# ON THE ACTIVE FLUX SCHEME FOR HYPERBOLIC PDES WITH SOURCE TERMS[*]

WASILIJ BARSUKOW[†], JONAS P. BERBERICH[‡], AND CHRISTIAN KLINGENBERG[‡]

**Abstract.** The active flux scheme is a finite volume scheme with additional point values distributed along the cell boundary. It is third order accurate and does not require a Riemann solver: the initial value problem at the particular points is solved instead. The intercell flux is then obtained from the evolved values along the cell boundary by quadrature. This paper focuses on the conceptual extension of active flux to include source terms, and thus for simplicity assumes the homogeneous part of the equations linear. To a large part the treatment of the source terms is independent of the choice of the homogeneous part of the system. Additionally, only systems are considered which admit characteristics (instead of characteristic cones). This is the case for scalar equations in any number of spatial dimensions and systems in one spatial dimension. Here, we succeed to extend the active flux method to include (possibly nonlinear) source terms while maintaining third order accuracy of the method. This requires a novel (approximate) operator for the evolution of point values and a modified update procedure of the cell average. For linear acoustics with gravity, it is shown how to achieve a well-balanced / stationarity preserving numerical method.

**Key words.** finite volume methods, active flux, source terms, balance laws, well-balanced methods, gravity

**AMS subject classifications.** 35L65, 35L45, 65M08

**1. Introduction.** Numerous phenomena of the physical world are modeled by hyperbolic balance laws (conservation laws augmented by source terms). This includes gas dynamics, the motion of water waves, plasma physics and even general relativity. Often physical modeling requires to include source terms, and conservation is modified due to creation or annihilation of some of the evolved quantities. Chemical reactions, for example, change the number density of a species and produce or absorb heat (i.e. internal energy). Gravity accelerates matter downwards and creates momentum. In the shallow water model describing the motion of a free water surface the bottom topography enters the equations through a source term. Rewriting the hydrodynamic equations in a different coordinate system (e.g. in polar coordinates) makes geometric source terms appear. All these applications require reliable numerical methods which are able to deal with source terms.

Reliable numerical methods for hyperbolic conservation laws with source terms first need to perform well in the homogeneous case. This means for example that they need to cope with discontinuities / weak solutions and with phenomena arising in multiple spatial dimensions, such as involutions and non-trivial stationary states. This requirement has led [ER13, FR15] to suggest *active flux*, an extension of the finite volume method. Additionally to the cell average, this scheme evolves point values located at the cell boundary. The update of the point values is achieved by using an evolution operator that includes multi-dimensional information. The presence of the point values along the cell boundary then allows to compute the intercell flux

[†]Institute of Mathematics, Zurich University, 8057 Zurich, Switzerland (wasilij.barsukow@math.uzh.ch).
[‡]Wuerzburg University, Emil-Fischer-Strasse 40, 97074 Wuerzburg, Germany.

1

via quadrature. It has been shown in [BHKR19] that this scheme is stationarity preserving and vorticity preserving for linear acoustics without any fix. It is third order accurate. Extensions to nonlinear systems have been recently suggested e.g. in [Fan17, HKS19, Bar19a]. Active flux therefore seems to be promising for resolving many of the structure preservation problems that currently available methods are facing (an overview of existing methods for balance laws is given below).

In view of the many applications that involve source terms, this paper therefore aims at deriving the necessary modifications for active flux to be applicable to balance laws while retaining its third order accuracy. Including the source term requires a number of modifications. The homogeneous part of the equations therefore is for simplicity assumed to be a linear hyperbolic system for which characteristics are available. This is the case for scalar equations in any number of spatial dimensions and for systems in one spatial dimension. For multi-dimensional systems, the concept of characteristics needs to be replaced by characteristics cones. In the homogeneous case, active flux has been used for this situation as well ([ER13, BHKR19]), but an extension to inhomogeneous systems in multi-d, and to nonlinear systems remains subject of future work. To a large part, the strategies presented in this paper will, however, remain valid when the homogeneous part of the equations is nonlinear as well, and even for nonlinear multi-dimensional systems.

As soon as a source term is added to a hyperbolic system, new stationary states arise which often are of particular interest. The stationarity is due to the flux divergence being equal to the source term. Many areas of application of balance laws involve studies of dynamics on top of such an equilibrium (e.g. astrophysics, meteorology, tsunami modeling, ... ). This requires the numerical method to be very accurate on the stationary states in order to avoid spurious, artificial perturbations. Therefore the error of a numerical solution representing one of those stationary states should not increase with time, thus allowing the simulation to run for a long time.

Numerical methods which achieve this are called *well-balanced*, introduced in [GL96]. They make sure that the discretization of the flux divergence and the discretization of the source term match, and that the numerical method keeps the desired stationary state exactly stationary for any resolution of the grid. The concept of well-balanced methods has been extensively used in the context of shallow water equations with non-flat bottom topography (e.g. [ABB+04, BV94, LeV98] and references therein). Here, the balance is the so-called lake-at-rest solution, which amounts to an algebraic condition and can thus be given explicitly.

Another area in which well-balanced methods have high relevance is the simulation of hydrodynamic processes using compressible Euler equations with gravitational source term. The so-called hydrostatic state (stationary state with no velocity) is described by one PDE for two unknown functions. There are many hydrostatic states, depending on the additional thermodynamical relation that one chooses in order to close this PDE. The fact that the stationary state is itself given by a differential equation that cannot be integrated makes well-balancing much more delicate in this context. There are two different ways which are currently used to construct well-balanced methods for the Euler equations with gravity. The first and more traditional way is to restrict the class of hydrostatic solutions which are balanced exactly or to choose a particular, but arbitrary hydrostatic state (e.g. [CL94, LGB11, DZBK16, CK15, BCK16, CCK+18, BCKR19, BCK19]). This is advantageous in all those applications where the stationary state is known, and the evolution of perturbations around it shall be studied. If no information on the stationary state can be assumed, then the only way to proceed is to make sure that the stationary states of the

numerical method are fulfilling some *discretization* of the corresponding PDE (e.g. [DZBK14, KM16, BKCK20]).

For linear numerical methods a theory of such *stationarity preserving* methods was given in [Bar19b], with a particular emphasis on this latter, more complicated, situation of the stationary states given by PDEs, and not by algebraic relations. It turns out that many standard numerical methods add diffusion even to those states that should remain stationary. The set of states that are actually kept stationary by such methods is very small (e.g. uniform constants). Stationarity preserving methods do not apply diffusion to certain discrete data. These data are described by a discrete version of the PDE governing the stationary states. Stationarity preserving methods thus keep stationary a much larger set of initial data. Independently of how these discrete equations actually look like, it is their existence that makes a qualitative difference. In a non-stationarity-preserving method, initial data sampled from an analytic stationary state will decay due to the diffusion and become unrecognizable in the end. In a stationarity preserving method, these initial data will evolve towards one of the many discrete stationary states approximating the steady PDE, and will remain there forever (up to machine precision). The long-time numerical solution will then indeed approximate the analytic stationary state. For more details, see [Bar19b]. In this paper we understand the concept of well-balancing in this sense of stationarity preservation.

In this paper, after extending the active flux scheme to include source terms, we construct a well-balanced active flux method for the equations of acoustics with gravity. The hydrostatic solutions of acoustics with gravity are comparable to those of the compressible Euler equations with gravity, since they are given via the same under-determined differential equation. We show that the active flux scheme endowed with an exact evolution operator is intrinsically well-balanced in this way. In practice, an approximate evolution operator needs to be used. Hence we introduce a modification of the approximate evolution operator which makes the scheme well-balanced even upon usage of an approximate evolution operator.

The paper is organized as follows: After the active flux scheme for homogeneous problems is introduced in section 2, the modifications necessary for including source terms are discussed. Section 3 discusses the evolution operators necessary for the update of the point values. Section 4 is devoted to the modifications in the update of the average. Here, the focus lies on linear systems of equations with possibly nonlinear source terms in one spatial dimension and on linear advection in multiple spatial dimensions. Section 5 discusses well-balancing of active flux for linear acoustics with gravity. Section 6 finally demonstrates numerically that the new method attains third order accuracy with linear and nonlinear source terms, can be used to compute Riemann problems, and displays well-balanced behavior for stationary states.

This work can be seen in the larger context of the quest for structure preserving numerical methods, of which well-balanced methods form an example. Extending these results to nonlinear hyperbolic equations with source terms and thus combining the structure preserving properties of active flux remains subject of future work. However, the procedures suggested in this paper are formulated with as little reference to the linearity of the equations as possible.
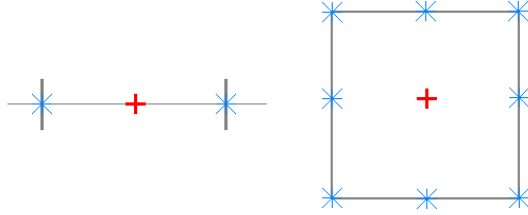
FIG. 1. *The degrees of freedom used for active flux. Stars indicate the location of point values, and the cross (placed in the center symbolically) refers to the cell average.* Left*: One spatial dimension.* Right*: Two spatial dimensions.*

**2. The active flux scheme.** Consider the initial value problem for an $m \times m$ system of hyperbolic balance laws in $d$ spatial dimensions[1]

(2.1)               $\partial_t q + \nabla \cdot \mathbf{f}(q) = s(q)$              $q : \mathbb{R}_0^+ \times \mathbb{R}^d \to \mathbb{R}^m, \ f, s : \mathbb{R}^m \to \mathbb{R}^m$

(2.2)                  $q(0, \mathbf{x}) := q_0(\mathbf{x})$

This section reviews the general idea of the active flux scheme. Some of the details then depend on the particular equation that is to be solved. After the general concept is outlined, the details that make it applicable to hyperbolic balance laws are discussed in sections 3 and 4.

**2.1. Degrees of freedom in the active flux scheme.** The active flux scheme ([ER13, BHKR19], first introduced in [VL77]) is an extension of the finite volume scheme. The active flux scheme evolves both the cell average and point values which are distributed along the cell boundary. In particular, here the following two choices are considered (see Figure 1):

- In one spatial dimension, there is a point value $q_{i+\frac{1}{2}}$ located at each cell interface $x_{i+\frac{1}{2}}$. Thus every cell has access to one cell average $\bar{q}_i$ and two point values at its interfaces.
- On Cartesian grids in two spatial dimensions, there is a point value $q_{i+\frac{1}{2},j}$, $q_{i,j+\frac{1}{2}}$ at each edge midpoint and one at each node $q_{i+\frac{1}{2},j+\frac{1}{2}}$. Every cell has access to one cell average $\bar{q}_{ij}$ and 8 point values distributed along the cell interface.

Note that the point values at cell interfaces are shared by the adjacent cells. Thus, in one spatial dimension, on average there are 2 degrees of freedom per cell: 1 cell average and 2 interface values shared each by 2 cells. In two spatial dimensions in the setup as described above there are 4 degrees of freedom per cell: 1 cell average, 4 edge values, each shared by two cells and 4 node values each shared by 4 cells.

Note also that active flux does not use a staggered grid. The degrees of freedom at the cell boundaries are not averages over staggered volumes, but point values. This also explains why there is no notion of a conservative update for these, because this concept only applies to averages. The update of the cell average in the active flux method is, of course, conservative (see below).

**2.2. Update of the cell average.** As the active flux scheme is an extension of the finite volume scheme, given a numerical flux, the update of the average happens in

---

[1]In this paper, indices never denote derivatives. Boldface symbols denote vectors that have the same dimension as the space.

the same way as for finite volume schemes. In this section, this finite volume aspect of active flux is described in an arbitrary number of spatial dimensions. The numerical flux, however, is obtained very differently in the active flux scheme ([ER13, FR15]). This is then described in detail in section 4.

Consider the computational domain to be subdivided into polygonal computational cells. Upon integration of (2.1) over one time step $[t^n, t^n + \Delta t]$ and over one computational cell $\mathcal{C}$ one obtains an evolution equation for the cell average $\bar{q}_{\mathcal{C}} := \frac{1}{|\mathcal{C}|} \int_{\mathcal{C}} \mathrm{d}\mathbf{x}\, q(t, \mathbf{x})$:

$$\frac{\bar{q}_{\mathcal{C}}^{n+1} - \bar{q}_{\mathcal{C}}^{n}}{\Delta t} + \frac{1}{|\mathcal{C}|}\frac{1}{\Delta t} \int_{t^n}^{t^n + \Delta t} \mathrm{d}t \int_{\partial\mathcal{C}} \mathrm{d}\sigma\, \mathbf{n} \cdot \mathbf{f}(q(t, \mathbf{x})) =$$

$$\frac{1}{\Delta t} \int_{t^n}^{t^n + \Delta t} \mathrm{d}t \frac{1}{|\mathcal{C}|} \int_{\mathcal{C}} \mathrm{d}\mathbf{x}\, s(q(t, \mathbf{x}))$$

Here, as usual, the index of the time step is denoted as a superscript and $q_{\mathcal{C}}^{n}$ denotes the average in cell $\mathcal{C}$ at time $t^n$. The boundary $\partial\mathcal{C}$ consists of edges $e$, such that one can rewrite

$$\frac{\bar{q}_{\mathcal{C}}^{n+1} - \bar{q}_{\mathcal{C}}^{n}}{\Delta t} + \frac{1}{|\mathcal{C}|}\frac{1}{\Delta t} \int_{t^n}^{t^n + \Delta t} \mathrm{d}t \sum_{e \subset \partial\mathcal{C}} \int_{e} \mathrm{d}\sigma\, \mathbf{n}_e \cdot \mathbf{f}(q(t, \mathbf{x})) =$$

$$\frac{1}{\Delta t} \int_{t^n}^{t^n + \Delta t} \mathrm{d}t \frac{1}{|\mathcal{C}|} \int_{\mathcal{C}} \mathrm{d}\mathbf{x}\, s(q(t, \mathbf{x}))$$

The vector $\mathbf{n}_e$ is the outward unit normal of edge $e$. This expression, so far exact, becomes a finite volume scheme upon replacing the exact normal flux and source averages by suitable approximations $\hat{f}_e$ and $\hat{s}_{\mathcal{C}}$:

(2.3)
$$\frac{\bar{q}_{\mathcal{C}}^{n+1} - \bar{q}_{\mathcal{C}}^{n}}{\Delta t} + \frac{1}{|\mathcal{C}|} \sum_{e \subset \partial\mathcal{C}} |e| \hat{f}_e = \hat{s}_{\mathcal{C}}$$

with

(2.4)
$$\hat{f}_e \simeq \frac{1}{\Delta t} \int_{t^n}^{t^n + \Delta t} \mathrm{d}t \frac{1}{|e|} \int_{e} \mathrm{d}\sigma\, \mathbf{n}_e \cdot \mathbf{f}(q(t, \mathbf{x}))$$

(2.5)
$$\hat{s}_{\mathcal{C}} \simeq \frac{1}{\Delta t} \int_{t^n}^{t^n + \Delta t} \mathrm{d}t \frac{1}{|\mathcal{C}|} \int_{\mathcal{C}} \mathrm{d}\mathbf{x}\, s(q(t, \mathbf{x}))$$

Usual finite volume schemes introduce a (piecewise continuous) reconstruction of the averages, and obtain the numerical flux by an exact or approximate short-time evolution of this reconstruction. For example, introducing a piecewise constant function whose averages match the given cell averages, and solving the Riemann problems at the cell interfaces allows to compute a numerical flux.

The active flux scheme does not need this. Indeed, the point values along the boundary can be used to immediately approximate (2.4)–(2.5) by quadrature. The

202  desired properties (most importantly the desired order of accuracy) of the resulting
203  scheme dictate the number of point values along each edge and also the points in time
204  at which these point values need to be available.
205      The source term also contributes to the update of the cell average. The quadrature
206  necessary to approximate the source term average (2.5) to sufficient order in space
207  and time is suggested in this paper for the first time and discussed in section 4.

**2.3. Update of the point values.** The cell average update, and in particular
209  the computation of the intercell fluxes, requires accurate point values at the cell
210  boundary to be available.
211      First consider the case where the source term vanishes: $s = 0$. For third order of
212  accuracy, the integrals in (2.4) need to be approximated by Simpson's rule. For the
213  integration in space this can easily be achieved using the available point values at each
214  cell interface as described in section 2.1. For the integration in time all point values
215  need to be available at $t^n, t^n + \frac{\Delta t}{2}$ and $t^n + \Delta t$. Altogether this yields a space-time
216  Simpson rule.
217      In order to obtain sufficiently accurate time evolved point values, in [VL77] it has
218  been suggested to reconstruct the data and to use an exact evolution operator. An
219  exact evolution operator generally is unavailable for nonlinear problems, and there-
220  fore in [Fan17, HKS19, Bar19a] approximate evolution operators have been proposed.
221  Even for linear systems of hyperbolic balance laws it is generally very difficult to ob-
222  tain closed-form exact evolution operators, as is shown in section 3.2. Therefore the
223  point values in the active flux scheme shall be evolved using a sufficiently high order
224  *approximate* evolution operator applied to a reconstruction of the discrete data. An
225  exact evolution operator provides the necessary upwinding in order to guarantee sta-
226  bility, and an approximate evolution operator needs to do the same. The approximate
227  evolution operator is introduced in section 3.3.

**2.4. Reconstruction.** The reconstruction shall interpolate the point values and
229  its average over the computational cell shall match the given cell average. In the
230  following, to simplify notation, in one spatial dimension a uniform grid is assumed,
231  although the reconstruction can immediately be generalized to nonuniform grids. In
232  two spatial dimensions, a Cartesian grid is used. As mentioned in section 2.1, in
233  one spatial dimension every cell has access to 3 degrees of freedom which makes a
234  parabolic reconstruction natural. With the above-mentioned setup it is unique and
235  reads ([VL77, FR15])

$$(2.6) \quad q_{\text{recon},i}(x) = -3(2\bar{q}_i - q_{i-\frac{1}{2}} - q_{i+\frac{1}{2}})\frac{(x - x_i)^2}{\Delta x^2}$$

$$(2.7) \qquad\qquad + (q_{i+\frac{1}{2}} - q_{i-\frac{1}{2}})\frac{x - x_i}{\Delta x} + \frac{6\bar{q}_i - q_{i-\frac{1}{2}} - q_{i+\frac{1}{2}}}{4} \qquad x \in [x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}]$$

In two spatial dimensions as described above, every cell has access to 9 degrees of freedom, and there is a unique biparabolic reconstruction, which reads

$$
\begin{aligned}
q_{\text{recon},ij}(\xi\Delta x, \eta\Delta y) := {} & \frac{9}{4}\bar{q}_{ij}\left(-1+4\xi^2\right)\left(-1+4\eta^2\right) \\
& -\frac{1}{4}q_{\text{W}}\left(-1-4\xi+12\xi^2\right)\left(-1+4\eta^2\right) \\
& -\frac{1}{4}q_{\text{E}}\left(-1+4\xi+12\xi^2\right)\left(-1+4\eta^2\right) \\
& -\frac{1}{4}q_{\text{S}}\left(-1+4\xi^2\right)\left(-1-4\eta+12\eta^2\right) \\
& -\frac{1}{4}q_{\text{N}}\left(-1+4\xi^2\right)\left(-1+4\eta+12\eta^2\right) \\
& +\frac{1}{16}q_{\text{SW}}(-1+2\xi)(-1+2\eta)(-1-2\eta+2\xi(-1+6\eta)) \\
& +\frac{1}{16}q_{\text{SE}}(1+2\xi)(-1+2\eta)(1+2\eta+2\xi(-1+6\eta)) \\
& +\frac{1}{16}q_{\text{NW}}(-1+2\xi)(1+2\eta)(1-2\eta+2\xi(1+6\eta)) \\
& +\frac{1}{16}q_{\text{NE}}(1+2\xi)(1+2\eta)(-1+2\eta+2\xi(1+6\eta))
\end{aligned}
\tag{2.8}
$$

with $\xi := x/\Delta x$, $\eta := y/\Delta y$ and

$$
q_{\text{NE}} = q_{i+\frac{1}{2},j+\frac{1}{2}} \qquad q_{\text{NW}} = q_{i-\frac{1}{2},j+\frac{1}{2}} \qquad q_{\text{SW}} = q_{i-\frac{1}{2},j-\frac{1}{2}} \qquad q_{\text{SE}} = q_{i+\frac{1}{2},j-\frac{1}{2}}
\tag{2.9}
$$

$$
q_{\text{N}} = q_{i,j+\frac{1}{2}} \qquad q_{\text{S}} = q_{i,j-\frac{1}{2}} \qquad q_{\text{E}} = q_{i+\frac{1}{2},j} \qquad q_{\text{W}} = q_{i-\frac{1}{2},j}
\tag{2.10}
$$

Note that both reconstructions are globally continuous, but generally not continuously differentiable at the cell interfaces.

**2.5. Overview of the algorithm.** The overall algorithm of active flux is as follows:

1. Given cell averages and point values, compute a reconstruction according to section 2.4.
2. Use the reconstruction as initial data in the update of the point values. The choices of evolution operators considered so far are discussed in section 2.3 and evolution operators in presence of source terms are suggested in section 3.3 below.
3. Given the updated point values along the cell interfaces, compute the inter-cell fluxes via quadrature (sections 2.2 and 4 for the homogeneous and the inhomogeneous cases, respectively).
4. Update the cell averages via (2.3).

A CFL-type condition arises in the update of the point values: the domain of dependence of the evolution operator needs to be contained in the neighbouring cells. Denoting by $\lambda_{\max}$ the maximum speed of propagation, the time step needs to be chosen as

$$
\Delta t \leq \frac{L_{\min}}{\lambda_{\max}}
\tag{2.11}
$$

where $L_{\min} = \Delta x$ in one spatial dimension, and $L_{\min} = \frac{1}{2}\min(\Delta x, \Delta y)$ in two spatial dimensions, when the point values are distributed as described in section 2.1.

**3. Evolution of the point values in presence of a source term.** The evolution of the point values needs to account for the source term. Additionally, in this paper a special focus shall lie on structure preservation properties of the resulting scheme. In the homogeneous case such properties have been observed upon usage of an exact evolution operator ([BHKR19]). In presence of a source term, one needs to use an approximate evolution operator (section 3.3), but should nevertheless aim at making it such that it does not spoil structure preservation (see section 5).

For certain equations, the inhomogeneous problem admits an exact solution (sections 3.1–3.2). This is valuable in order to assess specific properties of the numerical method later.

**3.1. Linear advection with a source term in multiple spatial dimensions.** Consider a scalar equation ($m = 1$) and $\mathbf{f}(q) = \mathbf{U}q$ with $\mathbf{U} \in \mathbb{R}^d$. Then

$$\partial_t q + \mathbf{U} \cdot \nabla q = s(q) \tag{3.1}$$

amounts to the ODE

$$\frac{\mathrm{d}}{\mathrm{d}t} q = s(q) \tag{3.2}$$

along the straight characteristic of velocity $\mathbf{U}$. This ODE can be easily solved analytically:

$$\int_{q_0(\mathbf{x}-\mathbf{U}t)}^{q(t,\mathbf{x})} \frac{\mathrm{d}p}{s(p)} = t \tag{3.3}$$

E.g. for $s(q) = \kappa q$ this yields $\ln \frac{q(t,\mathbf{x})}{q_0(\mathbf{x}-\mathbf{U}t)} = \kappa t$, or

$$q(t,\mathbf{x}) = q_0(\mathbf{x} - \mathbf{U}t) \exp(\kappa t) \tag{3.4}$$

and for $s(q) = \kappa q^B$, $B \neq 1$

$$q(t,\mathbf{x}) = \left( (q_0(\mathbf{x} - \mathbf{U}t))^{1-B} + (1 - B)\kappa t \right)^{\frac{1}{1-B}} \tag{3.5}$$

**3.2. Linear acoustics with gravity in one spatial dimension.** This section has threefold purpose. First, it introduces the acoustic equations with a gravity source term, which form a very useful system for the study of structure preservation of numerical methods. This is the set of equations for which a well-balanced method is derived in 5. This section also demonstrates the difficulties of finding an exact solution to an inhomogeneous system even if it is linear. Finally, the exact solution derived here is used later in order to assess the accuracy of the numerical method.

The equations of linear acoustics in one spatial dimension endowed with a gravity source term read:

$$\partial_t \rho + \partial_x v = 0 \tag{3.6}$$

$$\partial_t v + \partial_x p = \rho g \qquad g \in \mathbb{R} \tag{3.7}$$

$$\partial_t p + c^2 \partial_x v = 0 \tag{3.8}$$

The corresponding homogeneous problem (linear acoustics) is the linearization of the Euler equations around the background state of constant density $\rho_{\mathrm{bg}} = 1$, constant

pressure $p_{\mathrm{bg}}$ and vanishing velocity. Then the speed of sound $c = \sqrt{\frac{\gamma p_{\mathrm{bg}}}{\rho_{\mathrm{bg}}}}$ is a constant ($\mathbb{R} \ni \gamma > 1$). The full system (3.6)–(3.8) can be understood as a particular kind of a linearization of the Euler equations with gravity[2]

$$\text{(3.9)} \qquad \partial_t \rho + \partial_x(\rho v) = 0$$

$$\text{(3.10)} \qquad \partial_t(\rho v) + \partial_x(\rho v^2 + p) = \rho g$$

$$\text{(3.11)} \qquad \partial_t e + \partial_x(v(e + p)) = 0$$

$$\text{(3.12)} \qquad e = \frac{p}{\gamma - 1} + \frac{1}{2}\rho v^2 - \rho g x$$

The static (stationary and $v = 0$) states of (3.9)–(3.11) are governed by $\partial_x p = \rho g$. This equation can only be solved if e.g. $\rho$ is given as a function of $x$, or if another relation is provided between any two of the variables $p, \rho, e$. This multitude of possible stationary states is reflected in the linearization (3.6)–(3.8). (This is the reason for this particular choice of a linearization.) Observe that stationary states of (3.6)–(3.8) also are governed by $\partial_x p = \rho g$ and that $p$ can only be computed if $\rho$ is given as a function of $x$, or if an additional relation is provided that links $\rho$ and $p$. This is an example of a so-called non-trivial stationary state as introduced in [Bar19b]. Examples of stationarity preserving schemes for (3.6)–(3.8) have been discussed in [Bar18].

The exact solution of (3.6)–(3.8) is studied in the Appendix A. This solution is not part of the suggested method but only serves auxiliary purposes, such as accuracy checks. However it illustrates the difficulties encountered when solving linear systems with sources. To the authors' knowledge the exact solution to (3.6)–(3.8) is not available in the literature so far.

**3.3. Runge-Kutta method for linear systems with a source.** Consider an $m \times m$ linear system in characteristic variables:

$$\text{(3.13)} \qquad (\partial_t + \lambda_\ell \partial_x)Q_\ell = S_\ell(Q_1, \ldots, Q_m) \qquad \ell = 1, \ldots, m$$

From now on, the capital letter $Q$ denotes the characteristic variables of this particular system, whereas $q$ continues to denote a generic variable.

Recall the following theorem from [Bar19a]:

THEOREM 3.1. *Assume a hyperbolic CFL condition $\Delta x / \Delta t \to$ const as $\Delta t \to 0$. If the approximate evolution $Q^{\mathrm{approx}}(t, x)$ approximates the exact solution $Q(t, x)$ for fixed $x$ at least as*

$$\text{(3.14)} \qquad Q^{\mathrm{approx}}(t, x) = Q(t, x) + \mathcal{O}(t^3)$$

*and the quadrature rules used to approximate (2.4)–(2.5) yield the exact value up to an error of $\mathcal{O}(\Delta t^\alpha \Delta x^\beta)$, $\alpha + \beta \geq 3$ then active flux formally achieves third order accuracy.*

Note that the simple approach of evolving each component of the source term along its associated characteristic

(3.15)
$$Q_\ell(t, x) \simeq Q_{\ell,0}(x - \lambda_\ell t) + t S_\ell(Q_{1,0}(x - \lambda_\ell t), \ldots, Q_{m,0}(x - \lambda_\ell t)) \qquad \ell = 1, \ldots, m$$

---

[2]Note that often the energy equation is written with a source term $\rho g v$ appearing. This source term is unnecessary, as it can be removed by redefining the notion of total energy. When the total energy includes the potential energy $-\rho g x$ due to gravity, the conservation form of the energy equation is restored. The source term in the momentum equation remains.
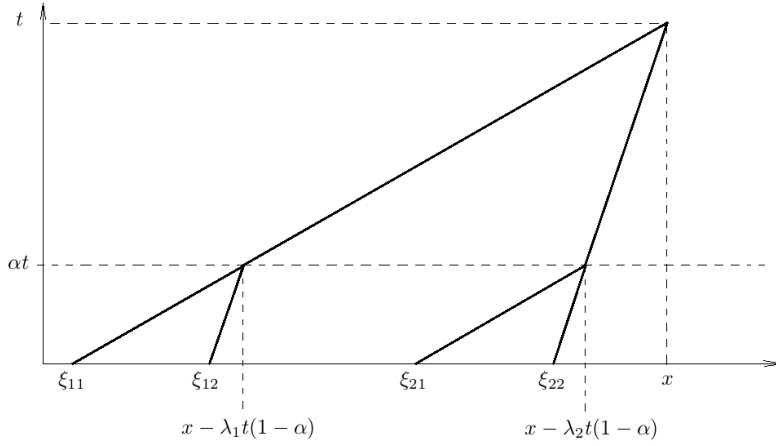
FIG. 2. *Illustration of the intermediate solutions and the involved characteristics for the first step in the Runge-Kutta scheme.*

353    fails to be accurate enough (the error is $\mathcal{O}(t^2)$ instead of $\mathcal{O}(t^3)$).

354        Recall the second order Runge-Kutta method for the ordinary differential equation

$$(3.16) \qquad \dot{q}(t) = s(t, q(t)) \qquad\qquad q : \mathbb{R}_0^+ \to \mathbb{R}$$

$$(3.17) \qquad q^{(1)}(\alpha t) = q(0) + \alpha t s(0, q(0))$$

$$(3.18) \qquad q(t) = q(0) + t\left(1 - \frac{1}{2\alpha}\right) s(0, q(0)) + t\frac{1}{2\alpha} s(\alpha t, q^{(1)}(\alpha t)) + \mathcal{O}(t^3)$$

361    for any $\alpha \in (0,1)$. In particular choosing $\alpha = \frac{1}{2}$ (midpoint method) involves a
362    predictor value at half time step. This can be taken as inspiration for constructing a
363    sufficiently accurate approximate evolution operator:

364        THEOREM 3.2 (RK2 evolution operator). *Choose (see Figure 2)*

$$(3.19) \qquad \xi_{\ell k} := x - \lambda_\ell t(1 - \alpha) - \lambda_k \alpha t$$

$$(3.20) \qquad Q^*_{k\ell} := Q_{k,0}(\xi_{\ell k}) + \alpha t S_k(Q_{1,0}(\xi_{\ell k}), \dots, Q_{m,0}(\xi_{\ell k})) \qquad k, \ell = 1, \dots, m$$

368    *and*

$$(3.21)$$

$$Q_\ell^{(1)}(t,x) := Q_{\ell,0}(x - \lambda_\ell t) + \left(1 - \frac{1}{2\alpha}\right) S_\ell(Q_{1,0}(x - \lambda_\ell t), \dots, Q_{m,0}(x - \lambda_\ell t)) t$$

$$(3.22) \qquad\qquad\qquad + \frac{t}{2\alpha} S_\ell\left(Q^*_{1\ell}, \dots, Q^*_{m\ell}\right) \qquad \ell = 1, \dots, m$$

372        *Then, for all $\alpha \in (0,1)$*

$$(3.23) \qquad Q_\ell^{(1)}(t,x) = Q_\ell(t,x) + \mathcal{O}(t^3) \qquad \ell = 1, \dots, m$$

375    Note that $Q^*_{\ell j}$ approximates $Q_\ell(\alpha t, x - \lambda_j t(1 - \alpha))$.

376  *Proof.* By explicitly computing the first three terms of the Taylor series in $t$ one
377  confirms the statement. The exact solution is

378  (3.24)  $Q_\ell(t,x) = Q_{\ell,0}(x) + t\partial_t Q_\ell\big|_{t=0} + \dfrac{t^2}{2}\partial_t^2 Q_\ell\big|_{t=0} + \mathcal{O}(t^3)$

379  (3.25)  $= Q_{\ell,0}(x) + t(S_{\ell,0} - \lambda_\ell \partial_x Q_{\ell,0})$

380  (3.26)  $+ \dfrac{t^2}{2}\left(\sum_k \dfrac{\partial S_\ell}{\partial Q_k}\Big(S_{k,0} - (\lambda_k + \lambda_\ell)\partial_x Q_{k,0}\Big) + \lambda_\ell^2 \partial_x^2 Q_{\ell,0}\right) + \mathcal{O}(t^3)$

381

382  where $S_{\ell,0}$ denotes

383
384  (3.27)  $S_{\ell,0} := S_\ell(Q_{1,0}(x),\ldots,Q_{m,0}(x))$

385  and $\dfrac{\partial S_\ell}{\partial Q_k}$ also is evaluated at $x$. Note that it has been used that $\partial_x \lambda_\ell = 0$ (i.e. that the
386  homogeneous system is linear), but the source $S$ can be any differentiable function of
387  $Q$.

388  Expand now (3.22) ($\ell = 1,\ldots,m$):

389  (3.28)  $\partial_t Q_{k\ell}^*\big|_{t=0} = -(\lambda_\ell(1-\alpha) + \lambda_k \alpha)\partial_x Q_{k,0} + \alpha S_{k,0}$

390  (3.29)  $\partial_t Q_\ell^{(1)}(t,x) = -\lambda_\ell \partial_x Q_{\ell,0}(x - \lambda_\ell t)$

391  (3.30)  $+ \left(1 - \dfrac{1}{2\alpha}\right)\left(t\sum_k \dfrac{\partial S_\ell}{\partial Q_k}\partial_x Q_{k,0}(x - \lambda_\ell t)(-\lambda_\ell)\right.$

392  (3.31)  $\left. + S_\ell(Q_{1,0}(x - \lambda_\ell t),\ldots,Q_{m,0}(x - \lambda_\ell t))\right)$

393  (3.32)  $+ \dfrac{1}{2\alpha}\left(t\sum_k \dfrac{\partial S_\ell}{\partial Q_k}\partial_t Q_{k\ell}^* + S_\ell\Big(Q_{1\ell}^*,\ldots,Q_{m\ell}^*\Big)\right)$

394  (3.33)  $\overset{t=0}{=} -\lambda_\ell \partial_x Q_{\ell,0} + S_{\ell,0}$

395  (3.34)  $\partial_t^2 Q_\ell^{(1)}(t,x)\big|_{t=0} = \lambda_\ell^2 \partial_x^2 Q_{\ell,0} + \left(1 - \dfrac{1}{2\alpha}\right)\left(2\sum_k \dfrac{\partial S_\ell}{\partial Q_k}\partial_x Q_{k,0}(-\lambda_\ell)\right)$

396  (3.35)  $+ \dfrac{1}{2\alpha}\left(2\sum_k \dfrac{\partial S_\ell}{\partial Q_k}\partial_t Q_{k\ell}^*\big|_{t=0}\right)$

397  (3.36)  $= \lambda_\ell^2 \partial_x^2 Q_{\ell,0} - \sum_k \dfrac{\partial S_\ell}{\partial Q_k}\Big(\partial_x Q_{k,0}(\lambda_\ell + \lambda_k) - S_{k,0}\Big)$  $\square$

398

399  Obviously the two Taylor series agree up to terms $\mathcal{O}(t^3)$, which proves the statement.

400  COROLLARY 3.3 (Midpoint method). *If* $\alpha = \frac{1}{2}$, *then for* $\ell, k = 1,\ldots m$

401  (3.37)  $\xi_{\ell,j} := x - (\lambda_\ell + \lambda_j)\dfrac{t}{2}$

402  (3.38)  $Q_{k\ell}^* := Q_{k,0}(\xi_{\ell k}) + \dfrac{t}{2}S_k(Q_{1,0}(\xi_{k\ell}),\ldots,Q_{m,0}(\xi_{k\ell}))$

403
404  (3.39)  $Q_\ell^{(1)}(t,x) := Q_{\ell,0}(x - \lambda_\ell t) + tS_\ell\Big(Q_{1\ell}^*,\ldots,Q_{m\ell}^*\Big)$

COROLLARY 3.4 (RK2 evolution operator for a scalar equation).   *For a scalar equation*

(3.40)                           $$(\partial_t + \lambda\partial_x)Q = S(Q)$$

*the algorithm reads*

(3.41)                           $$\xi := x - \lambda t$$

*and*

(3.42)      $$Q^{(1)}(t,x) := Q_0(x - \lambda t) + \left(1 - \frac{1}{2\alpha}\right)S(Q_0(x - \lambda t))t$$

(3.43)                    $$+ \frac{t}{2\alpha}S\Big(Q_0(\xi) + \alpha t S(Q_0(\xi))\Big)$$

For the equations (3.6)–(3.8) of linear acoustics with gravity, $\lambda_1 = c = -\lambda_2, \lambda_3 = 0$. The characteristic variables are

(3.44)           $$Q_1 = \frac{p + cv}{2} \qquad Q_2 = \frac{p - cv}{2} \qquad Q_3 = -\frac{p}{c^2} + \rho$$

and the gravity source term then is

(3.45)       $$S_1 = -S_2 = \frac{g}{2c}(Q_1 + Q_2) + \frac{cg}{2}Q_3 \qquad\qquad S_3 = 0$$

**4. Update of the cell average in presence of a source term.** The update of the cell average needs to include the space-time average of the source term according to (2.3) of section 2.2. This space-time average needs to be approximated by a suitable quadrature / approximation with sufficient order of accuracy. Active flux has a strong focus on providing discrete degrees of freedom along the boundary which allow to perform a quadrature along the boundary. However, the evaluation of the source term for the update of the cell average involves an averaging over the cell volume. It is more difficult to achieve the desired order of accuracy here, as the setup lacks the quadrature points that would have been natural for this task. A quadrature formula adapted to the geometry of the active flux method is derived here.

Active flux for equations with a source term is considered in [NR16] for stationary problems, and for parabolic problems with slowly varying boundary conditions. In these cases there is no need to use high order quadrature in time. Therefore the method suggested there cannot be used here.

**4.1. One spatial dimension.** The numerical discretization (2.5)

(4.1)                    $$\hat{s}_\mathcal{C} \simeq \frac{1}{\Delta t}\int\limits_{t^n}^{t^n + \Delta t} dt \, \frac{1}{|\mathcal{C}|}\int_\mathcal{C} d\mathbf{x}\, s(q(t, \mathbf{x}))$$

of the source term in (2.3) requires a space-time quadrature that is exact for parabolic functions. The natural candidate would be Simpson's rule in both space and time (as used for the numerical flux), but there are not enough quadrature points for it. For example in one spatial dimension, the available information is

<div align="center">

| $t^{n+1}$ | $q_{i-\frac{1}{2}}^{n+1}$ | | $q_{i+\frac{1}{2}}^{n+1}$ |
|---|---|---|---|
| $t^{n+\frac{1}{2}}$ | $q_{i-\frac{1}{2}}^{n+\frac{1}{2}}$ | | $q_{i+\frac{1}{2}}^{n+\frac{1}{2}}$ |
| $t^n$ | $q_{i-\frac{1}{2}}^{n+1}$ | $\boxed{\bar{q}_i^n}$ | $q_{i+\frac{1}{2}}^n$ |
| | $x_{i-\frac{1}{2}}$ | | $x_{i+\frac{1}{2}}$ |

</div>

These are only 7 values (the box emphasizes that one of the values is a cell average, whereas the others are point values).

**4.1.1. Linear source term.** Consider first a linear source term, i.e. $s'' = 0$. Such source terms are relevant in practice (e.g. compressible Euler equations with gravity), and therefore it is worth dealing with them specifically as they allow for a simpler approach. For linear source it is possible to first find a quadrature for $q$ and to apply $s$ to the result. In order to find a quadrature formula for $q$, one needs to find a space-time polynomial $p(t,x)$ of at least second degree which interpolates the available 7 data. Integrating this polynomial would yield a quadrature formula for $q$. Here we suggest to use

$$(4.2) \qquad \mathscr{P}(t,x) = (a_0 + a_1 x + a_2 t + a_3 x^2 + a_4 xt + a_5 t^2) + a_6 xt^2$$

There is a unique set of coefficients $a_0, \ldots, a_6$ which makes polynomial (4.2) fulfill

$$(4.3) \qquad \mathscr{P}(t^{n+1}, x_{i-\frac{1}{2}}) = q_{i-\frac{1}{2}}^{n+1} \qquad\qquad \mathscr{P}(t^{n+1}, x_{i+\frac{1}{2}}) = q_{i+\frac{1}{2}}^{n+1}$$

$$(4.4) \qquad \mathscr{P}(t^{n+\frac{1}{2}}, x_{i-\frac{1}{2}}) = q_{i-\frac{1}{2}}^{n+\frac{1}{2}} \qquad\qquad \mathscr{P}(t^{n+\frac{1}{2}}, x_{i+\frac{1}{2}}) = q_{i+\frac{1}{2}}^{n+\frac{1}{2}}$$

$$(4.5) \qquad \mathscr{P}(t^n, x_{i-\frac{1}{2}}) = q_{i-\frac{1}{2}}^n \qquad \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \mathrm{d}x\, \mathscr{P}(t^n, x) = q_i^n \qquad \mathscr{P}(t^n, x_{i+\frac{1}{2}}) = q_{i+\frac{1}{2}}^n$$

Inserting this polynomial in (2.5) and integrating it instead of the source yields the following quadrature formula:

$$(4.6) \qquad \frac{1}{\Delta t} \int_0^{\Delta t} \mathrm{d}t\, \frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} \mathrm{d}x\, q(t^n + t, x_i + x) =$$

$$\bar{q}_i^n + \frac{1}{12}\left( -5(q_{i-\frac{1}{2}}^n + q_{i+\frac{1}{2}}^n) + q_{i-\frac{1}{2}}^{n+1} + q_{i+\frac{1}{2}}^{n+1} + 4(q_{i-\frac{1}{2}}^{n+\frac{1}{2}} + q_{i+\frac{1}{2}}^{n+\frac{1}{2}}) \right)$$

The weights can be depicted as

<div align="center">

| $t^{n+1}$ | $\frac{1}{12}$ | | $\frac{1}{12}$ |
|---|---|---|---|
| $t^{n+\frac{1}{2}}$ | $\frac{4}{12}$ | | $\frac{4}{12}$ |
| $t^n$ | $-\frac{5}{12}$ | $\boxed{1}$ | $-\frac{5}{12}$ |
| | $x_{i-\frac{1}{2}}$ | | $x_{i+\frac{1}{2}}$ |

</div>

Again, the box indicates that the corresponding weight refers to the cell average, whereas the others multiply point values.

The time levels $(n, n+\frac{1}{2}, n+1)$ contribute with weights $\left(\frac{1}{6}, \frac{2}{3}, \frac{1}{6}\right)$, such that this quadrature formula is a modification of Simpson's rule in time. Note that it is not

possible to use terms proportional to $x^3$, $x^2 t$ or $t^3$ instead of the term $xt^2$ in the polynomial ansatz, as then the system (4.3)–(4.5) does not admit a solution. In a sense this is therefore the only choice of a simple quadrature formula.

Quadrature formula (4.6) can be used immediately in order to approximate (2.5) for linear source terms.

**4.1.2. Nonlinear source term.** For nonlinear $s$, the average

$$(4.7) \qquad \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \mathrm{d}x\, s(q(t^n, x))$$

in general is different from

$$(4.8) \qquad s\left( \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \mathrm{d}x\, q(t^n, x) \right)$$

Point values, however, do not present any difficulties: one can just evaluate $s$ on them. Therefore we suggest to consider a reconstruction $q_{\mathrm{recon},i}(x)$ that interpolates $q_{i-\frac{1}{2}}^n$ and $q_{i+\frac{1}{2}}^n$ and whose average agrees with $\bar{q}_i^n$. It is computed anyway in order to update the point values in time, see equation (2.7). This reconstruction can be easily evaluated at the midpoint of the cell. Then, instead of the cell averages, one works with a seventh point value $q_{\mathrm{recon},i}(0) = \frac{1}{4}(6\bar{q}_i^n - q_{i-\frac{1}{2}}^n - q_{i+\frac{1}{2}}^n)$. Of course, this is equivalent to replacing the average by a Simpson's rule in the quadrature, and thus the order of the quadrature is not reduced. Therefore when using only point values (the 6 pointwise degrees of freedom and one value at the cell midpoint) the weights of the quadrature formula read

| | | | |
|---|---|---|---|
| $t^{n+1}$ | $\frac{1}{12}$ | | $\frac{1}{12}$ |
| $t^{n+\frac{1}{2}}$ | $\frac{4}{12}$ | | $\frac{4}{12}$ |
| $t^n$ | $-\frac{3}{12}$ | $\frac{8}{12}$ | $-\frac{3}{12}$ |
| | $x_{i-\frac{1}{2}}$ | | $x_{i+\frac{1}{2}}$ |

Equation (2.5) then is replaced by the quadrature

$$(4.9) \quad \hat{s}_i = \frac{s(q_{i-\frac{1}{2}}^{n+1}) + s(q_{i+\frac{1}{2}}^{n+1}) + 4\left(s(q_{i-\frac{1}{2}}^{n+1}) + s(q_{i+\frac{1}{2}}^{n+1})\right) - 3\left(s(q_{i-\frac{1}{2}}^{n+1}) + s(q_{i+\frac{1}{2}}^{n+1})\right) + 8q_{\mathrm{recon},i}(0)}{12}$$

This quadrature can now be used for nonlinear $s$. As (4.9) uses a Simpson quadrature instead of the average, upon usage of a linear source $s$, it reduces to the expression (4.6) because of the quadratic reconstruction.

If the source term vanishes, the scheme becomes conservative in the sense that averages are updated using numerical fluxes.

**4.2. Two spatial dimensions.**

**4.2.1. Linear source term.** Similarly consider the setup of the active flux method on two-dimensional Cartesian grids as described in 2.1. The available de-
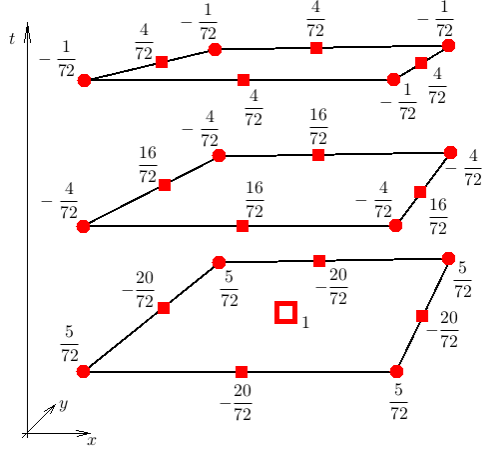
FIG. 3. *Illustration of the weights of the space time quadrature formula* (4.15).

grees of freedom are

(4.10) $$\text{3} \times \text{4 nodes: } q^n_{i\pm\frac{1}{2},j\pm\frac{1}{2}}, q^{n+\frac{1}{2}}_{i\pm\frac{1}{2},j\pm\frac{1}{2}}, q^{n+1}_{i\pm\frac{1}{2},j\pm\frac{1}{2}}$$

(4.11) $$\text{3} \times \text{2 vertical edges: } q^n_{i\pm\frac{1}{2},j}, q^{n+\frac{1}{2}}_{i\pm\frac{1}{2},j}, q^{n+1}_{i\pm\frac{1}{2},j}$$

(4.12) $$\text{3} \times \text{2 horizontal edges: } q^n_{i,j\pm\frac{1}{2}}, q^{n+\frac{1}{2}}_{i,j\pm\frac{1}{2}}, q^{n+1}_{i,j\pm\frac{1}{2}}$$

(4.13) $$\text{1 average: } \bar{q}^n_{ij}$$

The ansatz for a space-time polynomial is

(4.14) $$\mathscr{P}(t,x,y) = \left( \sum_{\zeta+\eta+\vartheta \leq 4} a_{\zeta\eta\vartheta} \cdot x^\zeta y^\eta t^\vartheta \right) + a_{212}x^2yt^2 + a_{122}xy^2t^2$$

It admits a unique solution to the interpolation problem given the available de-grees of freedom and yields the following quadrature formula (see also figure 3):

(4.15)
$$\frac{1}{\Delta x}\int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}}\mathrm{d}x\,\frac{1}{\Delta y}\int_{-\frac{\Delta y}{2}}^{\frac{\Delta y}{2}}\mathrm{d}y\,\frac{1}{\Delta t}\int_0^{\Delta t}\mathrm{d}t\,q(t,x,y) = \bar{q}^n_{ij}$$
$$-\frac{20}{72}\left(q^n_E + q^n_N + q^n_S + q^n_W\right) + \frac{5}{72}\left(q^n_{NE} + q^n_{NW} + q^n_{SE} + q^n_{SW}\right)$$
$$+\frac{16}{72}\left(q^{n+\frac{1}{2}}_E + q^{n+\frac{1}{2}}_N + q^{n+\frac{1}{2}}_S + q^{n+\frac{1}{2}}_W\right) - \frac{4}{72}\left(q^{n+\frac{1}{2}}_{NE} + q^{n+\frac{1}{2}}_{NW} + q^{n+\frac{1}{2}}_{SE} + q^{n+\frac{1}{2}}_{SW}\right)$$
$$+\frac{4}{72}\left(q^{n+1}_E + q^{n+1}_N + q^{n+1}_S + q^{n+1}_W\right) - \frac{1}{72}\left(q^{n+1}_{NE} + q^{n+1}_{NW} + q^{n+1}_{SE} + q^{n+1}_{SW}\right)$$

The time levels $(n, n+\frac{1}{2}, n+1)$ contribute again with weights $(\frac{1}{6}, \frac{2}{3}, \frac{1}{6})$, and the edges always contribute $-4$ times the nodes.

**4.2.2. Nonlinear source term.** Again, for nonlinear source instead of the average it is necessary to use the evaluation of the reconstruction at the cell midpoint.

This amounts to an approximation of the average by a two-dimensional Simpson rule.
Then the source term is approximating as follows:

$$
\begin{aligned}
\frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} \mathrm{d}x \, \frac{1}{\Delta y} \int_{-\frac{\Delta y}{2}}^{\frac{\Delta y}{2}} \mathrm{d}y \, \frac{1}{\Delta t} \int_{0}^{\Delta t} \mathrm{d}t \, s(q(t,x,y)) = \; & \frac{32}{72} s(q_{\mathrm{recon},ij}(0,0)) \\
& - \frac{12}{72} \left( s(q_{\mathrm{E}}^n) + s(q_{\mathrm{N}}^n) + s(q_{\mathrm{S}}^n) + s(q_{\mathrm{W}}^n) \right) \\
& + \frac{3}{72} \left( s(q_{NE}^n) + s(q_{NW}^n) + s(q_{SE}^n) + s(q_{SW}^n) \right) \\
& + \frac{16}{72} \left( s(q_{\mathrm{E}}^{n+\frac{1}{2}}) + s(q_{\mathrm{N}}^{n+\frac{1}{2}}) + s(q_{\mathrm{S}}^{n+\frac{1}{2}}) + s(q_{\mathrm{W}}^{n+\frac{1}{2}}) \right) \\
& - \frac{4}{72} \left( s(q_{NE}^{n+\frac{1}{2}}) + s(q_{NW}^{n+\frac{1}{2}}) + s(q_{SE}^{n+\frac{1}{2}}) + s(q_{SW}^{n+\frac{1}{2}}) \right) \\
& + \frac{4}{72} \left( s(q_{\mathrm{E}}^{n+1}) + s(q_{\mathrm{N}}^{n+1}) + s(q_{\mathrm{S}}^{n+1}) + s(q_{\mathrm{W}}^{n+1}) \right) \\
& - \frac{1}{72} \left( s(q_{NE}^{n+1}) + s(q_{NW}^{n+1}) + s(q_{SE}^{n+1}) + s(q_{SW}^{n+1}) \right)
\end{aligned}
\tag{4.16}
$$

In case that the data only depend on one of the variables, the two-dimensional quadratures (4.15) and (4.16) do *not* exactly reduce to the one dimensional quadratures (4.6) and (4.9). This is because (cf. Figure 3) the point values on edge midpoints $\left(0, \pm\frac{\Delta y}{2}\right)$ do not disappear even if the data depend only on $x$, and therefore the available degrees of freedom remain different from the one-dimensional case.

## 5. Well-balanced property for acoustics with gravity.

**5.1. Exact evolution operator.** As described in 3.2 a closed-form exact evolution operator for acoustics with gravity is very difficult to obtain. Nevertheless, it is still possible to show that a scheme endowed with such an operator would be well-balanced / stationarity preserving; i.e. that there exists a discretization of the stationary states of the PDE which remain exactly stationary. This proof does not require the evolution operator to be known explicitly, but only relies on the fact that it is exact. Besides its fundamental importance, this result is used in section 5.2 to analyze the situation for the approximate evolution operator and to restore the well-balanced property for it.

The numerical stationary states are best studied upon the (discrete) Fourier transform. Define $t_x := \exp(\hat{\imath} k_x \Delta x)$, $t_y := \exp(\hat{\imath} k_y \Delta y)$. Here $\hat{\imath}$ is the imaginary unit and $\mathbf{k} = (k_x, k_y) \in \mathbb{R}^2$ is the wave vector characterizing the spatial frequency of the Fourier mode. Applying the Fourier transform introduces one mode $\bar{q}$ for the averages and one mode $q$ for the point values; this implies writing $q_i := \bar{q} t_x^i t_y^j$, $q_{i+\frac{1}{2}} := q t_x^i t_y^j$.

THEOREM 5.1 (Stationarity preservation with exact evolution). *If the discrete data fulfill*

$$
\bar{\rho}_i = \frac{\rho_{i+\frac{1}{2}} + \rho_{i-\frac{1}{2}}}{2}
\tag{5.1}
$$

$$
\frac{p_{i+\frac{1}{2}} - p_{i-\frac{1}{2}}}{\Delta x} = g \frac{\rho_{i-\frac{1}{2}} + \rho_{i+\frac{1}{2}}}{2}
\tag{5.2}
$$

$$
\frac{\bar{p}_{i+\frac{3}{2}} - \bar{p}_{i+\frac{1}{2}}}{\Delta x} = g \frac{\rho_{i+\frac{3}{2}} + 4\rho_{i+\frac{1}{2}} + \rho_{i-\frac{1}{2}}}{6}
\tag{5.3}
$$

*and the exact evolution operator for* (3.6)–(3.8) *is used, then the numerical solution remains stationary.*

*Proof.* The proof consists of two parts.

i) Consider first the evolution of the point values. When the exact evolution opera-
tor is used to update the point values, they remain stationary if the reconstruction
fulfills

$$(5.4) \qquad v_{\text{recon}}(x) = \text{const} \qquad\qquad \partial_x p_{\text{recon}}(x) = \rho_{\text{recon}}(x)g$$

Upon the Fourier transform this becomes (w.l.o.g. $x_i = 0$)

$$(5.5) \qquad -3\left(2\bar{p} - p\left(1 + \frac{1}{t_x}\right)\right)\frac{2x}{\Delta x^2} + p\left(1 - \frac{1}{t_x}\right)\frac{1}{\Delta x} =$$

$$(5.6) \quad -3g\left(2\bar{\rho} - \rho\left(1 + \frac{1}{t_x}\right)\right)\frac{x^2}{\Delta x^2} + g\rho\left(1 - \frac{1}{t_x}\right)\frac{x}{\Delta x} + g\frac{6\bar{\rho} - \rho\left(1 + \frac{1}{t_x}\right)}{4}$$

This shall be valid for all $x$:

$$(5.7) \qquad\qquad 2\bar{\rho} - \rho(1 + 1/t_x) = 0$$

$$(5.8) \qquad\qquad -2\bar{p}t_x + p(t_x + 1) = \frac{\Delta x g\rho(t_x - 1)}{6}$$

$$(5.9) \qquad\qquad p(t_x - 1) = \Delta x g\frac{6\bar{\rho}t_x - \rho(t_x + 1)}{4}$$

These are three equations for four variables. In particular

$$(5.10) \qquad\qquad \bar{\rho} = \frac{\rho(1 + 1/t_x)}{2}$$

$$(5.11) \qquad\qquad p = \Delta x g\rho\frac{t_x + 1}{2(t_x - 1)}$$

$$(5.12) \qquad\qquad \bar{p} = \Delta x g\rho\frac{t_x^2 + 4t_x + 1}{6t_x(t_x - 1)}$$

These statements can be rewritten as finite difference formulae by inverting the
Fourier transform:

$$(5.13) \qquad\qquad \bar{\rho} = \frac{\rho_{i+\frac{1}{2}} + \rho_{i-\frac{1}{2}}}{2}$$

$$(5.14) \qquad\qquad \frac{p_{i+\frac{1}{2}} - p_{i-\frac{1}{2}}}{\Delta x} = g\frac{\rho_{i-\frac{1}{2}} + \rho_{i+\frac{1}{2}}}{2}$$

$$(5.15) \qquad\qquad \frac{\bar{p}_{i+1} - \bar{p}_i}{\Delta x} = g\frac{\rho_{i+\frac{3}{2}} + 4\rho_{i+\frac{1}{2}} + \rho_{i-\frac{1}{2}}}{6}$$

ii) Assume now (5.10)–(5.12) to be true. Simpson's rule in time for the flux average
is trivial, and thus the update of the cell average amounts to

$$(5.16) \qquad \frac{\bar{v}^{n+1} - \bar{v}^n}{\Delta t} + \frac{p(1 - 1/t_x)}{\Delta x} = \frac{\bar{v}^{n+1} - \bar{v}^n}{\Delta t} + g\rho\frac{t_x + 1}{2t_x}$$

$$(5.17) \qquad\qquad\qquad = \frac{\bar{v}^{n+1} - \bar{v}^n}{\Delta t} + g\bar{\rho} \qquad\qquad \square$$

The quadrature formula (4.6) for the source reduces to $g\bar{\rho}$ if the point values are
stationary, which implies $\bar{v}^{n+1} = \bar{v}^n$. This completes the proof.

The equations (5.10)–(5.12) contain $\rho$ as a free variable. One can rewrite the system making $p$ the free variable:

(5.18) $$\bar{\rho} = \frac{p(t_x - 1)}{t_x \Delta x g} \qquad \rho = \frac{2p(t_x - 1)}{\Delta x g(t_x + 1)} \qquad \bar{p} = p\frac{t_x^2 + 4t_x + 1}{3t_x(t_x + 1)}$$

This form will be useful later.

Equations (5.2)–(5.3) are finite difference approximations of $\partial_x p = \rho g$. Equation (5.1) implies that the reconstructed $\rho$ of the discrete stationary state is linear, which is clear: for quadratic reconstructions to fulfill (5.4), $\rho_{\text{recon}}$ has to be linear in each cell. The slope of the linear function can vary from cell to cell and is given by (5.2).

**5.2. Approximate evolution operator.** The above section identifies conditions (5.1)–(5.3) on the discrete data for them to remain stationary upon usage of the *exact* evolution operator. Unfortunately, such an operator is unavailable in practice. Having identified an approximate solution operator, which agrees with the exact solution up to terms $\mathcal{O}(t^3)$ in section 3.3, here we study whether it keeps the same data (5.1)–(5.3) stationary as well.

THEOREM 5.2. *If the discrete data fulfill* (5.1)–(5.3) *and the approximate evolution operator of theorem 3.2 for* (3.6)–(3.8) *is used, then both the pressure $p$ and the density $\rho$ remain stationary over one time step, but the velocity undergoes the time evolution*

(5.19) $$v_{i+\frac{1}{2}}(t) = -\frac{\alpha g^2}{4}\frac{\rho_{i+\frac{1}{2}} - \rho_{i-\frac{1}{2}}}{\Delta x}t^3$$

*Proof.* Assume the initial data to fulfill (5.1)–(5.3), or equivalently (5.4). Using (2.7) (and applying the discrete Fourier transform straight away) (5.4) implies

(5.20)
$$p_{\text{recon}}(x) = \frac{1}{4}\left(6\bar{p} - p\left(1 + \frac{1}{t_x}\right)\right) + \frac{x}{\Delta x}\left(1 - \frac{1}{t_x}\right)p - 3\frac{x^2}{\Delta x^2}\left(2\bar{p} - p\left(1 + \frac{1}{t_x}\right)\right)$$

(5.21)
$$\rho_{\text{recon}}(x) = \frac{1}{g\Delta x}\left(p\left(1 - \frac{1}{t_x}\right) - 6\frac{x}{\Delta x}\left(2\bar{p} - p\left(1 + \frac{1}{t_x}\right)\right)\right)$$

(5.22)
$$v_{\text{recon}}(x) = 0$$

and using (3.44) therefore

(5.23)
$$Q_{1,0}(x) = Q_{2,0}(x) = -\frac{p(1 + t_x) - 6\bar{p}t_x}{8t_x} + \frac{p(t_x - 1)x}{2\Delta x t_x} + \frac{3(p(1 + t_x) - 2\bar{p}t_x)x^2}{2\Delta x^2 t_x}$$

(5.24)
$$Q_{3,0}(x) = \frac{p(-1 + t_x)}{\Delta x g t_x} + \frac{p - 6\bar{p}t_x + pt_x}{4c^2 t_x}$$
$$+ \frac{\left(-\Delta x g p(t_x - 1) + 6c^2(p(1 + t_x) - 2\bar{p}t_x)\right)x}{c^2 \Delta x^2 g t_x} - \frac{3(p(1 + t_x) - 2\bar{p}t_x)x^2}{c^2 \Delta x^2 t_x}$$

Evaluating the Runge-Kutta algorithm of section 3.3 on these initial data (at

620   $x = \frac{\Delta x}{2}$) yields

621   (5.25)
$$\begin{pmatrix} \rho \\ v^* \\ p \end{pmatrix} \qquad \text{with} \qquad v^* = -\frac{\alpha g(t_x - 1)^2}{2\Delta x^2 t_x(t_x + 1)} p t^3$$

623   ($\alpha$ is the parameter appearing in the RK2 method.)

624       Recall that $\rho$ and $p$ are the Fourier coefficients of the point values of the density
625   and the pressure. Obviously $\rho$ and $p$ remain stationary, but the velocity does not.
626   Using (5.18) $v^*$ can be rewritten as

627   (5.26)
$$v^* = -\frac{\alpha g^2}{4\Delta x}\left(1 - \frac{1}{t_x}\right)\rho t^3 = -\frac{\alpha g^2}{4}\frac{\rho_{i+\frac{1}{2}} - \rho_{i-\frac{1}{2}}}{\Delta x} t^3$$

629   having applied the inverse Fourier transform in the last step.       □

630       Observe that the time evolution of the velocity is consistent with the accuracy of
631   the algorithm ($\mathcal{O}(t^3)$).

632       COROLLARY 5.3 (Stationarity preservation with approximate evolution). *If the*
633   *algorithm of section 3.3 is modified by adding the term*

634   (5.27)
$$\frac{\alpha g^2}{4}\frac{\rho_{i+\frac{1}{2}} - \rho_{i-\frac{1}{2}}}{\Delta x} t^3$$

636   *to the velocity evolution, then*
637    *i) its accuracy is not changed*
638   *ii) it becomes stationarity preserving / well-balanced with the same discrete station-*
639     *ary states as the exact evolution operator*

640       The two forms (5.25) and (5.19) of $v^*$ are equivalent, because the initial data
641   have been chosen to be stationary, and thus additionally fulfill (5.18). The proposed
642   modification is to *always* add $-v^*$ to the velocity evolution, irrespective of whether
643   the data fulfill (5.18) or not. At this point the Fourier coefficients of $\rho$ and $p$ are
644   independent and it matters whether the correction is used in the form (5.25) or (5.19).
645   Of course, also the inverse Fourier transform has to be applied to the expression first
646   in order for the correction to attain the form of a finite difference formula. Compact
647   finite difference formulae are in one-to-one-correspondence with Laurent polynomials
648   in $t_x$. An expression such as $\frac{1}{t_x+1} = 1 - t_x + t_x^2 \mp \ldots$ is an expression involving an
649   unbounded stencil and cannot be implemented in usual codes. Therefore (5.19) cannot
650   be used as a correction because the correction would have a non-compact stencil (just
651   as the equivalent expressions involving only $\bar{\rho}$ or $\bar{p}$). This is why the form (5.25) which
652   involves point values of $\rho$ is preferred.

653       Being always present in the velocity evolution (and not only at stationary states),
654   the modification (5.27) might in general affect the stability of the algorithm, but it
655   has not been found to have any effect on the stability in practice.

656   **6. Numerical examples.** The numerical examples of this section serve to il-
657   lustrate the performance of the new method. The equations discussed are linear
658   advection with different source terms (in one and two spatial dimensions, as intro-
659   duced in section 3.1) and linear acoustics with gravity (introduced in section 3.2). In
660   both cases it is demonstrated that the method achieves third order of accuracy in the
661   experiments. For acoustics with gravity additionally the discrete stationary states are
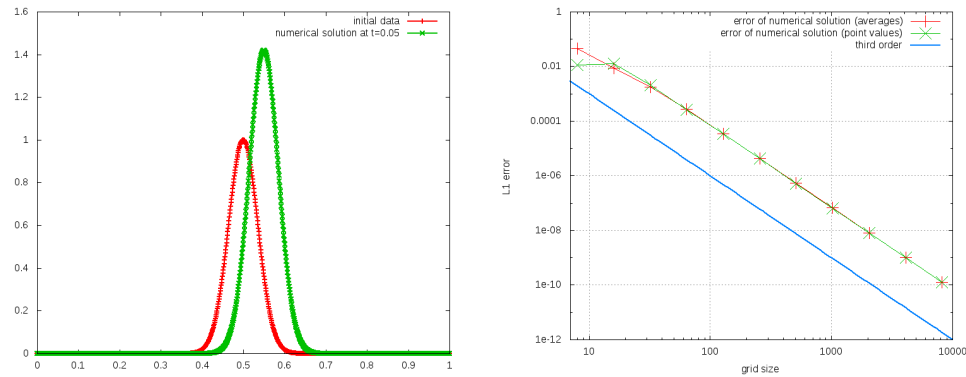662   studied and shown to agree with the prediction of section 5.

FIG. 4. *Gaussian initial data for* (6.1) *with* $\mathbf{U} = \mathbf{e}_x$, $\kappa = 7$. *Note that due to the source term, the Gaussian is advected and also changes shape. Exact evolution operator* (3.4) *and quadrature formula* (4.6) *have been used with CFL = 0.45.* Left: *Initial data and solution at* $t = 0.05$ *(cell averages) on a grid with 1000 cells.* Right: *Error of the numerical solution as a function of the grid size shows third order convergence.*
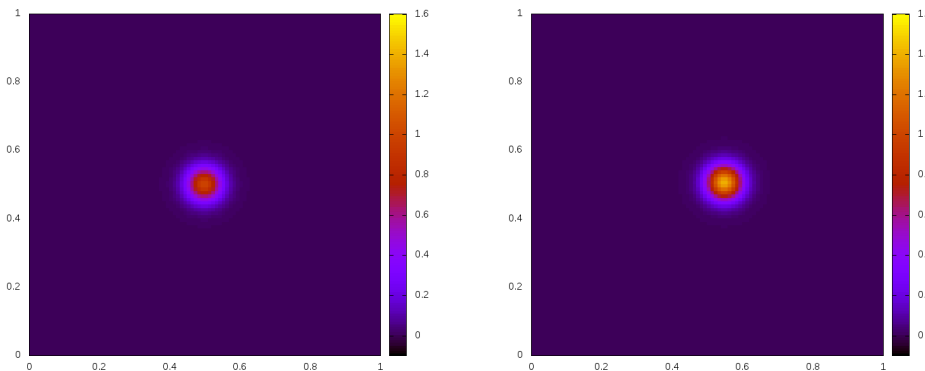


FIG. 5. *Gaussian initial data for* (6.1) *with* $\mathbf{U} = (1, 0.1)$, $\kappa = 7$. *Note that due to the source term, the Gaussian is advected and also changes shape. Exact evolution operator* (3.4) *and quadrature formula* (4.15) *have been used with CFL = 0.45.* Left: *Initial setup.* Right: *Numerical solution at* $t = 0.05$ *on a* $100 \times 100$ *Cartesian grid.*

## 6.1. Linear advection. Consider first

$$\partial_t q + \mathbf{U} \cdot \nabla q = \kappa q \tag{6.1}$$

with the exact solution given by (3.4). In Figures 4–6 the exact solution operator is used for the evolution of the point values and third order convergence is observed. This shows that the quadrature formulae (4.6) and (4.15) used to evolve the cell averages indeed yield a third order scheme. Figure 4 shows the setup for a one-dimensional situation together with a convergence study, Figure 5 shows the setup in two spatial dimensions and Figure 6 shows the corresponding convergence study.

Consider now

$$\partial_t q + \mathbf{U} \cdot \nabla q = \kappa q^B \qquad B \neq 1 \tag{6.2}$$

with the exact solution (3.5) and $\kappa = 7$, $B = 3$. Figure 7 (left) shows the initial data and the numerical solution, and Figure 7 (right) shows a convergence study for
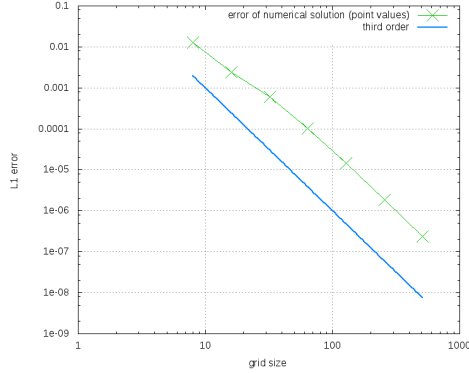
FIG. 6. *Convergence study for the setup shown in Figure 5. One observes third order accuracy.*
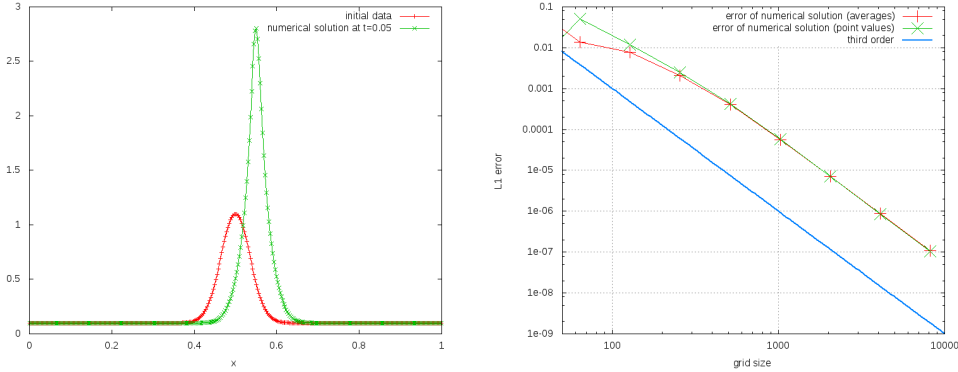


FIG. 7. *Gaussian initial data for (6.2) with $s(q) = \kappa q^B$ and $\mathbf{U} = \mathbf{e}_x$, $\kappa = 7$, $B = 3$. Runge-Kutta approximate evolution operator from Corollary 3.4 and quadrature formula (4.9) have been used with CFL = 0.45. The solution has been computed on a grid covering $[-1 : 2]$, but the error is only computed inside $[0, 1]$ to exclude any boundary influence. Left: Initial setup and solution at $t = 0.05$ (point values) on a grid with 1000 cells. Right: Error of the numerical solution as a function of the grid size shows third order convergence. The exact solution is given by (3.5).*

the approximate evolution operator from Corollary (3.4). One observes third order accuracy, as expected.

**6.2. Acoustics with gravity.** Consider now the equations of linear acoustics with a gravity source term (3.6)–(3.8). The exact solution operator is only partly available in closed form, and therefore the approximate Runge-Kutta evolution operator of section 3.3 is used in combination with the well-balancing fix (5.27). The parameter $\alpha$ in the Runge-Kutta method is chosen to $\alpha = \frac{1}{2}$.

Figure 8 shows a stationary setup given by

$$(6.3) \qquad p = A_1 x^2 + A_2 x + A_3 \qquad \rho = 2A_1 x/g + A_2/g \qquad v = 0$$

with $A_1 = 17, A_2 = -3, A_3 = 1$. This parabola is exactly recovered by the reconstruction, and thus remains stationary up to machine precision. This experiment shows that the well-balancing fix works as it should.
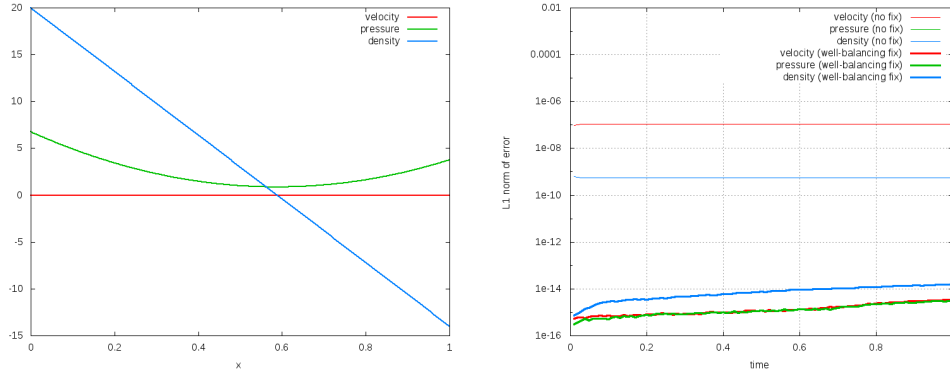
Fig. 8. *Setup of a stationary parabola* (6.3) *for* (3.6)–(3.8), *solved using the Runge-Kutta approximate evolution operator of section* 3.3 *with and without well-balancing* (5.27). *Here* $g = -1$, *and the setup is solved on a grid covering* $[-1.5, 2.5]$, *but the error is only measured inside* $[0, 1]$ $(\Delta x = 10^{-2})$ *to exclude the influence of the boundaries.* Left*: Setup.* Right*: Error of numerical solution as a function of time. Thin lines: without the well-balancing* (5.27). *Thick lines: including the well-balancing* (5.27). *In the latter case one only observes an evolution due to machine error.*
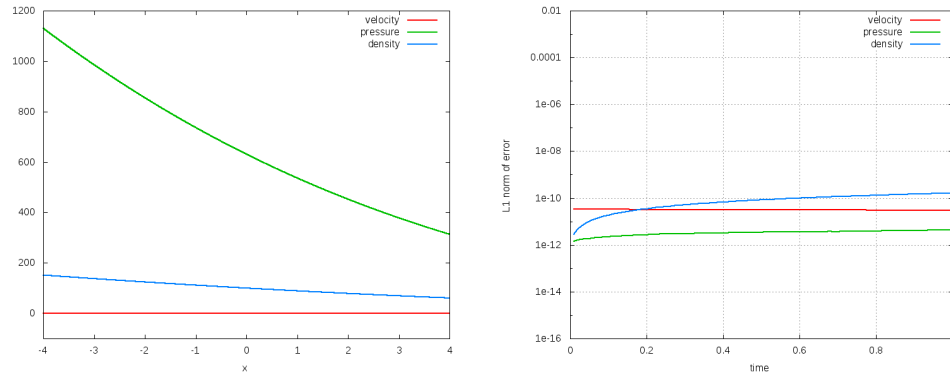


Fig. 9. *Stationary setup* (6.4) *for* (3.6)–(3.8), *solved using the Runge-Kutta approximate evolution operator of section* 3.3 *with well-balancing* (5.27). *Here* $g = -1$, *and the setup is solved on a grid covering* $[-5.5, 5.5]$, *but the error is only measured inside* $[-4, 4]$ $(\Delta x = 10^{-3})$ *to exclude the influence of the boundaries.* Left*: Setup (cell averages).* Right*: Error of numerical solution as a function of time. One observes a transition towards a numerical stationary state which then persists forever.*

Consider next (Figure 9) the stationary setup fulfilling $p = K\rho^\gamma$, i.e.

$$(6.4) \qquad \rho = \left( \frac{g(\gamma - 1)}{K\gamma} x + \rho_0^{\gamma - 1} \right)^{\frac{1}{\gamma - 1}}$$

with $K = 1, \gamma = 1.4$, $\rho_0 = 100$. This is reminiscent of an isentropic atmosphere in the context of the Euler equations. This setup is not recovered exactly by the reconstruction, but one observes a numerical evolution towards a discrete stationary state which then persists forever.

Next, a perturbation

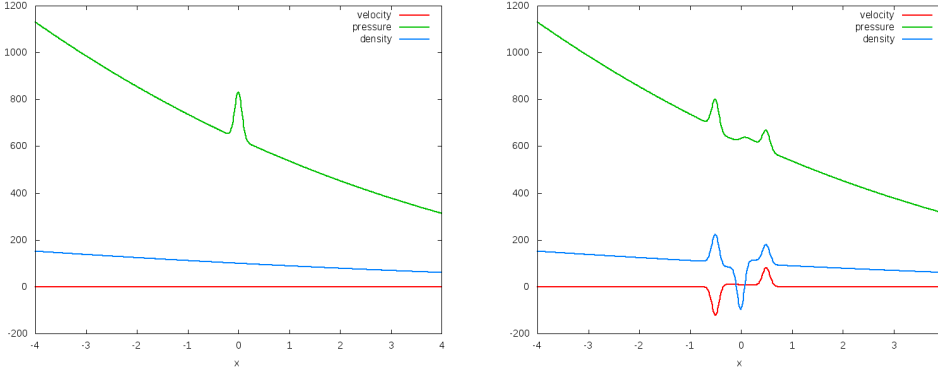$$(6.5) \qquad\qquad\qquad\qquad 200 \exp(-100x^2)$$

FIG. 10. *Setup* (6.4) *endowed with the pressure perturbation* (6.5) *solved using the Runge-Kutta approximate evolution operator of section* 3.3 *with well-balancing* (5.27). *Left: Initial data (cell averages). Right: Numerical solution (cell averages) at* $t = 0.5$ *on a grid covering* $[-5.5, 5.5]$, *but only the subinterval* $[-4, 4]$ *is considered in order to exclude the influence of the boundaries.* $\Delta x = 0.01$, *CFL = 0.45*.
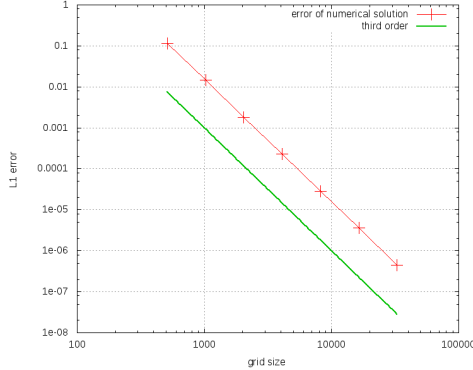


FIG. 11. *Setup of Figure* 10. *The error of the numerical solution is measured on the point values. One observes third order accuracy.*

in the pressure is added onto the setup (6.4). In order to study the accuracy of the scheme on this setup, it is solved on a grid of $131072 = 2^{18}$ cells and the solution is used as reference. Again, $g = -1, K = 1, \gamma = 1.4$. Figure 10 shows the setup and the numerical solution at $t = 0.5$, and Figure 11 shows a convergence study which displays third order convergence.

Consider finally a Riemann problem:

$$(6.6) \qquad \rho = 3.5 \qquad\qquad p = 1.5 \qquad\qquad v = \begin{cases} 1 & 0.25 \le x \le 0.75 \\ 3 & \text{else} \end{cases}$$

This Riemann problem can be solved exactly using the formula (A.18)–(A.22). Note that if all quantities are constant in space, then they solve

$$(6.7) \qquad\qquad\qquad\qquad \partial_t \rho = 0$$

$$(6.8) \qquad\qquad\qquad\qquad \partial_t p = 0$$

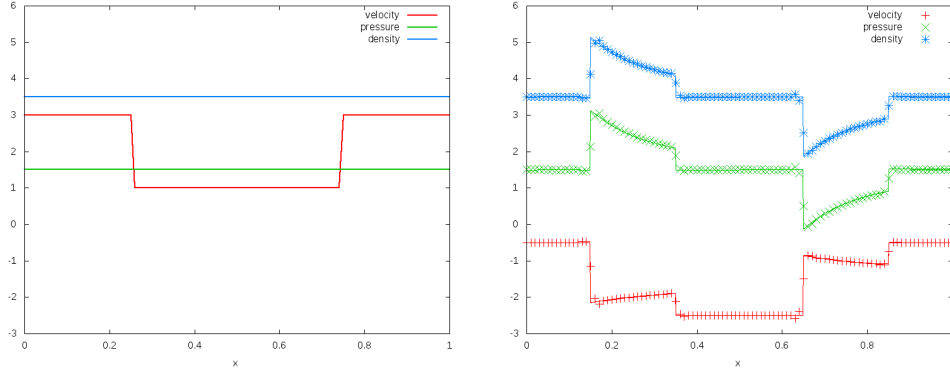$$(6.9) \qquad\qquad\qquad\qquad \partial_t v = \rho g$$

FIG. 12. *Riemann problem setup* (6.6) *solved using the Runge-Kutta approximate evolution operator of section* 3.3 *with well-balancing* (5.27). *Here,* $g = -10$. *Left: Initial data.* *Right: Numerical solution (dots) and exact solution (solid line) at* $t = 0.1$. *$\Delta x = 0.01$, CFL = 0.45. Point values of the numerical solution are shown are shown.*

which means that $\rho$ and $p$ remain stationary, but that $v = v(t = 0) + \rho g t$. The solution to the initial data (6.6) therefore can be obtained by adding the time evolution of

$$\begin{pmatrix} 0 \\ v_0(x) \\ 0 \end{pmatrix}$$ (via numerical quadrature of (A.18)–(A.22)) and the time evolution of

$$\begin{pmatrix} \rho \\ 0 \\ p \end{pmatrix}$$ which is just $$\begin{pmatrix} \rho \\ \rho g t \\ p \end{pmatrix}$$. Figure 12 shows the numerical and the exact solution.

**7. Conclusions and outlook.** Active flux is a novel kind of numerical method for hyperbolic problems, extending the finite volume method. Instead of computing the intercell flux via a Riemann problem it relies on a continuous reconstruction and on accurately evolved point values along the cell boundary. They then immediately serve as quadrature values for the computation of the intercell flux. The extension of active flux to time dependent balance laws presented in this paper requires a modification in both these aspects: the evolution of the point values and the average update need to account for the source term. Here, an approximate evolution operator is suggested for the point value update; this is done for linear systems with possibly nonlinear source terms in one spatial dimension, and linear scalar equations with source terms in multiple spatial dimensions. A suitable quadrature is suggested in order to approximate the contribution of the source term to the cell average. This quadrature can be applied to any system of (nonlinear) balance laws.

We aim at combining the strategy presented in this paper with an approximate evolution operator for a nonlinear homogeneous problem (such as those suggested in [Bar19a]) in future. Multi-dimensional systems of hyperbolic conservation laws are very different from their one-dimensional counterparts because in general characteristics are unavailable and need to be conceptually replaced by characteristic cones. Examples of evolution operators that make use of such cones can be found in [ER13, FR15, Fan17, BHKR19]. Combining these with an approximate evolution of

738 the source term shall pave the way towards the extension of active flux to nonlinear
739 multi-dimensional balance laws and the derivation of accurate structure preserving
740 (in particular well-balanced) methods for them.

741 **Appendix A. Exact solution of linear acoustics with gravity.**
742 System (3.6)–(3.8) can in principle be immediately solved exactly via Fourier
743 transform by inserting the ansatz

744 (A.1)
$$\begin{pmatrix} \rho \\ v \\ p \end{pmatrix} = \begin{pmatrix} \hat{\rho} \\ \hat{v} \\ \hat{p} \end{pmatrix} \exp(\mathring{\imath} k \cdot x - \mathring{\imath} \omega t)$$

746 into (3.6)–(3.8):

747 (A.2)
$$\omega \begin{pmatrix} \hat{\rho} \\ \hat{v} \\ \hat{p} \end{pmatrix} = \begin{pmatrix} 0 & k & 0 \\ \mathring{\imath} g & 0 & k \\ 0 & c^2 k & 0 \end{pmatrix} \begin{pmatrix} \hat{\rho} \\ \hat{v} \\ \hat{p} \end{pmatrix}$$

749 Therefore $\omega = 0$, or $\omega = \pm\sqrt{c^2 k^2 + \mathring{\imath} g k}$. The complex eigenvalue can be removed
750 upon transforming

751
752 (A.3)
$$\rho = \tilde{\rho} e^{\mu x} \qquad\qquad v = \tilde{v} e^{\mu x} \qquad\qquad p = \tilde{p} e^{\mu x}$$

753 with

754
755 (A.4)
$$\mu := \frac{g}{2c^2}$$

756 System (3.6)–(3.8) then reads

757 (A.5) $$\partial_t \tilde{\rho} + \partial_x \tilde{v} = -\mu \tilde{v}$$
758 (A.6) $$\partial_t \tilde{v} + \partial_x \tilde{p} = \tilde{\rho} g - \mu \tilde{p}$$
759
760 (A.7) $$\partial_t \tilde{p} + c^2 \partial_x \tilde{v} = -c^2 \mu \tilde{v}$$

761 Now, a solution of (A.5)–(A.7) shall be found. For better readability, drop the tilde.
762 Upon the Fourier transform (A.5)–(A.7) becomes

763 (A.8)
$$\omega \begin{pmatrix} \hat{\rho} \\ \hat{v} \\ \hat{p} \end{pmatrix} = \mathcal{E} \begin{pmatrix} \hat{\rho} \\ \hat{v} \\ \hat{p} \end{pmatrix} \qquad \mathcal{E} = \begin{pmatrix} 0 & k - \mathring{\imath}\mu & 0 \\ \mathring{\imath} g & 0 & k - \mathring{\imath}\mu \\ 0 & c^2 k - \mathring{\imath} c^2 \mu & 0 \end{pmatrix}$$

765 The eigenvalues of $\mathcal{E}$ are now real: $\omega_1 = 0$, $\omega_{2,3} = \pm c\sqrt{k^2 + \mu^2}$. Although this
766 transformation brings the endeavour of finding the exact solution to (3.6)–(3.8) into
767 the realm of the possible, technical difficulties prevent one from actually computing
768 all Green's functions in closed form.
769 Assume therefore that the only non-vanishing initial data are in the velocity.
770 Then the Fourier mode at initial time reads

771 (A.9)
$$\begin{pmatrix} 0 \\ \hat{v} \\ 0 \end{pmatrix} \exp(\mathring{\imath} k x)$$

772

773   and at a later time it becomes

774   (A.10)
$$\sum_{m=1}^{3} v_m \exp(\mathring{\imath}kx - \mathring{\imath}\omega_m t)$$

775

776   where the decomposition of $\begin{pmatrix} 0 \\ \hat{v} \\ 0 \end{pmatrix}$ in the eigenbasis of $\mathcal{E}$ is used, i.e.

777   (A.11)
$$\begin{pmatrix} 0 \\ \hat{v} \\ 0 \end{pmatrix} = \sum_{m=1}^{3} v_m \qquad\qquad \mathcal{E}v_m = \omega_m v_m$$

778

779   Such a basis is given e.g. by

780   (A.12)
$$e_1 = \begin{pmatrix} \mu + \mathring{\imath}k \\ 0 \\ g \end{pmatrix} \qquad\qquad e_{2,3} = \begin{pmatrix} \mu + \mathring{\imath}k \\ \pm \mathring{\imath}c\sqrt{k^2 + \mu^2} \\ c^2(\mu + \mathring{\imath}k) \end{pmatrix}$$

781

782   Collecting the terms yields the time evolution of the Fourier mode (A.9):

783   (A.13)
$$\hat{v} \exp(\mathring{\imath}kx) \begin{pmatrix} -\dfrac{(\mu + \mathring{\imath}k)\sin\left(ct\sqrt{k^2 + \mu^2}\right)}{c\sqrt{k^2 + \mu^2}} \\ \cos\left(ct\sqrt{k^2 + \mu^2}\right) \\ -\dfrac{c^2(\mu + \mathring{\imath}k)\sin\left(ct\sqrt{k^2 + \mu^2}\right)}{c\sqrt{k^2 + \mu^2}} \end{pmatrix}$$

784   (A.14)
$$= \hat{v} \begin{pmatrix} -(\mu + \partial_x) \\ \partial_t \\ -c^2(\mu + \partial_x) \end{pmatrix} \exp(\mathring{\imath}kx) \frac{\sin\left(ct\sqrt{k^2 + \mu^2}\right)}{c\sqrt{k^2 + \mu^2}}$$

785

786   Green's function is obtained by inserting the Fourier transform of a Dirac $\delta_{x'}$ at
787   $x'$, i.e. taking $\hat{v} = \frac{\exp(-\mathring{\imath}kx')}{\sqrt{2\pi}}$ and performing the inverse Fourier transform with the
788   help of formula 1.7 (30) in [Bat54]. This yields, wherever defined,

789   (A.15)
$$\begin{pmatrix} G_\rho(t, x; x') \\ G_v(t, x; x') \\ G_p(t, x; x') \end{pmatrix} = \begin{pmatrix} -(\mu + \partial_x) \\ \partial_t \\ -c^2(\mu + \partial_x) \end{pmatrix} \frac{1}{2c} J_0\left(\mu\sqrt{(ct)^2 - (x - x')^2}\right)$$

790   (A.16)
$$+ \begin{pmatrix} -\dfrac{\delta_{x+ct} - \delta_{x-ct}}{2c} \\ \dfrac{\delta_{x+ct} + \delta_{x-ct}}{2} \\ c\left(\delta_{x+ct} - \delta_{x-ct}\right) \end{pmatrix}$$

791

where $J_0$ is the 0-th order Bessel function of the first kind, and $J_0' = -J_1$. Then the solution is obtained by performing a convolution with the initial data. Reinstalling the tilde one has

$$\text{(A.17)} \qquad \tilde{v}(t,x) = \int \mathrm{d}x'\, G_v(t,x;x')\tilde{v}_0(x')$$

$$\text{(A.18)} \qquad v(t,x) = \int \mathrm{d}x'\, G_v(t,x;x')\mathrm{e}^{\mu(x-x')}v_0(x')$$

$$\text{(A.19)} \qquad = \frac{1}{2}\int \mathrm{d}x'\, \mathrm{e}^{\mu(x-x')}\partial_{ct}J_0\left(\mu\sqrt{(ct)^2 - (x-x')^2}\right)v_0(x')$$

$$\text{(A.20)} \qquad + \frac{1}{2}\left(\mathrm{e}^{-\mu ct}v_0(x+ct) + \mathrm{e}^{\mu ct}v_0(x-ct)\right)$$

$$\text{(A.21)} \qquad \rho(t,x) = -\frac{1}{2c}\int \mathrm{d}x'\, \mathrm{e}^{\mu(x-x')}\left(\mu + \partial_x\right)J_0\left(\mu\sqrt{(ct)^2 - (x-x')^2}\right)v_0(x')$$

$$\text{(A.22)} \qquad - \frac{1}{2c}\left(\mathrm{e}^{-\mu ct}v_0(x+ct) - \mathrm{e}^{\mu ct}v_0(x-ct)\right)$$

and analogously for $p$. However, it is easier to note that

$$\text{(A.23)} \qquad \partial_t(c^2\rho - p) = 0$$

such that

$$\text{(A.24)} \qquad p(t,x) = p_0(x) + c^2\left(\rho(t,x) - \rho_0(x)\right)$$

## REFERENCES

[ABB+04] Emmanuel Audusse, François Bouchut, Marie-Odile Bristeau, Rupert Klein, and Benoit Perthame. A fast and stable well-balanced scheme with hydrostatic reconstruction for shallow water flows. *SIAM Journal on Scientific Computing*, 25(6):2050–2065, 2004.

[Bar18] Wasilij Barsukow. *Low Mach number finite volume methods for the acoustic and Euler equations*. Doctoral thesis, University of Wuerzburg, 2018.

[Bar19a] W Barsukow. The active flux scheme for nonlinear problems which admit characteristic variables. *submitted to J. Sci. Comp*, 2019.

[Bar19b] Wasilij Barsukow. Stationarity preserving schemes for multi-dimensional linear systems. *Mathematics of Computation*, 88(318):1621–1645, 2019.

[Bat54] Harry Bateman. *Tables of integral transforms (volume 1)*, volume 1. McGraw-Hill Book Company, 1954.

[BCK16] Jonas P Berberich, Praveen Chandrashekar, and Christian Klingenberg. A general well-balanced finite volume scheme for euler equations with gravity. In *XVI International Conference on Hyperbolic Problems: Theory, Numerics, Applications*, pages 151–163. Springer, 2016.

[BCK19] Jonas P Berberich, Praveen Chandrashekar, and Christian Klingenberg. High order well-balanced finite volume methods for multi-dimensional systems of hyperbolic balance laws. *arXiv preprint arXiv:1903.05154*, 2019.

[BCKR19] Jonas P Berberich, Praveen Chandrashekar, Christian Klingenberg, and Friedrich K Röpke. Second order finite volume scheme for Euler equations with gravity which is well-balanced for general equations of state and grid systems. *Communications in Computational Physics*, 26:599–630, 2019.

[BHKR19] Wasilij Barsukow, Jonathan Hohm, Christian Klingenberg, and Philip L Roe. The active flux scheme on Cartesian grids and its low Mach number limit. *Journal of Scientific Computing*, 81(1):594–622, 2019.

[BKCK20] Jonas P Berberich, Roger Käppeli, Praveen Chandrashekar, and Christian Klingenberg. High order discretely well-balanced finite volume methods for Euler equations with

838　　　　　gravity – without any a priori information about the hydrostatic solution. *arXiv*
839　　　　　*preprint arXiv:2005.01811*, 2020.
840　[BV94]　Alfredo Bermudez and Ma Elena Vázquez. Upwind methods for hyperbolic conservation
841　　　　　laws with source terms. *Computers & Fluids*, 23(8):1049–1071, 1994.
842　[CCK+18]　Alina Chertock, Shumo Cui, Alexander Kurganov, Şeyma Nur Özcan, and Eitan Tad-
843　　　　　mor. Well-balanced schemes for the Euler equations with gravitation: Conservative
844　　　　　formulation using global fluxes. *Journal of Computational Physics*, 2018.
845　[CK15]　Praveen Chandrashekar and Christian Klingenberg. A second order well-balanced fi-
846　　　　　nite volume scheme for Euler equations with gravity. *SIAM Journal on Scientific*
847　　　　　*Computing*, 37(3):B382–B402, 2015.
848　[CL94]　P Cargo and AY LeRoux. A well balanced scheme for a model of atmosphere with
849　　　　　gravity. *COMPTES RENDUS DE L ACADEMIE DES SCIENCES SERIE I-*
850　　　　　*MATHEMATIQUE*, 318(1):73–76, 1994.
851　[DZBK14]　Vivien Desveaux, Markus Zenk, Christophe Berthon, and Christian Klingenberg. A
852　　　　　well-balanced scheme for the Euler equation with a gravitational potential. In
853　　　　　*Finite Volumes for Complex Applications VII-Methods and Theoretical Aspects*,
854　　　　　pages 217–226. Springer, 2014.
855　[DZBK16]　Vivien Desveaux, Markus Zenk, Christophe Berthon, and Christian Klingenberg. A
856　　　　　well-balanced scheme to capture non-explicit steady states in the Euler equations
857　　　　　with gravity. *International Journal for Numerical Methods in Fluids*, 81(2):104–
858　　　　　127, 2016.
859　[ER13]　Timothy A Eymann and Philip L Roe. Multidimensional active flux schemes. In *21st*
860　　　　　*AIAA computational fluid dynamics conference*, 2013.
861　[Fan17]　Duoming Fan. *On the acoustic component of active flux schemes for nonlinear hyper-*
862　　　　　*bolic conservation laws*. PhD thesis, University of Michigan, Dissertation, 2017.
863　[FR15]　Doreen Fan and Philip L Roe. Investigations of a new scheme for wave propagation. In
864　　　　　*22nd AIAA Computational Fluid Dynamics Conference*, page 2449, 2015.
865　[GL96]　Joshua M Greenberg and Alain-Yves LeRoux. A well-balanced scheme for the numerical
866　　　　　processing of source terms in hyperbolic equations. *SIAM Journal on Numerical*
867　　　　　*Analysis*, 33(1):1–16, 1996.
868　[HKS19]　Christiane Helzel, David Kerkmann, and Leonardo Scandurra. A new ADER method
869　　　　　inspired by the active flux method. *Journal of Scientific Computing*, 80(3):1463–
870　　　　　1497, 2019.
871　[KM16]　R Käppeli and S Mishra. A well-balanced finite volume scheme for the Euler equations
872　　　　　with gravitation-the exact preservation of hydrostatic equilibrium with arbitrary
873　　　　　entropy stratification. *Astronomy & Astrophysics*, 587:A94, 2016.
874　[LeV98]　Randall J LeVeque. Balancing source terms and flux gradients in high-resolution Go-
875　　　　　dunov methods: the quasi-steady wave-propagation algorithm. *Journal of compu-*
876　　　　　*tational physics*, 146(1):346–365, 1998.
877　[LGB11]　Randall J LeVeque, David L George, and Marsha J Berger. Tsunami modelling with
878　　　　　adaptively refined finite volume methods. *Acta Numerica*, 20:211–289, 2011.
879　[NR16]　Hiroaki Nishikawa and Philip L Roe. Third-order active-flux scheme for advection
880　　　　　diffusion: hyperbolic diffusion, boundary condition, and Newton solver. *Computers*
881　　　　　*& Fluids*, 125:71–81, 2016.
882　[VL77]　Bram Van Leer. Towards the ultimate conservative difference scheme. IV. A new ap-
883　　　　　proach to numerical convection. *Journal of computational physics*, 23(3):276–299,
884　　　　　1977.