begin{pmatrix} dv \\ duv \\ dv^2 + \frac{gd^2}{2} \end{pmatrix} , \qquad (3.23)
$$

where $v$ is the flow velocity in the $y$-direction in addition to the variables defined for (2.31), and

$$
\underline{S} = \begin{pmatrix} 0 \\ gdh_x \\ gdh_y \end{pmatrix} \qquad (3.24)
$$

into (3.1). The matrix $\mathbf{C}_n$ can be calculated simply from these for any edge orientation.

When $d \equiv h$ and $u \equiv v \equiv 0$ (quiescent flow in two dimensions) the desired balance is given by the equations

$$
\left( \frac{gd^2}{2} \right)_x = gdh_x , \quad \left( \frac{gd^2}{2} \right)_y = gdh_y . \qquad (3.25)
$$

The discretisation should satisfy (3.25) exactly in this special case.

The characteristic decomposition is now carried out on the eigenvectors of the matrix $(\tilde{\mathbf{A}}, \tilde{\mathbf{B}}) \cdot \hat{\vec{n}}$, which are

$$
\underline{\tilde{r}}_1 = \begin{pmatrix} 1 \\ \tilde{u} + \tilde{c}n^x \\ \tilde{v} + \tilde{c}n^y \end{pmatrix} , \quad \underline{\tilde{r}}_2 = \begin{pmatrix} 0 \\ -\tilde{c}n^y \\ \tilde{c}n^x \end{pmatrix} , \quad \underline{\tilde{r}}_3 = \begin{pmatrix} 1 \\ \tilde{u} - \tilde{c}n^x \\ \tilde{v} - \tilde{c}n^y \end{pmatrix} , \qquad (3.26)
$$

in which $(n^x, n^y) = \hat{\vec{n}}$ and

$$
\tilde{c} = \sqrt{\frac{g(d^{\mathrm{R}} + d^{\mathrm{L}})}{2}} ,
$$

$$
\tilde{u} = \frac{\sqrt{d^{\mathrm{R}}}u^{\mathrm{R}} + \sqrt{d^{\mathrm{L}}}u^{\mathrm{L}}}{\sqrt{d^{\mathrm{R}}} + \sqrt{d^{\mathrm{L}}}} , \quad \tilde{v} = \frac{\sqrt{d^{\mathrm{R}}}v^{\mathrm{R}} + \sqrt{d^{\mathrm{L}}}v^{\mathrm{L}}}{\sqrt{d^{\mathrm{R}}} + \sqrt{d^{\mathrm{L}}}} . \qquad (3.27)
$$

The superscripts $\cdot^{\mathrm{R}}$ and $\cdot^{\mathrm{L}}$ indicate here the evaluation of a quantity on either side of a cell edge, at its midpoint. The corresponding expressions for the eigenvalues (wave speeds) are

$$
\lambda_1 = \tilde{u}n^x + \tilde{v}n^y + \tilde{c} , \quad \lambda_2 = \tilde{u}n^x + \tilde{v}n^y , \quad \lambda_3 = \tilde{u}n^x + \tilde{v}n^y - \tilde{c} , \qquad (3.28)
$$

and the wave strengths,

$$\tilde{\alpha}_1 = \frac{\Delta d}{2} + \frac{1}{2\tilde{c}} \left( \Delta(du)n^x + \Delta(dv)n^y - (\tilde{u}n^x + \tilde{v}n^y)\Delta d \right)$$

$$\tilde{\alpha}_2 = \frac{1}{\tilde{c}} \left( (\Delta(dv) - \tilde{v}\Delta d)\, n^x - (\Delta(du) - \tilde{u}\Delta d)\, n^y \right)$$

$$\tilde{\alpha}_3 = \frac{\Delta d}{2} - \frac{1}{2\tilde{c}} \left( \Delta(du)n^x + \Delta(dv)n^y - (\tilde{u}n^x + \tilde{v}n^y)\Delta d \right) , \qquad (3.29)$$

complete the decomposition (3.5).

In this case, in order to provide the desired balance, the source term is written in the form (3.7), giving

$$\underline{S} = \begin{pmatrix} 0 \\ gd \\ 0 \end{pmatrix} \vec{\nabla} \cdot (h, 0) + \begin{pmatrix} 0 \\ 0 \\ gd \end{pmatrix} \vec{\nabla} \cdot (0, h) . \qquad (3.30)$$

At first glance this seems counterproductive, but it immediately allows the source term integral over an edge-cell to be approximated in a manner which will allow the discrete balance with the flux integral, $i.e.$ it can be approximated in the form (3.9) via (3.8). This leads to

$$\underline{\tilde{\mathbf{S}}}_{il} = L_{il} \begin{pmatrix} 0 \\ g\tilde{d}\Delta h n^x \\ g\tilde{d}\Delta h n^y \end{pmatrix} , \qquad (3.31)$$

which is used to obtain the coefficients which are used in the characteristic decomposition (3.11). In this case these are

$$\tilde{\beta}_1 = \frac{1}{2}\tilde{c}\Delta h , \quad \tilde{\beta}_2 = 0 , \quad \tilde{\beta}_3 = -\frac{1}{2}\tilde{c}\Delta h . \qquad (3.32)$$

By construction, it follows that $\tilde{\alpha}_k \tilde{\lambda}_k - \tilde{\beta}_k = 0$ for each $k$, $i.e.$

$$\tilde{\mathbf{R}} \left( L\tilde{\boldsymbol{\Lambda}}\tilde{\mathbf{R}}^{-1} \Delta\underline{U} - \tilde{\mathbf{R}}^{-1}\underline{\tilde{\mathbf{S}}} \right) = \underline{0} , \qquad (3.33)$$

when the flow is quiescent, and the numerical balance is assured.

## 3.4  Numerical results

The test cases presented in this section are essentially a subset of those described in Section 2.3.1 for the one-dimensional schemes, but applied to the

two-dimensional shallow water equations. For the purposes of presentation, comparisons will be made between breadth-averaged solutions for channel flows and exact solutions to the corresponding one-dimensional problem. These will obviously differ slightly in the non-quiescent cases due to the simplifications inherent in the one-dimensional model but still provide an accurate guide when the cross-flow velocity is small, as it is in the results presented.
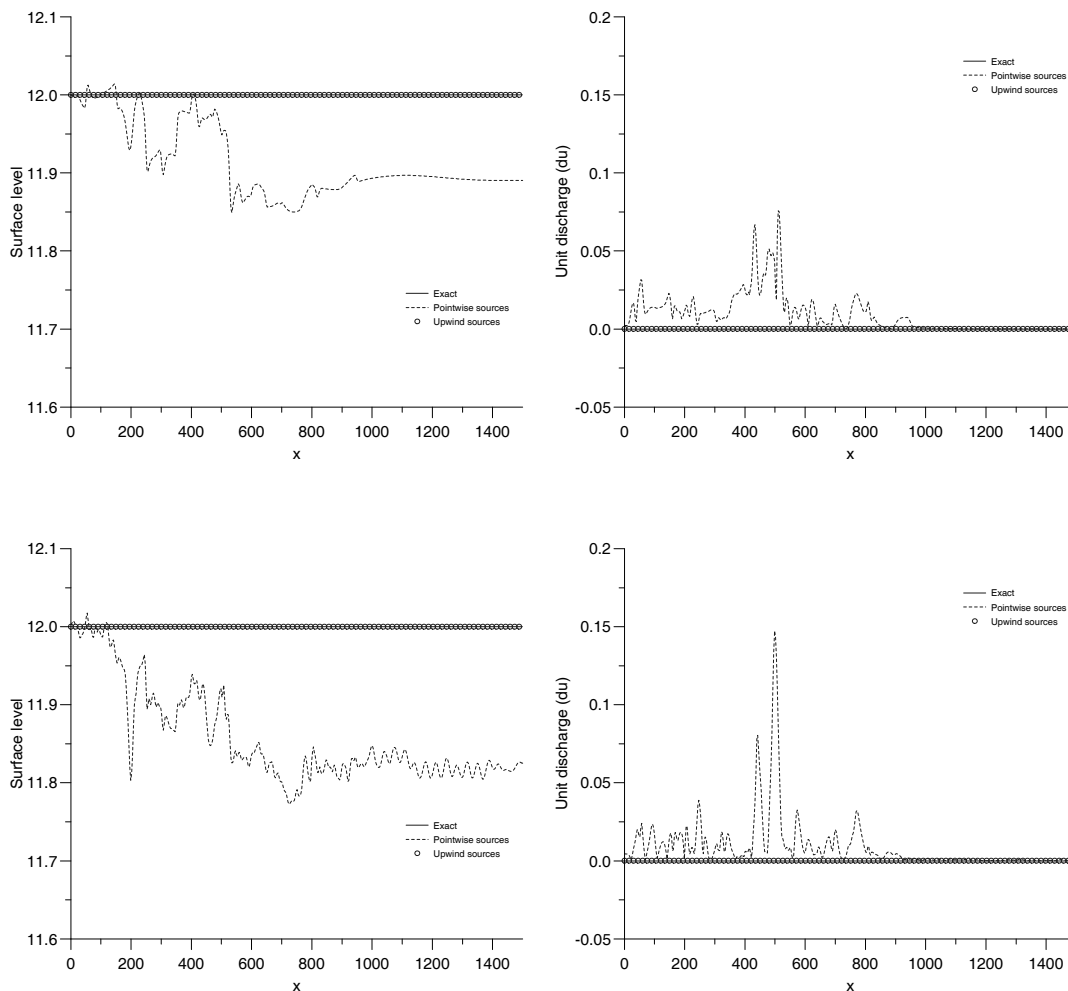


Figure 3.4: Water surface level and unit discharge for the still water test case for first order (top) and high resolution slope limited (bottom) schemes ($t = 1000$).

The ability of the new techniques to maintain the still water steady state is illustrated using the geometry of Figure 2.6 and a triangular grid with 4854 cells and 2738 nodes (giving about 300 cells along the channel, roughly half the one-dimensional grid resolution). As in one dimension, the upwind source term
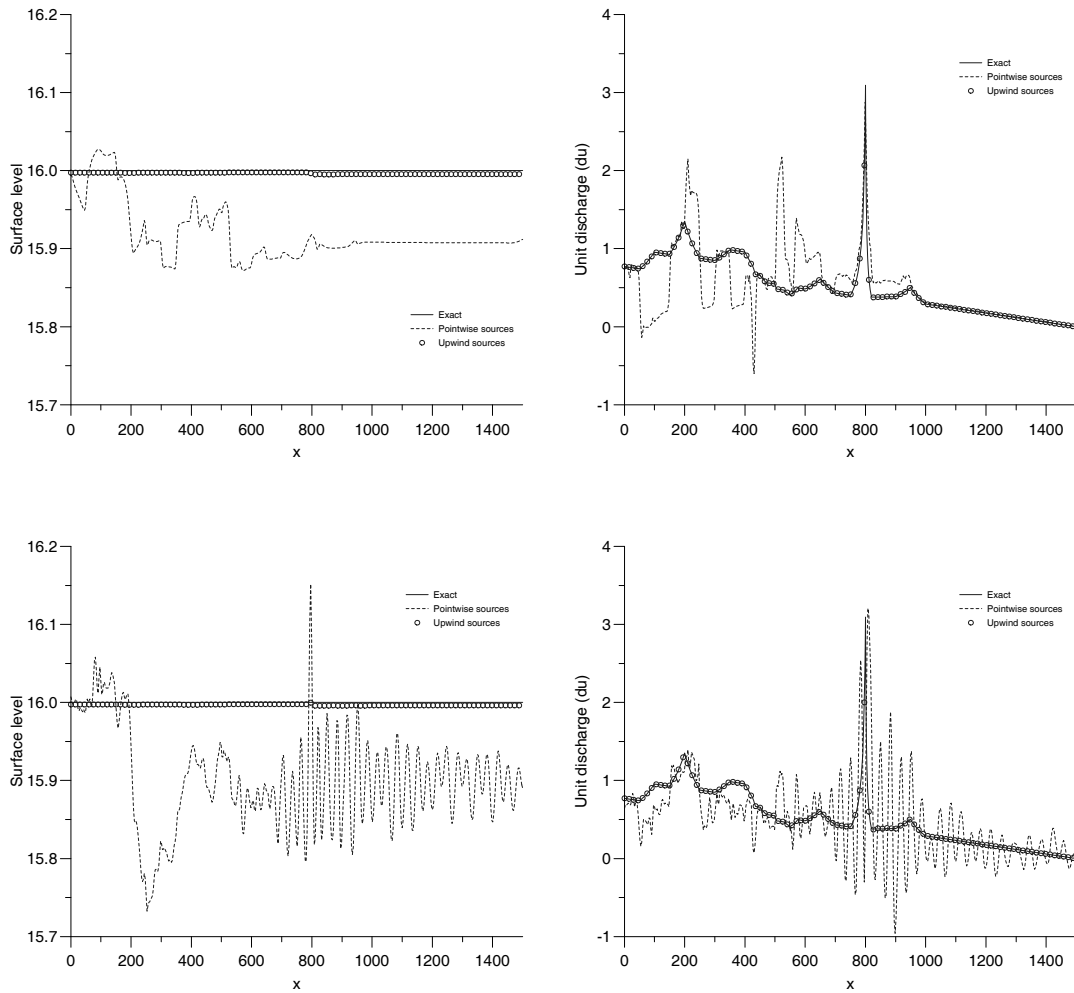
Figure 3.5: Water surface level and unit discharge for the tidal flow test case for first order (top) and high resolution slope limited (bottom) schemes.

discretisation maintains still, flat water indefinitely to machine accuracy in both the first order and the high resolution cases, see Figure 3.4. This is true of all of the channel shapes which were tested and each of the schemes described earlier in the text. The pointwise evaluation of the source term is clearly unable to match this.

Results for the tidal flow test case described in Section 2.3.1 are shown in Figure 3.5 for the same triangular grid. Again, the advantage of using the upwind source term discretisation is clearly visible and here, unlike in one dimension, the CFL number used to obtain the results is still 0.8. When the source terms are upwinded the results from the high resolution scheme are almost oscillation-free (although it must be remembered that the averaging across the channel breadth does produce a small amount of smoothing). Generally, it has been seen that the properties exhibited by the schemes in one dimension are carried over into higher dimensions.

# 4    Conclusions

In this paper a new method has been presented for the discretisation of source terms when they appear as part of a nonlinear system of conservation laws. Specifically, the correct approximation to the source terms is sought, given that a particular finite volume scheme has been used for the discretisation of the flux terms. Roe's scheme has been chosen here as the underlying numerical scheme, but the philosophy underlying the source term approximation (that the source terms must, in some sense, be discretised in the same manner as the flux derivatives) may also be applied to other finite volume methods. The discretisation builds on the work of many previous authors [5, 6, 1, 4], who approximated their source terms in a manner which took into account the flux discretisation and, as a consequence, allowed the numerical model to maintain specific equilibria which are satisfied by the mathematical model. The new aspect of this work is the generalisation of these techniques to high order TVD versions of Roe's scheme (using both flux limiters and slope limiters) and to arbitrary polygonal meshes in any number of dimensions. The methods have been designed specifically for

source terms which provide some sort of balance with the flux derivatives. Even so, the same techniques can easily be applied to other source terms (such as those which model bed friction in the shallow water equations) which do not exhibit a precise balance, but the advantages over the simple pointwise discretisation are less obvious.

The effectiveness of these techniques has been illustrated using the one- and two-dimensional shallow water equations (the extension to three-dimensional systems of equations is straightforward, though not described here in detail), in which source terms are used to model variations in the bed topography and (in one dimension) channel breadth. Particular attention has been paid to the special case of still water, and the schemes have been constructed so that they maintain this state. In fact, the improved accuracy of the new 'upwind' discretisation of the source terms is also shown in the approximation of other steady state solutions, particularly in one dimension when flux limiters have been used, and to a great extent by time-dependent test cases as well. The improvement is less marked for slope limited schemes, indicating that a more sophisticated approximation to the source term may be necessary away from the still water steady state. This has been shown by comparison with a selection of test cases for which exact solutions are available. The advantages over the commonly-used pointwise discretisations are particularly apparent when quantities depending on the flow velocity are compared. At this stage of the research, the main problem with the new technique (a problem which also applies to the old methods) is in the modelling of time-dependent problems. Here, in order to avoid spurious oscillations in the high resolution results a low CFL number has to be imposed (0.1 in the cases tested here), and in some cases the unphysical oscillations cannot be removed completely. This is because the TVD condition which is satisfied by the scheme is only valid for the homogeneous equations. The possible construction of a TVD condition in the presence of source terms is a topic for future research. In the meantime it may prove beneficial to apply a Flux-Corrected Transport approach since it is clear from the techniques presented in this paper how the source terms should be treated for both upwind and Lax-Wendroff schemes, and the first order upwind scheme appears to be robust enough to eradicate the unwanted oscillations.

An alternative method has also been proposed for the discretisation of the flux term in the case where it varies spatially but independently of the flow variables (as with one-dimensional models of shallow water flow through a channel of variable breadth). It has been shown that, in combination with the source term approximation, the method produces accurate solutions for a wide variety of steady state and time-dependent test cases.

## Acknowledgements

## References

[1] A.Bermúdez and M.E.Vázquez, 'Upwind methods for hyperbolic conservation laws with source terms', *Computers Fluids*, **23(8)**:1049–1071, 1994.

[2] A.Bermúdez, A.Dervieux, J-A.Desideri and M.E.Vázquez, 'Upwind schemes for the two-dimensional shallow water equations with variable depth using unstructured meshes' *Comput. Methods Appl. Mech. Engrg.*, **155**:49–72, 1998.

[3] N.Goutal and F.Maurel, Proceedings of the 2nd Workshop on Dam-Break Wave Simulation, Technical Report HE-43/97/016/A, Electricité de France, Département Laboratoire National d'Hydraulique, Groupe Hydraulique Fluviale, 1997.

[4] P.Garcia-Navarro and M.E.Vazquez-Cendon, 'Some considerations and improvements on the performance of Roe's scheme for 1d irregular geometries', Internal Report 23, Departamento de Matemática Aplicada, Universidade de Santiago do Compostela, 1997.

[5]  P.Glaister, 'Difference schemes for the shallow water equations', Numerical Analysis Report 9/87, Department of Mathematics, University of Reading, 1987.

[6]  P.Glaister, 'Prediction of supercritical flow in open channels', *Comput. Math. Applic.*, **24(7)**:69–75, 1992.

[7]  M.E.Hubbard, 'On the accuracy of one-dimensional models of steady converging/diverging open channel flows', Numerical Analysis Report 1/99, Department of Mathematics, University of Reading, 1999 (submitted to *Int. J. Numer. Methods Fluids*).

[8]  R.J.LeVeque, *Numerical methods for conservation laws*, Birkhäuser, Basel, 1992.

[9]  R.J.LeVeque, 'Balancing source terms and flux gradients in high-resolution Godunov methods: the quasi-steady wave-propagation algorithm', *J. Comput. Phys.*, **146(1)**:346–365, 1998.

[10]  A.Priestley, 'Roe-type schemes for super-critical flows in rivers', Numerical Analysis Report 13/89, Department of Mathematics, University of Reading, 1989.

[11]  P.L.Roe, 'Approximate Riemann solvers, parameter vectors, and difference schemes', *J. Comput. Phys.*, **43(2)**:357–372, 1981.

[12]  P.L.Roe, 'Fluctuations and signals - a framework for numerical evolution problems', in Numerical Methods for Fluid Dynamics, pp.219–257, ed. K.W.Morton, OUP, 1982.

[13]  P.L.Roe, 'Characteristic-based schemes for the Euler equations', *Ann. Rev. Fluid Mech.*, **18**:337–365, 1986.

[14]  P.K.Smolarkiewicz and L.G.Margolin, 'MPDATA: a finite-difference solver for geophysical flows', *J. Comput. Phys.*, **140**:1–22, 1998.

[15]  P.K.Sweby, 'High resolution schemes using flux limiters for hyperbolic conservation laws', *SIAM J. Numer. Anal.*, **21**:995–1011, 1984.

[16] M.E.Vázquez-Cendón, 'Improved treatment of source terms in upwind schemes for the shallow water equations in channels with irregular geometry', *J. Comput. Phys.*, **148(2)**:497–526, 1999.

[17] B.van Leer, 'Towards the ultimate conservative difference scheme V. A second order sequel to Godunov's method', *J. Comput. Phys.*, **32**:101–136, 1979.