

1 **ON THE ACTIVE FLUX SCHEME FOR HYPERBOLIC PDES WITH**
2 **SOURCE TERMS***

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

4 **Abstract.** The Active Flux scheme is a Finite Volume scheme with additional point values
5 distributed along the cell boundary. It is third order accurate and does not require a Riemann
6 solver: the continuous reconstruction serves as initial data for the evolution of the points values.
7 The intercell flux is then obtained from the evolved values along the cell boundary by quadrature.
8 This paper focuses on the conceptual extension of Active Flux to include source terms, and thus for
9 simplicity assumes the homogeneous part of the equations to be linear. To a large part, the treatment
10 of the source terms is independent of the choice of the homogeneous part of the system. Additionally,
11 only systems are considered which admit characteristics (instead of characteristic cones). This is the
12 case for scalar equations in any number of spatial dimensions and systems in one spatial dimension.
13 Here, we succeed to extend the Active Flux method to include (possibly nonlinear) source terms
14 while maintaining third order accuracy of the method. This requires a novel (approximate) operator
15 for the evolution of point values and a modified update procedure of the cell average. For linear
16 acoustics with gravity, it is shown how to achieve a well-balanced / stationarity preserving numerical
17 method.

18 **Key words.** finite volume methods, Active Flux, source terms, balance laws, well-balanced
19 methods, gravity

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

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

33 Reliable numerical methods for hyperbolic conservation laws with source terms
34 first need to perform well in the homogeneous case. This means for example that
35 they need to cope with discontinuities / weak solutions and with phenomena arising in
36 multiple spatial dimensions, such as involutions and non-trivial stationary states. This
37 requirement has led [ER13, FR15] to suggest *Active Flux*, an extension of the finite
38 volume method. Additionally to the cell average, this scheme evolves point values
39 located at the cell boundary. These are shared among neighbouring cells, which gives
40 rise to a continuous reconstruction. The update of the point values is achieved by using
41 an evolution operator that includes multi-dimensional information. The presence of

*Submitted to the editors DATE.

Funding: WB was supported by the German Academic Exchange Service (DAAD) with funds from the German Federal Ministry of Education and Research (BMBF) and the European Union (FP7-PEOPLE-2013-COFUND – grant agreement no. 605728) as well as by the Deutsche Forschungsgemeinschaft (DFG) through project 429491391 (BA 6878/1-1).

[†]Institute of Mathematics, Zurich University, 8057 Zurich, Switzerland
(wasilij.barsukow@math.uzh.ch).

[‡]Wuerzburg University, Emil-Fischer-Strasse 40, 97074 Wuerzburg, Germany.

42 the point values along the cell boundary then allows to compute the intercell flux via
 43 quadrature. Thus, Active Flux does not use Riemann solvers, while still evolving the
 44 cell average as one of the discrete degrees of freedom just as Finite Volume methods
 45 do. The additional (pointwise) degrees of freedom allow for the scheme to be of high
 46 order of accuracy on a compact stencil. It has been shown in [BHKR19] that this
 47 scheme is stationarity preserving and vorticity preserving for linear acoustics without
 48 any fix. It is third order accurate. Extensions to nonlinear systems have been recently
 49 suggested e.g. in [Fan17, HKS19, Bar21]. Active flux therefore seems to be promising
 50 for resolving many of the structure preservation problems that currently available
 51 methods are facing (an overview of existing methods for balance laws is given below).

52 In view of the many applications that involve source terms, this paper there-
 53 fore aims at deriving the necessary modifications for Active Flux to be applicable
 54 to balance laws while retaining its third order accuracy. Active flux for equations
 55 with a source term was considered in [NR16], where for stationary problems the
 56 necessary quadratures could be chosen of lower order of accuracy (trapezoidal rule)
 57 than in the original Active Flux method from [ER13] (Simpson’s rule) (see e.g. Eqn.
 58 (32) in [NR16]). For time-dependent problems, in [NR16] the reduced order of ac-
 59 curacy of these quadratures is remedied by using a high-order implicit time stepping
 60 method. The approach of the present work avoids sub-iterations and multi-step time
 61 integrators, and the high order in time is achieved through the choice of high order
 62 quadratures, that hardly entail any computational cost. Contrary to [NR16], this
 63 paper presents a fully explicit method for hyperbolic problems with source terms that
 64 reverts to the original Active Flux scheme of [ER11] when the source term vanishes.
 65 As we aim at resolving the acoustic time scale, explicit time stepping is very efficient.

66 Including the source term requires a number of modifications. The homogeneous
 67 part of the equations therefore is for simplicity assumed to be a linear hyperbolic
 68 system for which characteristics are available. This is the case for scalar equations
 69 in any number of spatial dimensions and for systems in one spatial dimension. For
 70 multi-dimensional systems, the concept of characteristics needs to be replaced by
 71 characteristics cones. In the homogeneous case, Active Flux has been used for this
 72 situation as well ([ER13, BHKR19]), but an extension to inhomogeneous systems in
 73 multi-d, and to nonlinear systems remains subject of future work. To a large part, the
 74 strategies presented in this paper will, however, remain valid when the homogeneous
 75 part of the equations is nonlinear as well, and even for nonlinear multi-dimensional
 76 systems.

77 As soon as a source term is added to a hyperbolic system, new stationary states
 78 arise which often are of particular interest. The stationarity is due to the flux di-
 79 vergence being equal to the source term. Many areas of application of balance laws
 80 involve studies of dynamics on top of such an equilibrium (e.g. astrophysics, meteorol-
 81 ogy, tsunami modeling, ...). This requires the numerical method to be very accurate
 82 on the stationary states in order to avoid spurious, artificial perturbations. Therefore
 83 the error of a numerical solution representing one of those stationary states should
 84 not increase with time, thus allowing the simulation to run for a long time (see e.g.
 85 the review [EHB⁺21]).

86 Numerical methods which achieve this are called *well-balanced*, introduced in
 87 [GL96]. They make sure that the discretization of the flux divergence and the dis-
 88 cretization of the source term match, and that the numerical method keeps the de-
 89 sired stationary state exactly stationary for any resolution of the grid. The concept
 90 of well-balanced methods has been extensively used in the context of shallow water
 91 equations with non-flat bottom topography (e.g. [ABB⁺04, BV94, LeV98] and refer-

ences 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 [Bar19], 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 [Bar19]. 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 underdetermined 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

142 the average. Here, the focus lies on linear systems of equations with possibly nonlinear
 143 source terms in one spatial dimension and on linear advection in multiple spatial
 144 dimensions. Section 5 discusses well-balancing of Active Flux for linear acoustics
 145 with gravity. Section 6 finally demonstrates numerically that the new method attains
 146 third order accuracy with linear and nonlinear source terms, can be used to compute
 147 Riemann problems, and displays well-balanced behavior for stationary states.

148 This work can be seen in the larger context of the quest for structure preserving
 149 numerical methods, of which well-balanced methods form an example. Extending
 150 these results to nonlinear hyperbolic equations with source terms and thus combining
 151 the structure preserving properties of Active Flux remains subject of future work.
 152 However, the procedures suggested in this paper are formulated with as little reference
 153 to the linearity of the equations as possible.

154 **2. The Active Flux scheme.** Consider the initial value problem for an $m \times m$
 155 system of hyperbolic balance laws in d spatial dimensions¹

$$156 \quad (2.1) \quad \partial_t q + \nabla \cdot \mathbf{f}(q) = s(q) \quad q : \mathbb{R}_0^+ \times \mathbb{R}^d \rightarrow \mathbb{R}^m, f, s : \mathbb{R}^m \rightarrow \mathbb{R}^m$$

$$157 \quad (2.2) \quad q(0, \mathbf{x}) := q_0(\mathbf{x})$$

159 This section reviews the general idea of the Active Flux scheme. Instead of introducing
 160 jumps at every cell interface, as is customary for finite volume schemes, Active Flux
 161 employs a continuous reconstruction and evolves point values at the cell interfaces
 162 independently. Given these point values, the update of the cell average is immediately
 163 possible by performing flux quadrature in time and along the cell interface. The
 164 distribution of degrees of freedom is discussed in section 2.1, and the update of the
 165 average in section 2.2. What remains, is the update of the point values. To this end,
 166 an IVP is solved (approximately) with the initial data given by the globally continuous
 167 reconstruction. This is very different from the usual approach of finite volume schemes
 168 and is described in sections 2.3–2.4. Some of the details of the (approximate) evolution
 169 operator then depend on the particular equation that is to be solved. After the general
 170 concept is outlined, the details that make it applicable to hyperbolic balance laws are
 171 discussed in sections 3 and 4.

172 **2.1. Degrees of freedom in the Active Flux scheme.** The Active Flux
 173 scheme ([ER13, BHKR19], first introduced in [vL77]) is an extension of the finite
 174 volume scheme. The Active Flux scheme evolves both the cell average and point
 175 values which are distributed along the cell boundary. In particular, here the following
 176 two choices are considered (see Figure 1):

- 177 • In one spatial dimension, there is a point value $q_{i+\frac{1}{2}}$ located at each cell
 178 interface $x_{i+\frac{1}{2}}$. Thus every cell has access to one cell average \bar{q}_i and two
 179 point values at its interfaces.
- 180 • On Cartesian grids in two spatial dimensions, there is a point value $q_{i+\frac{1}{2},j}$,
 181 $q_{i,j+\frac{1}{2}}$ at each edge midpoint and one at each node $q_{i+\frac{1}{2},j+\frac{1}{2}}$. Every cell has
 182 access to one cell average \bar{q}_{ij} and 8 point values distributed along the cell
 183 interface.

184 Note that the point values at cell interfaces are shared by the adjacent cells. As
 185 will be seen in the following, the reconstruction is globally continuous and no Riemann
 186 Problems arise. In one spatial dimension, on average there are 2 degrees of freedom
 187 per cell: 1 cell average and 2 interface values shared each by 2 cells. In two spatial

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

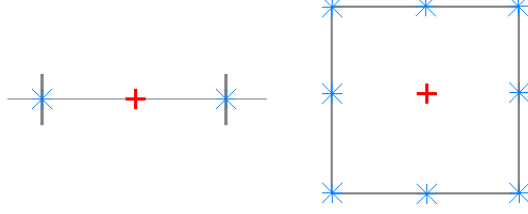


FIG. 1. *The degrees of freedom used for Active Flux. Stars indicate the location of point values, and the cross (placed in the center symbolically) refers to the cell average. Left: One spatial dimension. Right: Two spatial dimensions.*

188 dimensions in the setup as described above there are 4 degrees of freedom per cell: 1
 189 cell average, 4 edge values, each shared by two cells and 4 node values each shared by
 190 4 cells.

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

196 **2.2. Update of the cell average.** As the Active Flux scheme is an extension
 197 of the finite volume scheme, given a time-step-average of the flux through the cell
 198 interface, the update of the average happens in the same way as for finite volume
 199 schemes. As there is a point value located at the cell interface, a Riemann Solver is
 200 not required to obtain the flux. In this section, this finite volume aspect of Active
 201 Flux is described in an arbitrary number of spatial dimensions.

202 Consider the computational domain to be subdivided into polygonal computa-
 203 tional cells. Upon integration of (2.1) over one time step $[t^n, t^n + \Delta t]$ and over
 204 one computational cell \mathcal{C} one obtains an evolution equation for the cell average
 205 $\bar{q}_{\mathcal{C}} := \frac{1}{|\mathcal{C}|} \int_{\mathcal{C}} d\mathbf{x} q(t, \mathbf{x})$:

$$206 \quad \frac{\bar{q}_{\mathcal{C}}^{n+1} - \bar{q}_{\mathcal{C}}^n}{\Delta t} + \frac{1}{|\mathcal{C}|} \frac{1}{\Delta t} \int_{t^n}^{t^n + \Delta t} dt \int_{\partial \mathcal{C}} d\sigma \mathbf{n} \cdot \mathbf{f}(q(t, \mathbf{x})) =$$

$$207 \quad \frac{1}{\Delta t} \int_{t^n}^{t^n + \Delta t} dt \frac{1}{|\mathcal{C}|} \int_{\mathcal{C}} d\mathbf{x} s(q(t, \mathbf{x}))$$

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

$$212 \quad \frac{\bar{q}_{\mathcal{C}}^{n+1} - \bar{q}_{\mathcal{C}}^n}{\Delta t} + \frac{1}{|\mathcal{C}|} \frac{1}{\Delta t} \int_{t^n}^{t^n + \Delta t} dt \sum_{e \subset \partial \mathcal{C}} \int_e d\sigma \mathbf{n}_e \cdot \mathbf{f}(q(t, \mathbf{x})) =$$

$$213 \quad \frac{1}{\Delta t} \int_{t^n}^{t^n + \Delta t} dt \frac{1}{|\mathcal{C}|} \int_{\mathcal{C}} d\mathbf{x} s(q(t, \mathbf{x}))$$

214

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

$$218 \quad (2.3) \quad \frac{\bar{q}_C^{n+1} - \bar{q}_C^n}{\Delta t} + \frac{1}{|\mathcal{C}|} \sum_{e \in \partial \mathcal{C}} |e| \hat{f}_e = \hat{s}_C$$

219
 220 with

$$221 \quad (2.4) \quad \hat{f}_e \simeq \frac{1}{\Delta t} \int_{t^n}^{t^n + \Delta t} dt \frac{1}{|e|} \int_e d\sigma \mathbf{n}_e \cdot \mathbf{f}(q(t, \mathbf{x}))$$

$$222 \quad (2.5) \quad \hat{s}_C \simeq \frac{1}{\Delta t} \int_{t^n}^{t^n + \Delta t} dt \frac{1}{|\mathcal{C}|} \int_C d\mathbf{x} s(q(t, \mathbf{x}))$$

223

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

229 The Active Flux scheme does not need this. Indeed, the point values along the
 230 boundary can be used to immediately approximate (2.4)–(2.5) by quadrature. The
 231 desired properties (most importantly the desired order of accuracy) of the resulting
 232 scheme dictate the number of point values along each edge and also the points in time
 233 at which these point values need to be available.

234 The source term also contributes to the update of the cell average. The quadrature
 235 necessary to approximate the source term average (2.5) to sufficient order in space
 236 and time is suggested in this paper for the first time and discussed in section 4.

237 **2.3. Update of the point values.** The cell average update, and in particular
 238 the computation of the intercell fluxes, requires accurate point values at the cell
 239 boundary to be available.

240 First consider the case where the source term vanishes: $s = 0$. For third order of
 241 accuracy, the integrals in (2.4) need to be approximated by Simpson’s rule. For the
 242 integration in space this can easily be achieved using the available point values at each
 243 cell interface as described in section 2.1. For the integration in time all point values
 244 need to be available at $t^n, t^n + \frac{\Delta t}{2}$ and $t^n + \Delta t$. Altogether this yields a space-time
 245 Simpson rule.

246 In order to obtain sufficiently accurate time evolved point values, in [vL77] it has
 247 been suggested to reconstruct the data and to use an exact evolution operator. An
 248 exact evolution operator generally is unavailable for nonlinear problems, and there-
 249 fore in [Fan17, HKS19, Bar21] approximate evolution operators have been proposed.
 250 Even for linear systems of hyperbolic balance laws it is generally very difficult to ob-
 251 tain closed-form exact evolution operators, as is shown in section 3.2. Therefore the
 252 point values in the Active Flux scheme shall be evolved using a sufficiently high order
 253 *approximate* evolution operator applied to a reconstruction of the discrete data. An
 254 exact evolution operator provides the necessary unwinding in order to guarantee sta-
 255 bility, and an approximate evolution operator needs to do the same. The approximate
 256 evolution operator is introduced in section 3.3.

257 **2.4. Reconstruction.** The reconstruction shall interpolate the point values and
 258 its average over the computational cell shall match the given cell average. In the
 259 following, to simplify notation, in one spatial dimension a uniform grid is assumed,
 260 although the reconstruction can immediately be generalized to nonuniform grids. In
 261 two spatial dimensions, a Cartesian grid is used. See [ER13] for reconstruction on
 262 triangular grids. As mentioned in section 2.1, in one spatial dimension every cell has
 263 access to 3 degrees of freedom which makes a parabolic reconstruction natural. With
 264 the above-mentioned setup it is unique and reads ([vL77, FR15])

$$\begin{aligned}
 265 \quad (2.6) \quad q_{\text{recon},i}(x) &= -3(2\bar{q}_i - q_{i-\frac{1}{2}} - q_{i+\frac{1}{2}}) \frac{(x - x_i)^2}{\Delta x^2} \\
 266 \quad (2.7) \quad &+ (q_{i+\frac{1}{2}} - q_{i-\frac{1}{2}}) \frac{x - x_i}{\Delta x} + \frac{6\bar{q}_i - q_{i-\frac{1}{2}} - q_{i+\frac{1}{2}}}{4} \quad x \in [x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}] \\
 267
 \end{aligned}$$

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

$$\begin{aligned}
 270 \quad (2.8) \quad q_{\text{recon},ij}(\xi\Delta x, \eta\Delta y) &:= \frac{9}{4}\bar{q}_{ij}(-1 + 4\xi^2)(-1 + 4\eta^2) \\
 &- \frac{1}{4}q_W(-1 - 4\xi + 12\xi^2)(-1 + 4\eta^2) \\
 &- \frac{1}{4}q_E(-1 + 4\xi + 12\xi^2)(-1 + 4\eta^2) \\
 &- \frac{1}{4}q_S(-1 + 4\xi^2)(-1 - 4\eta + 12\eta^2) \\
 &- \frac{1}{4}q_N(-1 + 4\xi^2)(-1 + 4\eta + 12\eta^2) \\
 &+ \frac{1}{16}q_{SW}(-1 + 2\xi)(-1 + 2\eta)(-1 - 2\eta + 2\xi(-1 + 6\eta)) \\
 &+ \frac{1}{16}q_{SE}(1 + 2\xi)(-1 + 2\eta)(1 + 2\eta + 2\xi(-1 + 6\eta)) \\
 &+ \frac{1}{16}q_{NW}(-1 + 2\xi)(1 + 2\eta)(1 - 2\eta + 2\xi(1 + 6\eta)) \\
 271 \quad &+ \frac{1}{16}q_{NE}(1 + 2\xi)(1 + 2\eta)(-1 + 2\eta + 2\xi(1 + 6\eta))
 \end{aligned}$$

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

$$273 \quad (2.9) \quad q_{NE} = q_{i+\frac{1}{2},j+\frac{1}{2}} \quad q_{NW} = q_{i-\frac{1}{2},j+\frac{1}{2}} \quad q_{SW} = q_{i-\frac{1}{2},j-\frac{1}{2}} \quad q_{SE} = q_{i+\frac{1}{2},j-\frac{1}{2}}$$

$$274 \quad (2.10) \quad q_N = q_{i,j+\frac{1}{2}} \quad q_S = q_{i,j-\frac{1}{2}} \quad q_E = q_{i+\frac{1}{2},j} \quad q_W = q_{i-\frac{1}{2},j} \\
 275$$

276 Note that both reconstructions are globally continuous and no Riemann Problems
 277 are introduced. The reconstruction, however, is generally not continuously differen-
 278 tiable at the cell interfaces.

279 **2.5. Overview of the algorithm.** The overall algorithm of Active Flux is as
 280 follows:

- 281 1. Given cell averages and point values, compute a reconstruction according to
 282 section 2.4.
- 283 2. Use the reconstruction as initial data in the update of the point values. The
 284 choices of evolution operators considered so far are discussed in section 2.3
 285 and evolution operators in presence of source terms are suggested in section
 286 3.3 below.

- 287 3. Given the updated point values along the cell interfaces, compute the inter-
 288 cell fluxes via quadrature (sections 2.2 and 4 for the homogeneous and the
 289 inhomogeneous cases, respectively).
 290 4. Update the cell averages via (2.3).

291 The computations performed in the Active Flux algorithm are similar in structure
 292 and amount to high order Finite Volume methods, leading to similar time consump-
 293 tion in practice. The latter require a repeated evaluation of the reconstruction and
 294 of the numerical flux function for the individual steps of a time integrator (e.g. a
 295 Runge-Kutta method), while Active Flux performs several evaluations of the evolu-
 296 tion operator to compute values for the flux quadrature in time (without recomputing
 297 the reconstruction). The shared degrees of freedom lead to lower memory usage when
 298 compared to e.g. Discontinuous Galerkin (DG) methods.

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

$$303 \quad (2.11) \quad \Delta t \leq \frac{L_{\min}}{\lambda_{\max}}$$

304 where $L_{\min} = \Delta x$ in one spatial dimension, and $L_{\min} = \frac{1}{2} \min(\Delta x, \Delta y)$ in two spatial
 305 dimensions, when the point values are distributed as described in section 2.1. We
 306 introduce the CFL number as $\Delta t \lambda_{\max} / L_{\min}$.
 307

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

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

318 **3.1. Linear advection with a source term in multiple spatial dimen-**
 319 **sions.** Consider a scalar equation ($m = 1$) and $\mathbf{f}(q) = \mathbf{U}q$ with $\mathbf{U} \in \mathbb{R}^d$. Then

$$320 \quad (3.1) \quad \partial_t q + \mathbf{U} \cdot \nabla q = s(q)$$

322 amounts to the ODE

$$323 \quad (3.2) \quad \frac{d}{dt} q = s(q)$$

324 along the straight characteristic of velocity \mathbf{U} . This ODE can be easily solved ana-
 325 lytically:
 326

$$327 \quad (3.3) \quad \int_{q_0(\mathbf{x} - \mathbf{U}t)}^{q(t, \mathbf{x})} \frac{dp}{s(p)} = t$$

328 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

$$330 \quad (3.4) \quad q(t, \mathbf{x}) = q_0(\mathbf{x} - \mathbf{U}t) \exp(\kappa t)$$

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

$$333 \quad (3.5) \quad q(t, \mathbf{x}) = \left((q_0(\mathbf{x} - \mathbf{U}t))^{1-B} + (1-B)\kappa t \right)^{\frac{1}{1-B}}$$

335 **3.2. Linear acoustics with gravity in one spatial dimension.** This section
 336 has threefold purpose. First, it introduces the acoustic equations with a gravity
 337 source term, which form a very useful system for the study of structure preservation
 338 of numerical methods. This is the set of equations for which a well-balanced method
 339 is derived in 5. This section also demonstrates the difficulties of finding an exact
 340 solution to an inhomogeneous system even if it is linear. Finally, the exact solution
 341 derived here is used later in order to assess the accuracy of the numerical method.

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

$$344 \quad (3.6) \quad \partial_t \rho + \partial_x v = 0$$

$$345 \quad (3.7) \quad \partial_t v + \partial_x p = \rho g \quad g \in \mathbb{R}$$

$$346 \quad (3.8) \quad \partial_t p + c^2 \partial_x v = 0$$

348 The corresponding homogeneous problem (linear acoustics) is the linearization of
 349 the Euler equations around the background state of constant density $\rho_{\text{bg}} = 1$, constant
 350 pressure p_{bg} and vanishing velocity. Then the speed of sound $c = \sqrt{\frac{\gamma p_{\text{bg}}}{\rho_{\text{bg}}}}$ is a constant
 351 ($\mathbb{R} \ni \gamma > 1$). The full system (3.6)–(3.8) can be understood as a particular kind of a
 352 linearization of the Euler equations with gravity²

$$353 \quad (3.9) \quad \partial_t \rho + \partial_x(\rho v) = 0$$

$$354 \quad (3.10) \quad \partial_t(\rho v) + \partial_x(\rho v^2 + p) = \rho g$$

$$355 \quad (3.11) \quad \partial_t e + \partial_x(v(e + p)) = 0$$

$$356 \quad (3.12) \quad e = \frac{p}{\gamma - 1} + \frac{1}{2}\rho v^2 - \rho g x$$

358 The static (stationary and $v = 0$) states of (3.9)–(3.11) are governed by $\partial_x p = \rho g$.
 359 This equation can only be solved if e.g. ρ is given as a function of x , or if another
 360 relation is provided between any two of the variables p, ρ, e . This multitude of possible
 361 stationary states is reflected in the linearization (3.6)–(3.8). (This is the reason for
 362 this particular choice of a linearization.) Observe that stationary states of (3.6)–(3.8)
 363 also are governed by $\partial_x p = \rho g$ and that p can only be computed if ρ is given as a
 364 function of x , or if an additional relation is provided that links ρ and p . This is an
 365 example of a so-called non-trivial stationary state as introduced in [Bar19]. Examples
 366 of stationarity preserving schemes for (3.6)–(3.8) have been discussed in [Bar18].

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

²Note 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.

372 **3.3. Runge-Kutta method for linear systems with a source.** Consider an
 373 $m \times m$ linear system in characteristic variables:

$$374 \quad (3.13) \quad (\partial_t + \lambda_\ell \partial_x) Q_\ell = S_\ell(Q_1, \dots, Q_m) \quad \ell = 1, \dots, m$$

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

378 Recall the following theorem from [Bar21]:

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

$$382 \quad (3.14) \quad Q^{\text{approx}}(t, x) = Q(t, x) + \mathcal{O}(t^3)$$

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

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

$$389 \quad (3.15) \quad Q_\ell(t, x) \simeq Q_{\ell,0}(x - \lambda_\ell t) + t S_\ell(Q_{1,0}(x - \lambda_\ell t), \dots, Q_{m,0}(x - \lambda_\ell t)) \quad \ell = 1, \dots, m$$

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

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

$$393 \quad (3.16) \quad \dot{q}(t) = s(t, q(t)) \quad q : \mathbb{R}_0^+ \rightarrow \mathbb{R}$$

$$396 \quad (3.17) \quad q^{(1)}(\alpha t) = q(0) + \alpha t s(0, q(0))$$

$$397 \quad (3.18) \quad q(t) = q(0) + t \left(1 - \frac{1}{2\alpha}\right) s(0, q(0)) + t \frac{1}{2\alpha} s(\alpha t, q^{(1)}(\alpha t)) + \mathcal{O}(t^3)$$

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

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

$$403 \quad (3.19) \quad \xi_{\ell k} := x - \lambda_\ell t(1 - \alpha) - \lambda_k \alpha t$$

$$404 \quad (3.20) \quad Q_{k\ell}^* := Q_{k,0}(\xi_{\ell k}) + \alpha t S_k(Q_{1,0}(\xi_{\ell k}), \dots, Q_{m,0}(\xi_{\ell k})) \quad k, \ell = 1, \dots, m$$

406 and

$$407 \quad (3.21) \quad Q_\ell^{(1)}(t, x) := Q_{\ell,0}(x - \lambda_\ell t) + \left(1 - \frac{1}{2\alpha}\right) S_\ell(Q_{1,0}(x - \lambda_\ell t), \dots, Q_{m,0}(x - \lambda_\ell t))t$$

$$408 \quad (3.22) \quad + \frac{t}{2\alpha} S_\ell(Q_{1\ell}^*, \dots, Q_{m\ell}^*) \quad \ell = 1, \dots, m$$

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

$$411 \quad (3.23) \quad Q_\ell^{(1)}(t, x) = Q_\ell(t, x) + \mathcal{O}(t^3) \quad \ell = 1, \dots, m$$

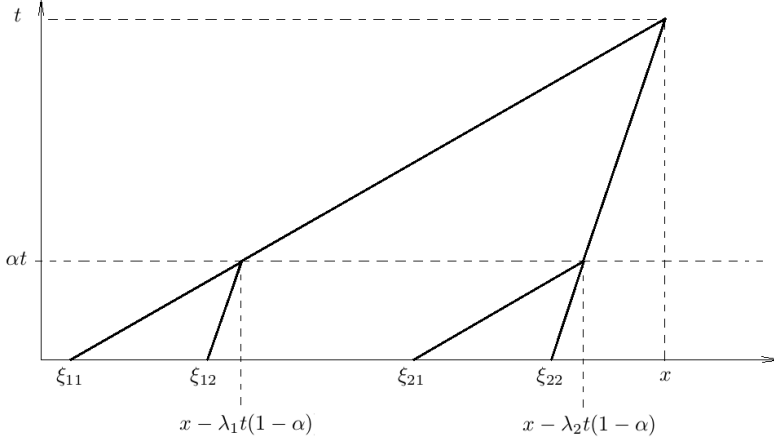


FIG. 2. Illustration of the intermediate solutions and the involved characteristics for the first step in the Runge-Kutta scheme.

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

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

$$416 \quad (3.24) \quad Q_{\ell}(t, x) = Q_{\ell,0}(x) + t \partial_t Q_{\ell} \Big|_{t=0} + \frac{t^2}{2} \partial_t^2 Q_{\ell} \Big|_{t=0} + \mathcal{O}(t^3)$$

$$417 \quad (3.25) \quad = Q_{\ell,0}(x) + t(S_{\ell,0} - \lambda_{\ell} \partial_x Q_{\ell,0})$$

$$418 \quad (3.26) \quad + \frac{t^2}{2} \left(\sum_k \frac{\partial S_{\ell}}{\partial Q_k} (S_{k,0} - (\lambda_k + \lambda_{\ell}) \partial_x Q_{k,0}) + \lambda_{\ell}^2 \partial_x^2 Q_{\ell,0} \right) + \mathcal{O}(t^3)$$

419

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

$$421 \quad (3.27) \quad S_{\ell,0} := S_{\ell}(Q_{1,0}(x), \dots, Q_{m,0}(x))$$

422

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

426 Expand now (3.22) ($\ell = 1, \dots, m$):

$$427 \quad (3.28) \quad \partial_t Q_{k\ell}^* \Big|_{t=0} = -(\lambda_\ell(1-\alpha) + \lambda_k\alpha)\partial_x Q_{k,0} + \alpha S_{k,0}$$

$$428 \quad (3.29) \quad \partial_t Q_\ell^{(1)}(t, x) = -\lambda_\ell \partial_x Q_{\ell,0}(x - \lambda_\ell t)$$

$$429 \quad (3.30) \quad + \left(1 - \frac{1}{2\alpha}\right) \left(t \sum_k \frac{\partial S_\ell}{\partial Q_k} \partial_x Q_{k,0}(x - \lambda_\ell t)(-\lambda_\ell) \right.$$

$$430 \quad (3.31) \quad \left. + S_\ell(Q_{1,0}(x - \lambda_\ell t), \dots, Q_{m,0}(x - \lambda_\ell t)) \right)$$

$$431 \quad (3.32) \quad + \frac{1}{2\alpha} \left(t \sum_k \frac{\partial S_\ell}{\partial Q_k} \partial_t Q_{k\ell}^* + S_\ell(Q_{1\ell}^*, \dots, Q_{m\ell}^*) \right)$$

$$432 \quad (3.33) \quad \stackrel{t=0}{=} -\lambda_\ell \partial_x Q_{\ell,0} + S_{\ell,0}$$

$$433 \quad (3.34) \quad \partial_t^2 Q_\ell^{(1)}(t, x) \Big|_{t=0} = \lambda_\ell^2 \partial_x^2 Q_{\ell,0} + \left(1 - \frac{1}{2\alpha}\right) \left(2 \sum_k \frac{\partial S_\ell}{\partial Q_k} \partial_x Q_{k,0}(-\lambda_\ell) \right)$$

$$434 \quad (3.35) \quad + \frac{1}{2\alpha} \left(2 \sum_k \frac{\partial S_\ell}{\partial Q_k} \partial_t Q_{k\ell}^* \Big|_{t=0} \right)$$

$$435 \quad (3.36) \quad = \lambda_\ell^2 \partial_x^2 Q_{\ell,0} - \sum_k \frac{\partial S_\ell}{\partial Q_k} \left(\partial_x Q_{k,0}(\lambda_\ell + \lambda_k) - S_{k,0} \right) \quad \square$$

436

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

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

$$439 \quad (3.37) \quad \xi_{\ell j} := x - (\lambda_\ell + \lambda_j) \frac{t}{2}$$

$$440 \quad (3.38) \quad Q_{k\ell}^* := Q_{k,0}(\xi_{k\ell}) + \frac{t}{2} S_k(Q_{1,0}(\xi_{k\ell}), \dots, Q_{m,0}(\xi_{k\ell}))$$

$$441 \quad (3.39) \quad Q_\ell^{(1)}(t, x) := Q_{\ell,0}(x - \lambda_\ell t) + t S_\ell(Q_{1\ell}^*, \dots, Q_{m\ell}^*)$$

442

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

$$445 \quad (3.40) \quad (\partial_t + \lambda \partial_x) Q = S(Q)$$

447 *the algorithm reads*

$$448 \quad (3.41) \quad \xi := x - \lambda t$$

450 *and*

$$451 \quad (3.42) \quad Q^{(1)}(t, x) := Q_0(x - \lambda t) + \left(1 - \frac{1}{2\alpha}\right) S(Q_0(x - \lambda t))t$$

$$452 \quad (3.43) \quad + \frac{t}{2\alpha} S\left(Q_0(\xi) + \alpha t S(Q_0(\xi))\right)$$

453

454 For the equations (3.6)–(3.8) of linear acoustics with gravity, $\lambda_1 = c = -\lambda_2$, $\lambda_3 =$
455 0 . The characteristic variables are

$$456 \quad (3.44) \quad Q_1 = \frac{p + cv}{2} \quad Q_2 = \frac{p - cv}{2} \quad Q_3 = -\frac{p}{c^2} + \rho$$

457

458 and the gravity source term then is

$$459 \quad (3.45) \quad S_1 = -S_2 = \frac{g}{2c}(Q_1 + Q_2) + \frac{cg}{2}Q_3 \quad S_3 = 0$$

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

471 **4.1. One spatial dimension.** The numerical discretization (2.5)

$$472 \quad (4.1) \quad \hat{s}_c \simeq \frac{1}{\Delta t} \int_{t^n}^{t^n + \Delta t} dt \frac{1}{|\mathcal{C}|} \int_{\mathcal{C}} dx s(q(t, \mathbf{x}))$$

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

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	$q_{i-\frac{1}{2}}^{n+1}$	$q_{i+\frac{1}{2}}^n$
	$x_{i-\frac{1}{2}}$	$x_{i+\frac{1}{2}}$

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

481 **4.1.1. Linear source term.** Consider first a linear source term, i.e. $s'' = 0$.
 482 Such source terms are relevant in practice (e.g. compressible Euler equations with
 483 gravity), and therefore it is worth dealing with them specifically as they allow for a
 484 simpler approach. For linear source it is possible to first find a quadrature for q
 485 and to apply s to the result. In order to find a quadrature formula for q , one needs to
 486 find a space-time polynomial $p(t, x)$ of at least second degree which interpolates the
 487 available 7 data. Integrating this polynomial would yield a quadrature formula for q .
 488 Here we suggest to use

$$489 \quad (4.2) \quad \mathcal{P}(t, x) = (a_0 + a_1x + a_2t + a_3x^2 + a_4xt + a_5t^2) + a_6xt^2$$

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

492 (4.3) $\mathcal{P}(t^{n+1}, x_{i-\frac{1}{2}}) = q_{i-\frac{1}{2}}^{n+1}$ $\mathcal{P}(t^{n+1}, x_{i+\frac{1}{2}}) = q_{i+\frac{1}{2}}^{n+1}$

493 (4.4) $\mathcal{P}(t^{n+\frac{1}{2}}, x_{i-\frac{1}{2}}) = q_{i-\frac{1}{2}}^{n+\frac{1}{2}}$ $\mathcal{P}(t^{n+\frac{1}{2}}, x_{i+\frac{1}{2}}) = q_{i+\frac{1}{2}}^{n+\frac{1}{2}}$

494 (4.5) $\mathcal{P}(t^n, x_{i-\frac{1}{2}}) = q_{i-\frac{1}{2}}^n$ $\int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} dx \mathcal{P}(t^n, x) = q_i^n$ $\mathcal{P}(t^n, x_{i+\frac{1}{2}}) = q_{i+\frac{1}{2}}^n$

495

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

498 (4.6)
$$\frac{1}{\Delta t} \int_0^{\Delta t} dt \frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} dx q(t^n + t, x_i + x) =$$

499
$$\bar{q}_i^n + \frac{1}{12} \left(-5(q_{i-\frac{1}{2}}^n + q_{i+\frac{1}{2}}^n) + q_{i-\frac{1}{2}}^{n+1} + q_{i+\frac{1}{2}}^{n+1} + 4(q_{i-\frac{1}{2}}^{n+\frac{1}{2}} + q_{i+\frac{1}{2}}^{n+\frac{1}{2}}) \right)$$

500 The weights can be depicted as

501

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

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

504 The time levels $(n, n + \frac{1}{2}, n + 1)$ contribute with weights $(\frac{1}{6}, \frac{2}{3}, \frac{1}{6})$, such that this
505 quadrature formula is a modification of Simpson's rule in time. Note that it is not
506 possible to use terms proportional to x^3 , x^2t or t^3 instead of the term xt^2 in the
507 polynomial ansatz, as then the system (4.3)–(4.5) does not admit a solution. In a
508 sense this is therefore the only choice of a simple quadrature formula.

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

511 **4.1.2. Nonlinear source term.** For nonlinear s , the average

512 (4.7)
$$\int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} dx s(q(t^n, x))$$

513

514 in general is different from

515 (4.8)
$$s \left(\int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} dx q(t^n, x) \right)$$

516

517 Point values, however, do not present any difficulties: one can just evaluate s on
518 them. Therefore we suggest to consider a reconstruction $q_{\text{recon},i}(x)$ that interpolates
519 $q_{i-\frac{1}{2}}^n$ and $q_{i+\frac{1}{2}}^n$ and whose average agrees with \bar{q}_i^n . It is computed anyway in order

520 to update the point values in time, see equation (2.7). This reconstruction can be
 521 easily evaluated at the midpoint of the cell. Then, instead of the cell averages, one
 522 works with a seventh point value $q_{\text{recon},i}(0) = \frac{1}{4}(6\bar{q}_i^n - q_{i-\frac{1}{2}}^n - q_{i+\frac{1}{2}}^n)$. Of course, this is
 523 equivalent to replacing the average by a Simpson's rule in the quadrature, and thus
 524 the order of the quadrature is not reduced. Therefore when using only point values
 525 (the 6 pointwise degrees of freedom and one value at the cell midpoint) the weights
 526 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	$-\frac{3}{12}$	$-\frac{3}{12}$
	$x_{i-\frac{1}{2}}$	$x_{i+\frac{1}{2}}$

528 Equation (2.5) then is replaced by the quadrature

$$\begin{aligned}
 529 \quad (4.9) \quad \hat{s}_i &= \frac{s(q_{i-\frac{1}{2}}^{n+1}) + s(q_{i+\frac{1}{2}}^{n+1}) + 4(s(q_{i-\frac{1}{2}}^{n+1}) + s(q_{i+\frac{1}{2}}^{n+1})) - 3(s(q_{i-\frac{1}{2}}^{n+1}) + s(q_{i+\frac{1}{2}}^{n+1})) + 8q_{\text{recon},i}(0)}{12} \\
 530
 \end{aligned}$$

531 This quadrature can now be used for nonlinear s . As (4.9) uses a Simpson quadrature
 532 instead of the average, upon usage of a linear source s , it reduces to the expression
 533 (4.6) because of the quadratic reconstruction.

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

536 4.2. Two spatial dimensions.

537 **4.2.1. Linear source term.** Similarly consider the setup of the Active Flux
 538 method on two-dimensional Cartesian grids as described in 2.1. The available degrees
 539 of freedom are

$$540 \quad (4.10) \quad 3 \times 4 \text{ nodes: } \bar{q}_{i\pm\frac{1}{2},j\pm\frac{1}{2}}^n, \bar{q}_{i\pm\frac{1}{2},j\pm\frac{1}{2}}^{n+\frac{1}{2}}, \bar{q}_{i\pm\frac{1}{2},j\pm\frac{1}{2}}^{n+1}$$

$$541 \quad (4.11) \quad 3 \times 2 \text{ vertical edges: } q_{i\pm\frac{1}{2},j}^n, q_{i\pm\frac{1}{2},j}^{n+\frac{1}{2}}, q_{i\pm\frac{1}{2},j}^{n+1}$$

$$542 \quad (4.12) \quad 3 \times 2 \text{ horizontal edges: } q_{i,j\pm\frac{1}{2}}^n, q_{i,j\pm\frac{1}{2}}^{n+\frac{1}{2}}, q_{i,j\pm\frac{1}{2}}^{n+1}$$

$$543 \quad (4.13) \quad 1 \text{ average: } \bar{q}_{ij}^n$$

545 The ansatz for a space-time polynomial is

$$546 \quad (4.14) \quad \mathcal{P}(t, x, y) = \left(\sum_{\zeta+\eta+\vartheta \leq 4} a_{\zeta\eta\vartheta} \cdot x^\zeta y^\eta t^\vartheta \right) + a_{212} x^2 y t^2 + a_{122} x y^2 t^2$$

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

$$\begin{aligned}
 550 \quad (4.15) \quad & \frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} dx \frac{1}{\Delta y} \int_{-\frac{\Delta y}{2}}^{\frac{\Delta y}{2}} dy \frac{1}{\Delta t} \int_0^{\Delta t} dt q(t, x, y) = \bar{q}_{ij}^n \\
 & - \frac{20}{72} (q_E^n + q_N^n + q_S^n + q_W^n) + \frac{5}{72} (q_{NE}^n + q_{NW}^n + q_{SE}^n + q_{SW}^n) \\
 & + \frac{16}{72} (q_E^{n+\frac{1}{2}} + q_N^{n+\frac{1}{2}} + q_S^{n+\frac{1}{2}} + q_W^{n+\frac{1}{2}}) - \frac{4}{72} (q_{NE}^{n+\frac{1}{2}} + q_{NW}^{n+\frac{1}{2}} + q_{SE}^{n+\frac{1}{2}} + q_{SW}^{n+\frac{1}{2}}) \\
 551 \quad & + \frac{4}{72} (q_E^{n+1} + q_N^{n+1} + q_S^{n+1} + q_W^{n+1}) - \frac{1}{72} (q_{NE}^{n+1} + q_{NW}^{n+1} + q_{SE}^{n+1} + q_{SW}^{n+1})
 \end{aligned}$$

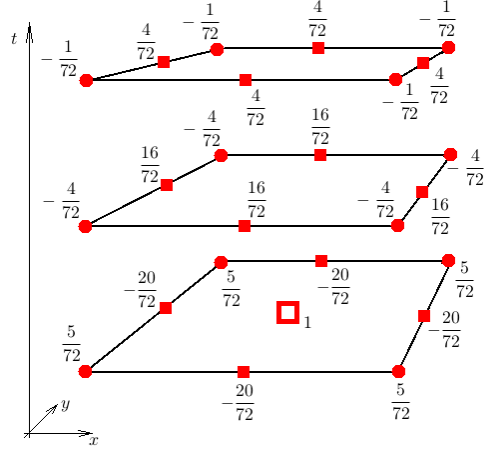


FIG. 3. Illustration of the weights of the space time quadrature formula (4.15).

552 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
 553 always contribute -4 times the nodes.

554 **4.2.2. Nonlinear source term.** Again, for nonlinear source instead of the av-
 555 erage it is necessary to use the evaluation of the reconstruction at the cell midpoint.
 556 This amounts to an approximation of the average by a two-dimensional Simpson rule.
 557 Then the source term is approximating as follows:

$$\begin{aligned}
 & \frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} dx \frac{1}{\Delta y} \int_{-\frac{\Delta y}{2}}^{\frac{\Delta y}{2}} dy \frac{1}{\Delta t} \int_0^{\Delta t} dt s(q(t, x, y)) = \frac{32}{72} s(q_{\text{recon}, ij}(0, 0)) \\
 & - \frac{12}{72} (s(q_E^n) + s(q_N^n) + s(q_S^n) + s(q_W^n)) \\
 & + \frac{7}{72} (s(q_{NE}^n) + s(q_{NW}^n) + s(q_{SE}^n) + s(q_{SW}^n)) \\
 558 \quad (4.16) \quad & + \frac{16}{72} (s(q_E^{n+\frac{1}{2}}) + s(q_N^{n+\frac{1}{2}}) + s(q_S^{n+\frac{1}{2}}) + s(q_W^{n+\frac{1}{2}})) \\
 & - \frac{4}{72} (s(q_{NE}^{n+\frac{1}{2}}) + s(q_{NW}^{n+\frac{1}{2}}) + s(q_{SE}^{n+\frac{1}{2}}) + s(q_{SW}^{n+\frac{1}{2}})) \\
 & + \frac{4}{72} (s(q_E^{n+1}) + s(q_N^{n+1}) + s(q_S^{n+1}) + s(q_W^{n+1})) \\
 559 \quad & - \frac{1}{72} (s(q_{NE}^{n+1}) + s(q_{NW}^{n+1}) + s(q_{SE}^{n+1}) + s(q_{SW}^{n+1}))
 \end{aligned}$$

560 In case that the data only depend on one of the variables, the two-dimensional
 561 quadratures (4.15) and (4.16) do *not* exactly reduce to the one dimensional quadra-
 562 tures (4.6) and (4.9). This is because (cf. Figure 3) the point values on edge midpoints
 563 $(0, \pm \frac{\Delta y}{2})$ do not disappear even if the data depend only on x , and therefore the avail-
 564 able degrees of freedom remain different from the one-dimensional case.

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

566 **5.1. Exact evolution operator.** As described in 3.2 a closed-form exact evo-
 567 lution operator for acoustics with gravity is very difficult to obtain. Nevertheless,

568 it is still possible to show that a scheme endowed with such an operator would be
 569 well-balanced / stationarity preserving; i.e. that there exists a discretization of the
 570 stationary states of the PDE which remain exactly stationary. This proof does not
 571 require the evolution operator to be known explicitly, but only relies on the fact that
 572 it is exact. Besides its fundamental importance, this result is used in section 5.2
 573 to analyze the situation for the approximate evolution operator and to restore the
 574 well-balanced property for it.

575 The numerical stationary states are best studied upon the (discrete) Fourier trans-
 576 form. Define $t_x := \exp(\mathfrak{i}k_x \Delta x)$, $t_y := \exp(\mathfrak{i}k_y \Delta y)$. Here \mathfrak{i} is the imaginary unit and
 577 $\mathbf{k} = (k_x, k_y) \in \mathbb{R}^2$ is the wave vector characterizing the spatial frequency of the Fourier
 578 mode. Applying the Fourier transform introduces one mode \bar{q} for the averages and
 579 one mode q for the point values; this implies writing $q_i := \bar{q} t_x^i t_y^j$, $q_{i+\frac{1}{2}} := q t_x^i t_y^j$.

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

$$582 \quad (5.1) \quad \bar{\rho}_i = \frac{\rho_{i+\frac{1}{2}} + \rho_{i-\frac{1}{2}}}{2}$$

$$583 \quad (5.2) \quad \frac{p_{i+\frac{1}{2}} - p_{i-\frac{1}{2}}}{\Delta x} = g \frac{\rho_{i-\frac{1}{2}} + \rho_{i+\frac{1}{2}}}{2}$$

$$584 \quad (5.3) \quad \frac{\bar{p}_{i+\frac{3}{2}} - \bar{p}_{i+\frac{1}{2}}}{\Delta x} = g \frac{\rho_{i+\frac{3}{2}} + 4\rho_{i+\frac{1}{2}} + \rho_{i-\frac{1}{2}}}{6}$$

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

588 *Proof.* The proof consists of two parts.

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

$$592 \quad (5.4) \quad v_{\text{recon}}(x) = \text{const} \quad \partial_x p_{\text{recon}}(x) = \rho_{\text{recon}}(x)g$$

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

$$595 \quad (5.5) \quad -3 \left(2\bar{p} - p \left(1 + \frac{1}{t_x} \right) \right) \frac{2x}{\Delta x^2} + p \left(1 - \frac{1}{t_x} \right) \frac{1}{\Delta x} =$$

$$596 \quad (5.6) \quad -3g \left(2\bar{p} - \rho \left(1 + \frac{1}{t_x} \right) \right) \frac{x^2}{\Delta x^2} + g\rho \left(1 - \frac{1}{t_x} \right) \frac{x}{\Delta x} + g \frac{6\bar{p} - \rho \left(1 + \frac{1}{t_x} \right)}{4}$$

598 This shall be valid for all x :

$$599 \quad (5.7) \quad 2\bar{p} - \rho(1 + 1/t_x) = 0$$

$$600 \quad (5.8) \quad -2\bar{p}t_x + p(t_x + 1) = \frac{\Delta x g \rho(t_x - 1)}{6}$$

$$601 \quad (5.9) \quad p(t_x - 1) = \Delta x g \frac{6\bar{p}t_x - \rho(t_x + 1)}{4}$$

602

603 These are three equations for four variables. In particular

$$604 \quad (5.10) \quad \bar{\rho} = \frac{\rho(1 + 1/t_x)}{2}$$

$$605 \quad (5.11) \quad p = \Delta x g \rho \frac{t_x + 1}{2(t_x - 1)}$$

$$606 \quad (5.12) \quad \bar{p} = \Delta x g \rho \frac{t_x^2 + 4t_x + 1}{6t_x(t_x - 1)}$$

607

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

$$610 \quad (5.13) \quad \bar{\rho} = \frac{\rho_{i+\frac{1}{2}} + \rho_{i-\frac{1}{2}}}{2}$$

$$611 \quad (5.14) \quad \frac{p_{i+\frac{1}{2}} - p_{i-\frac{1}{2}}}{\Delta x} = g \frac{\rho_{i-\frac{1}{2}} + \rho_{i+\frac{1}{2}}}{2}$$

$$612 \quad (5.15) \quad \frac{\bar{p}_{i+1} - \bar{p}_i}{\Delta x} = g \frac{\rho_{i+\frac{3}{2}} + 4\rho_{i+\frac{1}{2}} + \rho_{i-\frac{1}{2}}}{6}$$

613

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

$$616 \quad (5.16) \quad \frac{\bar{v}^{n+1} - \bar{v}^n}{\Delta t} + \frac{p(1 - 1/t_x)}{\Delta x} = \frac{\bar{v}^{n+1} - \bar{v}^n}{\Delta t} + g\rho \frac{t_x + 1}{2t_x}$$

$$617 \quad (5.17) \quad = \frac{\bar{v}^{n+1} - \bar{v}^n}{\Delta t} + g\bar{\rho} \quad \square$$

618

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

621 The equations (5.10)–(5.12) contain ρ as a free variable. One can rewrite the
622 system making p the free variable:

$$623 \quad (5.18) \quad \bar{\rho} = \frac{p(t_x - 1)}{t_x \Delta x g} \quad \rho = \frac{2p(t_x - 1)}{\Delta x g(t_x + 1)} \quad \bar{p} = p \frac{t_x^2 + 4t_x + 1}{3t_x(t_x + 1)}$$

624

625 This form will be useful later.

626 Equations (5.2)–(5.3) are finite difference approximations of $\partial_x p = \rho g$. By con-
627 struction, the discrete stationary states are those whose reconstruction fulfills (5.4) in
628 every cell. Equation (5.1) implies that the reconstructed ρ of the discrete stationary
629 state is linear, which is clear: for quadratic reconstructions to fulfill (5.4), ρ_{recon} has
630 to be linear in each cell. The slope of the linear function can vary from cell to cell
631 and is given by (5.2).

632 **5.2. Approximate evolution operator.** The above section identifies condi-
633 tions (5.1)–(5.3) on the discrete data for them to remain stationary upon usage of the
634 *exact* evolution operator. Unfortunately, such an operator is unavailable in practice.
635 Having identified an approximate solution operator, which agrees with the exact so-
636 lution up to terms $\mathcal{O}(t^3)$ in section 3.3, here we study whether it keeps the same data
637 (5.1)–(5.3) stationary as well.

638 **THEOREM 5.2.** *If the discrete data fulfill (5.1)–(5.3) and the approximate evolu-*
639 *tion operator of theorem 3.2 for (3.6)–(3.8) is used, then both the pressure p and the*

640 density ρ remain stationary over one time step, but the velocity undergoes the time
641 evolution

$$642 \quad (5.19) \quad v_{i+\frac{1}{2}}(t) = -\frac{\alpha g^2}{4} \frac{\rho_{i+\frac{1}{2}} - \rho_{i-\frac{1}{2}}}{\Delta x} t^3$$

644 *Proof.* Assume the initial data to fulfill (5.1)–(5.3), or equivalently (5.4). Using
645 (2.7) (and applying the discrete Fourier transform straight away) (5.4) implies

(5.20)

$$646 \quad p_{\text{recon}}(x) = \frac{1}{4} \left(6\bar{p} - p \left(1 + \frac{1}{t_x} \right) \right) + \frac{x}{\Delta x} \left(1 - \frac{1}{t_x} \right) p - 3 \frac{x^2}{\Delta x^2} \left(2\bar{p} - p \left(1 + \frac{1}{t_x} \right) \right)$$

(5.21)

$$647 \quad \rho_{\text{recon}}(x) = \frac{1}{g\Delta x} \left(p \left(1 - \frac{1}{t_x} \right) - 6 \frac{x}{\Delta x} \left(2\bar{p} - p \left(1 + \frac{1}{t_x} \right) \right) \right)$$

(5.22)

$$648 \quad v_{\text{recon}}(x) = 0$$

650 and using (3.44) therefore

(5.23)

$$651 \quad Q_{1,0}(x) = Q_{2,0}(x) = -\frac{p(1+t_x) - 6\bar{p}t_x}{8t_x} + \frac{p(t_x-1)x}{2\Delta x t_x} + \frac{3(p(1+t_x) - 2\bar{p}t_x)x^2}{2\Delta x^2 t_x}$$

(5.24)

$$652 \quad Q_{3,0}(x) = \frac{p(-1+t_x)}{\Delta x g t_x} + \frac{p - 6\bar{p}t_x + p t_x}{4c^2 t_x}$$

$$653 \quad + \frac{(-\Delta x g p(t_x-1) + 6c^2(p(1+t_x) - 2\bar{p}t_x))x}{c^2 \Delta x^2 g t_x} - \frac{3(p(1+t_x) - 2\bar{p}t_x)x^2}{c^2 \Delta x^2 t_x}$$

655 Evaluating the Runge-Kutta algorithm of section 3.3 on these initial data (at
656 $x = \frac{\Delta x}{2}$) yields

$$657 \quad (5.25) \quad (\rho, v^*, p)^T \quad \text{with} \quad v^* = -\frac{\alpha g(t_x-1)^2}{2\Delta x^2 t_x(t_x+1)} p t^3$$

659 (α is the parameter appearing in the RK2 method.)

660 Recall that ρ and p are the Fourier coefficients of the point values of the density
661 and the pressure. Obviously ρ and p remain stationary, but the velocity does not.
662 Using (5.18) v^* can be rewritten as

$$663 \quad (5.26) \quad v^* = -\frac{\alpha g^2}{4\Delta x} \left(1 - \frac{1}{t_x} \right) \rho t^3 = -\frac{\alpha g^2}{4} \frac{\rho_{i+\frac{1}{2}} - \rho_{i-\frac{1}{2}}}{\Delta x} t^3$$

665 having applied the inverse Fourier transform in the last step. \square

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

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

$$670 \quad (5.27) \quad \frac{\alpha g^2}{4} \frac{\rho_{i+\frac{1}{2}} - \rho_{i-\frac{1}{2}}}{\Delta x} t^3$$

671 *to the velocity evolution, then*

- 673 *i) its accuracy is not changed*
 674 *ii) it becomes stationarity preserving / well-balanced with the same discrete station-*
 675 *ary states as the exact evolution operator.*

676 The two forms (5.25) and (5.19) of v^* are equivalent, because the initial data
 677 have been chosen to be stationary, and thus additionally fulfill (5.18). The proposed
 678 modification is to *always* add $-v^*$ to the velocity evolution, irrespective of whether
 679 the data fulfill (5.18) or not. At this point the Fourier coefficients of ρ and p are
 680 independent and it matters whether the correction is used in the form (5.25) or (5.19).
 681 Of course, also the inverse Fourier transform has to be applied to the expression first
 682 in order for the correction to attain the form of a finite difference formula. Compact
 683 finite difference formulae are in one-to-one-correspondence with Laurent polynomials
 684 in t_x . An expression such as $\frac{1}{t_x+1} = 1 - t_x + t_x^2 \mp \dots$ is an expression involving an
 685 unbounded stencil and cannot be implemented in usual codes. Therefore (5.19) cannot
 686 be used as a correction because the correction would have a non-compact stencil (just
 687 as the equivalent expressions involving only $\bar{\rho}$ or \bar{p}). This is why the form (5.25) which
 688 involves point values of ρ is preferred.

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

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

699 **6.1. Linear advection.** Consider first

$$700 \quad (6.1) \quad \partial_t q + \mathbf{U} \cdot \nabla q = \kappa q$$

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

708 Consider now

$$709 \quad (6.2) \quad \partial_t q + \mathbf{U} \cdot \nabla q = \kappa q^B \quad B \neq 1$$

711 with the exact solution (3.5) and $\kappa = 7$, $B = 3$. Figure 7 (left) shows the initial
 712 data and the numerical solution, and Figure 7 (right) shows a convergence study for
 713 the approximate evolution operator from Corollary (3.4). One observes third order
 714 accuracy, as expected.

715 **6.2. Acoustics with gravity.** Consider now the equations of linear acoustics
 716 with a gravity source term (3.6)–(3.8). The exact solution operator is only partly avail-
 717 able in closed form, and therefore the approximate Runge-Kutta evolution operator of
 718 section 3.3 is used in combination with the well-balancing fix (5.27). The parameter
 719 α in the Runge-Kutta method is chosen to $\alpha = \frac{1}{2}$ and CFL = 0.9 everywhere.

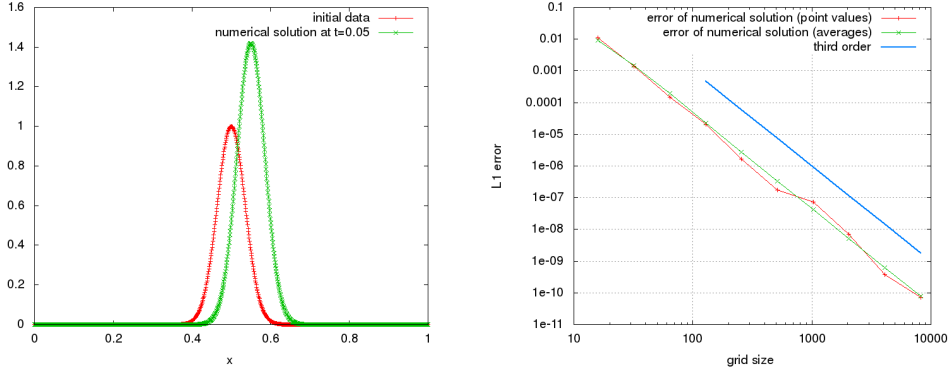


FIG. 4. Gaussian initial data for (6.1) with $\mathbf{U} = \mathbf{e}_x$, $\kappa = 7$. Note that due to the source term, the Gaussian is advected and also changes shape. Exact evolution operator (3.4) and quadrature formula (4.6) have been used with $CFL = 0.9$. Left: Initial data and solution at $t = 0.05$ (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.

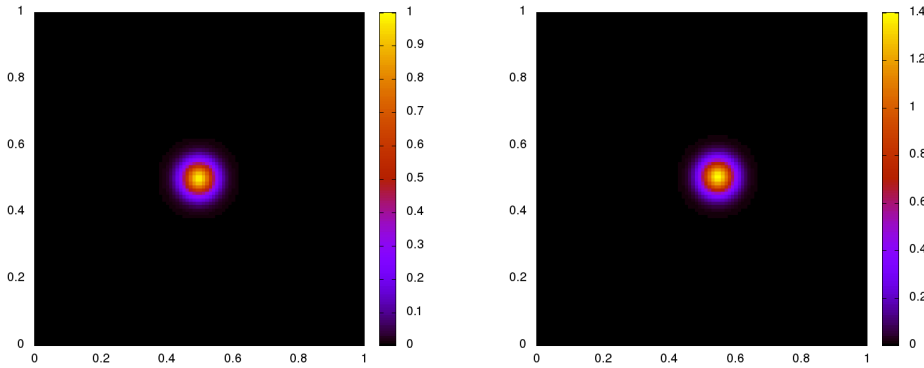


FIG. 5. Gaussian initial data for (6.1) with $\mathbf{U} = (1, 0.1)$, $\kappa = 7$. Note that due to the source term, the Gaussian is advected and also changes shape. Exact evolution operator (3.4) and quadrature formula (4.15) have been used with $CFL = 0.9$. Left: Initial setup. Right: Numerical solution at $t = 0.05$ on a 100×100 Cartesian grid.

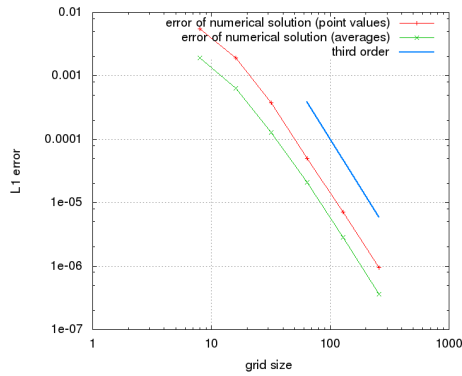


FIG. 6. Convergence study for the setup shown in Figure 5. One observes third order accuracy.

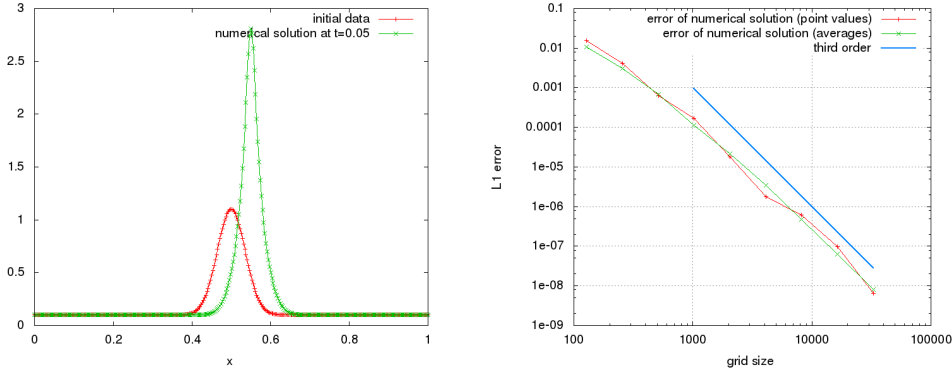


FIG. 7. Gaussian initial data for (6.2) with $s(q) = \kappa q^B$ and $\mathbf{U} = \mathbf{e}_x$, $\kappa = 7$, $B = 3$. Runge-Kutta approximate evolution operator from Corollary 3.4 (with $\alpha = \frac{1}{2}$) and quadrature formula (4.9) have been used with $CFL = 0.9$. The solution has been computed 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 $t = 0.05$ (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. The exact solution is given by (3.5).

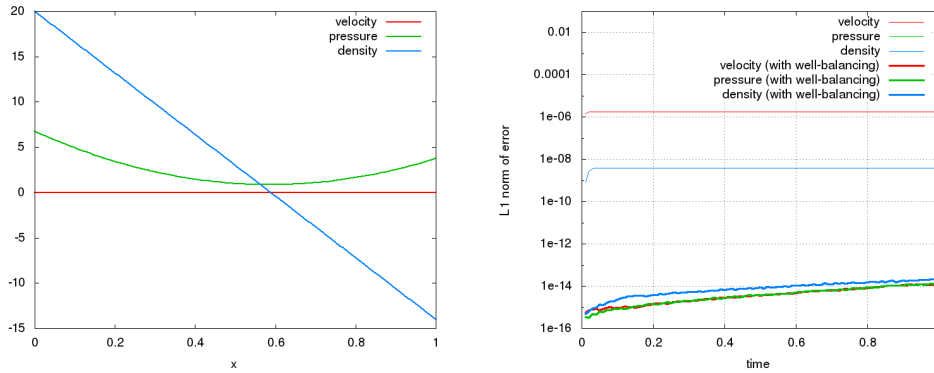


FIG. 8. Setup of a stationary parabola (6.3) for (3.6)–(3.8), solved using the Runge-Kutta approximate evolution operator of section 3.3 with and without well-balancing (5.27). Here $g = -1$, and the setup is solved on a grid covering $[-1.5, 2.5]$, but the error is only measured inside $[0, 1]$ ($\Delta x = 10^{-2}$) to exclude the influence of the boundaries. Left: Setup. Right: Error of numerical solution (point values) 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.

720 Figure 8 shows a stationary setup given by

$$721 \quad (6.3) \quad p = A_1 x^2 + A_2 x + A_3 \quad \rho = 2A_1 x/g + A_2/g \quad v = 0$$

723 with $A_1 = 17$, $A_2 = -3$, $A_3 = 1$. This parabola is exactly recovered by the reconstruc-
724 tion, and thus remains stationary up to machine precision. This experiment shows
725 that the well-balancing fix works as it should.

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

$$727 \quad (6.4) \quad \rho = \left(\frac{g(\gamma - 1)}{K\gamma} x + \rho_0^{\gamma-1} \right)^{\frac{1}{\gamma-1}}$$

728

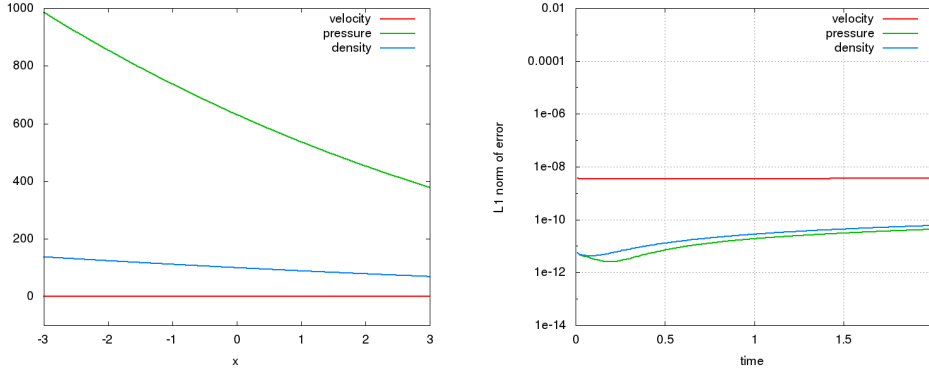


FIG. 9. Stationary setup (6.4) for (3.6)–(3.8), solved using the Runge-Kutta approximate evolution operator of section 3.3 with well-balancing (5.27). Here $g = -1$, and the setup is solved on a grid covering $[-5.5, 5.5]$, but the error is only measured inside $[-3, 3]$ ($\Delta x = 1/300$) to exclude the influence of the boundaries. Left: Setup (cell averages). Right: Error of numerical solution (point values) as a function of time. One observes a transition towards a numerical stationary state which then persists forever.

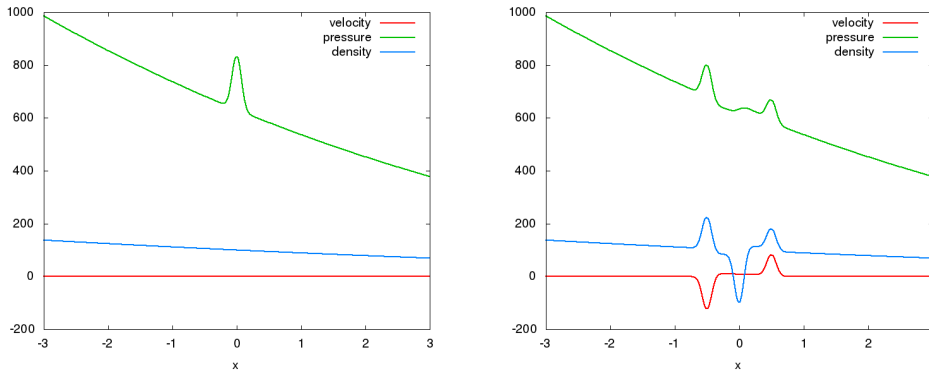


FIG. 10. Setup (6.4) endowed with the pressure perturbation (6.5) solved using the Runge-Kutta approximate evolution operator of section 3.3 with well-balancing (5.27). Left: Initial data (cell averages). Right: Numerical solution (cell averages) at $t = 0.5$ on a grid covering $[-5.5, 5.5]$, but only the subinterval $[-3, 3]$ is considered in order to exclude the influence of the boundaries. $\Delta x = 0.01$, $CFL = 0.9$.

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

733 Next, a perturbation

734 (6.5)
$$200 \exp(-100x^2)$$

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

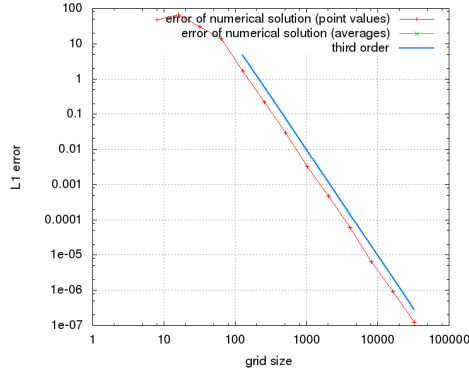


FIG. 11. Setup of Figure 10. The error of the numerical solution is measured on the point values. One observes third order accuracy.

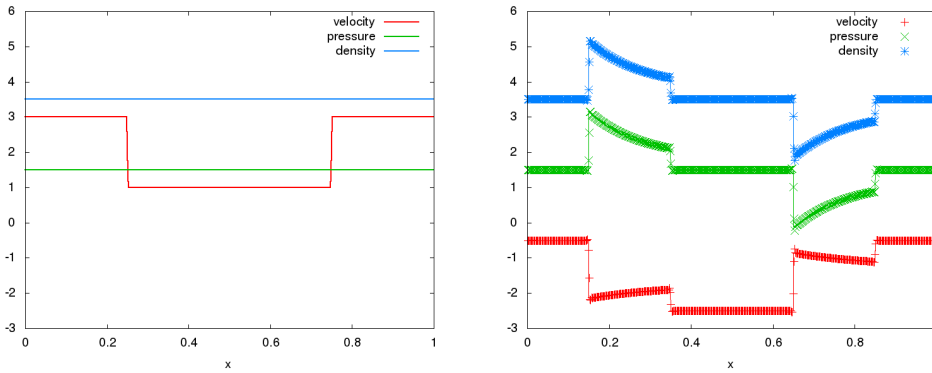


FIG. 12. Riemann problem setup (6.6) solved using the Runge-Kutta approximate evolution operator of section 3.3 with well-balancing (5.27). Here, $g = -10$. Left: Initial data. Right: Numerical solution (dots) and exact solution (solid line) at $t = 0.1$. $\Delta x = 0.01$, $CFL = 0.9$. Averages of the numerical solution are shown.

741 Consider finally a Riemann problem:

$$742 \quad (6.6) \quad \rho = 3.5 \quad p = 1.5 \quad v = \begin{cases} 1 & 0.25 \leq x \leq 0.75 \\ 3 & \text{else} \end{cases}$$

743

744 This Riemann problem can be solved exactly using the formula (A.17)–(A). Note
745 that if all quantities are constant in space, then they solve

$$746 \quad (6.7) \quad \partial_t \rho = 0 \quad \partial_t p = 0 \quad \partial_t v = \rho g$$

748 which means that ρ and p remain stationary, but that $v = v(t=0) + \rho g t$. The solution
749 to the initial data (6.6) therefore can be obtained by adding the time evolution of
750 $(0, v_0(x), 0)^T$ (via numerical quadrature of (A.17)–(A)) and the time evolution of
751 $(\rho, 0, p)^T$ which is just $(\rho, \rho g t, p)^T$. Figure 12 shows the numerical and the exact
752 solution.

753 **7. Conclusions and outlook.** Active flux is a novel kind of numerical method
754 for hyperbolic problems, extending the finite volume method. Instead of computing

755 the intercell flux via a Riemann problem it relies on a continuous reconstruction and
 756 on accurately evolved point values along the cell boundary. They then immediately
 757 serve as quadrature values for the computation of the intercell flux. The extension
 758 of Active Flux to time dependent balance laws presented in this paper requires a
 759 modification in both these aspects: the evolution of the point values and the average
 760 update need to account for the source term. Here, an approximate evolution operator
 761 is suggested for the point value update; this is done for linear systems with possibly
 762 nonlinear source terms in one spatial dimension, and linear scalar equations with
 763 source terms in multiple spatial dimensions. A suitable quadrature is suggested in
 764 order to approximate the contribution of the source term to the cell average. This
 765 quadrature can be applied to any system of (nonlinear) balance laws.

766 We aim at combining the strategy presented in this paper with an approximate
 767 evolution operator for a nonlinear homogeneous problem (such as those suggested
 768 in [Bar21]) in future. Multi-dimensional systems of hyperbolic conservation laws
 769 are very different from their one-dimensional counterparts because in general char-
 770 acteristics are unavailable and need to be conceptually replaced by characteristic
 771 cones. Examples of evolution operators that make use of such cones can be found in
 772 [ER13, FR15, Fan17, BHKR19]. Combining these with an approximate evolution of
 773 the source term shall pave the way towards the extension of Active Flux to nonlinear
 774 multi-dimensional balance laws and the derivation of accurate structure preserving
 775 (in particular well-balanced) methods for them.

776 **Appendix A. Exact solution of linear acoustics with gravity.**

777 System (3.6)–(3.8) can in principle be immediately solved exactly via Fourier
 778 transform by inserting the ansatz

$$\begin{aligned}
 & \begin{pmatrix} \rho \\ v \\ p \end{pmatrix} = \begin{pmatrix} \hat{\rho} \\ \hat{v} \\ \hat{p} \end{pmatrix} \exp(\mathfrak{i}k \cdot x - \mathfrak{i}\omega t) \\
 & \text{(A.1)}
 \end{aligned}$$

781 into (3.6)–(3.8):

$$\begin{aligned}
 & \omega \begin{pmatrix} \hat{\rho} \\ \hat{v} \\ \hat{p} \end{pmatrix} = \begin{pmatrix} 0 & k & 0 \\ \mathfrak{i}g & 0 & k \\ 0 & c^2k & 0 \end{pmatrix} \begin{pmatrix} \hat{\rho} \\ \hat{v} \\ \hat{p} \end{pmatrix} \\
 & \text{(A.2)}
 \end{aligned}$$

784 Therefore $\omega = 0$, or $\omega = \pm\sqrt{c^2k^2 + \mathfrak{i}gk}$. The complex eigenvalue can be removed
 785 upon transforming

$$\begin{aligned}
 & \rho = \tilde{\rho}e^{\mu x} \qquad v = \tilde{v}e^{\mu x} \qquad p = \tilde{p}e^{\mu x} \\
 & \text{(A.3)}
 \end{aligned}$$

788 with

$$\begin{aligned}
 & \mu := \frac{g}{2c^2} \\
 & \text{(A.4)}
 \end{aligned}$$

791 System (3.6)–(3.8) then reads

$$\begin{aligned}
 & \partial_t \tilde{\rho} + \partial_x \tilde{v} = -\mu \tilde{v} \\
 & \text{(A.5)}
 \end{aligned}$$

$$\begin{aligned}
 & \partial_t \tilde{v} + \partial_x \tilde{p} = \tilde{\rho}g - \mu \tilde{p} \\
 & \text{(A.6)}
 \end{aligned}$$

$$\begin{aligned}
 & \partial_t \tilde{p} + c^2 \partial_x \tilde{v} = -c^2 \mu \tilde{v} \\
 & \text{(A.7)}
 \end{aligned}$$

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

$$798 \quad (A.8) \quad \omega \begin{pmatrix} \hat{\rho} \\ \hat{v} \\ \hat{p} \end{pmatrix} = \mathcal{E} \begin{pmatrix} \hat{\rho} \\ \hat{v} \\ \hat{p} \end{pmatrix} \quad \mathcal{E} = \begin{pmatrix} 0 & k - \mathfrak{i}\mu & 0 \\ \mathfrak{i}g & 0 & k - \mathfrak{i}\mu \\ 0 & c^2k - \mathfrak{i}c^2\mu & 0 \end{pmatrix}$$

800 The eigenvalues of \mathcal{E} are now real: $\omega_1 = 0$, $\omega_{2,3} = \pm c\sqrt{k^2 + \mu^2}$. Although this
 801 transformation brings the endeavour of finding the exact solution to (3.6)–(3.8) into
 802 the realm of the possible, technical difficulties prevent one from actually computing
 803 all Green's functions in closed form.

804 Assume therefore that the only non-vanishing initial data are in the velocity.
 805 Then the Fourier mode at initial time reads

$$806 \quad (A.9) \quad (0, \hat{v}, 0)^T \exp(\mathfrak{i}kx)$$

808 and at a later time it becomes

$$809 \quad (A.10) \quad \sum_{m=1}^3 v_m \exp(\mathfrak{i}kx - \mathfrak{i}\omega_m t)$$

811 where the decomposition of $(0, \hat{v}, 0)^T$ in the eigenbasis of \mathcal{E} is used, i.e.

$$812 \quad (A.11) \quad (0, \hat{v}, 0)^T = \sum_{m=1}^3 v_m \quad \mathcal{E}v_m = \omega_m v_m$$

814 Such a basis is given e.g. by

$$815 \quad (A.12) \quad e_1 = \begin{pmatrix} \mu + \mathfrak{i}k \\ 0 \\ g \end{pmatrix} \quad e_{2,3} = \begin{pmatrix} \mu + \mathfrak{i}k \\ \pm \mathfrak{i}c\sqrt{k^2 + \mu^2} \\ c^2(\mu + \mathfrak{i}k) \end{pmatrix}$$

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

$$818 \quad (A.13) \quad \hat{v} \exp(\mathfrak{i}kx) \begin{pmatrix} \frac{(\mu + \mathfrak{i}k) \sin(ct\sqrt{k^2 + \mu^2})}{c\sqrt{k^2 + \mu^2}} \\ \cos(ct\sqrt{k^2 + \mu^2}) \\ \frac{c^2(\mu + \mathfrak{i}k) \sin(ct\sqrt{k^2 + \mu^2})}{c\sqrt{k^2 + \mu^2}} \end{pmatrix}$$

$$819 \quad (A.14) \quad = \hat{v} \begin{pmatrix} -(\mu + \partial_x) \\ \partial_t \\ -c^2(\mu + \partial_x) \end{pmatrix} \exp(\mathfrak{i}kx) \frac{\sin(ct\sqrt{k^2 + \mu^2})}{c\sqrt{k^2 + \mu^2}}$$

820

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

$$\begin{aligned}
 824 \quad (\text{A.15}) \quad \begin{pmatrix} G_\rho(t, x; x') \\ G_v(t, x; x') \\ G_p(t, x; x') \end{pmatrix} &= \begin{pmatrix} -(\mu + \partial_x) \\ \partial_t \\ -c^2(\mu + \partial_x) \end{pmatrix} \frac{1}{2c} J_0 \left(\mu \sqrt{(ct)^2 - (x - x')^2} \right) \\
 825 &+ \begin{pmatrix} -\frac{\delta_{x+ct} - \delta_{x-ct}}{2c} \\ \frac{\delta_{x+ct} + \delta_{x-ct}}{2} \\ c(\delta_{x+ct} - \delta_{x-ct}) \end{pmatrix} \\
 826 &
 \end{aligned}$$

827 where J_0 is the 0-th order Bessel function of the first kind, and $J'_0 = -J_1$. Then the
 828 solution is obtained by performing a convolution with the initial data. Reinstalling
 829 the tilde one has

$$830 \quad (\text{A.16}) \quad \tilde{v}(t, x) = \int dx' G_v(t, x; x') \tilde{v}_0(x')$$

$$\begin{aligned}
 831 \quad (\text{A.17}) \quad v(t, x) &= \int dx' G_v(t, x; x') e^{\mu(x-x')} v_0(x') \\
 832 &= \frac{1}{2} \int dx' e^{\mu(x-x')} \partial_{ct} J_0 \left(\mu \sqrt{(ct)^2 - (x - x')^2} \right) v_0(x') \\
 833 &+ \frac{1}{2} \left(e^{-\mu ct} v_0(x + ct) + e^{\mu ct} v_0(x - ct) \right)
 \end{aligned}$$

$$\begin{aligned}
 834 \quad (\text{A.18}) \quad \rho(t, x) &= -\frac{1}{2c} \int dx' e^{\mu(x-x')} (\mu + \partial_x) J_0 \left(\mu \sqrt{(ct)^2 - (x - x')^2} \right) v_0(x') \\
 835 &- \frac{1}{2c} \left(e^{-\mu ct} v_0(x + ct) - e^{\mu ct} v_0(x - ct) \right) \\
 836 &
 \end{aligned}$$

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

$$838 \quad (\text{A.19}) \quad \partial_t (c^2 \rho - p) = 0$$

840 such that

$$841 \quad (\text{A.20}) \quad p(t, x) = p_0(x) + c^2 \left(\rho(t, x) - \rho_0(x) \right)$$

843 **Acknowledgments.** We thank Philip L. Roe for valuable comments and advice.

844 REFERENCES

- 845 [ABB⁺04] Emmanuel Audusse, François Bouchut, Marie-Odile Bristeau, Rupert Klein, and Benoit
 846 Perthame. A fast and stable well-balanced scheme with hydrostatic reconstruction
 847 for shallow water flows. *SIAM Journal on Scientific Computing*, 25(6):2050–2065,
 848 2004.
- 849 [Bar18] Wasilij Barsukow. *Low Mach number finite volume methods for the acoustic and Euler*
 850 *equations*. Doctoral thesis, University of Wuerzburg, 2018.
- 851 [Bar19] Wasilij Barsukow. Stationarity preserving schemes for multi-dimensional linear systems.
 852 *Mathematics of Computation*, 88(318):1621–1645, 2019.

- 853 [Bar21] Wasilij Barsukow. The active flux scheme for nonlinear problems. *Journal of Scientific*
854 *Computing*, 86(1):1–34, 2021.
- 855 [Bat54] Harry Bateman. *Tables of integral transforms (volume 1)*, volume 1. McGraw-Hill
856 Book Company, 1954.
- 857 [BCK16] Jonas P Berberich, Praveen Chandrashekar, and Christian Klingenberg. A general
858 well-balanced finite volume scheme for euler equations with gravity. In *XVI In-*
859 *ternational Conference on Hyperbolic Problems: Theory, Numerics, Applications*,
860 pages 151–163. Springer, 2016.
- 861 [BCK19] Jonas P Berberich, Praveen Chandrashekar, and Christian Klingenberg. High order
862 well-balanced finite volume methods for multi-dimensional systems of hyperbolic
863 balance laws. *arXiv preprint arXiv:1903.05154*, 2019.
- 864 [BCKR19] Jonas P Berberich, Praveen Chandrashekar, Christian Klingenberg, and Friedrich K
865 Röpke. Second order finite volume scheme for Euler equations with gravity which
866 is well-balanced for general equations of state and grid systems. *Communications*
867 *in Computational Physics*, 26:599–630, 2019.
- 868 [BHKR19] Wasilij Barsukow, Jonathan Hohm, Christian Klingenberg, and Philip L Roe. The
869 active flux scheme on Cartesian grids and its low Mach number limit. *Journal of*
870 *Scientific Computing*, 81(1):594–622, 2019.
- 871 [BKCK20] Jonas P Berberich, Roger Käppeli, Praveen Chandrashekar, and Christian Klingenberg.
872 High order discretely well-balanced finite volume methods for Euler equations with
873 gravity – without any a priori information about the hydrostatic solution. *arXiv*
874 *preprint arXiv:2005.01811*, 2020.
- 875 [BV94] Alfredo Bermudez and Ma Elena Vázquez. Upwind methods for hyperbolic conservation
876 laws with source terms. *Computers & Fluids*, 23(8):1049–1071, 1994.
- 877 [CCK⁺18] Alina Chertock, Shumo Cui, Alexander Kurganov, Şeyma Nur Özcan, and Eitan Tad-
878 mor. Well-balanced schemes for the Euler equations with gravitation: Conservative
879 formulation using global fluxes. *Journal of Computational Physics*, 2018.
- 880 [CK15] Praveen Chandrashekar and Christian Klingenberg. A second order well-balanced fi-
881 nite volume scheme for Euler equations with gravity. *SIAM Journal on Scientific*
882 *Computing*, 37(3):B382–B402, 2015.
- 883 [CL94] P Cargo and AY LeRoux. A well balanced scheme for a model of atmosphere with
884 gravity. *COMPTES RENDUS DE L ACADEMIE DES SCIENCES SERIE I-*
885 *MATHEMATIQUE*, 318(1):73–76, 1994.
- 886 [DZBK14] Vivien Desveaux, Markus Zenk, Christophe Berthon, and Christian Klingenberg. A
887 well-balanced scheme for the Euler equation with a gravitational potential. In
888 *Finite Volumes for Complex Applications VII-Methods and Theoretical Aspects*,
889 pages 217–226. Springer, 2014.
- 890 [DZBK16] Vivien Desveaux, Markus Zenk, Christophe Berthon, and Christian Klingenberg. A
891 well-balanced scheme to capture non-explicit steady states in the Euler equations
892 with gravity. *International Journal for Numerical Methods in Fluids*, 81(2):104–
893 127, 2016.
- 894 [EHB⁺21] PVF Edelmann, L Horst, JP Berberich, R Andrassy, J Higl, C Klingenberg, and
895 FK Roepke. Well-balanced treatment of gravity in astrophysical fluid dynamics
896 simulations at low Mach numbers. *arXiv preprint arXiv:2102.13111*, 2021.
- 897 [ER11] Timothy A Eymann and Philip L Roe. Active flux schemes for systems. In *20th AIAA*
898 *computational fluid dynamics conference*, 2011.
- 899 [ER13] Timothy A Eymann and Philip L Roe. Multidimensional active flux schemes. In *21st*
900 *AIAA computational fluid dynamics conference*, 2013.
- 901 [Fan17] Duoming Fan. *On the acoustic component of active flux schemes for nonlinear hyper-*
902 *bolic conservation laws*. PhD thesis, University of Michigan, Dissertation, 2017.
- 903 [FR15] Doreen Fan and Philip L Roe. Investigations of a new scheme for wave propagation. In
904 *22nd AIAA Computational Fluid Dynamics Conference*, page 2449, 2015.
- 905 [GL96] Joshua M Greenberg and Alain-Yves LeRoux. A well-balanced scheme for the numerical
906 processing of source terms in hyperbolic equations. *SIAM Journal on Numerical*
907 *Analysis*, 33(1):1–16, 1996.
- 908 [HKS19] Christiane Helzel, David Kerkmann, and Leonardo Scandurra. A new ADER method
909 inspired by the active flux method. *Journal of Scientific Computing*, 80(3):1463–
910 1497, 2019.
- 911 [KM16] R Käppeli and S Mishra. A well-balanced finite volume scheme for the Euler equations
912 with gravitation-the exact preservation of hydrostatic equilibrium with arbitrary
913 entropy stratification. *Astronomy & Astrophysics*, 587:A94, 2016.
- 914 [LeV98] Randall J LeVeque. Balancing source terms and flux gradients in high-resolution Go-

- 915 dunov methods: the quasi-steady wave-propagation algorithm. *Journal of compu-*
916 *tational physics*, 146(1):346–365, 1998.
- 917 [LGB11] Randall J LeVeque, David L George, and Marsha J Berger. Tsunami modelling with
918 adaptively refined finite volume methods. *Acta Numerica*, 20:211–289, 2011.
- 919 [NR16] Hiroaki Nishikawa and Philip L Roe. Third-order active-flux scheme for advection
920 diffusion: hyperbolic diffusion, boundary condition, and Newton solver. *Computers*
921 *& Fluids*, 125:71–81, 2016.
- 922 [vL77] Bram van Leer. Towards the ultimate conservative difference scheme. IV. A new ap-
923 proach to numerical convection. *Journal of computational physics*, 23(3):276–299,
924 1977.