

High order well-balanced finite volume WENO schemes and discontinuous Galerkin methods for a class of hyperbolic systems with source terms

Yulong Xing^a, Chi-Wang Shu^{b,*}

^a *Department of Mathematics, Brown University, Providence, RI 02912, United States*

^b *Division of Applied Mathematics, Brown University, Box F, Providence, RI 02912, United States*

Received 31 March 2005; received in revised form 31 July 2005; accepted 6 October 2005

Available online 22 November 2005

Abstract

Hyperbolic balance laws have steady state solutions in which the flux gradients are nonzero but are exactly balanced by the source term. In our earlier work [J. Comput. Phys. 208 (2005) 206–227; J. Sci. Comput., accepted], we designed a well-balanced finite difference weighted essentially non-oscillatory (WENO) scheme, which at the same time maintains genuine high order accuracy for general solutions, to a class of hyperbolic systems with separable source terms including the shallow water equations, the elastic wave equation, the hyperbolic model for a chemosensitive movement, the nozzle flow and a two phase flow model. In this paper, we generalize high order finite volume WENO schemes and Runge–Kutta discontinuous Galerkin (RKDG) finite element methods to the same class of hyperbolic systems to maintain a well-balanced property. Finite volume and discontinuous Galerkin finite element schemes are more flexible than finite difference schemes to treat complicated geometry and adaptivity. However, because of a different computational framework, the maintenance of the well-balanced property requires different technical approaches. After the description of our well-balanced high order finite volume WENO and RKDG schemes, we perform extensive one and two dimensional simulations to verify the properties of these schemes such as the exact preservation of the balance laws for certain steady state solutions, the non-oscillatory property for general solutions with discontinuities, and the genuine high order accuracy in smooth regions. © 2005 Elsevier Inc. All rights reserved.

Keywords: Hyperbolic balance laws; WENO scheme; Discontinuous Galerkin method; High order accuracy; Source term; Conservation laws; Shallow water equation; Elastic wave equation; Chemosensitive movement; Nozzle flow; Two phase flow

1. Introduction

In this paper, we are interested in designing high order weighted essentially non-oscillatory (WENO) finite volume schemes and Runge–Kutta discontinuous Galerkin (RKDG) finite element methods for solving hyperbolic systems of conservation laws with source terms (also called balance laws)

* Corresponding author. Tel.: +1 401 863 2549; fax: +1 401 863 1355.

E-mail addresses: xing@dam.brown.edu (Y. Xing), shu@dam.brown.edu (C.-W. Shu).

$$u_t + f_1(u, x, y)_x + f_2(u, x, y)_y = g(u, x, y) \quad (1.1)$$

or in the one dimensional case

$$u_t + f(u, x)_x = g(u, x), \quad (1.2)$$

where u is the solution vector, $f_1(u, x, y)$ and $f_2(u, x, y)$ (or $f(u, x)$) are the fluxes and $g(u, x, y)$ (or $g(u, x)$) is the source term.

These balance laws often admit steady state solutions in which the source term is exactly balanced by the flux gradients. Such cases, along with their perturbations, are very difficult to capture numerically. A straightforward treatment of the source terms will fail to preserve this balance. The objective of well-balanced schemes is to preserve exactly some of these steady state solutions. This objective should be achieved without sacrificing the high order accuracy and non-oscillatory properties of the scheme when applied to general, non-steady state solutions.

A typical example considered extensively in the literature for balance laws is the shallow water equation with a non-flat bottom topology. Research on numerical methods for the solution of the shallow water system has attracted significant attention in the past two decades. An early, important result in computing such solutions was given by Bermudez and Vazquez [3]. They proposed the idea of the “exact C-property”, which means that the scheme is “exact” when applied to the stationary case $h + b = \text{constant}$ and $hu = 0$, where h , b and u are the water height, the given bottom topography, and the velocity of the fluid, respectively, see (6.1) in Section 6.1. A good scheme for the shallow water system should satisfy this property. Also, Bermudez and Vazquez introduced in [3] the first order Q-scheme and the idea of source term upwinding. After this pioneering work, many other schemes for the shallow water equations with such well-balanced property have been developed. LeVeque [21] introduced a quasi-steady wave propagation algorithm. A Riemann problem is introduced in the center of each grid cell such that the flux difference exactly cancels the source term. Zhou et al. [38] used the surface gradient method for the treatment of the source terms. They used $h + b$ for the reconstruction instead of using h . For more related work, see also [13,14,17,19,20,24,26,33,34,37]. In particular, the authors of [33,34] presented well-balanced ENO and WENO schemes for the shallow water equations and other equations.

Our development of well-balanced WENO finite volume schemes and discontinuous Galerkin methods is based on our recent work [35,36]. In [35], we developed a well-balanced high order finite difference WENO scheme for solving the shallow water equation, which is non-oscillatory, well balanced (satisfying the exact C-property) for still water, and genuinely high order in smooth regions. Different from [33,34], a key ingredient of the technique used in [35] is a special decomposition of the source term, allowing a discretization to the source term to be both high order accurate for general solutions and exactly well balanced with the flux gradient for still water. Extensive one and two dimensional numerical experiments were provided in [35] to demonstrate the good behavior of this scheme. In [36], we extended this idea of decomposition of source terms to a general class of balance laws with separable source terms, allowing the design of well-balanced high order finite difference WENO scheme for all balance laws falling into this category. This class is quite general, including, besides the shallow water equations, the elastic wave equation, the hyperbolic model for a chemosensitive movement, the nozzle flow and a two phase flow model.

In this paper, we consider finite volume WENO schemes first introduced by Liu et al. [22], see also [16,27], and Runge–Kutta discontinuous Galerkin (RKDG) finite element methods that were originally developed by Cockburn and Shu [8], see also [7,9]. We will generalize these schemes to obtain high order well-balanced schemes. The crucial difference between the finite volume and the finite difference WENO schemes is that the WENO reconstruction procedure for a finite volume scheme applies to the solution and not to the flux function values. As a consequence, finite volume schemes are more suitable for computations in complex geometry and for using adaptive meshes, however the maintenance of the well-balanced property requires different technical approaches. The RKDG methods can be considered as a generalization of finite volume schemes, even though they do not require a reconstruction and evolve the complete polynomial in each cell forward in time. The RKDG methods are therefore easier to use for multi-dimensional problems in complex geometry, than the finite volume schemes, as the complicated reconstruction procedure can be avoided. Even

though the detailed technical approaches are different, the framework of the algorithm construction in this paper follows that in [35,36].

This paper is organized as follows. In Sections 2 and 3, we give a brief review of finite volume WENO schemes and RKDG schemes for the homogeneous conservation laws. In Section 4, we describe the class of balance laws under consideration and develop well-balanced finite volume WENO schemes, which at the same time are genuinely high order accurate for the general solutions. The well-balanced generalization of the RKDG schemes is presented in Section 5. In Section 6, we give several examples in applications which fall into the category of balance laws discussed in Section 4, and show selective numerical results in one and two dimensions to demonstrate the behavior of our well-balanced finite volume WENO schemes and RKDG schemes, verifying high order accuracy, the well-balanced property, and good resolution for smooth and discontinuous solutions. Concluding remarks are given in Section 7.

2. A review of high order finite volume WENO schemes

In this section, we briefly review the basic ideas of finite volume WENO schemes. For further details, we refer to [2,16,18,22,27,29–31].

First, we consider a scalar hyperbolic conservation law equation in one dimension

$$u_t + f(u)_x = 0, \tag{2.1}$$

and discretize the computational domain with cells $I_i = [x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}]$, $i = 1, \dots, N$. We denote the size of the i th cell by Δx_i and the center of the cell by $x_i = \frac{1}{2}(x_{i-\frac{1}{2}} + x_{i+\frac{1}{2}})$. Let $\bar{u}(x_i, t) = \frac{1}{\Delta x_i} \int_{I_i} u(x, t) dx$ denote the cell average of $u(\cdot, t)$ over the cell I_i . In a finite volume scheme, our computational variables are $\bar{u}_i(t)$, which approximate the cell averages $\bar{u}(x_i, t)$.

For finite volume schemes, we solve an integrated version of (2.1):

$$\frac{d}{dt} \bar{u}(x_i, t) = -\frac{1}{\Delta x_i} \left(f(u(x_{i+\frac{1}{2}}), t) - f(u(x_{i-\frac{1}{2}}), t) \right).$$

This is approximated by the following conservative scheme:

$$\frac{d}{dt} \bar{u}_i(t) = -\frac{1}{\Delta x_i} \left(\hat{f}_{i+\frac{1}{2}} - \hat{f}_{i-\frac{1}{2}} \right) \tag{2.2}$$

with $\hat{f}_{i+\frac{1}{2}} = F(u_{i+\frac{1}{2}}^-, u_{i+\frac{1}{2}}^+)$ being the numerical flux. Here $u_{i+\frac{1}{2}}^-$ and $u_{i+\frac{1}{2}}^+$ are the high order pointwise approximations to $u(x_{i+\frac{1}{2}}, t)$, obtained from the cell averages by a high order WENO reconstruction procedure.

In order to obtain a stable scheme, the numerical flux $F(a, b)$ needs to be a monotone flux, namely F is a nondecreasing function of its first argument a and a nonincreasing function of its second argument b . There are many choices of these fluxes, such as the Godunov flux, the Engquist–Osher flux and the Lax–Friedrichs (LF) flux. The difference among these fluxes is significant for low order schemes but becomes less significant for higher order reconstructions. The simplest and most inexpensive monotone flux is the Lax–Friedrichs flux:

$$F(a, b) = \frac{1}{2} (f(a) + f(b) - \alpha(b - a)), \tag{2.3}$$

where $\alpha = \max_u |f'(u)|$. Depending on whether the maximum is taken globally (along the line of computation) or locally, this flux is referred to as the Lax–Friedrichs (LF) or the local Lax–Friedrichs (LLF) flux.

The approximations $u_{i+\frac{1}{2}}^-$ and $u_{i+\frac{1}{2}}^+$ are computed through the neighboring cell average values \bar{u}_j . For a $(2k - 1)$ th order WENO scheme, we first compute k reconstructed values

$$\hat{u}_{i+\frac{1}{2}}^{(r)} = \sum_{j=0}^{k-1} c_{rj} \bar{u}_{i-r+j}, \quad r = 0, \dots, k - 1,$$

corresponding to k different candidate stencils

$$S_r(i) = \{x_{i-r}, \dots, x_{i-r+k-1}\}, \quad r = 0, \dots, k - 1. \tag{2.4}$$

The coefficients c_{rj} are chosen such that each of these k reconstructed values is k th order accurate, see [29]. Also, we obtain the k reconstructed values $\tilde{u}_{i-\frac{1}{2}}^{(r)}$, of k th order accuracy, using

$$\tilde{u}_{i-\frac{1}{2}}^{(r)} = \sum_{j=0}^{k-1} \tilde{c}_{rj} \bar{u}_{i-r+j}, \quad r = 0, \dots, k-1,$$

with

$$\tilde{c}_{rj} = c_{r-1,j},$$

based on the same stencils (2.4). The $(2k - 1)$ th order WENO reconstruction is a convex combination of all these k reconstructed values

$$u_{i+\frac{1}{2}}^- = \sum_{r=0}^{k-1} w_r \hat{u}_{i+\frac{1}{2}}^{(r)}, \quad u_{i-\frac{1}{2}}^+ = \sum_{r=0}^{k-1} \tilde{w}_r \tilde{u}_{i-\frac{1}{2}}^{(r)}.$$

The nonlinear weights w_r satisfy $w_r \geq 0$, $\sum_{j=0}^{k-1} w_r = 1$, and are defined in the following way:

$$w_r = \frac{\alpha_r}{\sum_{s=0}^{k-1} \alpha_s}, \quad \alpha_r = \frac{d_r}{(\varepsilon + \beta_r)^2}. \tag{2.5}$$

Here d_r are the linear weights which yield $(2k - 1)$ th order accuracy, β_r are the so-called ‘‘smoothness indicators’’ of the stencil $S_r(i)$ which measure the smoothness of the function $u(x)$ in the stencil. ε is a small constant used to avoid the denominator to become zero and is typically taken as 10^{-6} . By symmetry, \tilde{w}_r is computed by:

$$\tilde{w}_r = \frac{\tilde{\alpha}_r}{\sum_{s=0}^{k-1} \tilde{\alpha}_s}, \quad \tilde{\alpha}_r = \frac{\tilde{d}_r}{(\varepsilon + \beta_r)^2}, \tag{2.6}$$

with

$$\tilde{d}_r = d_{k-1+r}. \tag{2.7}$$

The exact form of the smoothness indicators and other details about the WENO reconstruction can be found in [18,29].

For hyperbolic systems such as the shallow water equations, we use the local characteristic decomposition, which is more robust than a component by component version. First, we compute an average state $\bar{u}_{i+\frac{1}{2}}$ between \bar{u}_i and \bar{u}_{i+1} , using either the simple arithmetic mean or a Roe’s average [25]. The WENO procedure is used on

$$\bar{v}_j = R^{-1} \bar{u}_j, \quad j \text{ in a neighborhood of } i. \tag{2.8}$$

where $R = (r_1, \dots, r_n)$ is the matrix whose columns are the right eigenvectors of $f'(\bar{u}_{i+\frac{1}{2}})$. The reconstructed values $v_{i+\frac{1}{2}}^\pm$ thus computed are then projected back into the physical space by left multiplying with R , yielding finally the reconstructed values in the physical space.

With the reconstructed values $u_{i+\frac{1}{2}}^\pm$, the right-hand side of (2.2) can be computed through (2.3) to high order accuracy. Together with a TVD high order Runge–Kutta time discretization [30], this completes the description of a high order finite volume WENO scheme.

Finite volume WENO schemes in the two dimensional case have the same framework but are more complicated to implement. In this paper, we consider only rectangular cells for simplicity, although the technique also works for general triangulations. Consider the two dimensional homogeneous conservation law

$$u_t + f_1(u, x, y)_x + f_2(u, x, y)_y = 0, \tag{2.9}$$

together with a spatial discretization of the computational domain with cells $I_{ij} = [x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}] \times [y_{j-\frac{1}{2}}, y_{j+\frac{1}{2}}]$, $i = 1, \dots, N_x$, $j = 1, \dots, N_y$. As usual, we use the notations:

$$\Delta x_i = x_{i+\frac{1}{2}} - x_{i-\frac{1}{2}}, \quad \Delta y_j = y_{j+\frac{1}{2}} - y_{j-\frac{1}{2}}$$

to denote the grid sizes.

We integrate (2.9) over the interval I_{ij} to obtain:

$$\begin{aligned} \frac{d}{dt} \bar{u}(x_i, y_j, t) = & -\frac{1}{\Delta x_i \Delta y_j} \left(\int_{y_{j-\frac{1}{2}}}^{y_{j+\frac{1}{2}}} f_1(u(x_{i+\frac{1}{2}}, y, t)) dy - \int_{y_{j-\frac{1}{2}}}^{y_{j+\frac{1}{2}}} f_1(u(x_{i-\frac{1}{2}}, y, t)) dy \right. \\ & \left. + \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} f_2(u(x, y_{j+\frac{1}{2}}, t)) dx - \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} f_2(u(x, y_{j-\frac{1}{2}}, t)) dx \right), \end{aligned} \tag{2.10}$$

where

$$\bar{u}(x_i, y_j, t) = \frac{1}{\Delta x_i \Delta y_j} \int_{y_{j-\frac{1}{2}}}^{y_{j+\frac{1}{2}}} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} u(\xi, \eta, t) d\xi d\eta$$

is the cell average. We approximate (2.10) by the conservative scheme

$$\frac{d}{dt} \bar{u}_{ij}(t) = -\frac{1}{\Delta x_i} \left((\hat{f}_1)_{i+\frac{1}{2}j} - (\hat{f}_1)_{i-\frac{1}{2}j} \right) - \frac{1}{\Delta y_j} \left((\hat{f}_2)_{i,j+\frac{1}{2}} - (\hat{f}_2)_{i,j-\frac{1}{2}} \right), \tag{2.11}$$

where the numerical flux $(\hat{f}_1)_{i+\frac{1}{2}j}$ is defined by

$$(\hat{f}_1)_{i+\frac{1}{2}j} = \sum_{\alpha} w_{\alpha} F \left(u_{x_{i+\frac{1}{2}}, y_j + \beta_{\alpha} \Delta y_j}^{-}, u_{x_{i+\frac{1}{2}}, y_j + \beta_{\alpha} \Delta y_j}^{+} \right), \tag{2.12}$$

where β_{α} and w_{α} are the Gaussian quadrature nodes and weights, to approximate the integration in y :

$$\frac{1}{\Delta y_j} \int_{y_{j-\frac{1}{2}}}^{y_{j+\frac{1}{2}}} f_1(u(x_{i+\frac{1}{2}}, y, t)) dy.$$

The monotone flux $F(a, b)$ is the same as defined above (for example, Formula (2.3)). $u_{x_{i+\frac{1}{2}}, y_j + \beta_{\alpha} \Delta y_j}^{\pm}$ are the $(2k - 1)$ th order accurate reconstructed values obtained by a WENO reconstruction procedure. In this procedure, for rectangular meshes, if we use the tensor products of one dimensional polynomials, i.e., polynomials in Q^{k-1} , things can proceed as in one dimension. A practical way to perform the reconstruction in two dimensions is given as follows. We first perform a one dimensional reconstruction in one of the directions (e.g., the y -direction), obtaining one dimensional cell averages of the function u in the other direction (e.g., the x -direction). We then perform a reconstruction in the other direction to obtain the approximated point values, see [27,29].

Similarly, we can compute the flux $(\hat{f}_2)_{i,j+\frac{1}{2}}$ by

$$(\hat{f}_2)_{i,j+\frac{1}{2}} = \sum_{\alpha} w_{\alpha} F \left(u_{x_i + \beta_{\alpha} \Delta x_i, y_{j+\frac{1}{2}}}^{-}, u_{x_i + \beta_{\alpha} \Delta x_i, y_{j+\frac{1}{2}}}^{+} \right). \tag{2.13}$$

3. A review of high order discontinuous Galerkin methods

In this section, we give a short overview of another widely used high order scheme, namely the Runge–Kutta discontinuous Galerkin method, which was first introduced by Cockburn and Shu. We refer to [5,7–10] for more information.

Again, a scalar hyperbolic conservation law in one dimension is considered:

$$u_t + f(u)_x = 0, \quad u(x, 0) = u_0(x). \tag{3.1}$$

As before, we discretize the computational domain into cells $I_i = [x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}]$, and denote the size of the i th cell by Δx_i and the maximum mesh size by $h = \max_i \Delta x_i$.

First, we multiply Eq. (3.1) by an arbitrary smooth function v , integrate it over cell I_j and perform integration by parts to obtain

$$\int_{I_j} \partial_t u(x, t) v(x) \, dx - \int_{I_j} f(u(x, t)) \partial_x v(x) \, dx + f(u(x_{j+\frac{1}{2}}, t)) v(x_{j+\frac{1}{2}}) - f(u(x_{j-\frac{1}{2}}, t)) v(x_{j-\frac{1}{2}}) = 0, \quad (3.2)$$

$$\int_{I_j} u(x, 0) v(x) \, dx = \int_{I_j} u_0(x) v(x) \, dx.$$

The main difference between the DG method and a traditional finite element method lies in the choice of the test space and solution space. Here, we seek an approximation u_h to u which belongs to the finite dimensional space

$$V_h = V_h^k \equiv \{v : v|_{I_j} \in P^k(I_j), j = 1, \dots, N\}, \quad (3.3)$$

where $P^k(I)$ denotes the space of polynomials in I of degree at most k . Notice that u_h can be discontinuous at the cell boundary $x_{j+\frac{1}{2}}$. In Eq. (3.2), we replace the smooth functions v by test functions v_h from the test space V_h , and u by the numerical solution u_h . Together with the replacement of the nonlinear flux $f(u(x_{j+\frac{1}{2}}, t))$ by a numerical flux $\hat{f}_{j+\frac{1}{2}} = F(u_h(x_{j+\frac{1}{2}}^-, t), u_h(x_{j+\frac{1}{2}}^+, t))$, we obtain the numerical scheme denoted by

$$\int_{I_j} \partial_t u_h(x, t) v_h(x) \, dx - \int_{I_j} f(u_h(x, t)) \partial_x v_h(x) \, dx + \hat{f}_{j+\frac{1}{2}} v_h(x_{j+\frac{1}{2}}^-) - \hat{f}_{j-\frac{1}{2}} v_h(x_{j-\frac{1}{2}}^+) = 0, \quad (3.4)$$

$$\int_{I_j} u_h(x, 0) v_h(x) \, dx = \int_{I_j} u_0(x) v_h(x) \, dx.$$

As before, $F(a, b)$ is chosen as a monotone flux to recover a finite volume monotone scheme for the piecewise constant $k = 0$ case. We could, for example, again use the simple Lax–Friedrichs flux (2.3).

Another important ingredient for the RKDG method is that a slope limiter procedure should be performed after each inner stage in the Runge–Kutta time stepping. This is necessary for computing solutions with strong discontinuities. There are many choices for the slope limiters, see, e.g. [23]. In this paper, we use the total variation bounded (TVB) limiter in [28,8,6,9]; we refer to these references for the details of this limiter.

Together with a TVD high order Runge–Kutta time discretization [30], we have then finished the description of the RKDG method.

Multi-dimensional problems can be handled in the same fashion. We also perform an integration by parts (Green’s formula) first, and then replace the boundary values by numerical fluxes. The main difference is that the fluxes are now integrals along the cell boundary, which can be calculated by Gauss–quadrature rules. For more details, we refer to [5,6,9].

4. Construction of well-balanced finite volume WENO schemes

In this section, we design a genuine high order finite volume WENO scheme for a class of general balance laws (1.1). We will concentrate our discussion on the one dimensional case (1.2). Generalization to the multi-dimensional case (1.1) can be done in some situations, for example the cases discussed in [35,36]; we present the details for the two dimensional shallow water equations in Section 6.2.

Our main objective is to preserve certain steady state solutions while maintaining high order accuracy for general solutions. The main idea in [35,36] to design a well-balanced high order finite difference WENO scheme is to decompose the source term into a sum of several terms, each of which is discretized independently using a finite difference formula consistent with that of approximating the flux derivative terms in the conservation law. We follow a similar idea here and decompose the integral of the source term into a sum of several terms, then compute each of them in a way consistent with that of computing the corresponding flux terms. We first consider the case that (1.2) is a scalar balance law. The case of systems will be explored later.

We are interested in preserving exactly certain steady state solutions u of (1.2):

$$f(u, x)_x = g(u, x). \quad (4.1)$$

As in [36], we make some assumptions on Eq. (1.2) and the steady state solution u of (4.1) that we are interested to preserve exactly:

Assumption 4.1. The steady state solution u of (4.1) that we are interested to preserve satisfies

$$a(u, x) \equiv \frac{u + p(x)}{q(x)} = \text{constant} \tag{4.2}$$

for some known functions $p(x)$ and $q(x)$.

Assumption 4.2. The source term $g(u, x)$ in (1.2) can be decomposed as

$$g(u, x) = \sum_j s_j(a(u, x))t'_j(x) \tag{4.3}$$

for some known functions s_j and t_j .

Note that Assumption 4.1 given here is more restrictive than that in [36]. This is due to the additional difficulties related to the finite volume formulation.

We would like to preserve exactly the steady state solutions u which satisfy Assumption 4.1, for a balance law (1.2) with a source term satisfying Assumption 4.2.

Now let us describe the details of the algorithm. We consider the semi-discrete formulation of the balance law

$$\frac{d}{dt} \bar{u}_i(t) = -\frac{1}{\Delta x_i} (f(u(x_{i+\frac{1}{2}}), t) - f(u(x_{i-\frac{1}{2}}), t)) + \frac{1}{\Delta x_i} \int_{I_i} g(u, x) dx. \tag{4.4}$$

The time discretization is usually performed by the classical high order Runge–Kutta method. Before stating our numerical scheme, we first present the procedure to reconstruct the pointwise values by the WENO reconstruction procedure, and then decompose the integral of the source term into several terms, with the objective of keeping the exact balance property without reducing the high order accuracy of the scheme. The scheme is then finally introduced with a minor change on the flux term, compared with the original WENO scheme.

The first step in building the algorithm is to reconstruct $u_{i+\frac{1}{2}}^\pm$ from the given cell averages \bar{u}_i , by the WENO reconstruction procedure explained in Section 2, which are high order accurate approximations to the exact value $u(x_{i+\frac{1}{2}})$. We use the smoothness indicators β_r to measure the smoothness of the variable u . The WENO reconstruction can be eventually written out as

$$u_{i+\frac{1}{2}}^+ = \sum_{k=-r+1}^r w_k \bar{u}_{i+k} \equiv S_u^+(\bar{u})_i, \quad u_{i+\frac{1}{2}}^- = \sum_{k=-r}^{r-1} \tilde{w}_k \bar{u}_{i+k} \equiv S_u^-(\bar{u})_i. \tag{4.5}$$

where $r = 3$ for the fifth order WENO approximation and the coefficients w_k and \tilde{w}_k depend nonlinearly on the smoothness indicators involving the cell average \bar{u} , following (2.5) and (2.6). Here we obtain a linear operator $S_u^\pm(v)$ (linear in v) which is obtained from a WENO reconstruction with fixed coefficients w_k calculated from the cell averages \bar{u} . Once again, our purpose is to find a high order finite volume scheme for a class of conservation laws which can preserve the steady state solution (4.2). The key idea here is to use the *linear* operators $S_u^\pm(v)$ and apply them to reconstruct the functions \bar{p}_i and \bar{q}_i . Thus

$$\begin{aligned} p_{i+\frac{1}{2}}^+ &= S_u^+(\bar{p})_i = \sum_{k=-r+1}^r w_k \bar{p}_{i+k}, & p_{i+\frac{1}{2}}^- &= S_u^-(\bar{p})_i = \sum_{k=-r}^{r-1} \tilde{w}_k \bar{p}_{i+k}, \\ q_{i+\frac{1}{2}}^+ &= S_u^+(\bar{q})_i = \sum_{k=-r+1}^r w_k \bar{q}_{i+k}, & q_{i+\frac{1}{2}}^- &= S_u^-(\bar{q})_i = \sum_{k=-r}^{r-1} \tilde{w}_k \bar{q}_{i+k}. \end{aligned} \tag{4.6}$$

With the reconstructed values $p_{i+\frac{1}{2}}^\pm$ and $q_{i+\frac{1}{2}}^\pm$, we obtain the pointwise value of $a(u, x)$ by $a(u, x)_{i+\frac{1}{2}}^\pm = \frac{u_{i+\frac{1}{2}}^\pm + p_{i+\frac{1}{2}}^\pm}{q_{i+\frac{1}{2}}^\pm}$.

Clearly, $p_{i+\frac{1}{2}}^\pm$ and $q_{i+\frac{1}{2}}^\pm$ are high order accurate pointwise approximation to the function of $p(x)$ and $q(x)$ at the cell boundary $x_{i+\frac{1}{2}}$. Hence, $a(u, x)_{i+\frac{1}{2}}^\pm$ is a high order approximation to $a(u(x_{i+\frac{1}{2}}), x_{i+\frac{1}{2}})$.

Now assume that u is the steady state solution satisfying (4.2), namely

$$u + p(x) = cq(x)$$

for some constant c . If the cell averages \bar{u}_i , \bar{p}_i and \bar{q}_i are computed in the same fashion (e.g., all computed exactly, or all computed with the same numerical quadrature) from u , $p(x)$ and $q(x)$, then we clearly also have

$$\bar{u}_i + \bar{p}_i = c\bar{q}_i$$

for the same constant c . Since the reconstructed values $u_{i+\frac{1}{2}}^\pm, p_{i+\frac{1}{2}}^\pm$ and $q_{i+\frac{1}{2}}^\pm$ are computed from the cell averages \bar{u}_j, \bar{p}_j and \bar{q}_j with the same linear operators $S_u^\pm(v)$, we clearly have

$$u_{i+\frac{1}{2}}^\pm + p_{i+\frac{1}{2}}^\pm = cq_{i+\frac{1}{2}}^\pm$$

for the same constant c , that is,

$$a(u, x)_{i+\frac{1}{2}}^\pm = c \tag{4.7}$$

for the same constant c .

Clearly, for a steady state solution u satisfying Assumptions 4.1 and 4.2,

$$\frac{d}{dx} \left(f(u, x) - \sum_j s_j(a(u, x))t_j(x) \right) = f(u, x)_x - \sum_j s_j(a(u, x))t'_j(x) = f(u, x)_x - g(u, x) = 0.$$

Therefore, $f(u, x) - \sum_j s_j(a(u, x))t_j(x)$ is a constant. We would need to choose suitably $(t_j)_{i+\frac{1}{2}}^\pm$, which should be high order approximations to $t_j(x_{i+\frac{1}{2}})$ such that

$$f(u_{i+\frac{1}{2}}^\pm) - \sum_j s_j(a(u, x)_{i+\frac{1}{2}}^\pm)(t_j)_{i+\frac{1}{2}}^\pm = \text{constant} \tag{4.8}$$

for a steady state solution u satisfying Assumptions 4.1 and 4.2. In the applications stated later in Section 6, we will specify the choices of $(t_j)_{i+\frac{1}{2}}^\pm$ in each case.

The integral of the source term takes the form

$$\int_{I_i} g(u, x) \, dx = \sum_j \int_{I_i} s_j(a(u, x))t'_j(x) \, dx.$$

We need to decompose it further in the following way in order to obtain a well-balanced scheme

$$\begin{aligned} \sum_j \int_{I_i} s_j(a(u, x))t'_j(x) \, dx &= \sum_j \left(\frac{1}{2} \left(s_j(a(u, x)_{i-\frac{1}{2}}^+) + s_j(a(u, x)_{i+\frac{1}{2}}^-) \right) \int_{I_i} t'_j(x) \, dx \right. \\ &\quad \left. + \int_{I_i} \left(s_j(a(u, x)) - \frac{1}{2} \left(s_j(a(u, x)_{i-\frac{1}{2}}^+) + s_j(a(u, x)_{i+\frac{1}{2}}^-) \right) \right) t'_j(x) \, dx \right) \\ &= \sum_j \left(\frac{1}{2} \left(s_j(a(u, x)_{i-\frac{1}{2}}^+) + s_j(a(u, x)_{i+\frac{1}{2}}^-) \right) (t_j(x_{i+\frac{1}{2}}) - t_j(x_{i-\frac{1}{2}})) \right. \\ &\quad \left. + \int_{I_i} \left(s_j(a(u, x)) - \frac{1}{2} \left(s_j(a(u, x)_{i-\frac{1}{2}}^+) + s_j(a(u, x)_{i+\frac{1}{2}}^-) \right) \right) t'_j(x) \, dx \right). \end{aligned} \tag{4.9}$$

The purpose of this decomposition is to ensure the balance with the flux difference term on the right-hand side of (4.4), see the proof of Proposition 4.3. We remark that $\frac{1}{2}(s_j(a(u, x)_{i-\frac{1}{2}}^+) + s_j(a(u, x)_{i+\frac{1}{2}}^-))$ can also be replaced by $s_j(\frac{\bar{u}+p(x)}{q(x)})$ where as usual the overbar denotes the cell average over the cell I_i , which could be used when there is a singularity at the boundary, for example, in the application in Section 6.5.

Now we are ready to describe the final form of the algorithm

$$\frac{d}{dt} \bar{u}_i(t) = -\frac{1}{\Delta x_i} (\hat{f}_{i+\frac{1}{2}} - \hat{f}_{i-\frac{1}{2}}) + \frac{1}{\Delta x_i} \hat{g}_i, \tag{4.10}$$

with

$$\hat{g}_i = \sum_j \left(\frac{1}{2} \left(s_j(a(u, x)_{i-\frac{1}{2}}^+) + s_j(a(u, x)_{i+\frac{1}{2}}^-) \right) \left((\hat{t}_j)_{i+\frac{1}{2}} - (\hat{t}_j)_{i-\frac{1}{2}} \right) + g_{i,j} \right), \tag{4.11}$$

where $(\hat{t}_j)_{i+\frac{1}{2}}$ is a high order approximation to $t_j(x_{i+\frac{1}{2}})$, whose definition will be described below, and $g_{i,j}$ is any high order approximation to the integral

$$\int_{I_i} \left(s_j(a(u, x)) - \frac{1}{2} \left(s_j(a(u, x)_{i-\frac{1}{2}}^+) + s_j(a(u, x)_{i+\frac{1}{2}}^-) \right) \right) t'_j(x) dx. \tag{4.12}$$

Comparing with (4.9), it is clear that \hat{g}_i is a high order approximation to the source term in (4.4).

The numerical flux $\hat{f}_{i+\frac{1}{2}}$ is defined by a monotone flux such as the Lax–Friedrichs flux (2.3)

$$F(u_{i+\frac{1}{2}}^-, u_{i+\frac{1}{2}}^+) = \frac{1}{2} \left[f(u_{i+\frac{1}{2}}^-) + f(u_{i+\frac{1}{2}}^+) - \alpha(u_{i+\frac{1}{2}}^+ - u_{i+\frac{1}{2}}^-) \right]. \tag{4.13}$$

We need to make a modification to this flux, by replacing $\alpha(u_{i+\frac{1}{2}}^+ - u_{i+\frac{1}{2}}^-)$ in (4.13) with $\alpha \text{sign}(q(x))(a(u, x)_{i+\frac{1}{2}}^+ - a(u, x)_{i+\frac{1}{2}}^-)$. The numerical flux now becomes

$$\hat{f}_{i+\frac{1}{2}} = \frac{1}{2} \left[f(u_{i+\frac{1}{2}}^-) + f(u_{i+\frac{1}{2}}^+) - \alpha \text{sign}(q(x))(a(u, x)_{i+\frac{1}{2}}^+ - a(u, x)_{i+\frac{1}{2}}^-) \right]. \tag{4.14}$$

We would need to assume here that $q(x)$ in (4.2) does not change sign. The constant α should be suitably adjusted by the size of $\frac{1}{q(x)}$ in order to maintain enough artificial viscosity. This modification does not affect accuracy. For the steady state solution (4.2),

$$\alpha \text{sign}(q(x))(a(u, x)_{i+\frac{1}{2}}^+ - a(u, x)_{i+\frac{1}{2}}^-) = 0$$

because of (4.7). Hence, the effect of these viscosity terms becomes zero and the numerical flux turns out to be in a simple form

$$\hat{f}_{i+\frac{1}{2}} = \frac{1}{2} \left[f(u_{i+\frac{1}{2}}^-) + f(u_{i+\frac{1}{2}}^+) \right]. \tag{4.15}$$

Following this, we treat the approximation $(\hat{t}_j)_{i+\frac{1}{2}}$ in (4.11) in a similar way:

$$(\hat{t}_j)_{i+\frac{1}{2}} = \frac{1}{2} \left[(t_j)_{i+\frac{1}{2}}^- + (t_j)_{i+\frac{1}{2}}^+ \right] \tag{4.16}$$

where, as mentioned before, $(t_j)_{i+\frac{1}{2}}^\pm$ are high order approximations to $t_j(x_{i+\frac{1}{2}})$ satisfying (4.8). Note that we implement (4.16) for the general case, not only for the steady solution. There is no viscosity term in the source term, compared with the numerical flux (4.14).

For the remaining source term $g_{i,j}$, we simply use a suitable high order Gauss quadrature to evaluate the integral. The approximation of the values at those Gauss points are obtained by the WENO reconstruction procedure. It is easy to observe that high order accuracy is guaranteed for our scheme, and even if discontinuities exist in the solution, non-oscillatory property is maintained.

We now formulate the preservation of the steady state solution (4.2) by our numerical scheme.

Proposition 4.3. *The WENO-LF schemes as implemented above with (4.10), (4.11), (4.14) and (4.16) are exact for steady state solutions satisfying (4.2) and can maintain the original high order accuracy for general solutions.*

Proof: The high order accuracy is straightforward to observe. We only prove the well-balanced property here. First, for the steady state solution $a(u, x) = c$ for some constant c , the reconstructed values $a(u, x)_{i-\frac{1}{2}}^\pm$ are also equal to the same constant c , see (4.7). Hence, we notice that the source term $g_{i,j}$, which is a high order numerical approximation of the integral in (4.12) by a Gauss quadrature, is simply zero since $a(u, x)$ is equal to $a(u, x)_{i-\frac{1}{2}}^\pm$ at each Gauss point. Furthermore, in this case the flux terms take the form (4.15) and (4.16). Therefore, the truncation error reduces to

$$\begin{aligned} & -\hat{f}_{i+\frac{1}{2}} + \hat{f}_{i-\frac{1}{2}} + \sum_j \frac{1}{2} \left(s_j(a(u, x)_{i-\frac{1}{2}}^+) + s_j(a(u, x)_{i+\frac{1}{2}}^-) \right) \left((\hat{t}_j)_{i+\frac{1}{2}} - (\hat{t}_j)_{i-\frac{1}{2}} \right) \\ & = -\hat{f}_{i+\frac{1}{2}} + \sum_j s_j(c)(\hat{t}_j)_{i+\frac{1}{2}} + \hat{f}_{i-\frac{1}{2}} - \sum_j s_j(c)(\hat{t}_j)_{i-\frac{1}{2}} = 0, \end{aligned}$$

where we have used (4.7) for the first equality, and (4.8), (4.15) and (4.16) for the second equality. This finishes the proof. \square

We now discuss the system case. The framework described for the scalar case can be applied to systems provided that we have certain knowledge about the steady state solutions to be preserved in the form of (4.2). Typically, for a system with m equations, we would have m relationships in the form of (4.2):

$$a_1(u, x) = \text{constant}, \quad \dots \quad a_m(u, x) = \text{constant} \tag{4.17}$$

for the steady state solutions that we would like to preserve exactly. Here we require that, for the steady state solution (4.17), $a_j(u, x) = \frac{\sum_k b_k u_k + p_j(x)}{q_j(x)}$, where $u = (u_1, \dots, u_m)$, b_k are arbitrary constants, and $p_j(x)$ and $q_j(x)$ are arbitrary known functions of x . We would then still aim for decomposing each component of the source term in the form of (4.3), where s_j could be arbitrary functions of $a_1(u, x), \dots, a_m(u, x)$, and the functions s_j and t_j could be different for different components of the source vector. The remaining procedure is then the same as that for the scalar case and we again obtain well-balanced high order WENO schemes. Examples of such systems will be given in Section 6. We should also mention that local characteristic decomposition is typically used in high order WENO schemes in order to obtain better non-oscillatory property for strong discontinuities. When reconstructing the point value at $x_{i+\frac{1}{2}}$, the local characteristic matrix R , consisting of the right eigenvectors of the Jacobian at $u_{i+\frac{1}{2}}$, is a constant matrix for fixed i . Hence this characteristic decomposition procedure does not alter the argument presented above for the scalar case.

5. Construction of well-balanced discontinuous Galerkin schemes

In this section, we generalize the idea used in Section 4 to RKDG schemes. A well-balanced high order RKDG scheme will be designed for a class of conservation laws satisfying Assumptions 4.1 and 4.2. The basic idea is the same as that for the finite volume schemes, such as the technique of decomposing the source term and replacing the viscosity term in the numerical fluxes. We start with the description in the scalar case.

Consider now Eq. (1.2). Following the description in Section 3, the semi-discrete DG scheme for (1.2) is

$$\int_{I_j} \partial_t u_h(x, t) v_h(x) \, dx - \int_{I_j} f(u_h(x, t)) \partial_x v_h(x) \, dx + \hat{f}_{j+\frac{1}{2}} v_h(x_{j+\frac{1}{2}}^-) - \hat{f}_{j-\frac{1}{2}} v_h(x_{j-\frac{1}{2}}^+) = \int_{I_j} g(u_h(x, t), t) v_h(x) \, dx, \tag{5.1}$$

$$\int_{I_j} u_h(x, 0) v_h(x) \, dx = \int_{I_j} u_0(x) v_h(x) \, dx. \tag{5.2}$$

First, we define a high order approximation $a_h(u_h, x) = \frac{u_h + p_h}{q_h}$ to $a(u_h, x)$, where p_h and q_h are L^2 projections of p and q into V_h , see (5.2) for such a projection. Now assume that u is the steady state solution satisfying (4.2), namely

$$u(x) + p(x) = cq(x)$$

for some constant c , and u_h is the L^2 projection of this steady state solution. Clearly, since the L^2 projection is a linear operator,

$$u_h(x) + p_h(x) = cq_h(x)$$

for the same constant c at every point x . This implies

$$a_h(u_h, x) = \frac{u_h(x) + p_h(x)}{q_h(x)} = c$$

for the same constant c .

For such steady state solution u satisfying Assumptions 4.1 and 4.2, we have

$$\frac{d}{dx} \left(f(u, x) - \sum_j s_j(a(u, x)) t_j(x) \right) = 0.$$

We would need to suitably choose a function $(t_j)_h$, which should be a high order approximation to t_j and should satisfy the condition

$$f(u_h(x)) - \sum_j s_j(a_h(u_h(x), x)) (t_j)_h(x) = \text{constant} \tag{5.3}$$

for all x . The construction of $(t_j)_h$ follows a similar procedure as that for the construction of $(t_j)_{i+\frac{1}{2}}^\pm$ for the finite volume well-balanced scheme in Section 4. We will describe in detail the construction of $(t_j)_h$ for each application case in Section 6.

Similar to the decomposition of the source term in the well-balanced finite volume schemes (4.9), we decompose the integral of the source term on the right-hand side of (5.1) as:

$$\begin{aligned} \int_{I_i} g(u_h, x) v_h \, dx &= \sum_j \left(\frac{1}{2} \left(s_j(a(u_h, x)_{i-\frac{1}{2}}^+) + s_j(a(u_h, x)_{i+\frac{1}{2}}^-) \right) \int_{I_i} t_j'(x) v_h \, dx \right. \\ &\quad \left. + \int_{I_i} \left(s_j(a(u_h, x)) - \frac{1}{2} \left(s_j(a(u_h, x)_{i-\frac{1}{2}}^+) + s_j(a(u_h, x)_{i+\frac{1}{2}}^-) \right) \right) t_j'(x) v_h \, dx \right) \\ &= \sum_j \left(\frac{1}{2} \left(s_j(a(u_h, x)_{i-\frac{1}{2}}^+) + s_j(a(u_h, x)_{i+\frac{1}{2}}^-) \right) \left(t_j(x_{i+\frac{1}{2}}) v_h(x_{i+\frac{1}{2}}^-) - t_j(x_{i-\frac{1}{2}}) v_h(x_{i-\frac{1}{2}}^+) \right) \right. \\ &\quad \left. - \int_{I_i} t_j(x) v_h'(x) \, dx \right) + \int_{I_i} \left(s_j(a(u_h, x)) - \frac{1}{2} \left(s_j(a(u_h, x)_{i-\frac{1}{2}}^+) + s_j(a(u_h, x)_{i+\frac{1}{2}}^-) \right) \right) t_j'(x) v_h \, dx \end{aligned}$$

We then replace this source term with a high order approximation of it given by

$$\begin{aligned} &\sum_j \left(\frac{1}{2} \left(s_j(a_h(u_h, x)_{i-\frac{1}{2}}^+) + s_j(a_h(u_h, x)_{i+\frac{1}{2}}^-) \right) \left((\hat{t}_j)_{h,i+\frac{1}{2}} v_h(x_{i+\frac{1}{2}}^-) - (\hat{t}_j)_{h,i-\frac{1}{2}} v_h(x_{i-\frac{1}{2}}^+) - \int_{I_i} (t_j)_h(x) v_h'(x) \, dx \right) \right. \\ &\quad \left. + \int_{I_i} \left(s_j(a_h(u_h, x)) - \frac{1}{2} \left(s_j(a_h(u_h, x)_{i-\frac{1}{2}}^+) + s_j(a_h(u_h, x)_{i+\frac{1}{2}}^-) \right) \right) (t_j)_h(x) v_h \, dx \right), \end{aligned} \tag{5.4}$$

where $(\hat{t}_j)_{h,i+\frac{1}{2}}$ is a high order approximation to $t_j(x_{i+\frac{1}{2}})$, whose definition will be described below.

To deal with the ‘‘hat’’ terms (numerical fluxes and approximations to $t_j(x_{i+\frac{1}{2}})$), we use the relation between the finite volume schemes and the DG schemes. If we simply take the test function v_h in the DG scheme as the constant function 1, we obtain the evolution of the cell averages similar to that for a finite volume scheme. We have already explained the construction of the ‘‘hat’’ terms for well-balanced finite volume schemes. Here we simply copy those definitions (4.14) and (4.16) from Section 4 without further explanation

$$\begin{aligned} \hat{f}_{i+\frac{1}{2}} &= \frac{1}{2} \left[f((u_h)_{i+\frac{1}{2}}^-) + f((u_h)_{i+\frac{1}{2}}^+) - \alpha \operatorname{sign}(q(x)) \left(a_h(u_h, x)_{i+\frac{1}{2}}^+ - a_h(u_h, x)_{i+\frac{1}{2}}^- \right) \right], \\ (\hat{t}_j)_{h,i+\frac{1}{2}} &= \frac{1}{2} \left[(t_j)_h(x_{i+\frac{1}{2}}^-) + (t_j)_h(x_{i+\frac{1}{2}}^+) \right]. \end{aligned}$$

A combination of the above equations gives the final version of our well-balanced high order RKDG schemes if one more modification on the slope limiter procedure is provided. Usually, we perform the limiter on the function u_h after each Runge–Kutta stage. Now, our purpose is to maintain the steady state solution u which satisfies $a(u, x) = \text{constant}$. The above limiter procedure could destroy the preservation of such steady state, since if the limiter is enacted, the resulting modified solution u_h may no longer satisfy $a_h(u_h, x) = \text{constant}$. We therefore propose to first check whether any limiting is needed based on the function $a_h(u_h, x)$ in each Runge–Kutta stage, where the cell averages of $a_h(u_h, x)$ (needed to implement the TVB limiter) are computed by a suitable Gauss quadrature. If a certain cell is flagged by this procedure needing limiting, then the actual limiter is implemented on u_h , not on $a_h(u_h, x)$. When the limiting procedure is implemented this way, if the steady state u satisfying $a(u, x) = \text{constant}$ is reached, no cell will be flagged as requiring limiting since $a_h(u_h, x)$ is equal to the same constant, hence u_h will not be limited and therefore the steady state is preserved.

It is easy to compute the remaining integrals because u_h , $(t_j)_h$ and v_h are all piecewise polynomials in the space V_h . This finishes the description of the RKDG schemes. We can clearly observe that the accuracy is maintained. We also state below the proposition claiming the exact preservation of the steady state solution (4.2). The proof is similar to that of Proposition 4.3 for the finite volume schemes, and is therefore omitted.

Proposition 5.1. *The RKDG schemes as stated above are exact for steady state solutions satisfying (4.2) and can maintain the original high order accuracy for general solutions.*

The extension of the well-balanced high order RKDG schemes to the system case follows the same idea as that for the well-balanced finite volume schemes.

6. Applications

In this section, we give several examples from applications which fall into the category of balance laws considered in the previous sections, and present well-balanced high order finite volume WENO and discontinuous Galerkin schemes for them. Due to page limitation, only selected numerical results are shown to give a glimpse of how these methods work. Fifth order finite volume WENO scheme and third order finite element RKDG scheme are implemented as examples. In all numerical tests, time discretization is by the third order TVD Runge–Kutta method in [30]. For finite volume WENO schemes, the CFL number is taken as 0.6, except for the accuracy tests where smaller time steps are taken to ensure that spatial errors dominate. For the third order RKDG scheme, the CFL number is 0.18. For the TVB limiter implemented in the RKDG scheme, the TVB constant M (see [28,8] for its definition) is taken as 0 in most numerical examples, unless otherwise stated.

6.1. One dimensional shallow water equations

The shallow water equations have wide applications in ocean and hydraulic engineering and river, reservoir, and open channel flows, among others. We consider the system with a geometrical source term due to the bottom topology. In one space dimension, the equations take the form

$$\begin{cases} h_t + (hu)_x = 0, \\ (hu)_t + (hu^2 + \frac{1}{2}gh^2)_x = -ghb_x, \end{cases} \tag{6.1}$$

where h denotes the water height, u is the velocity of the fluid, b represents the given bottom topography and g is the gravitational constant.

The steady state solution we are interested in preserving satisfies (4.17) in the form

$$a_1 \equiv h + b = \text{constant}, \quad a_2 \equiv u = 0.$$

The first component of the source term is 0. A decomposition of the second component of the source term in the form of (4.3) is

$$-ghb_x = -g(h + b)b_x + \frac{1}{2}g(b^2)_x,$$

i.e., $s_1 = s_1(a_1) = -g(h + b)$, $s_2 = \frac{1}{2}g$, $t_1(x) = b(x)$, and $t_2(x) = b^2(x)$. For the finite volume schemes, we apply the WENO reconstruction to the function $(b(x), 0)^T$, with coefficients computed from $(h, hu)^T$, to obtain $b_{i+\frac{1}{2}}^\pm$. We define

$$(t_1)_{i+\frac{1}{2}}^\pm = b_{i+\frac{1}{2}}^\pm, \quad (t_2)_{i+\frac{1}{2}}^\pm = (b_{i+\frac{1}{2}}^\pm)^2.$$

Under these definitions and if the steady state $h + b = c$, $u = 0$ is reached, we have

$$\begin{aligned} f(u_{i+\frac{1}{2}}^-) - \sum_j s_j (a(u, x)_{i+\frac{1}{2}}^-) (t_j)_{i+\frac{1}{2}}^- &= \frac{1}{2}g (h_{i+\frac{1}{2}}^-)^2 - \frac{1}{2}g (b_{i+\frac{1}{2}}^-)^2 + g \frac{1}{2} (h_{i+\frac{1}{2}}^- + b_{i+\frac{1}{2}}^- + h_{i-\frac{1}{2}}^+ + b_{i-\frac{1}{2}}^+) b_{i+\frac{1}{2}}^- \\ &= \frac{1}{2}g (h_{i+\frac{1}{2}}^- + b_{i+\frac{1}{2}}^-) (h_{i+\frac{1}{2}}^- - b_{i+\frac{1}{2}}^-) + gcb_{i+\frac{1}{2}}^- = \frac{1}{2}gc (h_{i+\frac{1}{2}}^- - b_{i+\frac{1}{2}}^- + 2b_{i+\frac{1}{2}}^-) \\ &= \frac{1}{2}gc^2, \end{aligned}$$

which is a constant. A similar manipulation leads to

$$f(u_{i+\frac{1}{2}}^+) - \sum_j s_j (a(u, x)_{i+\frac{1}{2}}^+) (t_j)_{i+\frac{1}{2}}^+ = \frac{1}{2}gc^2.$$

For the RKDG method, we define

$$(t_1)_h(x) = b_h(x), \quad (t_2)_h(x) = (b_h(x))^2,$$

where $b_h(x)$ is the L^2 projection of $b(x)$ to the finite element space V_h . A similar manipulation as in the finite volume case leads to

$$f(u_h) - \sum_j s_j(a_h(u_h, x))(t_j)_h = \frac{1}{2}gc^2$$

when the steady state $h + b = c, u = 0$ is reached, satisfying our requirement.

Next, we provide numerical results to demonstrate the good properties of the well-balanced finite volume WENO and finite element RKDG schemes when applied to the one dimensional shallow water equations. The gravitation constant g is taken as 9.812 m/s^2 during the computation.

6.1.1. Test for the exact C-property

The purpose of the first test problem is to verify that the schemes indeed maintain the exact C-property over a non-flat bottom. We choose two different functions for the bottom topography given by ($0 \leq x \leq 10$):

$$b(x) = 5e^{-\frac{2}{3}(x-5)^2}, \tag{6.2}$$

which is smooth, and

$$b(x) = \begin{cases} 4 & \text{if } 4 \leq x \leq 8, \\ 0 & \text{otherwise,} \end{cases} \tag{6.3}$$

which is discontinuous. The initial data are the stationary solution:

$$h + b = 10, \quad hu = 0.$$

This steady state should be exactly preserved. We compute the solution until $t = 0.5$ using $N = 200$ uniform cells. In order to demonstrate that the exact C-property is indeed maintained up to round-off error, we use single precision, double precision and quadruple precision to perform the computation, and show the L^1 and L^∞ errors for the water height h (note: h in this case is not a constant function!) and the discharge hu in Tables 1 and 2 for the two bottom functions (6.2) and (6.3) and different precisions. For the RKDG method, the errors are computed based on the numerical solutions at cell centers. We can clearly see that the L^1 and L^∞ errors are at the level of round-off errors for different precisions, verifying the exact C-property.

We have also computed stationary solutions using initial conditions which are not the steady state solutions and letting time evolve into a steady state, obtaining similar results with the exact C-property.

6.1.2. Testing the orders of accuracy

In this example we will test the high order accuracy of our schemes for a smooth solution. There are some known exact solutions to the shallow water equation with non-flat bottom in the literature, such as some stationary solutions, but they are not generic test cases for accuracy. We have therefore chosen to use the following bottom function and initial conditions:

$$b(x) = \sin^2(\pi x), \quad h(x, 0) = 5 + e^{\cos(2\pi x)}, \quad (hu)(x, 0) = \sin(\cos(2\pi x)), \quad x \in [0, 1]$$

Table 1
 L^1 and L^∞ errors for different precisions for the stationary solution with a smooth bottom (6.2)

	Precision	L^1 error		L^∞ error	
		h	hu	h	hu
FV	Single	4.07E – 06	3.75E – 05	1.33E – 05	1.33E – 04
	Double	2.50E – 14	2.23E – 13	7.64E – 14	7.97E – 13
	Quadruple	3.49E – 33	2.90E – 32	1.39E – 32	9.62E – 32
RKDG	Single	6.44E – 06	2.44E – 05	2.57E – 05	1.75E – 04
	Double	6.82E – 15	2.90E – 14	2.84E – 14	2.14E – 13
	Quadruple	9.06E – 31	3.92E – 33	8.05E – 29	1.12E – 31

Table 2
 L^1 and L^∞ errors for different precisions for the stationary solution with a nonsmooth bottom (6.3)

	Precision	L^1 error		L^∞ error	
		h	hu	h	hu
FV	Single	6.50E – 06	2.61E – 05	1.91E – 05	1.53E – 04
	Double	1.73E – 14	5.88E – 14	4.62E – 14	2.43E – 13
	Quadruple	2.69E – 32	9.30E – 32	5.85E – 32	3.04E – 31
RKDG	Single	5.76E – 07	3.54E – 07	9.54E – 07	1.18E – 06
	Double	1.41E – 15	8.90E – 16	3.55E – 15	2.83E – 15
	Quadruple	2.69E – 31	1.62E – 35	8.06E – 29	8.18E – 34

with periodic boundary conditions, see [35]. Since the exact solution is not known explicitly for this case, we use the fifth order finite volume WENO scheme with $N = 12,800$ cells to compute a reference solution, and treat this reference solution as the exact solution in computing the numerical errors. We compute up to $t = 0.1$ when the solution is still smooth (shocks develop later in time for this problem). Table 3 contains the L^1 errors for the cell averages and numerical orders of accuracy for the finite volume and RKDG schemes, respectively. We can clearly see that fifth order accuracy is achieved for the WENO scheme, and third order accuracy is achieved for the RKDG scheme. For the RKDG scheme, the TVB constant M is taken as 32. Notice that the CFL number we have used for the finite volume scheme decreases with the mesh size and is recorded in Table 3. For the RKDG method, the CFL number is fixed at 0.18.

6.1.3. A small perturbation of a steady state water

The following quasi-stationary test case was proposed by LeVeque [21]. It was chosen to demonstrate the capability of the proposed scheme for computations on a rapidly varying flow over a smooth bed, and the perturbation of a stationary state.

The bottom topography consists of one hump:

$$b(x) = \begin{cases} 0.25(\cos(10\pi(x - 1.5)) + 1) & \text{if } 1.4 \leq x \leq 1.6, \\ 0 & \text{otherwise.} \end{cases} \tag{6.4}$$

The initial conditions are given with

$$(hu)(x, 0) = 0 \quad \text{and} \quad h(x, 0) = \begin{cases} 1 - b(x) + \epsilon & \text{if } 1.1 \leq x \leq 1.2, \\ 1 - b(x) & \text{otherwise,} \end{cases} \tag{6.5}$$

where ϵ is a non-zero perturbation constant. Two cases have been run: $\epsilon = 0.2$ (big pulse) and $\epsilon = 0.001$ (small pulse). Theoretically, for small ϵ , this disturbance should split into two waves, propagating left and right at the characteristic speeds $\pm\sqrt{gh}$. Many numerical methods have difficulty with the calculations involving such small perturbations of the water surface [21]. Both sets of initial conditions are shown in Fig. 1. The solution at time $t = 0.2$ s for the big pulse $\epsilon = 0.2$, obtained on a 200 cell uniform grid with simple transmissive boundary conditions, and compared with a 3000 cell solution, is shown in Fig. 2 for the FV scheme and in Fig. 4 for

Table 3
 L^1 errors and numerical orders of accuracy for the example in Section 6.1.2

No. of cells	FV schemes				RKDG schemes				
	CFL	h		hu		h		hu	
		L^1 error	Order	L^1 error	Order	L^1 error	Order	L^1 error	Order
25	0.6	1.48E – 02		9.45E – 02		2.35E – 03		2.12E – 02	
50	0.6	2.40E – 03	2.63	1.98E – 02	2.26	1.15E – 04	4.36	1.01E – 03	4.39
100	0.4	2.97E – 04	3.01	2.58E – 03	2.93	1.24E – 05	3.20	1.09E – 04	3.21
200	0.3	2.43E – 05	3.61	2.13E – 04	3.60	1.02E – 06	3.59	8.97E – 06	3.60
400	0.2	1.02E – 06	4.57	8.96E – 06	4.57	1.11E – 07	3.19	9.79E – 07	3.19
800	0.1	3.26E – 08	4.97	2.85E – 07	4.97	1.30E – 08	3.09	1.14E – 07	3.08

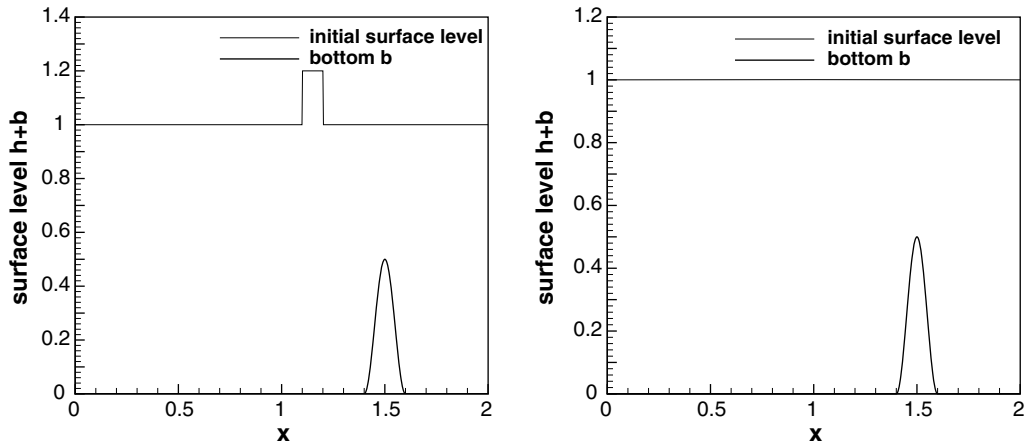


Fig. 1. The initial surface level $h + b$ and the bottom b for a small perturbation of a steady state water. Left: a big pulse $\epsilon = 0.2$; right: a small pulse $\epsilon = 0.001$.

the RKDG scheme. The results for the small pulse $\epsilon = 0.001$ are shown in Figs. 3 and 5. For this small pulse problem, we take $\epsilon = 10^{-9}$ in the WENO weight formula (2.5), such that it is smaller than the square of the perturbation. At this time, the downstream-traveling water pulse has already passed the bump. We can clearly see that there are no spurious numerical oscillations.

6.1.4. The dam breaking problem over a rectangular bump

In this example we simulate the dam breaking problem over a rectangular bump, which involves a rapidly varying flow over a discontinuous bottom topography. This example was used in [33].

The bottom topography takes the form:

$$b(x) = \begin{cases} 8 & \text{if } |x - 750| \leq 1500/8, \\ 0 & \text{otherwise} \end{cases} \tag{6.6}$$

for $x \in [0, 1500]$. The initial conditions are

$$(hu)(x, 0) = 0 \quad \text{and} \quad h(x, 0) = \begin{cases} 20 - b(x) & \text{if } x \leq 750, \\ 15 - b(x) & \text{otherwise.} \end{cases} \tag{6.7}$$

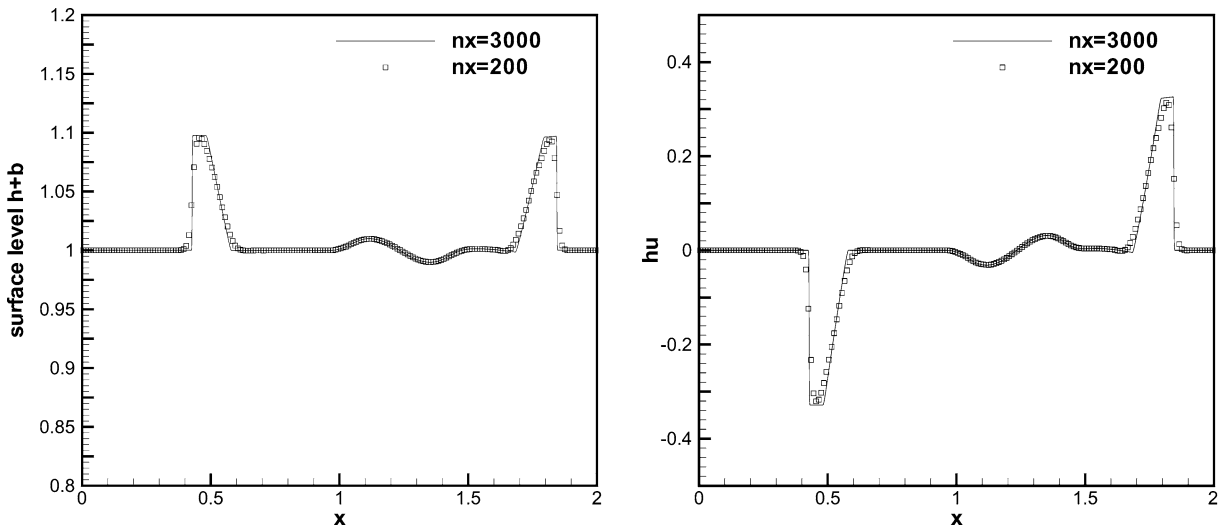


Fig. 2. FV scheme: small perturbation of a steady state water with a big pulse. $t = 0.2$ s. Left: surface level $h + b$; right: the discharge hu .

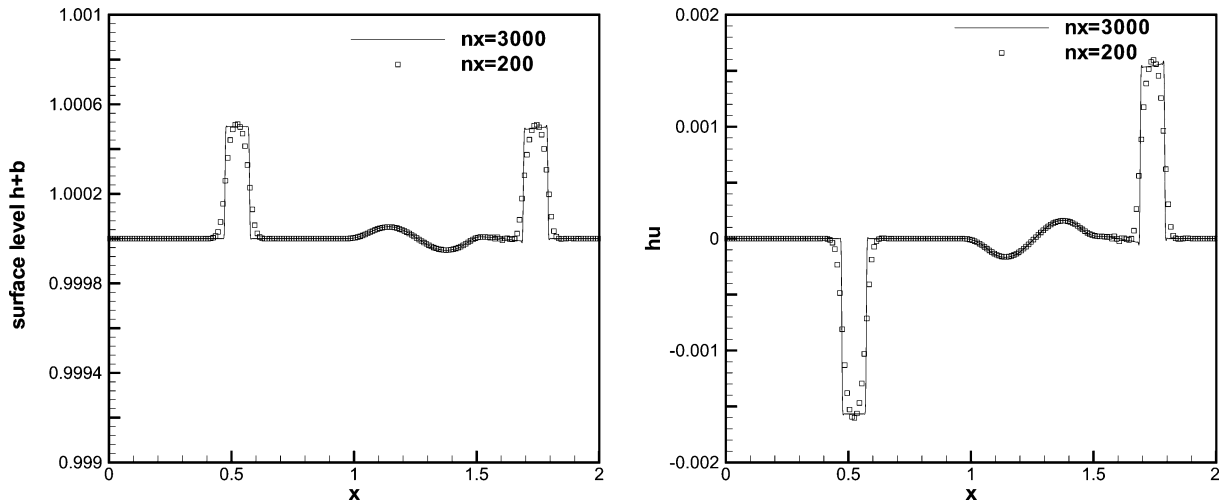


Fig. 3. FV scheme: small perturbation of a steady state water with a small pulse. $t = 0.2$ s. Left: surface level $h + b$; right: the discharge hu .

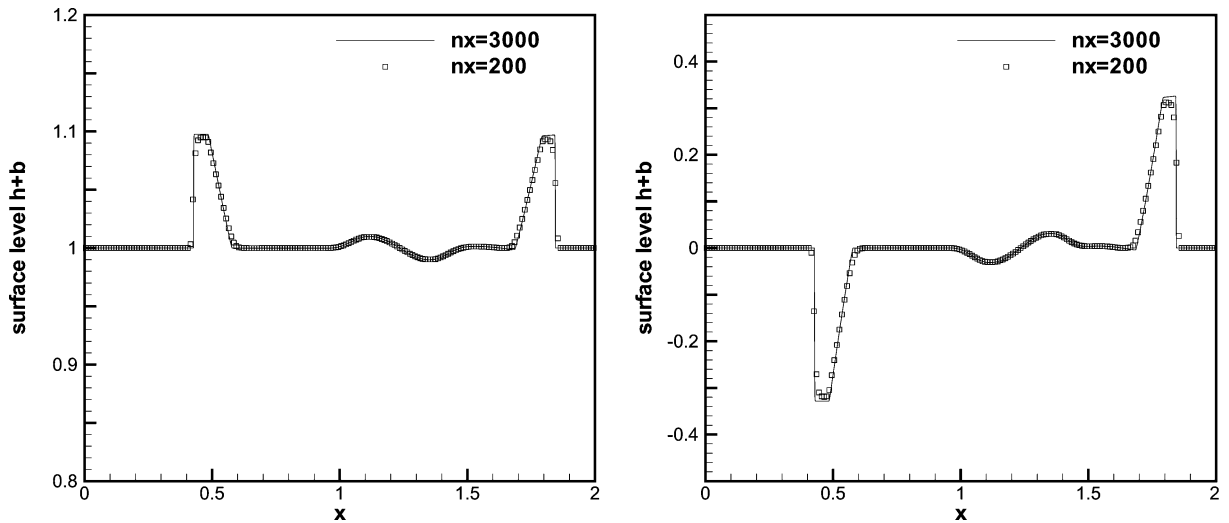


Fig. 4. RKDG scheme: small perturbation of a steady state water with a big pulse. $t = 0.2$ s. Left: surface level $h + b$; right: the discharge hu .

The numerical results obtained by the FV scheme with 400 uniform cells (and a comparison with the results using 4000 uniform cells) are shown in Figs. 6 and 7, with two different ending time $t = 15$ s and $t = 60$ s. Figs. 8 and 9 demonstrate the numerical results by the RKDG scheme, with the same number of uniform cells. In this example, the water height $h(x)$ is discontinuous at the points $x = 562.5$ and $x = 937.5$, while the surface level $h(x) + b(x)$ is smooth there. Both schemes work well for this example, giving well resolved, non-oscillatory solutions using 400 cells which agree with the converged results using 4000 cells.

6.1.5. Steady flow over a hump

The purpose of this test case is to study the convergence in time towards steady flow over a bump. These are classical test problems for transcritical and subcritical flows, and they are widely used to test numerical schemes for shallow water equations. For example, they have been considered by the *working group on dam break modelling* [12], and have been used as test cases in, e.g. [32].

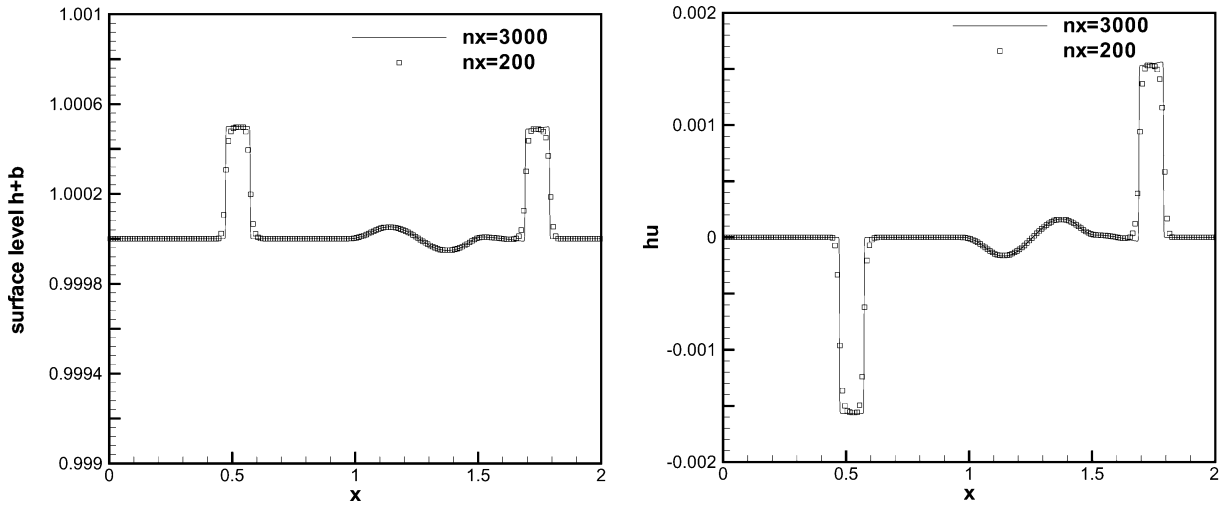


Fig. 5. RKDG scheme: small perturbation of a steady state water with a small pulse. $t = 0.2$ s. Left: surface level $h + b$; right: the discharge hu .

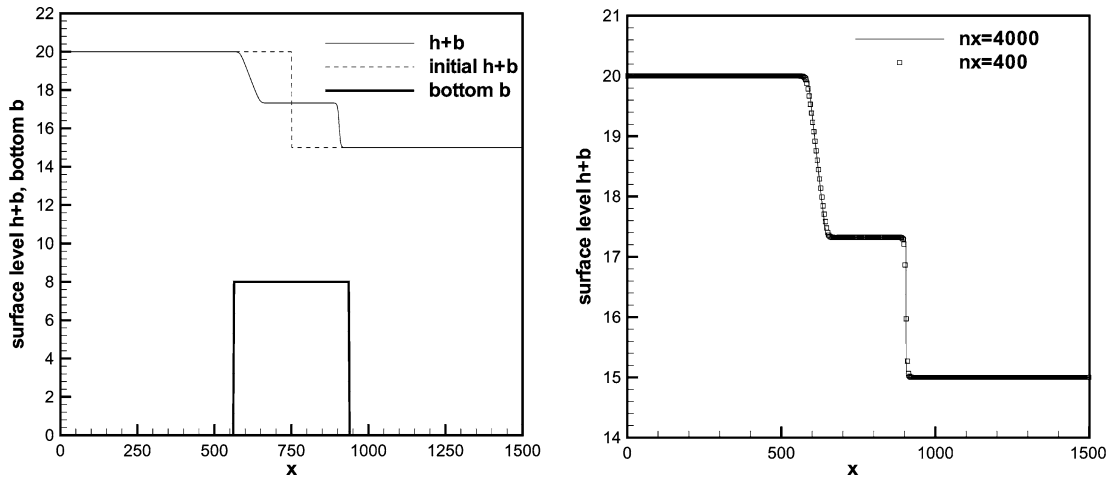


Fig. 6. FV scheme: The surface level $h + b$ for the dam breaking problem at time $t = 15$ s. Left: the numerical solution using 400 grid cells, plotted with the initial condition and the bottom topography; right: the numerical solution using 400 and 4000 grid cells.

The bottom function is given by:

$$b(x) = \begin{cases} 0.2 - 0.05(x - 10)^2 & \text{if } 8 \leq x \leq 12, \\ 0 & \text{otherwise} \end{cases} \tag{6.8}$$

for a channel of length 25 m. The initial conditions are taken as

$$h(x, 0) = 0.5 - b(x) \quad \text{and} \quad u(x, 0) = 0.$$

Depending on different boundary conditions, the flow can be subcritical or transcritical with or without a steady shock. The computational parameters common for all three cases are: uniform mesh size $\Delta x = 0.125$ m, ending time $t = 200$ s. Analytical solutions for the various cases are given in [12].

(a) Transcritical flow without a shock.

- upstream: The discharge $hu = 1.53 \text{ m}^2/\text{s}$ is imposed.
- downstream: The water height $h = 0.66$ m is imposed when the flow is subcritical.

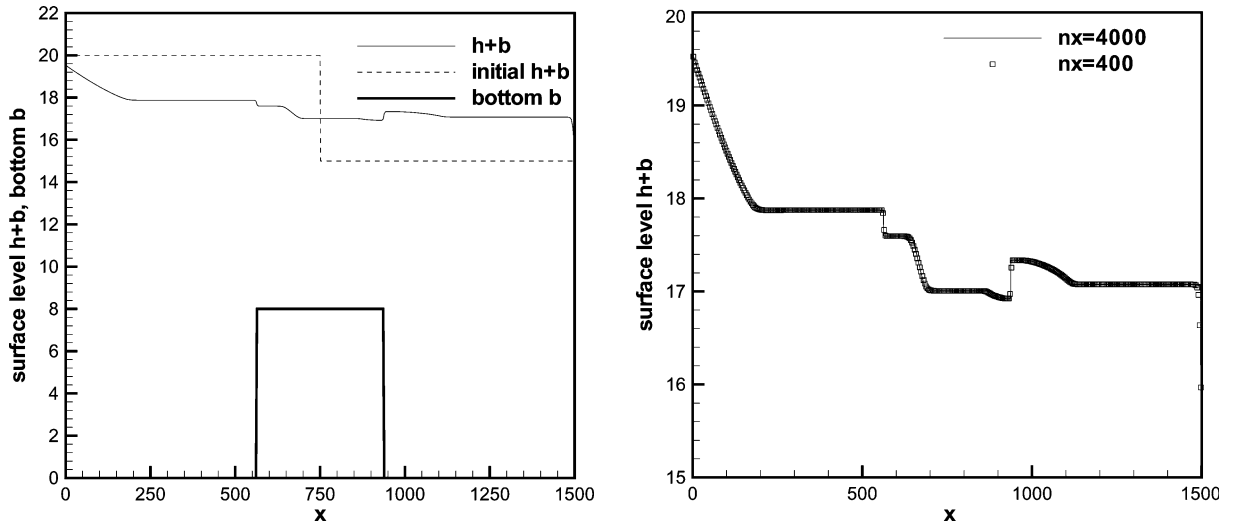


Fig. 7. FV scheme: The surface level $h + b$ for the dam breaking problem at time $t = 60$ s. Left: the numerical solution using 400 grid cells, plotted with the initial condition and the bottom topography; right: the numerical solution using 400 and 4000 grid cells.

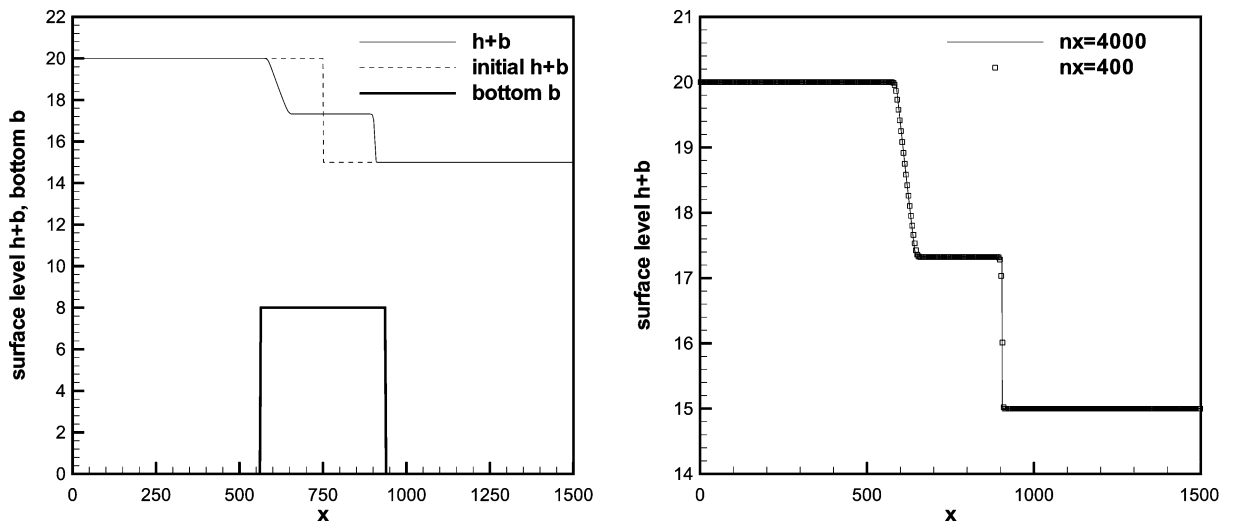


Fig. 8. RKDG scheme: The surface level $h + b$ for the dam breaking problem at time $t = 15$ s. Left: the numerical solution using 400 grid cells, plotted with the initial condition and the bottom topography; right: the numerical solution using 400 and 4000 grid cells.

The surface level $h + b$ and the discharge hu , as the numerical flux for the water height h in Eq. (6.1), are plotted in Figs. 10 and 11, which show very good agreement with the analytical solution. The correct capturing of the discharge hu is usually more difficult than the surface level $h + b$, as noticed by many authors. The numerical errors for the discharge hu of our well-balanced finite volume WENO and RKDG schemes are both very small.

(b) Transcritical flow with a shock.

- upstream: The discharge $hu = 0.18 \text{ m}^2/\text{s}$ is imposed.
- downstream: The water height $h = 0.33 \text{ m}$ is imposed.

In this case, the Froude number $Fr = u/\sqrt{gh}$ increases to a value larger than one above the bump, and then decreases to less than one. A stationary shock can appear on the surface. The surface level $h + b$ and the discharge hu , as the numerical flux for the water height h in Eq. (6.1), are plotted in Fig. 12 and 14, which

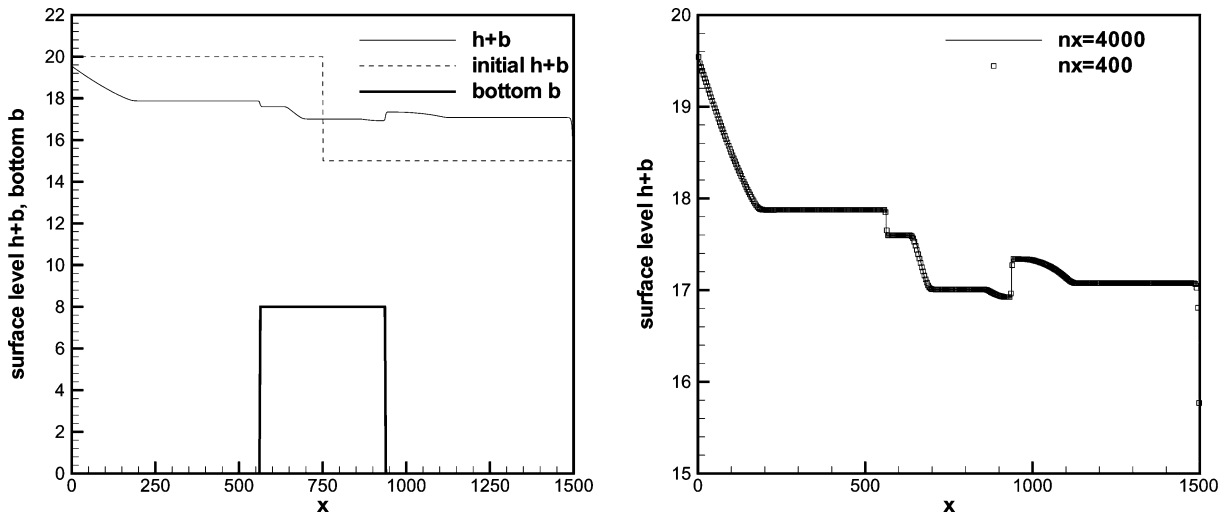


Fig. 9. RKDG scheme: The surface level $h + b$ for the dam breaking problem at time $t = 60$ s. Left: the numerical solution using 400 grid cells, plotted with the initial condition and the bottom topography; right: the numerical solution using 400 and 4000 grid cells.

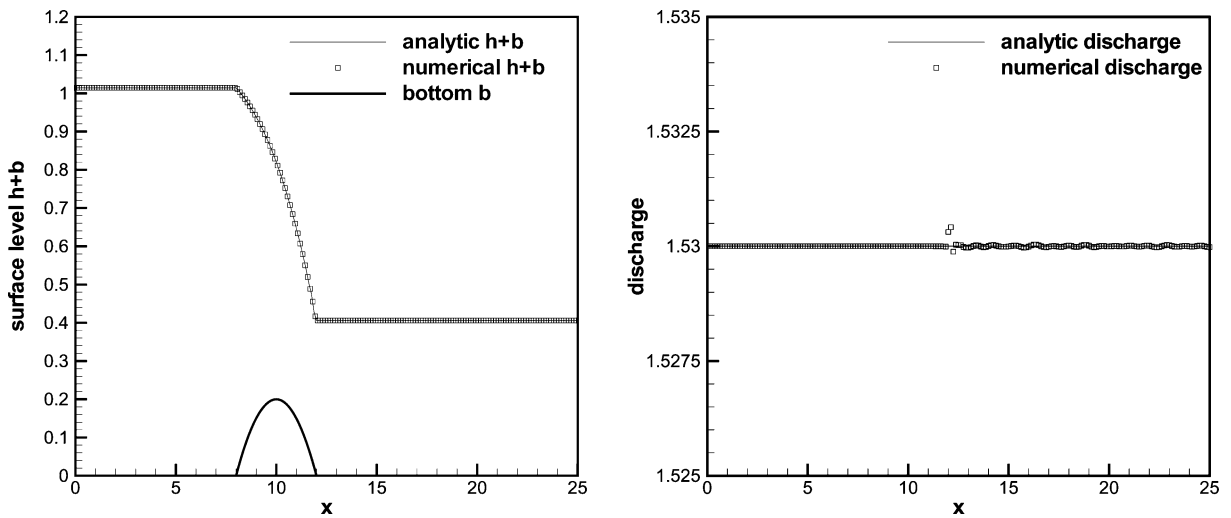


Fig. 10. FV scheme: steady transcritical flow over a bump without a shock. Left: the surface level $h + b$; right: the discharge hu as the numerical flux for the water height h .

show non-oscillatory results in good agreement with the analytical solution. In Figs. 13 and 15, we compare the pointwise errors of the numerical solutions obtained with 200 and 400 uniform cells. We have also performed such error comparisons for the cases of the transcritical flow without a shock and of the subcritical flow, obtaining qualitatively similar results. We have therefore omitted them to save space.

(c) Subcritical flow.

- upstream: The discharge $hu = 4.42 \text{ m}^2/\text{s}$ is imposed.
- downstream: The water height $h = 2 \text{ m}$ is imposed.

This is a subcritical flow. The surface level $h + b$ and the discharge hu , as the numerical flux for the water height h in Eq. (6.1), are plotted in Figs. 16 and 17, which are in good agreement with the analytical solution.

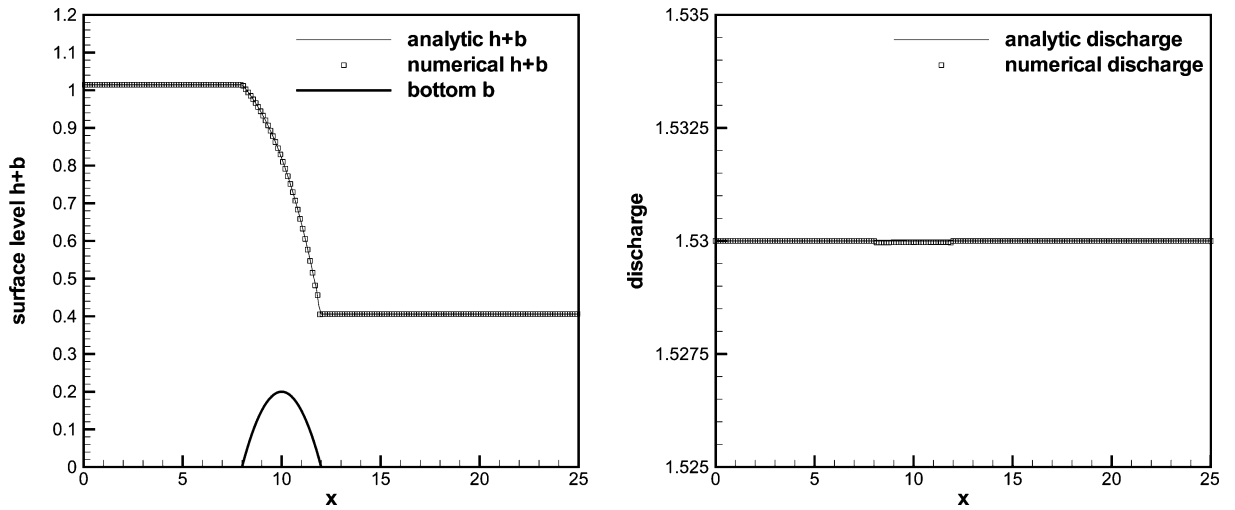


Fig. 11. RKDG scheme: steady transcritical flow over a bump without a shock. Left: the surface level $h + b$; right: the discharge hu as the numerical flux for the water height h .

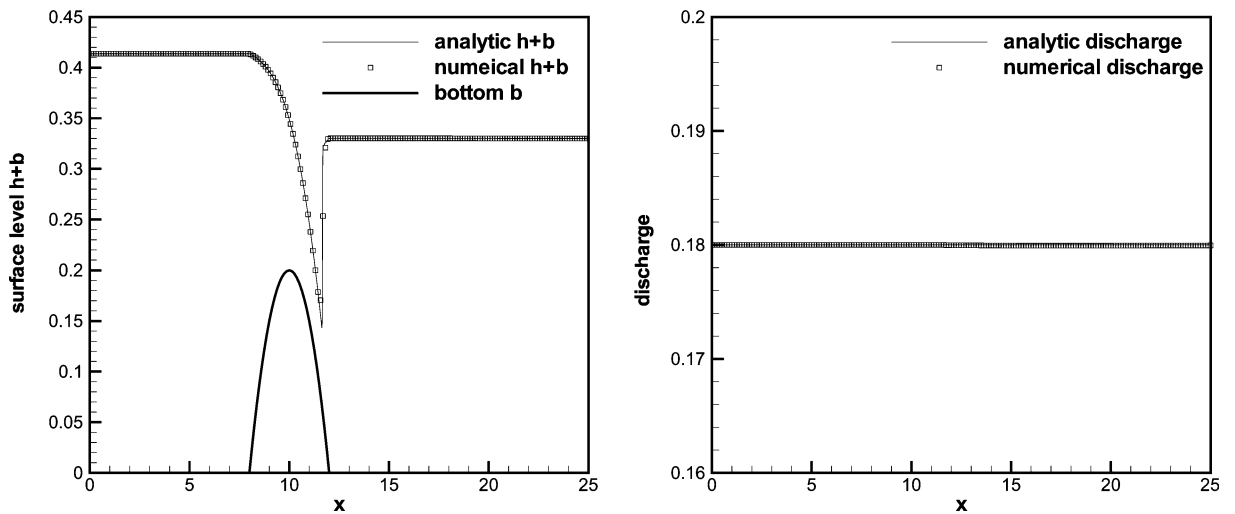


Fig. 12. FV scheme: steady transcritical flow over a bump with a shock. Left: the surface level $h + b$; right: the discharge hu as the numerical flux for the water height h .

6.2. Two dimensional shallow water equations

The shallow water system in two space dimensions takes the form:

$$\begin{cases} h_t + (hu)_x + (hv)_y = 0, \\ (hu)_t + (hu^2 + \frac{1}{2}gh^2)_x + (huv)_y = -ghb_x, \\ (hv)_t + (huv)_x + (hv^2 + \frac{1}{2}gh^2)_y = -ghb_y, \end{cases} \tag{6.9}$$

where again h is the water height, (u, v) is the velocity of the fluid, b represents the bottom topography and g is the gravitational constant.

We are interested in preserving the still water solution, which takes the form (satisfying (4.17))

$$a_1 \equiv h + b = \text{constant}, \quad a_2 \equiv u = 0, \quad a_3 \equiv v = 0.$$

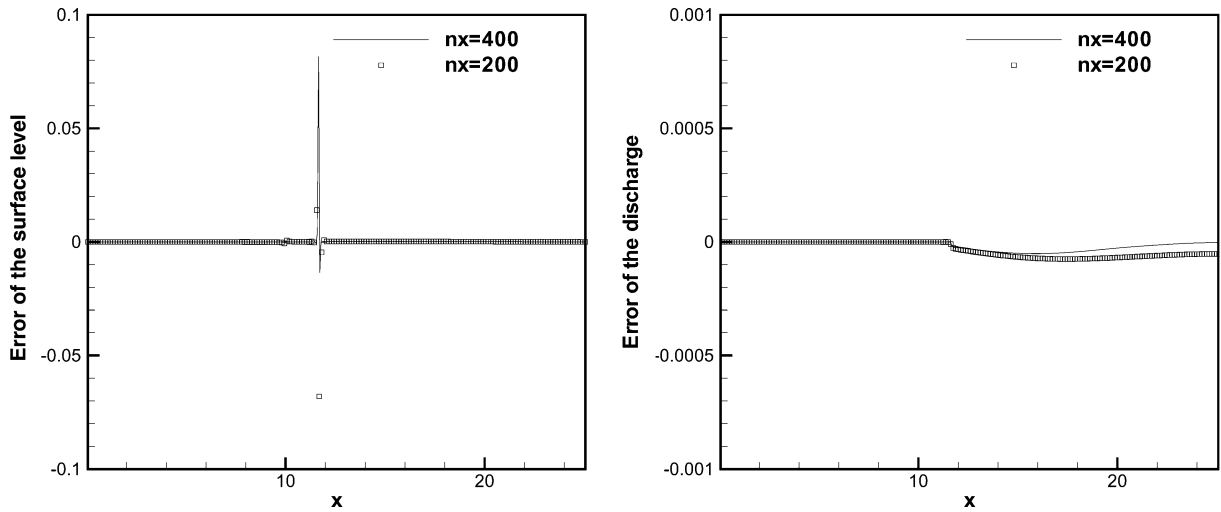


Fig. 13. FV scheme: steady transcritical flow over a bump with a shock. Pointwise error comparison between numerical solutions using 200 and 400 cells. Left: the surface level $h + b$; right: the discharge hu as the numerical flux for the water height h .

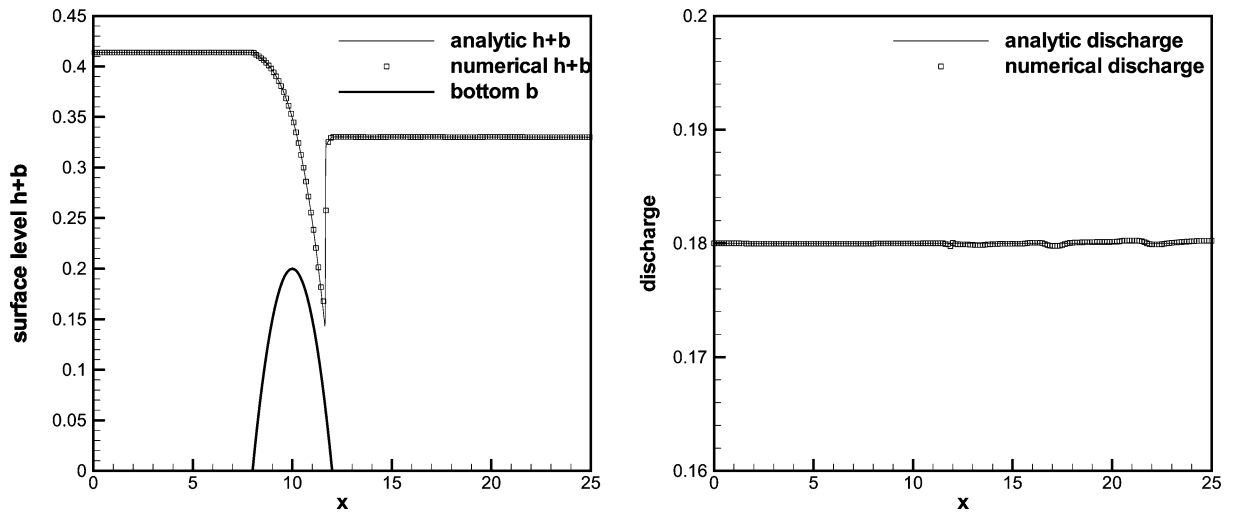


Fig. 14. RKDG scheme: steady transcritical flow over a bump with a shock. Left: the surface level $h + b$; right: the discharge hu as the numerical flux for the water height h .

The first component of the source term is 0. Similarly as in one dimensional case, we decompose the second and third components of the source term as

$$-ghb_x = -g(h + b)b_x + \frac{1}{2}g(b^2)_x, \quad -ghb_y = -g(h + b)b_y + \frac{1}{2}g(b^2)_y,$$

i.e., $s_1 = s_1(a_1) = -g(h + b)$, $s_2 = \frac{1}{2}g$, $t_1(x) = b(x)$, $t_2(x) = b^2(x)$ for the second component, and $s_1 = s_1(a_1) = -g(h + b)$, $s_2 = \frac{1}{2}g$, $t_1(x) = b(x)$, $t_2(x) = b^2(x)$ for the third component.

For the finite volume scheme, we apply the WENO reconstruction to the function $(b(x), 0, 0)^T$, with coefficients computed from $(h, hu, hv)^T$, to obtain $b_{i+\frac{1}{2}j}^\pm$ and $b_{i,j+\frac{1}{2}}^\pm$. We define, for the source term of the second equation,

$$(t_1)_{i+\frac{1}{2}j}^\pm = b_{i+\frac{1}{2}j}^\pm, \quad (t_2)_{i+\frac{1}{2}j}^\pm = \left(b_{i+\frac{1}{2}j}^\pm\right)^2,$$

and, for the source term of the third equation,

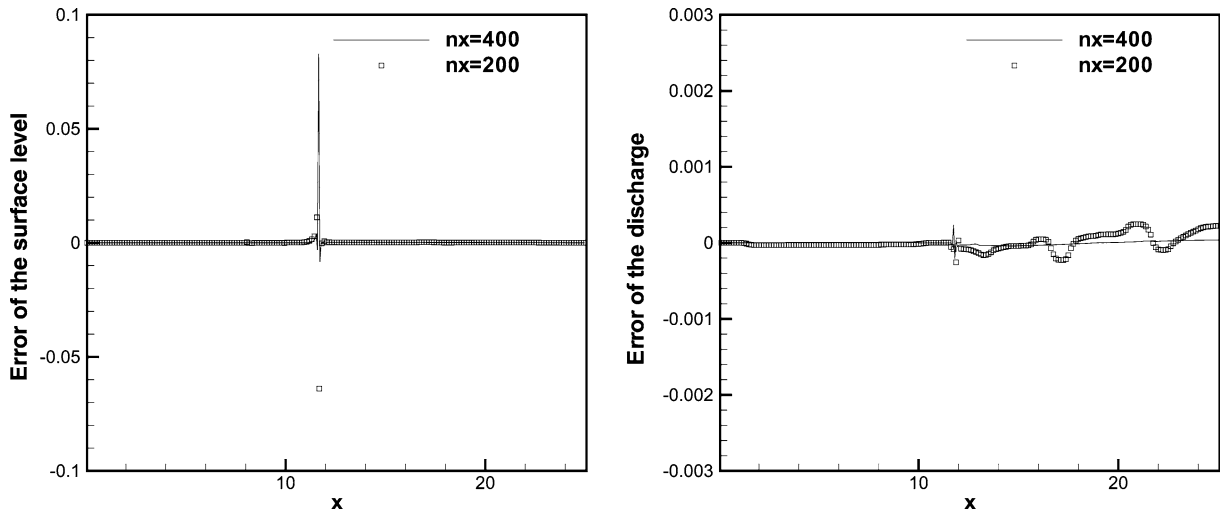


Fig. 15. RKDG scheme: steady transcritical flow over a bump with a shock. Pointwise error comparison between numerical solutions using 200 and 400 cells. Left: the surface level $h + b$; right: the discharge hu as the numerical flux for the water height h .

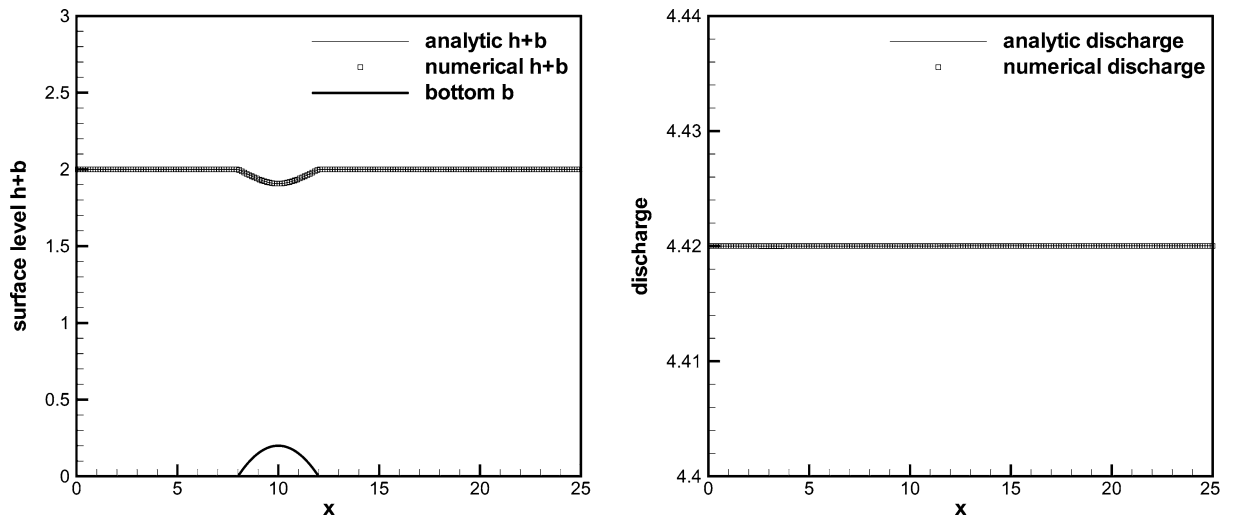


Fig. 16. FV scheme: steady subcritical flow over a bump. Left: the surface level $h + b$; right: the discharge hu as the numerical flux for the water height h .

$$(t_1)_{i,j+\frac{1}{2}}^\pm = b_{i,j+\frac{1}{2}}^\pm, \quad (t_2)_{i,j+\frac{1}{2}}^\pm = (b_{i,j+\frac{1}{2}}^\pm)^2.$$

We can verify, similar to the one dimensional case, that these choices of t_j^\pm will maintain the requirement for the steady state solution satisfying $h + b = c$, $u = v = 0$.

For the RKDG method, we define

$$(t_1)_h(x, y) = b_h(x, y), \quad (t_2)_h(x, y) = (b_h(x, y))^2,$$

where $b_h(x, y)$ is the L^2 projection of $b(x, y)$ to the finite element space V_h , for the source terms of both the second and the third equations.

We now show numerical examples to demonstrate the behavior of our well-balanced schemes for the two dimensional shallow water equations.

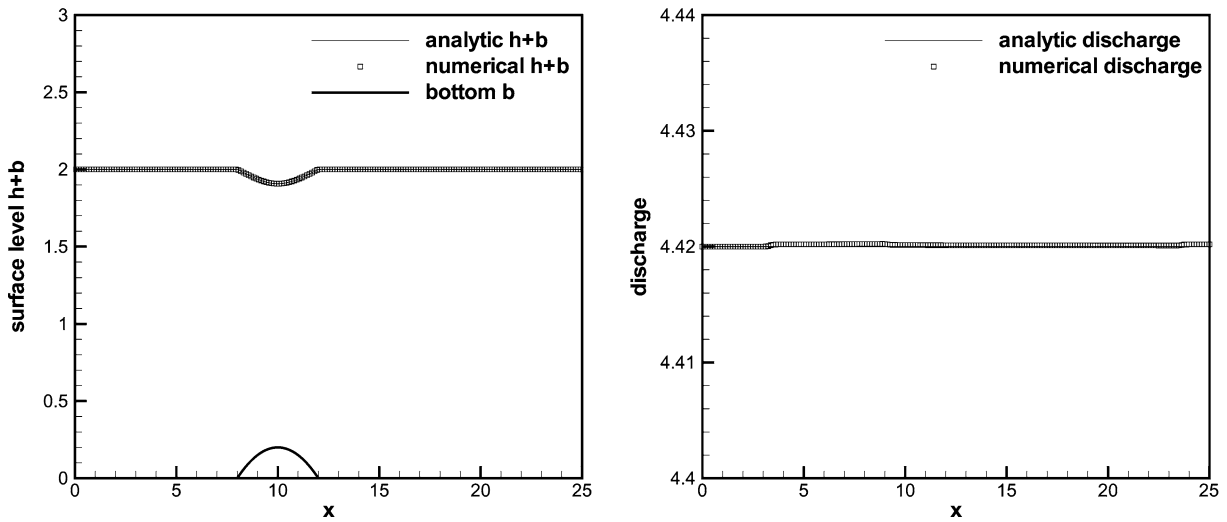


Fig. 17. RKDG scheme: steady subcritical flow over a bump. Left: the surface level $h + b$; right: the discharge hu as the numerical flux for the water height h .

6.2.1. Test for the exact C-property in two dimensions

This example is used to check that our schemes indeed maintain the exact C-property over a non-flat bottom. The two dimensional hump

$$b(x, y) = 0.8e^{-50((x-0.5)^2+(y-0.5)^2)}, \quad x, y \in [0, 1] \tag{6.10}$$

is chosen to be the bottom. $h(x, y, 0) = 1 - b(x, y)$ is the initial depth of the water. Initial velocity is set to be zero. This surface should remain flat. The computation is performed to $t = 0.1$ using single, double and quadruple precisions with a 100×100 uniform mesh. Table 4 contains the L^1 errors for the water height h (which is not a constant function) and the discharges hu and hv for both schemes. We can clearly see that the L^1 errors are at the level of round-off errors for different precisions, verifying the exact C-property.

6.2.2. Testing the orders of accuracy

In this example we check the numerical orders of accuracy when the schemes are applied to the following two dimensional problem. The bottom topography and the initial data are given by:

$$b(x, y) = \sin(2\pi x) + \cos(2\pi y), \quad h(x, y, 0) = 10 + e^{\sin(2\pi x)} \cos(2\pi y),$$

$$(hu)(x, y, 0) = \sin(\cos(2\pi x)) \sin(2\pi y), \quad (hv)(x, y, 0) = \cos(2\pi x) \cos(\sin(2\pi y))$$

defined over a unit square, with periodic boundary conditions. The terminal time is taken as $t = 0.05$ to avoid the appearance of shocks in the solution. Since the exact solution is also not known explicitly for this case, we use the same fifth order WENO scheme with an extremely refined mesh consisting of 1600×1600 cells to

Table 4
 L^1 errors for different precisions for the stationary solution in Section 6.2.1

	Precision	L^1 error			
		h	hu	hv	
FV	Single	1.09E - 06	8.87E - 07	8.87E - 07	
	Double	8.16E - 16	9.31E - 16	8.47E - 16	
	Quadruple	7.30E - 34	7.31E - 34	7.34E - 34	
RKDG	Single	9.40E - 08	3.58E - 07	3.60E - 07	
	Double	6.20E - 17	1.14E - 15	1.16E - 15	
	Quadruple	5.87E - 34	8.35E - 34	8.36E - 34	

compute a reference solution, and treat this reference solution as the exact solution in computing the numerical errors. The TVB constant M in the limiter for the RKDG scheme is taken as 40 here. Tables 5 and 6 contain the L^1 errors and orders of accuracy for the cell averages. We can clearly see that, in this two dimensional test case, fifth order accuracy is achieved for the finite volume WENO scheme and third order accuracy is achieved for the RKDG scheme.

6.2.3. A small perturbation of a two dimensional steady state water

This is a classical example to show the capability of the proposed scheme for the perturbation of the stationary state, given by LeVeque [21]. It is analogous to the test done previously in Section 6.1.3 in one dimension.

We solve the system in the rectangular domain $[0, 2] \times [0, 1]$. The bottom topography is an isolated elliptical shaped hump:

$$b(x, y) = 0.8e^{-5(x-0.9)^2 - 50(y-0.5)^2}. \tag{6.11}$$

The surface is initially given by:

$$h(x, y, 0) = \begin{cases} 1 - b(x, y) + 0.01 & \text{if } 0.05 \leq x \leq 0.15, \\ 1 - b(x, y) & \text{otherwise,} \end{cases} \tag{6.12}$$

$$hu(x, y, 0) = hv(x, y, 0) = 0.$$

So the surface is almost flat except for $0.05 \leq x \leq 0.15$, where h is perturbed upward by 0.01. Figs. 18 and 19 display the right-going disturbance as it propagates past the hump, on two different uniform meshes with 200×100 cells and 600×300 cells for comparison. The surface level $h + b$ is presented at different times. The results indicate that both schemes can resolve the complex small features of the flow very well.

6.3. Elastic wave equation

We consider the propagation of compressional waves [1,34] in an one dimensional elastic rod with a given media density $\rho(x)$. The equations of motion in a Lagrangian frame are given by the balance laws:

$$\begin{cases} (\rho\varepsilon)_t + (-\rho u)_x = -u \frac{d\rho}{dx}, \\ (\rho u)_t + (-\sigma)_x = 0, \end{cases} \tag{6.13}$$

Table 5
FV scheme: L^1 errors and numerical orders of accuracy for the example in Section 6.2.2

Number of cells	CFL	h		hu		hv	
		L^1 error	Order	L^1 error	Order	L^1 error	Order
25×25	0.6	7.91E – 03		2.12E – 02		6.52E – 02	
50×50	0.6	1.13E – 03	2.81	2.01E – 03	3.40	9.23E – 03	2.82
100×100	0.6	8.89E – 05	3.66	1.25E – 04	4.00	7.19E – 04	3.68
200×200	0.4	4.07E – 06	4.45	5.19E – 06	4.59	3.30E – 05	4.45
400×400	0.3	1.42E – 07	4.84	1.84E – 07	4.82	1.16E – 06	4.84
800×800	0.2	4.38E – 09	5.02	5.99E – 09	4.94	3.63E – 08	4.99

Table 6
RKDG scheme: L^1 errors and numerical orders of accuracy for the example in Section 6.2.2

Number of cells	h		hu		hv	
	L^1 error	Order	L^1 error	Order	L^1 error	Order
25×25	2.45E – 03		1.36E – 02		2.05E – 02	
50×50	5.73E – 04	2.10	2.92E – 03	2.22	4.75E – 03	2.11
100×100	1.06E – 04	2.43	5.31E – 04	2.46	8.51E – 04	2.48
200×200	1.71E – 05	2.63	8.82E – 05	2.59	1.39E – 04	2.61
400×400	2.52E – 06	2.76	1.32E – 05	2.74	2.10E – 05	2.73
800×800	3.52E – 07	2.84	1.89E – 06	2.80	3.01E – 06	2.81

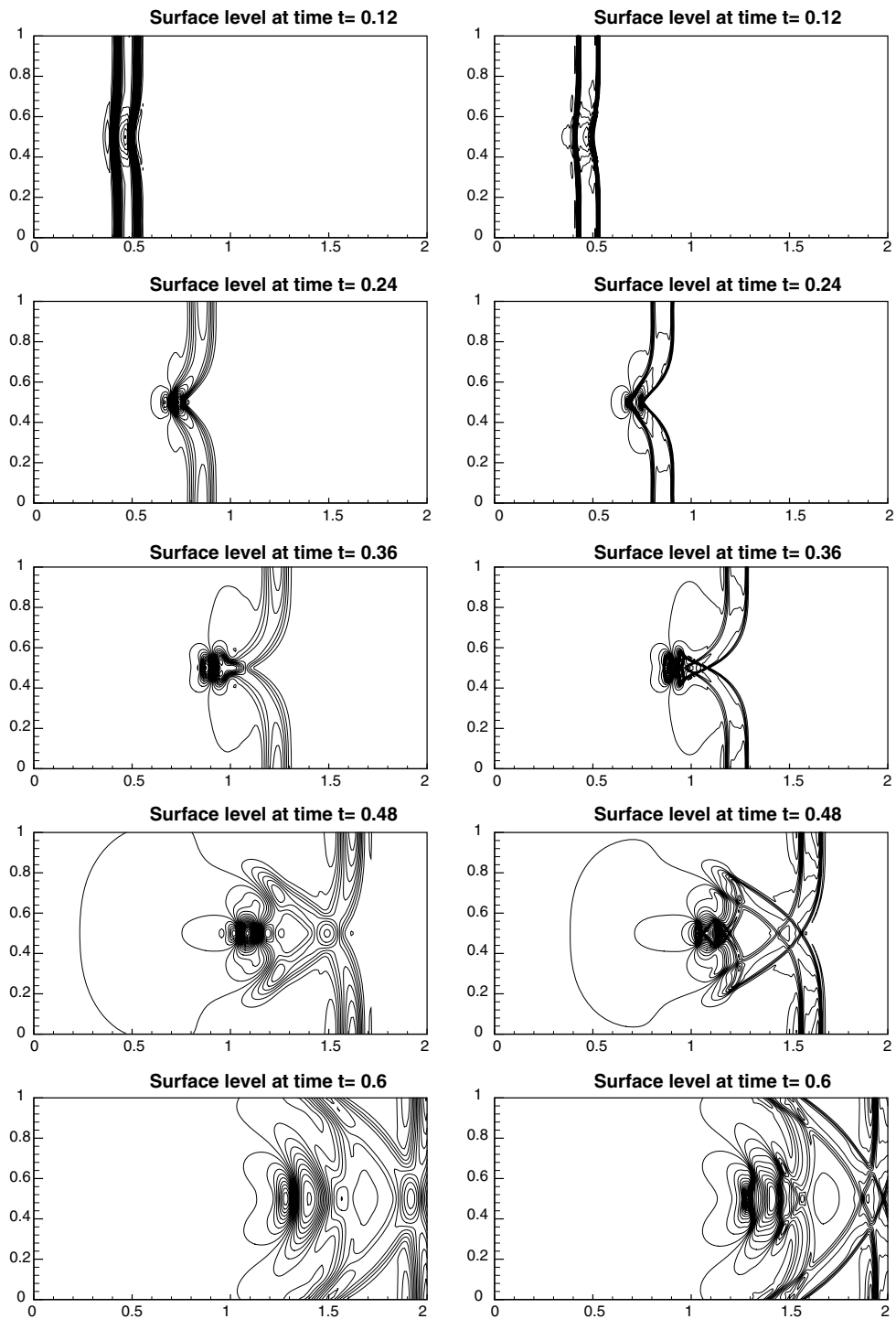


Fig. 18. FV scheme: The contours of the surface level $h + b$ for the problem in Section 6.2.3. 30 uniformly spaced contour lines. From top to bottom: at time $t = 0.12$ from 0.99942 to 1.00656; at time $t = 0.24$ from 0.99318 to 1.01659; at time $t = 0.36$ from 0.98814 to 1.01161; at time $t = 0.48$ from 0.99023 to 1.00508; and at time $t = 0.6$ from 0.99514 to 1.00629. Left: results with a 200×100 uniform mesh. Right: results with a 600×300 uniform mesh.

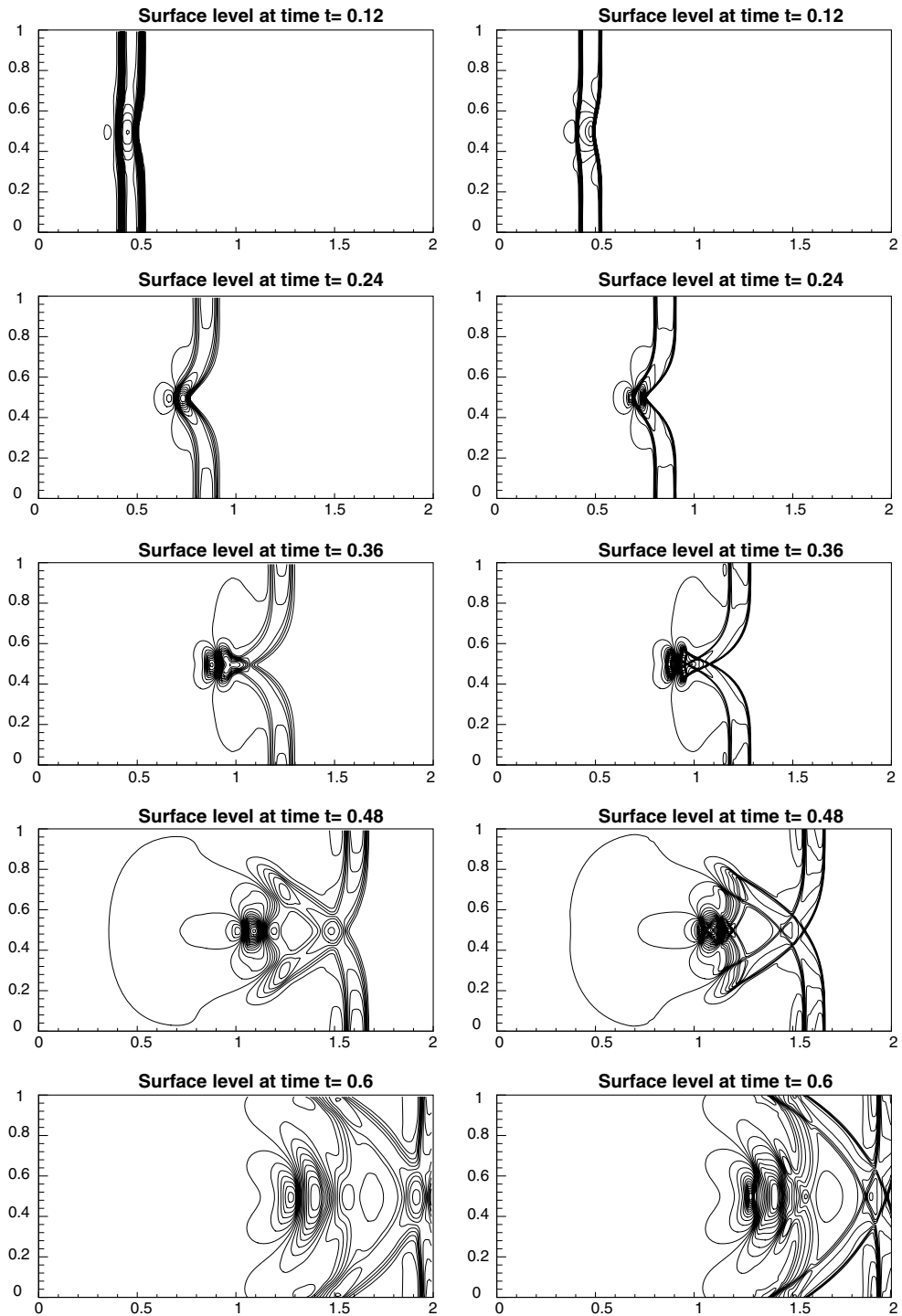


Fig. 19. RKDG scheme: The contours of the surface level $h + b$ for the problem in Section 6.2.3. 30 uniformly spaced contour lines. From top to bottom: at time $t = 0.12$ from 0.99942 to 1.00656; at time $t = 0.24$ from 0.99318 to 1.01659; at time $t = 0.36$ from 0.98814 to 1.01161; at time $t = 0.48$ from 0.99023 to 1.00508; and at time $t = 0.6$ from 0.99514 to 1.00629. Left: results with a 200×100 uniform mesh. Right: results with a 600×300 uniform mesh.

where $\varepsilon = \varepsilon(x, t)$ is the strain, $u = u(x, t)$ is the velocity and σ is a given stress–strain relationship $\sigma(\varepsilon, x)$. The equation of linear acoustics can be obtained from the elasticity problem if the stress–strain relationship is linear,

$$\sigma(\varepsilon, x) = K(x)\varepsilon,$$

where $K(x)$ is the given bulk modulus of compressibility.

The steady state we are interested to preserve for this problem is characterized by

$$a_1 \equiv \sigma(\varepsilon, x) = \text{constant}, \quad a_2 \equiv u = \text{constant}$$

which is of the form (4.17). The second component of the source term is 0. The first component of the source term is already in the form of (4.3) with $s_1 = s_1(a_2) = -u = -\frac{\rho u}{\rho}$ and $t_1 = \rho(x)$.

For finite volume schemes, we apply the WENO reconstruction to the function $(0, \rho(x))^T$, with coefficients computed from $(\rho\varepsilon, \rho u)^T$, to obtain $\rho_{i+\frac{1}{2}}^\pm$. We then define $(t_1)_{i+\frac{1}{2}}^\pm = \rho_{i+\frac{1}{2}}^\pm$, which leads to

$$f(u_{i+\frac{1}{2}}^\pm) - \sum_j s_j(a(u, x)_{i+\frac{1}{2}}^\pm)(t_j)_{i+\frac{1}{2}}^\pm = -(\rho u)_{i+\frac{1}{2}}^\pm + \frac{(\rho u)_{i+\frac{1}{2}}^\pm}{\rho_{i+\frac{1}{2}}^\pm} \rho_{i+\frac{1}{2}}^\pm = 0,$$

satisfying our requirement. For the RKDG scheme, we define

$$(t_1)_h(x) = \rho_h(x),$$

where $\rho_h(x)$ is the L^2 projection of $\rho(x)$ to the finite element space V_h . We can then easily verify the requirement

$$f(u_h) - \sum_j s_j(a_h(u_h, x))(t_j)_h = 0$$

for the steady state solution.

Next, we present the numerical result for a linear acoustic test [1]. The properties of the media are given by

$$c(x) = \sqrt{\frac{K(x)}{\rho(x)}} = 1 + 0.5 \sin(10\pi x), \quad Z(x) = \rho(x)c(x) = 1 + 0.25 \cos(10\pi x).$$

The initial conditions are given by

$$\rho\varepsilon(x, 0) = \begin{cases} \frac{-1.75+0.75\cos(10\pi x)}{c^2(x)} & \text{if } 0.4 < x < 0.6, \\ \frac{-1}{c^2(x)} & \text{otherwise,} \end{cases} \quad u(x, 0) = 0.$$

It is a test case where the impedance $Z(x)$ and hence the eigenvectors are both spatially varying. We perform the computation with 200 uniform cells, with the ending time $t = 0.4$ s. An “exact” reference solution is computed with the same scheme over a 2000 uniform cells. The simulation results are shown in Fig. 20. The numerical resolution shows very good agreement with the “exact” reference solution.

6.4. Chemosensitive movement

Originated from biology, chemosensitive movement [11,15] is a process by which cells change their direction reacting to the presence of a chemical substance, approaching chemically favorable environments and avoiding unfavorable ones. Hyperbolic models for chemotaxis are recently introduced [15] and take the form

$$\begin{cases} n_t + (nu)_x = 0, \\ (nu)_t + (nu^2 + n)_x = n\chi'(c) \frac{\partial c}{\partial x} - \sigma nu, \end{cases} \quad (6.14)$$

where the chemical concentration $c = c(x, t)$ is given by the parabolic equation

$$\frac{\partial c}{\partial t} - D_c \Delta c = n - c.$$

Here, $n(x, t)$ is the cell density, $nu(x, t)$ is the population flux and σ is the friction coefficient.

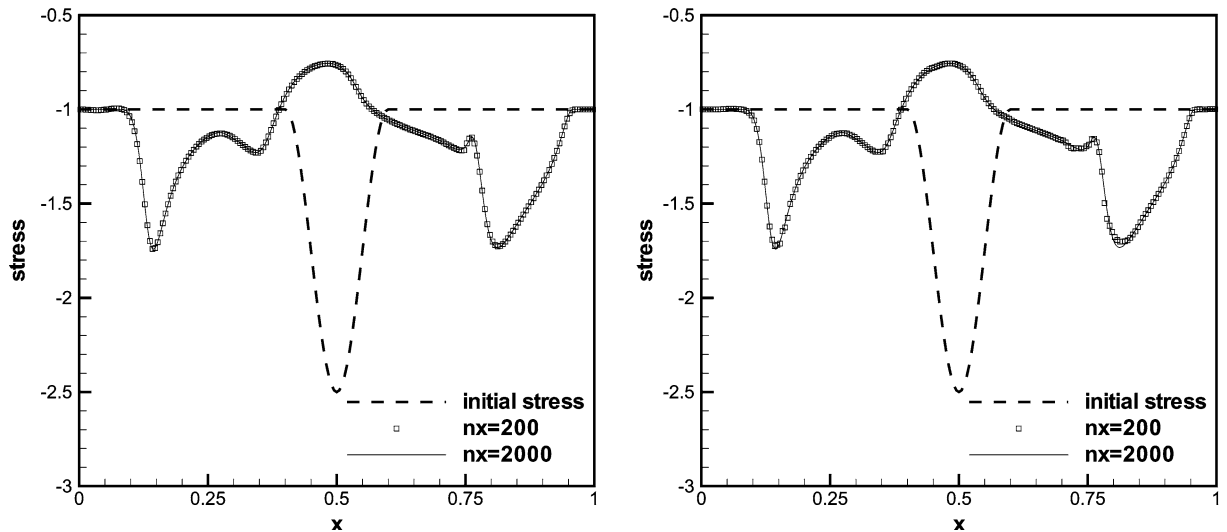


Fig. 20. The numerical (symbols) and the “exact” reference (solid line) stress $\sigma(x)$ at time $t = 0.4$ s. Left: FV schemes; right: RKDG schemes.

We would like to preserve the steady state solution to (6.14) with a zero population flux, which satisfies

$$n\chi'(c)c_x - n_x = 0, \quad nu = 0, \tag{6.15}$$

where $c = c(x)$ does not depend on t in steady state. The first equality above does not seem to be of the form (4.17). However, (6.15) is equivalent to

$$a_1 \equiv \frac{n}{e^{\chi(c)}} = \text{constant}, \quad a_2 \equiv nu = 0,$$

which is clearly in the form of (4.17). The first component of the source term is 0. A decomposition of the second component of the source term in the form of (4.3) is

$$n\chi'(c) \frac{\partial c}{\partial x} - \sigma nu = \frac{n}{e^{\chi(c)}} \frac{d}{dx} e^{\chi(c)} - \sigma nu,$$

i.e., $s_1 = s_1(a_1) = \frac{n}{e^{\chi(c)}}$, $s_2 = s_2(a_2) = \sigma nu$, $t_1(x) = e^{\chi(c(x))}$, and $t_2(x) = x$.

For the finite volume scheme, we apply the WENO reconstruction to the function $(e^{\chi(c(x))}, 0)^T$, with coefficients computed from $(n, nu)^T$, to obtain $(e^{\chi(c(x))})_{i+\frac{1}{2}}^\pm$. We define

$$(t_1)_{i+\frac{1}{2}}^\pm = (e^{\chi(c(x))})_{i+\frac{1}{2}}^\pm, \quad (t_2)_{i+\frac{1}{2}}^\pm = x_{i+\frac{1}{2}}.$$

In the case of steady state,

$$f(u_{i+\frac{1}{2}}^\pm) - \sum_j s_j(a(u, x)_{i+\frac{1}{2}}^\pm)(t_j)_{i+\frac{1}{2}}^\pm = n_{i+\frac{1}{2}}^\pm - \frac{n_{i+\frac{1}{2}}^\pm}{(e^{\chi(c(x))})_{i+\frac{1}{2}}^\pm} (e^{\chi(c(x))})_{i+\frac{1}{2}}^\pm = 0,$$

which satisfies our requirement. For the RKDG scheme, we define

$$(t_1)_h(x) = (e^{\chi(c(x))})_h, \quad (t_2)_h(x) = x,$$

where $(e^{\chi(c(x))})_h$ is the L^2 projection of $e^{\chi(c(x))}$ to the finite element space V_h . A similar manipulation as in the finite volume case leads to

$$f(u_h) - \sum_j s_j(a_h(u_h, x))(t_j)_h = 0.$$

Our technique can also be applied to the two dimensional case of this application.

We give an numerical example here to test the high order accuracy for smooth solutions for our schemes. The initial conditions are taken as

$$n(x, 0) = 1 + 0.2 \cos(\pi x), \quad u(x, 0) = 0, \quad x \in [-1, 1]$$

with

$$c(x) = e^{-16x^2}, \quad \chi(c) = \log(1 + c), \quad \sigma = 0$$

with a periodic boundary condition. Since the exact solution is not known explicitly for this problem, we use the same fifth order WENO scheme with $N = 5120$ cells to compute a reference solution and treat it as the exact solution when computing the numerical errors for the cell averages. Final time $t = 0.5$ s is used to avoid the development of shocks. The TVB constant M in the limiter for the RKDG scheme is taken as 13 for this example. Table 7 contains the L^1 errors and numerical orders of accuracy. We can clearly see that expected order accuracy is achieved for this example.

6.5. A model in fluid mechanics with spherical symmetry

A classical singularity arising in fluid mechanics in case of spherical symmetry leads to the following model equation:

$$u_t + \left(\frac{u^2}{2}\right)_x = \frac{1}{x}u^2, \tag{6.16}$$

which has been considered in [4]. Notice that the source term is a nonlinear function of u . The steady state for this problem is given by

$$\frac{du}{dx} = \frac{u}{x} \Rightarrow a(u, x) \equiv \frac{u}{x} = \text{constant}$$

which is of the form (4.2) with $p(x) = 0$ and $q(x) = x$. The source term can be rewritten as

$$\frac{u^2}{x} = \left(\frac{u}{x}\right)^2 x = \left(\frac{u}{x}\right)^2 \left(\frac{x^2}{2}\right)_x$$

which is in the form of (4.3) with $s_1(a) = a^2 = \left(\frac{u}{x}\right)^2$ and $t_1(x) = \frac{x^2}{2}$. Note that here s_1 is a nonlinear function of a .

For finite volume schemes, we apply the WENO reconstruction to the function $q(x) = x$, with coefficients computed from u , to obtain $q_{i+\frac{1}{2}}^\pm$. Since x is a polynomial with degree 1, the reconstructed $q_{i+\frac{1}{2}}^\pm$ should be exactly $x_{i+\frac{1}{2}}$ no matter how we compute the WENO coefficients. Hence, we can use $x_{i+\frac{1}{2}}$ directly, without applying WENO reconstruction on it. We then define $(t_1)_{i+\frac{1}{2}}^\pm = \frac{x_{i+\frac{1}{2}}^2}{2}$, which leads to

$$f(u_{i+\frac{1}{2}}^\pm) - \sum_j s_j(a(u, x)_{i+\frac{1}{2}}^\pm)(t_j)_{i+\frac{1}{2}}^\pm = \frac{(u_{i+\frac{1}{2}}^\pm)^2}{2} - \left(\frac{u_{i+\frac{1}{2}}^\pm}{x_{i+\frac{1}{2}}^\pm}\right)^2 \frac{x_{i+\frac{1}{2}}^2}{2} = \frac{(u_{i+\frac{1}{2}}^\pm)^2}{2} - \frac{(u_{i+\frac{1}{2}}^\pm)^2}{2} = 0,$$

satisfying our requirement. For the RKDG scheme, we define

Table 7
 L^1 errors and numerical orders of accuracy for the example in Section 6.4

No. of cells	FV schemes					RKDG schemes			
	CFL	$\rho\epsilon$		ρu		$\rho\epsilon$		ρu	
		L^1 error	Order	L^1 error	Order	L^1 error	Order	L^1 error	Order
20	0.6	9.70E-03		7.41E-03		1.27E-04		1.46E-04	
40	0.6	1.03E-03	3.24	8.85E-04	3.07	1.75E-05	2.85	2.07E-05	2.82
80	0.5	1.07E-04	3.26	8.80E-05	3.33	1.32E-06	3.73	1.89E-06	3.46
160	0.4	5.63E-06	4.25	5.63E-06	3.97	1.21E-07	3.45	1.97E-07	3.26
320	0.3	2.21E-07	4.67	1.89E-07	4.89	1.29E-08	3.23	2.27E-08	3.12
640	0.1	7.18E-09	4.94	6.07E-08	4.96	1.57E-09	3.03	2.76E-09	3.04

$$(t_1)_h(x) = \frac{x^2}{2}$$

and we can then easily verify the requirement

$$f(u_h) - \sum_j s_j(a_h(u_h, x))(t_j)_h = 0$$

for the steady state solution.

Next, we present a numerical result to demonstrate the well-balanced property. The initial and boundary conditions are given by

$$u(x, 0) = 0, \quad x \in [-5, 5], \tag{6.17}$$

$$u(x = -5, t) = 10, \quad u(x = 5, t) = -10. \tag{6.18}$$

The choice of these information allows us to compute the steady state, which is $u = -2x$. Numerical computations are performed by the well-balanced version of finite volume WENO schemes and RKDG methods. To see the benefit of well-balanced schemes, we also use a non-well-balanced finite volume WENO schemes and RKDG methods, and compare the results. We use 100 uniform cells here. The comparison of the convergence history, measured by the L^1 norm of the difference with the steady state, is given in Fig. 21. The advantage of the well-balanced schemes can be easily observed. Also, we compute the L^1 and L^∞ errors at time $t = 10$, with single precision and double precision. The results are shown in Table 8. We can clearly see that the errors are at the level of round-off errors for different precisions, verifying the well-balanced property.

6.6. Other applications

There are many other application problems which admit steady states that can be approximated by our well-balanced schemes. These include the nozzle flow problem, a two phase flow model and a typical example with a stiff source term. We refer to [36] for more details of the first two models. The model with a stiff source term takes the form:

$$u_t + u_x = -\frac{1}{\epsilon}u(u - 1). \tag{6.19}$$

We can easily check that our well-balanced schemes can be applied to these models. Due to page limitation, we do not include computational results for these models here.

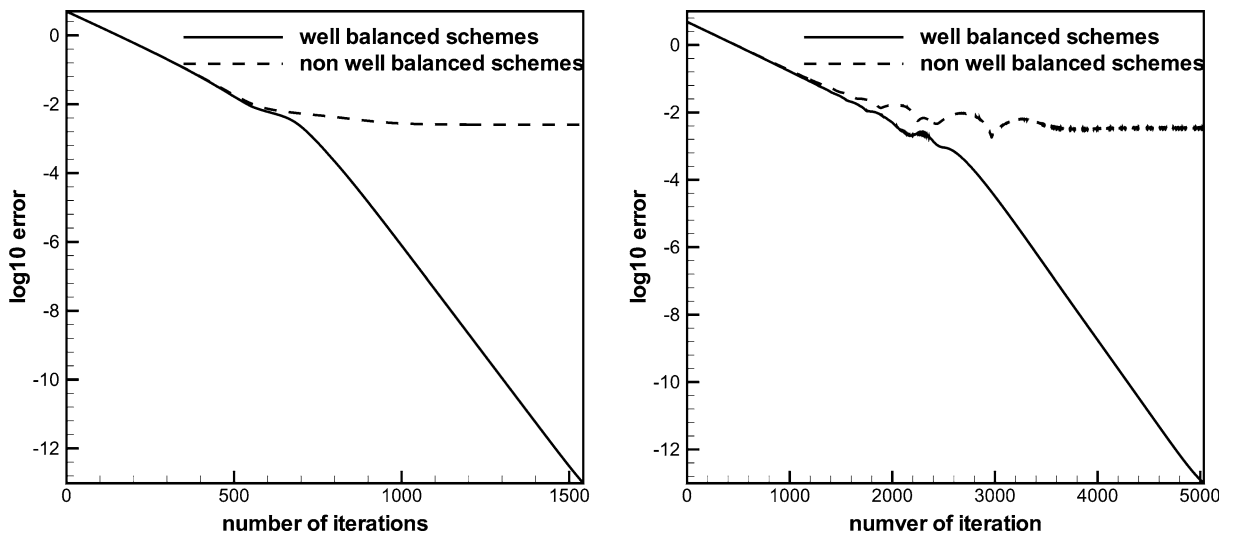


Fig. 21. Comparison of the convergence history in L^1 error. Left: FV WENO schemes; right: RKDG schemes.

Table 8
 L^1 and L^∞ errors for different precisions for the steady state (6.17)–(6.18)

Precision	FV		DG	
	L^1 error	L^∞ error	L^1 error	L^∞ error
Single	6.06E – 06	2.24E – 05	2.63E – 05	9.87E – 05
Double	1.60E – 14	7.42E – 14	3.25E – 14	2.16E – 13

7. Concluding remarks

In this paper, we have extended the high order finite volume WENO and finite element discontinuous Galerkin schemes to solve a class of conservation laws with separable source terms including the shallow water equations, the elastic wave equation, the hyperbolic model for a chemosensitive movement and a model of fluid mechanics in case of spherical symmetry. Our technique can also be applied to other application problems such as the nozzle flow problem and a two phase flow model [36], but we have not included them in this paper to save space. A special decomposition of the source terms allows us to design specific approximations such that the resulting schemes maintain properties of the exact preservation of the balance laws for certain steady state solutions, their original high order accuracy and essentially non-oscillatory property for general solutions. Extensive numerical examples are given to demonstrate the exactness property, accuracy, and non-oscillatory shock resolution of the proposed numerical method.

Acknowledgments

Research supported by ARO Grant W911NF-04-1-0291, NSF Grants DMS-0207451 and DMS-0510345, and AFOSR Grant F49620-02-1-0113.

References

- [1] D.S. Bale, R.J. LeVeque, S. Mitran, J.A. Rossmann, A wave propagation method for conservation laws and balance laws with spatially varying flux functions, *SIAM Journal on Scientific Computing* 24 (2002) 955–978.
- [2] D.S. Balsara, C.-W. Shu, Monotonicity preserving weighted essentially non-oscillatory schemes with increasingly high order of accuracy, *Journal of Computational Physics* 160 (2000) 405–452.
- [3] A. Bermudez, M.E. Vazquez, Upwind methods for hyperbolic conservation laws with source terms, *Computers and Fluids* 23 (1994) 1049–1071.
- [4] R. Botchorishvili, B. Perthame, A. Vasseur, Equilibrium schemes for scalar conservation laws with stiff sources, *Mathematics of Computation* 72 (2003) 131–157. Also, an extended version containing more numerical examples is located at <http://www.inria.fr/rrrt/rr-3891.html>.
- [5] B. Cockburn, Discontinuous Galerkin methods for convection-dominated problems, in: T.J. Barth, H. Deconinck (Eds.), *High-Order Methods for Computational Physics*, Lecture Notes in Computational Science and Engineering, vol. 9, Springer, Berlin, 1999, pp. 69–224.
- [6] B. Cockburn, S. Hou, C.-W. Shu, The Runge–Kutta local projection discontinuous Galerkin finite element method for conservation laws IV: the multidimensional case, *Mathematics of Computation* 54 (1990) 545–581.
- [7] B. Cockburn, S.-Y. Lin, C.-W. Shu, TVB Runge–Kutta local projection discontinuous Galerkin finite element method for conservation laws III: one dimensional systems, *Journal of Computational Physics* 84 (1989) 90–113.
- [8] B. Cockburn, C.-W. Shu, TVB Runge–Kutta local projection discontinuous Galerkin finite element method for conservation laws II: general framework, *Mathematics of Computation* 52 (1989) 411–435.
- [9] B. Cockburn, C.-W. Shu, The Runge–Kutta discontinuous Galerkin method for conservation laws V: multidimensional systems, *Journal of Computational Physics* 141 (1998) 199–224.
- [10] B. Cockburn, C.-W. Shu, Runge–Kutta discontinuous Galerkin methods for convection-dominated problems, *Journal of Scientific Computing* 16 (2001) 173–261.
- [11] F. Filbet, C.-W. Shu, Approximation of hyperbolic models for chemosensitive movement, *SIAM Journal on Scientific Computing* (to appear).
- [12] N. Goutal, F. Maurel, in: *Proceedings of the Second 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.
- [13] J.M. Greenberg, A.Y. LeRoux, A well-balanced scheme for the numerical processing of source terms in hyperbolic equations, *SIAM Journal on Numerical Analysis* 33 (1996) 1–16.

- [14] A. Harten, P.D. Lax, B. Van Leer, On upstream differencing and Godunov-type schemes for hyperbolic conservation laws, *SIAM Review* 25 (1983) 35–61.
- [15] T. Hillen, Hyperbolic models for chemosensitive movement, *Mathematical Models and Methods in Applied Sciences* 12 (2002) 1007–1034.
- [16] C. Hu, C.-W. Shu, Weighted essentially non-oscillatory schemes on triangular meshes, *Journal of Computational Physics* 150 (1999) 97–127.
- [17] M.E. Hubbard, P. Garcia-Navarro, Flux difference splitting and the balancing of source terms and flux gradients, *Journal of Computational Physics* 165 (2000) 89–125.
- [18] G. Jiang, C.-W. Shu, Efficient implementation of weighted ENO schemes, *Journal of Computational Physics* 126 (1996) 202–228.
- [19] S. Jin, A steady state capturing method for hyperbolic systems with geometrical source terms, *Mathematical Modelling and Numerical Analysis* 35 (2001) 631–646.
- [20] A. Kurganov, D. Levy, Central-upwind schemes for the Saint–Venant system, *Mathematical Modelling and Numerical Analysis* 36 (2002) 397–425.
- [21] R.J. LeVeque, Balancing source terms and flux gradients on high-resolution Godunov methods: the quasi-steady wave-propagation algorithm, *Journal of Computational Physics* 146 (1998) 346–365.
- [22] X.-D. Liu, S. Osher, T. Chan, Weighted essentially nonoscillatory schemes, *Journal of Computational Physics* 115 (1994) 200–212.
- [23] J. Qiu, C.-W. Shu, Runge–Kutta discontinuous Galerkin method using WENO limiters, *SIAM Journal on Scientific Computing* 26 (2005) 907–929.
- [24] T.C. Rebollo, A.D. Delgado, E.D.F. Nieto, A family of stable numerical solvers for the shallow water equations with source terms, *Computer Methods in Applied Mechanics and Engineering* 192 (2003) 203–225.
- [25] P.L. Roe, Approximate Riemann solvers, parameter vectors, and difference schemes, *Journal of Computational Physics* 43 (1981) 357–372.
- [26] G. Russo, Central schemes for balance laws, in: *Proceedings of the VIII International Conference on Nonlinear Hyperbolic Problems*, Magdeburg, 2000.
- [27] J. Shi, C. Hu, C.-W. Shu, A technique of treating negative weights in WENO schemes, *Journal of Computational Physics* 175 (2002) 108–127.
- [28] C.-W. Shu, TVB uniformly high-order schemes for conservation laws, *Mathematics of Computation* 49 (1987) 105–121.
- [29] C.-W. Shu, Essentially non-oscillatory and weighted essentially non-oscillatory schemes for hyperbolic conservation laws, in: B. Cockburn, C. Johnson, C.-W. Shu, E. Tadmor (Eds.), *Advanced Numerical Approximation of Nonlinear Hyperbolic Equations*, in: A. Quarteroni (Ed.), *Lecture Notes in Mathematics*, vol. 1697, Springer, Berlin, 1998, pp. 325–432.
- [30] C.-W. Shu, S. Osher, Efficient implementation of essentially non-oscillatory shock-capturing schemes, *Journal of Computational Physics* 77 (1988) 439–471.
- [31] C.-W. Shu, S. Osher, Efficient implementation of essentially non-oscillatory shock-capturing schemes, II, *Journal of Computational Physics* 83 (1989) 32–78.
- [32] M.E. Vazquez-Cendon, Improved treatment of source terms in upwind schemes for the shallow water equations in channels with irregular geometry, *Journal of Computational Physics* 148 (1999) 497–526.
- [33] S. Vukovic, L. Sopta, ENO and WENO schemes with the exact conservation property for one-dimensional shallow water equations, *Journal of Computational Physics* 179 (2002) 593–621.
- [34] S. Vukovic, N. Crnjacic-Zic, L. Sopta, WENO schemes for balance laws with spatially varying flux, *Journal of Computational Physics* 199 (2004) 87–109.
- [35] Y. Xing, C.-W. Shu, High order finite difference WENO schemes with the exact conservation property for the shallow water equations, *Journal of Computational Physics* 208 (2005) 206–227.
- [36] Y. Xing, C.-W. Shu, High order well-balanced finite difference WENO schemes for a class of hyperbolic systems with source terms, *Journal of Scientific Computing* (accepted).
- [37] K. Xu, A well-balanced gas-kinetic scheme for the shallow-water equations with source terms, *Journal of Computational Physics* 178 (2002) 533–562.
- [38] J.G. Zhou, D.M. Causon, C.G. Mingham, D.M. Ingram, The surface gradient method for the treatment of source terms in the shallow-water equations, *Journal of Computational Physics* 168 (2001) 1–25.