

# An Eulerian–Lagrangian method for optimization problems governed by multidimensional nonlinear hyperbolic PDEs

Alina Chertock · Michael Herty · Alexander Kurganov

Received: 1 April 2012 / Published online: 2 April 2014  
© Springer Science+Business Media New York 2014

**Abstract** We present a numerical method for solving tracking-type optimal control problems subject to scalar nonlinear hyperbolic balance laws in one and two space dimensions. Our approach is based on the formal optimality system and requires numerical solutions of the hyperbolic balance law forward in time and its nonconservative adjoint equation backward in time. To this end, we develop a hybrid method, which utilizes advantages of both the Eulerian finite-volume central-upwind scheme (for solving the balance law) and the Lagrangian discrete characteristics method (for solving the adjoint transport equation). Experimental convergence rates as well as numerical results for optimization problems with both linear and nonlinear constraints and a duct design problem are presented.

**Keywords** Optimal control · Multidimensional hyperbolic partial differential equations · Numerical methods

## 1 Introduction

We are concerned with numerical approaches for optimization problems governed by hyperbolic PDEs in both one and two space dimensions. As a prototype, we consider

---

A. Chertock  
Department of Mathematics, North Carolina State University, Raleigh, NC 27695, USA  
e-mail: chertock@math.ncsu.edu

M. Herty (✉)  
Department of Mathematics, RWTH Aachen University, Templergraben 55, 52056 Aachen, Germany  
e-mail: herty@mathc.rwth-aachen.de; herty@igpm.rwth-aachen.de

A. Kurganov  
Mathematics Department, Tulane University, New Orleans, LA 70118, USA  
e-mail: kurganov@math.tulane.edu

a tracking type problem for a terminal state  $u_d$  prescribed at some given time  $t = T$  and the control acts as initial condition  $u_0$ . A mathematical formulation of this optimal control problem is reduced to minimizing a functional, and, for instance, in the two-dimensional (2-D) case it can be stated as follows:

$$\min_{u_0} J(u(\cdot, \cdot, T); u_d(\cdot, \cdot)), \quad (1.1a)$$

where  $J$  is a given functional and  $u$  is the unique entropy solution of the nonlinear scalar balance law

$$\begin{aligned} u_t + f(u)_x + g(u)_y &= h(u, x, y, t), & (x, y) \in \Omega \subseteq \mathbb{R}^2, & \quad t > 0, \\ u(x, y, 0) &= u_0(x, y), & (x, y) \in \Omega \subseteq \mathbb{R}^2. & \end{aligned} \quad (1.1b)$$

If  $\Omega \neq \mathbb{R}^2$ , then (1.1b) is supplemented by appropriate boundary conditions.

In recent years, there has been tremendous progress in both analytical and numerical studies of problems of type (1.1a), (1.1b), see, e.g., [1–3, 8–10, 13, 18, 19, 21–24, 28, 40, 44, 45]. Its solution relies on the property of the evolution operator  $\mathcal{S}_t : u_0(\cdot, \cdot) \rightarrow u(\cdot, \cdot, t) = \mathcal{S}_t u_0(\cdot, \cdot)$  for (1.1b). It is known that the semi-group  $\mathcal{S}_t$  generated by a nonlinear hyperbolic conservation/balance law is generically nondifferentiable in  $L^1$  even in the scalar one-dimensional (1-D) case (see, e.g., [10, Example 1]). A calculus for the first-order variations of  $\mathcal{S}_t u$  with respect to  $u_0$  has been established in [10, Theorems 2.2 and 2.3] for general 1-D systems of conservation laws with a piecewise Lipschitz continuous  $u_0$  that contains finitely many discontinuities. Therein, the concept of generalized first order tangent vectors has been introduced to characterize the evolution of variations with respect to  $u_0$ , see [10, equations (2.16)–(2.18)]. This result has been extended to BV initial data in [3, 8] and lead to the introduction of a differential structure for  $u_0 \rightarrow \mathcal{S}_t u_0$ , called shift-differentiability, see e.g. [3, Definition 5.1]. Related to that equations for the generalized cotangent vectors have been introduced for 1-D systems in [11, Proposition 4]. These equations (also called adjoint equations) consists of a nonconservative transport equation [11, equation (4.2)] and an ordinary differential equation [11, equations (4.3)–(4.5)] for the tangent vector and shift in the positions of possible shocks in  $u(x, t)$ , respectively. Necessary conditions for a general optimal control problem have been established in [11, Theorem 1]. However, this result was obtained using strong assumptions on  $u_0$  (see [11, Remark 4] and [3, Example 5.5]), which in the 1-D scalar case can be relaxed as shown for example in [13, 44, 46]. We note that the nonconservative transport part of the adjoint equation has been intensively studied also independently from the optimal control context. In the scalar case we refer to [4–6, 13, 36, 44, 46] for a notion of solutions and properties of solutions to those equations. The multidimensional nonconservative transport equation was studied in [7], but without a discussion of optimization issues. Analytical results for optimal control problems in the case of a scalar hyperbolic conservation law with a convex flux have also been developed using a different approach in [44, 46].

Numerical methods for the optimal control problems have been discussed in [2, 20, 22, 44, 46]. In [18, 19], the adjoint equation has been discretized using a

Lax–Friedrichs-type scheme, obtained by including conditions along shocks and modifying the Lax–Friedrichs numerical viscosity. Convergence of the modified Lax–Friedrichs scheme has been rigorously proved in the case of a smooth convex flux function. Convergence results have also been obtained in [44] for the class of schemes satisfying the one-sided Lipschitz condition (OSLC) and in [2] for implicit-explicit finite-volume methods.

In [13], analytical and numerical results for the optimal control problem (1.1a) coupled with the 1-D inviscid Burgers equation have been presented in the particular case of a least-square cost functional  $J$ . Therein, existence of a minimizer  $u_0$  was proven, however, uniqueness could not be obtained for discontinuous  $u$ . This result was also extended to the discretized optimization problem provided that the numerical schemes satisfy either the OSLC or discrete Oleinik’s entropy condition. Furthermore, convergence of numerical schemes was investigated in the case of convex flux functions and with a-priori known shock positions, and numerical resolution of the adjoint equations in both the smooth and nonsmooth cases was studied.

In this paper, we consider the problem (1.1a), (1.1b) with the least-square cost functional,

$$J(u(\cdot, \cdot, T); u_d(\cdot, \cdot)) := \frac{1}{2} \iint_{\Omega} (u(x, y, T) - u_d(x, y))^2 dx dy, \tag{1.1c}$$

and a general nonlinear scalar hyperbolic balance law. To the best of our knowledge, there is no analytical calculus available for the multidimensional problem (1.1a)–(1.1c). We therefore study the problem numerically and focus on designing highly *accurate* and *robust* numerical approach for both the forward equation and the non-conservative part of the adjoint equation (leaving aside possible additional conditions necessary to track the variations of shock positions). We treat the forward and adjoint equations separately: The hyperbolic balance law (1.1b) is numerically solved forward in time from  $t = 0$  to  $t = T$ , while the corresponding adjoint linear transport equation is integrated backward in time from  $t = T$  to  $t = 0$ . Since these two equations are of a different nature, they are attempted by different numerical methods in contrast to [13, 18, 19, 44].

The main source of difficulty one comes across while numerically solving the forward equation is the loss of smoothness, that is, the solution may develop shocks even for infinitely smooth initial data. To accurately resolve the shocks, we apply a high-resolution shock capturing Eulerian finite-volume method to (1.1b), in particular, we use a second-order semi-discrete central-upwind scheme introduced in [31–34], which is a reliable “black-box” solver for general multidimensional (systems of) hyperbolic conservation and balance laws.

The nonconservative part of the arising adjoint equation is a linear transport equation with generically discontinuous coefficients, and thus the adjoint equation is hard to accurately solve by conventional Eulerian methods. We therefore use a Lagrangian approach to numerically trace the solution of the adjoint transport equation along the backward characteristics. The resulting Lagrangian method achieves a superb numerical resolution thanks to its low numerical dissipation.

The paper is organized as follows. In Sect. 2, we briefly revise the formal adjoint calculus and additional interior conditions on the shock position in the 1-D case. Then, in Sect. 3 we present an iterative numerical optimization algorithm followed by the description of our hybrid Eulerian–Lagrangian numerical method, which is applied to a variety of 1-D and 2-D optimal control problems in Sect. 4. Convergence properties of the designed 1-D scheme are discussed in Sect. 5.

## 2 The adjoint equation

We are interested in solving the optimization problem (1.1a)–(1.1c). Formally, we proceed as follows: We introduce the Lagrangian for this problem as

$$L(u, p) = \frac{1}{2} \iint_{\mathbb{R}^2} (u(x, y, T) - u_d(x, y))^2 dx dy - \iint_{\mathbb{R}^2} p(u_t + f(u)_x + g(u)_y - h(u, x, y, t)) dx dy.$$

We integrate by parts and compute the variations with respect to  $u$  and  $p$ . In a strong formulation, the variation with respect to  $p$  leads to (1.1b) while the variation with respect to  $u$  results in the following adjoint equation:

$$-p_t - f'(u)p_x - g'(u)p_y = h_u(u, x, y, t) p, \tag{2.1a}$$

subject to the terminal condition

$$p(x, y, T) = u(x, y, T) - u_d(x, y). \tag{2.1b}$$

For sufficiently smooth solutions  $u$ , the above calculations are exact and the coupled systems (1.1b), (2.1a), (2.1b) together with

$$p(x, y, 0) = 0 \text{ a.e. } (x, y) \in \mathbb{R}^2 \tag{2.2}$$

represent the first-order optimality system for problem (1.1a)–(1.1c), in which (1.1b) should be solved forward in time while the adjoint problem (2.1a), (2.1b) is to be solved backward in time. These computations however have to be seriously modified once the solution  $u$  possesses discontinuities. This was demonstrated in [8, 10, 13, 18, 21, 44] and we will now briefly review the relevant results.

Consider a scalar 1-D conservation law  $u_t + f(u)_x = 0$  subject to the initial data  $u(x, 0) = u_0(x)$ . Denote its weak solution by  $u(\cdot, t) = \mathcal{S}_t u_0(\cdot)$ . Assume that both the solution  $u(x, t)$  and initial data  $u_0(x)$  contain a single discontinuity at  $x_s(t)$  and  $x_s(0)$ , respectively, and that  $u(x, t)$  is smooth elsewhere. As discussed above, the first-order variation of  $\mathcal{S}_t$  with respect to  $u_0$  can be computed using, for example, the concept of shift-differentiability, and in the assumed situation, this amounts to considering the corresponding linearized shock position  $\hat{x}_s$  given by  $\frac{d}{dt} (\hat{x}_s(t)[u(\cdot, t)]) = [(f'(u) -$

$\frac{d}{dt} x_s(t) \hat{u}(\cdot, t)$ ], where  $[u(\cdot, t)] = u(x_s(t)+, t) - u(x_s(t)-, t)$ . Here,  $\hat{u}$  is a solution of the linearized PDE

$$\hat{u}_t + (f'(u)\hat{u})_x = 0.$$

Using the variations  $(\hat{u}, \hat{x})$  shift-differentiability of the nonlinear cost functional  $J(u(\cdot, T), u_d(\cdot)) = J(\mathcal{S}_t u_0(\cdot), u_d(\cdot))$  is also obtained, see [13, 21, 25, 44]. It can be shown that the variation  $\delta J(\mathcal{S}_t u_0) \hat{u}$  of  $J$  with respect to  $u_0$  can be also expressed using an adjoint (or cotangent vector) formulation. More precisely, it has been shown in [18] that in the particular setting considered above,

$$\delta J(\mathcal{S}_t u_0)(\hat{u}) = \int_{\mathbb{R}} p(x, 0) \hat{u}(x, 0) dx, \tag{2.3}$$

$p(x, T) = u(x, T) - u_d(x)$  and  $p(x_s(t), t) = \frac{1}{2}[(u - u_d)^2]_T/[u]_T$ , where  $[u]_T$  represents the jump in  $u(x, T)$  across the shock at the final time. This value can be viewed as a finite difference approximation to the derivative  $u - u_d$  being the adjoint solution on either side of the shock. The latter condition requires that the adjoint solution is constant on all characteristics leading into the shock. Those rigorous results hold true once the shock position in the solution  $u$  is a-priori known. For  $u_0$  to be optimal we require the variation on the left-hand side of (2.3) to be equal to zero for all feasible variations  $\hat{u}$ .

In [19], the above result has been extended to the 1-D scalar convex case with smooth initial data that break down over time. Note that due to the nature of the adjoint equation, the region outside the extremal backwards characteristics emerging at  $x_s(T)$  are independent on the value along  $x_s(t)$ . Hence, neglecting this condition yields nonuniqueness of the solution  $p$  in the area leading into the shock, see [2, 10, 13, 44]. A recent numerical approach in [18] is based on an attempt to capture the behavior inside the region leading into the shock by adding a diffusive term to the scheme with a viscosity depending on the grid size. The results show a convergence rate of  $(\Delta x)^\alpha$  with  $\alpha < 1$ . However, to the best of our knowledge, no multidimensional extension of the above approach is available.

In this paper, we focus on the development of suitable numerical discretizations of both forward and adjoint equations. The forward equation is solved using a second-order Godunov-type central-upwind scheme [31–34], which is based on the evolution of a piecewise polynomial reconstruction of  $u$ . The latter is then used to solve the adjoint equation by the method of characteristics applied both outside and inside the region entering shock(s). Similar to [2] we therefore expect convergence of the proposed method outside the region of the backwards shock. We numerically study the behavior of  $J$  and demonstrate that even in the case of evolving discontinuities, the desired profile  $u_d$  can be recovered using the presented approach.

Since the systems (1.1b) and (2.1a), (2.1b) are fully coupled, an iterative procedure is to be imposed in order to solve the optimization problem. This approach can either be seen as a damped block Gauss-Seidel method for the system or as a gradient descent for the reduced cost functional. The reduced cost functional  $\tilde{J}(u_0) := J(\mathcal{S}_t u_0; u_d)$

is obtained when replacing  $u$  in (1.1a) by  $\mathcal{S}_t u_0$ . Then, a straightforward computation shows that  $\frac{\delta}{\delta u_0} \tilde{J}(u_0) = p(\cdot, \cdot, 0)$  if no shocks are present. As discussed above, rigorous results on the linearized cost in the presence of shocks are available only in the 1-D case, see [18, 21] and [13, Proposition 4.1, Remark 4.7 and Proposition 5.1].

### 3 Numerical method

In this section, we describe the iterative optimization algorithm along with the Eulerian–Lagrangian numerical method used in its implementation.

The underlying optimization problem can be formulated as follows: Given a terminal state  $u_d(x, y)$ , find an initial datum  $u_0(x, y)$  which by time  $t = T$  will either evolve into  $u(x, y, T) = u_d(x, y)$  or will be as close as possible to  $u_d$  in the  $L^2$ -norm. To solve the problem iteratively, we implement the following algorithm and generate a sequence  $\{u_0^{(m)}(x)\}$ ,  $m = 0, 1, 2, \dots$

#### *Iterative Optimization Algorithm*

1. Choose an initial guess  $u_0^{(0)}(x, y)$  and prescribed tolerance  $tol$ .
2. Solve the problem (1.1b) with  $u(x, y, 0) = u_0^{(0)}(x, y)$  forward in time from  $t = 0$  to  $t = T$  by an Eulerian finite-volume method (described in Sect. 3.1) to obtain  $u^{(0)}(x, y, T)$ .
3. **Iterations for  $m = 0, 1, 2, \dots$**   
**while**  $J(u^{(m)}; u_d) = \frac{1}{2} \iint_{\Omega} (u^{(m)}(x, y, T) - u_d(x, y))^2 dx dy > tol$  **or**  
**while**  $|J(u^{(m)}; u_d) - J(u^{(m-1)}; u_d)| > tol$ 
  - (a) Solve the linear transport equation (2.1a) subject to the terminal condition  $p^{(m)}(x, y, T) := u^{(m)}(x, y, T) - u_d(x, y)$  backward in time from  $t = T$  to  $t = 0$  using the Lagrangian numerical scheme (described in Sect. 3.2) to obtain  $p^{(m)}(x, y, 0)$ .
  - (b) Update the control  $u_0$  using either a gradient descent or quasi-Newton method [12, 29, 42].
  - (c) Solve the problem (1.1b) with  $u(x, y, 0) = u_0^{(m+1)}(x, y)$  forward in time from  $t = 0$  to  $t = T$  by an Eulerian finite-volume method (described in Sect. 3.1) to obtain  $u^{(m+1)}(x, y, T)$ .
  - (d) Set  $m := m + 1$ .

*Remark 3.1* Note that in the given approach the full solution  $u$  needs not to be stored during the iteration.

*Remark 3.2* The above algorithm is similar to the continuous approach from [13] though, unlike [13], we focus on the numerical methods in steps 2, 3(a) and 3(c) and thus do not use an approximation to the generalized tangent vectors to improve the gradient descent method.

*Remark 3.3* If we use a steepest descent update in step 3(b) for some stepsize  $\sigma > 0$  as

$$u_0^{(m+1)}(x, y) := u_0^{(m)}(x, y) - \sigma p^{(m)}(x, y, 0),$$

then, due to a global finite-volume approximation of  $u$  and an appropriate projection of  $p$  (see Sect. 3.1), we obtain a piecewise polynomial control  $u_0^{(m+1)}$  in this step of the algorithm. The fact that the control  $u_0^{(m+1)}$  is always piecewise polynomial prevents the accumulation of discontinuities in the forward solution in our algorithm. Clearly, other (higher-order) gradient-based optimization methods can be used to speed up the convergence, especially in the advance stages of the above iterative procedure, see, e.g., [29,42] for more details.

*Remark 3.4* In [13], a higher number of iterations is reported when the adjoints to the Rankine–Hugoniot condition are neglected. In view of 2-D examples, we opt in this paper for the higher number of optimization steps  $m$  compared with the necessity to recompute the shock location and solve for the adjoint linearized Rankine–Hugoniot condition.

### 3.1 Finite-volume method for (1.1b)

In this section, we briefly describe a second-order semi-discrete central-upwind scheme from [34] (see also [31–33]), which has been applied in steps 2 and 3(c) in our iterative optimization algorithm on page 5.

We start by introducing a uniform spatial grid, which is obtained by dividing the computational domain into finite-volume cells  $C_{j,k} := [x_{j-\frac{1}{2}}, x_{j+\frac{1}{2}}] \times [y_{k-\frac{1}{2}}, y_{k+\frac{1}{2}}]$  with  $x_{j+\frac{1}{2}} - x_{j-\frac{1}{2}} = \Delta x, \forall j$  and  $y_{k+\frac{1}{2}} - y_{k-\frac{1}{2}} = \Delta y, \forall k$ . We assume that at a certain time  $t$ , the computed solution is available and realized in terms of its cell averages

$$\bar{u}_{j,k}(t) \approx \frac{1}{\Delta x \Delta y} \iint_{C_{j,k}} u(x, y, t) dx dy.$$

according to the central-upwind approach, the cell averages are evolved in time by solving the following system of ODEs:

$$\frac{d}{dt} \bar{u}_{j,k}(t) = -\frac{F_{j+\frac{1}{2},k} - F_{j-\frac{1}{2},k}}{\Delta x} - \frac{G_{j,k+\frac{1}{2}} - G_{j,k-\frac{1}{2}}}{\Delta y} + h(\bar{u}_{j,k}, x_j, y_k, t), \quad (3.1)$$

using an appropriate ODE solver. In this paper, we implement the third-order strong stability preserving Runge–Kutta (SSP RK) method from [26].

The numerical fluxes  $F$  and  $G$  in (3.1) are given by (see [31,34] for details):

$$\begin{aligned}
 F_{j+\frac{1}{2},k} &= \frac{a_{j+\frac{1}{2},k}^+ f(u_{j,k}^E) - a_{j+\frac{1}{2},k}^- f(u_{j+1,k}^W)}{a_{j+\frac{1}{2},k}^+ - a_{j+\frac{1}{2},k}^-} + \frac{a_{j+\frac{1}{2},k}^+ a_{j+\frac{1}{2},k}^-}{a_{j+\frac{1}{2},k}^+ - a_{j+\frac{1}{2},k}^-} \left[ u_{j+1,k}^W - u_{j,k}^E \right], \\
 G_{j,k+\frac{1}{2}} &= \frac{b_{j,k+\frac{1}{2}}^+ g(u_{j,k}^N) - b_{j,k+\frac{1}{2}}^- g(u_{j,k+1}^S)}{b_{j,k+\frac{1}{2}}^+ - b_{j,k+\frac{1}{2}}^-} + \frac{b_{j,k+\frac{1}{2}}^+ b_{j,k+\frac{1}{2}}^-}{b_{j,k+\frac{1}{2}}^+ - b_{j,k+\frac{1}{2}}^-} \left[ u_{j,k+1}^S - u_{j,k}^N \right].
 \end{aligned}
 \tag{3.2}$$

Here,  $u_{j,k}^{E,W,N,S}$  are the point values of the piecewise linear reconstruction for  $u$

$$\tilde{u}(x, y) := \bar{u}_{j,k} + (u_x)_{j,k}(x - x_j) + (u_y)_{j,k}(y - y_k), \quad (x, y) \in C_{j,k}, \tag{3.3}$$

at  $(x_{j+\frac{1}{2}}, y_k)$ ,  $(x_{j-\frac{1}{2}}, y_k)$ ,  $(x_j, y_{k+\frac{1}{2}})$ , and  $(x_j, y_{k-\frac{1}{2}})$ , respectively. Namely, we have:

$$\begin{aligned}
 u_{j,k}^E &:= \tilde{u}(x_{j+\frac{1}{2}} - 0, y_k) = \bar{u}_{j,k} + \frac{\Delta x}{2}(u_x)_{j,k}, & u_{j,k}^W &:= \tilde{u}(x_{j-\frac{1}{2}} + 0, y_k) = \bar{u}_{j,k} \\
 &\quad - \frac{\Delta x}{2}(u_x)_{j,k}, \\
 u_{j,k}^N &:= \tilde{u}(x_j, y_{k+\frac{1}{2}} - 0) = \bar{u}_{j,k} + \frac{\Delta y}{2}(u_y)_{j,k}, & u_{j,k}^S &:= \tilde{u}(x_j, y_{k-\frac{1}{2}} + 0) = \bar{u}_{j,k} \\
 &\quad - \frac{\Delta y}{2}(u_y)_{j,k}.
 \end{aligned}
 \tag{3.4}$$

The numerical derivatives  $(u_x)_{j,k}$  and  $(u_y)_{j,k}$  are (at least) first-order approximations of  $u_x(x_j, y_k, t)$  and  $u_y(x_j, y_k, t)$ , respectively, and are computed using a nonlinear limiter that would ensure a non-oscillatory nature of the reconstruction (3.3). In our numerical experiments, we have used a generalized minmod reconstruction [35,37, 43,47]:

$$\begin{aligned}
 (u_x)_{j,k} &= \text{minmod} \left( \theta \frac{\bar{u}_{j,k} - \bar{u}_{j-1,k}}{\Delta x}, \frac{\bar{u}_{j+1,k} - \bar{u}_{j-1,k}}{2\Delta x}, \theta \frac{\bar{u}_{j+1,k} - \bar{u}_{j,k}}{\Delta x} \right), \\
 (u_y)_{j,k} &= \text{minmod} \left( \theta \frac{\bar{u}_{j,k} - \bar{u}_{j,k-1}}{\Delta y}, \frac{\bar{u}_{j,k+1} - \bar{u}_{j,k-1}}{2\Delta y}, \theta \frac{\bar{u}_{j,k+1} - \bar{u}_{j,k}}{\Delta y} \right),
 \end{aligned} \quad \theta \in [1, 2].
 \tag{3.5}$$

Here, the minmod function is defined as

$$\text{minmod}(z_1, z_2, \dots) := \begin{cases} \min_j \{z_j\}, & \text{if } z_j > 0 \ \forall j, \\ \max_j \{z_j\}, & \text{if } z_j < 0 \ \forall j, \\ 0, & \text{otherwise.} \end{cases}$$

The parameter  $\theta$  is used to control the amount of numerical viscosity present in the resulting scheme: larger values of  $\theta$  correspond to less dissipative but, in general, more oscillatory reconstructions. In all of our numerical experiments, we have taken  $\theta = 1.5$ .

The one-sided local speeds in the  $x$ - and  $y$ -directions,  $a_{j+\frac{1}{2},k}^\pm$  and  $b_{j,k+\frac{1}{2}}^\pm$ , are determined as follows:

$$\begin{aligned}
 a_{j+\frac{1}{2},k}^+ &= \max \left\{ \max_{u \in [\min\{u_{j,k}^E, u_{j+1,k}^W\}, \max\{u_{j,k}^E, u_{j+1,k}^W\}]} f'(u), 0 \right\}, \\
 a_{j+\frac{1}{2},k}^- &= \min \left\{ \min_{u \in [\min\{u_{j,k}^E, u_{j+1,k}^W\}, \max\{u_{j,k}^E, u_{j+1,k}^W\}]} f'(u), 0 \right\}, \\
 b_{j,k+\frac{1}{2}}^+ &= \max \left\{ \max_{u \in [\min\{u_{j,k}^N, u_{j,k+1}^S\}, \max\{u_{j,k}^N, u_{j,k+1}^S\}]} g'(u), 0 \right\}, \\
 b_{j,k+\frac{1}{2}}^- &= \min \left\{ \min_{u \in [\min\{u_{j,k}^N, u_{j,k+1}^S\}, \max\{u_{j,k}^N, u_{j,k+1}^S\}]} g'(u), 0 \right\}.
 \end{aligned} \tag{3.6}$$

*Remark 3.5* In equations (3.1)–(3.6), we suppress the dependence of  $\bar{u}_{j,k}$ ,  $F_{j+\frac{1}{2},k}$ ,  $G_{j,k+\frac{1}{2}}$ ,  $u_{j,k}^{E,W,N,S}$ ,  $(u_x)_{j,k}$ ,  $(u_y)_{j,k}$ ,  $a_{j+\frac{1}{2},k}^\pm$ , and  $b_{j,k+\frac{1}{2}}^\pm$  on  $t$  to simplify the notation.

*Remark 3.6* We solve the balance law (1.1b) starting from time  $t^0 = 0$  and compute the solution at time levels  $t^n$ ,  $n = 1, \dots, N_T$ , where  $t^{N_T} = T$ . Since the obtained approximate solution is to be used for solving the adjoint problem backward in time, we store the values of  $u$  and its discrete derivatives,  $u_x$  and  $u_y$ , so the piecewise linear approximants (3.3) are available at all of the time levels.

### 3.2 Discrete method of characteristics for (2.1a), (2.1b)

We finish the description of the proposed Eulerian–Lagrangian numerical method by presenting the Lagrangian method for the backward transport equation (2.1a).

Since equation (2.1a) is linear (with possibly discontinuous coefficients), we follow [14, 15] and solve it using the Lagrangian approach, which is a numerical version of the method of characteristics. In this method, the solution is represented by a certain number of its point values prescribed at time  $T$  at the points  $(x_i^c(T), y_i^c(T))$ , which may (or may not) coincide with the uniform grid points  $(x_j, y_k)$  used in the numerical solution of the forward problem (1.1b). The location of these characteristics points is tracked backward in time by numerically solving the following system of ODEs:

$$\begin{cases} \frac{dx_i^c(t)}{dt} = f'(u(x_i^c(t), y_i^c(t), t)), \\ \frac{dy_i^c(t)}{dt} = g'(u(x_i^c(t), y_i^c(t), t)), \\ \frac{dp_i^c(t)}{dt} = -h_u(u(x_i^c(t), y_i^c(t), t), x_i^c(t), y_i^c(t), t)p_i^c(t), \end{cases} \tag{3.7}$$

where  $x_i^c(t), y_i^c(t)$  is a position of the  $i$ th characteristics point at time  $t$  and  $p_i^c(t)$  is a corresponding value of  $p$  at that point.

Notice that since the piecewise linear approximants  $\tilde{u}$  are only available at the discrete time levels  $t = t^n, n = 0, \dots, N_T$ , we solve the system (3.7) using the second-order SSP RK method (the Heun method), see [26]. Unlike the third-order SSP RK method, the Heun method only uses the data from the current time level and the previous one and does not require data from any intermediate time levels. This is the main advantage of the Heun method since the right-hand side of (3.7) can be easily computed at each discrete time level  $t = t^n$ . To this end, we simply check at which finite-volume cell the characteristics point is located at a given discrete time moment, and then set

$$u(x_i^c(t^n), y_i^c(t^n), t^n) = \tilde{u}(x_i^c(t^n), y_i^c(t^n), t^n) = \bar{u}_{j,k}^n + (u_x)_{j,k}^n(x_i^c(t^n) - x_j) + (u_y)_{j,k}^n(y_i^c(t^n) - y_k), \text{ if } (x_i^c(t^n), y_i^c(t^n)) \in C_{j,k}. \tag{3.8}$$

Thus, as the time reaches 0, we will have the set of the characteristics points  $(x_i^c(0), y_i^c(0))$  and the corresponding point values  $p_i^c(0)$ . We then obtain the updated initial data (step 3(b) in our iterative optimization algorithm on page 5) by first setting

$$u_0^{(m+1)}(x_i^c(0), y_i^c(0)) := u_0^{(m)}(x_i^c(0), y_i^c(0)) - \sigma p_i^c(0), \tag{3.9}$$

where  $\sigma = \mathcal{O}(\Delta x)$ , and then projecting the resulting set of discrete data onto the uniform grid. The latter is done as follows: For a given grid point  $(x_j, y_k)$ , we find four characteristics points  $(x_{i_1}^c(0), y_{k_1}^c(0)), (x_{i_2}^c(0), y_{k_2}^c(0)), (x_{i_3}^c(0), y_{k_3}^c(0))$ , and  $(x_{i_4}^c(0), y_{k_4}^c(0))$  such that

$$\begin{aligned} (x_{i_1}^c(0) - x_j)^2 + (y_{i_1}^c(0) - y_k)^2 &= \min_{i: x_{i_1}^c(0) > x_j, y_{i_1}^c(0) > y_k} \left[ (x_{i_1}^c(0) - x_j)^2 + (y_{i_1}^c(0) - y_k)^2 \right], \\ (x_{i_2}^c(0) - x_j)^2 + (y_{i_2}^c(0) - y_k)^2 &= \min_{i: x_{i_2}^c(0) < x_j, y_{i_2}^c(0) > y_k} \left[ (x_{i_2}^c(0) - x_j)^2 + (y_{i_2}^c(0) - y_k)^2 \right], \\ (x_{i_3}^c(0) - x_j)^2 + (y_{i_3}^c(0) - y_k)^2 &= \min_{i: x_{i_3}^c(0) < x_j, y_{i_3}^c(0) < y_k} \left[ (x_{i_3}^c(0) - x_j)^2 + (y_{i_3}^c(0) - y_k)^2 \right], \\ (x_{i_4}^c(0) - x_j)^2 + (y_{i_4}^c(0) - y_k)^2 &= \min_{i: x_{i_4}^c(0) > x_j, y_{i_4}^c(0) < y_k} \left[ (x_{i_4}^c(0) - x_j)^2 + (y_{i_4}^c(0) - y_k)^2 \right], \end{aligned} \tag{3.10}$$

and then use a bilinear interpolation between these four characteristics points to obtain  $u_0^{(m+1)}(x_j, y_k)$ .

We are then ready to proceed with the next,  $(m + 1)$ -st iteration.

*Remark 3.7* In the 1-D case, the projection procedure from the characteristics points onto the uniform grid simplifies significantly, as one just needs to locate the two characteristics points that are closest to the given grid point from the left and from the right, and then use a linear interpolation.

*Remark 3.8* It is well-known that one of the difficulties emerging when Lagrangian methods are applied to linear transport equations with discontinuous coefficients is

that the distances between characteristics points constantly change. As a result, characteristics points may either cluster or spread too far from each other. This may lead not only to a poor resolution of the computed solution, but also to an extremely low efficiency of the method. To overcome this difficulties, we either add characteristics points to fill appearing gaps between the points (the solution values at the added points can be obtained using a similar bilinear/linear interpolation) or remove some of the clustered characteristics points.

*Remark 3.9* Notice that if  $u_d(x, y)$  is discontinuous, the solution of the optimization problem is not unique: Due to the loss of information at the shock, there will be infinitely many different initial data  $u_0(x, y)$ , which would lead to the same solution of (1.1b).

*Remark 3.10* It should be observed that the solution of the adjoint transport equation (2.1a) can be computed back in time using different methods. We refer the reader, e.g., to [25], where several first-order discretizations of the 1-D transport equation without source terms have been presented and analyzed. In the numerical results below, we compare the proposed discrete method of characteristics with an upwind approach from [25] when applied to the 1-D equation (2.1a) (see Example 1 in Sect. 4.1).

### 4 Numerical results

In this section, we illustrate the performance of the proposed method on a number of numerical examples. We start with a grid convergence study by considering a linear advection equation in (1.1b) as a constraint. We then consider a control problem of the inviscid Burgers equation and demonstrate the non-uniqueness of optimal controls in the case of nonsmooth desired state. Next, we numerically solve a duct design problem modified from [44]. Finally, we apply the proposed method to the 2-D inviscid Burgers equation.

#### 4.1 The one-dimensional case

In this section, we consider the 1-D version of the optimization problem (1.1a)–(1.1c):

$$\min_{u_0} J(u(\cdot, T); u_d(\cdot)), \quad J(u(\cdot, T); u_d(\cdot)) := \frac{1}{2} \int_I (u(x, T) - u_d(x))^2 dx, \quad (4.1a)$$

where  $u$  is a solution of the scalar hyperbolic PDE

$$\begin{aligned} u_t + f(u)_x &= h(u, x, t), & x \in I \subseteq \mathbb{R}, & \quad t > 0, \\ u(x, 0) &= u_0(x), & x \in I \subseteq \mathbb{R}. & \end{aligned} \quad (4.1b)$$

If  $I \neq \mathbb{R}$ , then (4.1b) is augmented with appropriate boundary conditions.

The corresponding adjoint problem is

$$\begin{aligned} -p_t - f'(u)p_x &= h_u(u, x, t)p, & x \in I \subseteq \mathbb{R}, \quad t > 0, \\ p(x, T) &= p_T(x) := u(x, T) - u_d(x) & x \in I \subseteq \mathbb{R}. \end{aligned} \quad (4.2)$$

*Example 1: Linear constraints*

We numerically solve the optimization problem (4.1a)–(4.1b) with the terminal state  $u_d(x) = e^{-(x-\pi)^2}$  and subject to a linear advection equation as constraint:

$$u_t + u_x = 0, \quad x \in [0, 2\pi], \quad t > 0, \quad (4.3)$$

with the periodic boundary conditions. The corresponding adjoint equation is

$$-p_t - p_x = 0 \quad x \in [0, 2\pi], \quad t < T. \quad (4.4)$$

Since both (4.3) and (4.4) are linear advection equations with constant coefficients, the exact solution  $u_0$  of the studied optimization problem is unique and can be easily obtained:

$$u_0(x) = u_d(x - T). \quad (4.5)$$

In this example, we start the iterative optimization algorithm (page 5) with the constant initial condition,

$$u_0^{(0)}(x) \equiv 0.5,$$

and illustrate that the proposed Eulerian–Lagrangian method is second-order accurate. We use spatial grids with  $\Delta x = 2\pi/100, 2\pi/200, 2\pi/400$  and  $2\pi/800$  for the Eulerian part of the method and the corresponding number of the characteristics points (100, 200, 400 or 800) for the Lagrangian part. We set  $T = 2\pi$  and  $tol$ , and measure the  $L^1$ -errors for both the control (see (4.5)),

$$\|e_0\|_{L^1} := \Delta x \sum_j \left| u_0^{(m)}(x_j) - u_d(x_j - 2\pi) \right|,$$

and the terminal state,

$$\|e_T\|_{L^1} := \Delta x \sum_j \left| u^{(m)}(x_j, 2\pi) - u_d(x_j) \right|.$$

The results of the numerical grid convergence study displayed in the Table 1 for  $tol = (\Delta x)^4$ , clearly show that the expected second-order rate of convergence has been achieved.

It should be pointed out that problem (4.1a) does not contain any regularization with respect to  $u_0$ . In general, this may lead to a deterioration of the iteration numbers

**Table 1** Example 1:  $L^1$ -convergence study

$\Delta x$	No. of iterations	$\ e_0\ _{L^1}$	Rate	$\ e_T\ _{L^1}$	Rate
$2\pi/50$	29	$8.32 \times 10^{-2}$	–	$5.05 \times 10^{-2}$	–
$2\pi/100$	81	$2.13 \times 10^{-2}$	1.97	$1.26 \times 10^{-2}$	2.00
$2\pi/200$	207	$5.46 \times 10^{-3}$	1.96	$3.24 \times 10^{-3}$	1.96
$2\pi/400$	505	$1.36 \times 10^{-3}$	2.00	$8.09 \times 10^{-4}$	2.00
$2\pi/800$	1189	$3.44 \times 10^{-4}$	1.98	$2.04 \times 10^{-4}$	1.99

Optimization is terminated when  $J \leq tol = (\Delta x)^4$

**Table 2** Example 1:  $L^1$ -convergence study for the regularized cost functional

$\Delta x$	$\alpha$	No. of iterations	$\ e_0\ _{L^1}$	Rate	$\ e_T\ _{L^1}$	Rate
$2\pi/100$	$10^0$	70	$7.92 \times 10^{-3}$	-	$7.95 \times 10^{-3}$	-
	$10^{-1}$	75	$1.92 \times 10^{-2}$	-	$1.16 \times 10^{-2}$	-
	$10^{-2}$	80	$2.13 \times 10^{-2}$	-	$1.27 \times 10^{-2}$	-
$2\pi/200$	$10^0$	135	$2.05 \times 10^{-3}$	1.95	$2.06 \times 10^{-3}$	1.94
	$10^{-1}$	191	$4.93 \times 10^{-3}$	1.96	$2.97 \times 10^{-3}$	1.97
	$10^{-2}$	205	$5.42 \times 10^{-3}$	1.97	$3.24 \times 10^{-3}$	1.97
$2\pi/400$	$10^0$	313	$5.16 \times 10^{-4}$	1.99	$5.18 \times 10^{-4}$	1.99
	$10^{-1}$	464	$1.25 \times 10^{-3}$	1.98	$7.50 \times 10^{-4}$	1.98
	$10^{-2}$	500	$1.36 \times 10^{-3}$	1.99	$8.09 \times 10^{-4}$	2.00

Optimization is terminated when  $J_{reg} \leq tol = (\Delta x)^4$

within the optimization. In order to quantify this effect, we present, in Table 2, the results obtained for the regularized problem

$$\min_{u_0} J_{reg}(u(\cdot, T); u_d(\cdot)),$$

$$J_{reg}(u(\cdot, T); u_d(\cdot)) := J(u(\cdot, T); u_d(\cdot)) + \frac{\alpha}{2} \int_I (u_0(x) - u^*(x))^2 dx, \quad (4.6)$$

where  $u$  is a solution of (4.3),  $u^*(x) = e^{-(x-\pi)^2}$ , and  $\alpha$  is a regularization parameter. Obviously, the gradient computation in (3.9) has to be modified accordingly, but no further changes compared with the previous example are made. As expected, Table 2 shows that increasing value of parameter  $\alpha$  leads a decreasing number of iterations while preserving the second order of accuracy of the method.

Furthermore, we compare the performance of the proposed discrete method of characteristics with an upwind approach from [25] for solving the adjoint transport equation (4.4). To this end, we apply a first-order version of the central-upwind scheme (briefly discussed in ‘‘Appendix A’’) to the forward equation (4.3) and then use either the method of characteristics or the first-order upwind scheme for solving the adjoint equation (4.4). Since in this case, the overall accuracy of the method is one, we chose a much larger  $tol = (\Delta x)^2$ . Table 3 shows the number of iterations and  $L^1$ -errors for both cases. We observe that the iteration numbers are slightly better if the adjoint

**Table 3** Example 1: Comparison of the upwind scheme and the method of characteristics for solving (4.4)

$\Delta x$	Upwind			Characteristics		
	No. of Iterations	$\ e_0\ _{L^1}$	$\ e_T\ _{L^1}$	No. of Iterations	$\ e_0\ _{L^1}$	$\ e_T\ _{L^1}$
$2\pi/100$	46	$2.50 \times 10^{-1}$	$1.90 \times 10^{-1}$	38	$2.10 \times 10^{-1}$	$2.00 \times 10^{-1}$
$2\pi/200$	107	$1.40 \times 10^{-1}$	$1.00 \times 10^{-1}$	99	$1.10 \times 10^{-1}$	$1.00 \times 10^{-1}$
$2\pi/400$	251	$7.25 \times 10^{-2}$	$5.22 \times 10^{-2}$	243	$5.20 \times 10^{-2}$	$5.12 \times 10^{-2}$
$2\pi/800$	584	$3.71 \times 10^{-2}$	$2.61 \times 10^{-2}$	574	$2.62 \times 10^{-2}$	$2.60 \times 10^{-2}$

Optimization is terminated when  $J \leq tol = (\Delta x)^2$

equation is solved by the method of characteristics; they are also lower compared to the corresponding numbers in Table 1 due to the reduced tolerance. As expected, the rate of convergence is approximately one.

*Example 2: Nonlinear constraints*

In this section, we consider the optimization problem (4.1a)–(4.1b) subject to the nonlinear inviscid Burgers equation,

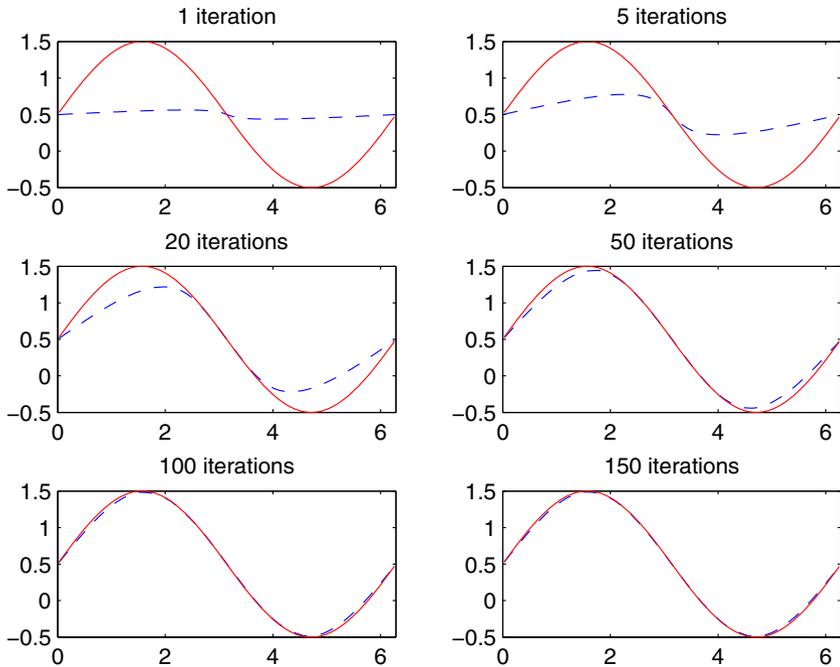
$$\begin{aligned}
 u_t + \left(\frac{u^2}{2}\right)_x &= 0, & x \in [0, 2\pi], \quad t > 0, \\
 u(x, 0) &= u_0(x), & x \in [0, 2\pi], \\
 u &\text{ is } 2\pi\text{-periodic,}
 \end{aligned}
 \tag{4.7}$$

as the constraint and use its solution at a certain time as a terminal state for the control problem. We generate this terminal state  $u_d$  by solving (4.7) with the initial data given by

$$u_0(x) = \frac{1}{2} + \sin x, \quad x \in [0, 2\pi].
 \tag{4.8}$$

It is easy to show that for  $t < \pi/2$  the solution of (4.7), (4.8) is smooth, whereas it breaks down and develops a shock wave at the critical time of  $t = \pi/2$  (later on, the shock travels to the right with a constant speed of  $s = 0.5$ ). In the following, we will consider both smooth,  $u(x, T = \pi/4)$ , and nonsmooth,  $u(x, T = 2)$ , solutions of (4.7), (4.8), computed by the 1-D version of the second-order semi-discrete central-upwind scheme (see Sect. 3.1), and use them as terminal states  $u_d$  in the optimization problem.

In Figs. 1 and 2, we plot the recovered smooth optimal initial data  $u_0^{(m)}(x)$  together with the exact initial data  $u_0(x)$  from (4.8) and the computed terminal state  $u^{(m)}(x, T = \pi/4)$  together with  $u_d(x)$ , respectively. The Eulerian part of the computations were performed on a uniform grid with  $\Delta x = 2\pi/100$ . At the beginning of each backward Lagrangian step, the characteristics points were placed at the center of each finite-volume cell. We have also solved the optimization problem on finer grids, namely, with  $\Delta x = 2\pi/200$  and  $2\pi/400$ . The plots look very similar to those shown in Figs. 1



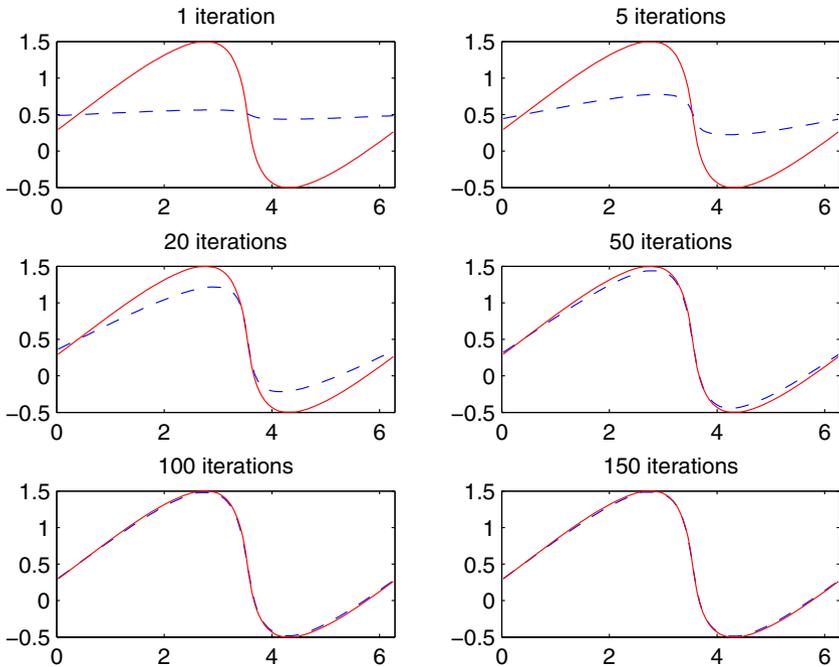
**Fig. 1** Example 2 (smooth case): Recovered initial data  $u_0^{(m)}(x)$  with  $m = 1, 5, 20, 50, 100$  and  $150$  (dashed line) and the exact initial data (4.8) (solid line)

and 2 and therefore are not presented here. In all of the computations, the iterations were started with the initial guess  $u_0^{(0)}(x) \equiv 0.5$ .

We also perform a numerical grid convergence study, whose results are presented in Table 4. One can clearly see that similarly to the case of linear constraint (Example 1), the expected second-order rate of convergence has been achieved in the smooth nonlinear case with  $tol = 25(\Delta x)^4$ .

We then proceed with the nonsmooth terminal state  $u_d(x) = u(x, T = 2)$ . In this case, the solution of the studied optimization problem is not unique. Our Eulerian–Lagrangian method recovers just one of the possible solutions, which is presented in Figs. 3 and 4. In Fig. 3, we plot the recovered initial data  $u_0^{(m)}(x)$  together with the initial data (4.8), used to generate  $u(x, T = 2)$  by the central-upwind scheme. In Fig. 4, we show the computed terminal state  $u^{(m)}(x, T = 2)$  together with  $u_d(x)$ . The plotted solutions are computed on a uniform grid with  $\Delta x = 2\pi/100$  and the corresponding number of characteristics points (as before, similar results are obtained on finer grids but not reported here).

Further, we conduct a comparison of the convergence behaviour for the smooth ( $T = \pi/4$ ) and nonsmooth ( $T = 2$ ) solutions of the optimization problem (4.1a), (4.7). In the nonsmooth case, the terminal state  $u_d$  is discontinuous. The discontinuity is located at  $\bar{x} = \pi + 1$  with the left  $u_d(\bar{x}-) = u_l = 3/2$  and right  $u_d(\bar{x}+) = u_r = -1/2$  states, respectively. The extremal backward characteristics emerging at  $\bar{x}$  reach the points  $x_l = \bar{x} - 3$  and  $x_r = 1 + \bar{x}$  at time  $t = 0$ . In order to avoid the problem of



**Fig. 2** Example 2 (smooth case): Recovered solution of the optimal control problem,  $u^{(m)}(x, T = \pi/4)$  with  $m = 1, 5, 20, 50, 100$  and  $150$  (dashed line) and the terminal state  $u_d$  (solid line)

**Table 4** Example 2:  $L^1$ -convergence study for smooth solutions

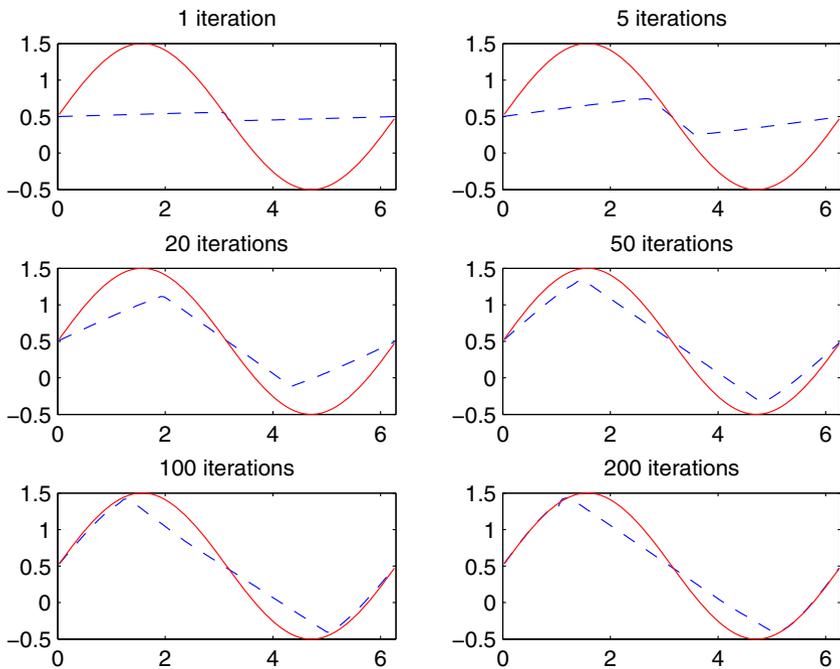
$\Delta x$	No. of iterations	$\ e_0\ _{L^1}$	Rate	$\ e_T\ _{L^1}$	Rate
$2\pi/100$	127	$6.44 \times 10^{-2}$	–	$6.10 \times 10^{-2}$	–
$2\pi/200$	404	$1.75 \times 10^{-2}$	1.88	$1.57 \times 10^{-2}$	1.94
$2\pi/400$	1365	$3.20 \times 10^{-3}$	2.45	$3.07 \times 10^{-3}$	2.35

nonuniqueness of the optimal solution in the region  $x_l \leq x \leq x_r$ , we compute in the nonsmooth case the  $L^1$ -error only on the intervals  $I_l := [0, x_l]$  and  $I_r := [x_r, 2\pi]$ , that is, in this case we define

$$\|e_0\|_{L^1_{loc}} := \Delta x \sum_j \chi_{I_l \cup I_r}(x_j) |u_0^{(m)}(x_j) - u_0(x_j)|.$$

Here,  $u_0$  is given by equation (4.8) and  $\chi_I$  is the characteristic function on the interval  $I$ . The tolerance for both cases is set to  $tol = (\Delta x)^2$  and the results are reported in Table 5.

We also check the numerical convergence of the computed adjoint solution  $p$ . Notice that inside the region of the extremal backwards characteristics,  $p$  is constant and according to [19], its value is equal to  $\frac{1}{2} \frac{(u(\bar{x}+, T) - u_d(\bar{x}))^2 - (u(\bar{x}-, T) - u_d(\bar{x}))^2}{u(\bar{x}+, T) - u(\bar{x}-, T)}$ . Therefore,  $p(x, 0) \equiv 0$  everywhere, and we measure the computed  $p$  in the maximum



**Fig. 3** Example 2 (nonsmooth case): Recovered initial data  $u_0^{(m)}(x)$  with  $m = 1, 5, 20, 50, 100$  and  $200$  (dashed line) and  $u_0(x)$  given by (4.8) (solid line)

norm at time  $t = 0$ . The results are shown in Table 6, where we demonstrate the behavior of  $\max |p_j(0)|$  both outside,  $\max_{j:x_j \in I_l \cup I_r} |p_j(0)|$ , and inside,  $\max_{j:x_j \in [x_l, x_r]} |p_j(0)|$ , the shock region. As expected, a pointwise convergence is only observed away from the region of the extremal backwards characteristics.

*Example 3: Duct design problem*

In this example, we consider a duct design problem from [16] (see also [44, p. 233]).

In the original model in [16], the flow within a duct area  $A(x)$  is described by the 1-D Euler equations. Under several simplifying assumptions, this problem can be reduced to an optimization problem for  $A$  subject to the following ODE constraint (see [17]):

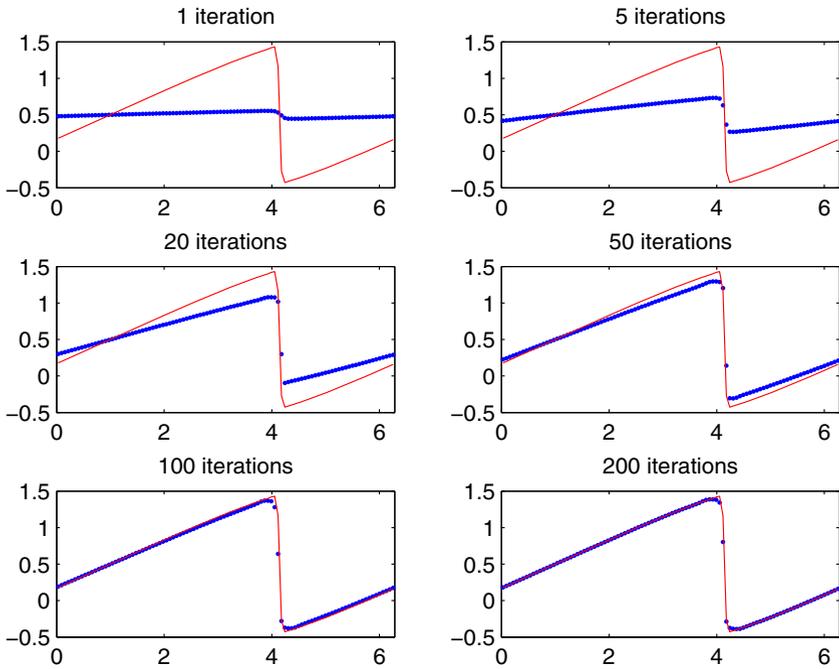
$$f(u)_x = h(u, A), \quad x \in (0, 1), \tag{4.9}$$

with

$$f(u) = u + \frac{H}{u}, \quad h(u, A) = -\frac{A'}{A} \left( \gamma u - \frac{H}{u} \right), \tag{4.10}$$

where  $\gamma$  and  $H$  are positive constants, and appropriate boundary conditions are supplied.

The boundary conditions for (4.9), (4.10) given in [16] are such that the solution has a discontinuity at some unknown point  $x_*$ , at which the Rankine–Hugoniot condition



**Fig. 4** Example 2 (nonsmooth case): Recovered solution of the optimal control problem,  $u^{(m)}(x, T = 2)$  with  $m = 1, 5, 20, 50, 100$  and  $200$  (plotted with *points*) and the terminal state  $u_d$  (*solid line*)

**Table 5** Example 2: Comparison of convergence behavior for  $J \leq tol = (\Delta x)^2$  and terminal times  $T = \pi/4$  and  $T = 2$ , for which the terminal states  $u_d$  is smooth and nonsmooth, respectively

$\Delta x$	$T = \pi/4$			$T = 2$		
	No. of iterations	$\ e_0\ _{L^1}$	$\ e_T\ _{L^1}$	No. of iterations	$\ e_0\ _{L^1_{loc}}$	$\ e_T\ _{L^1_{loc}}$
$2\pi/50$	21	$4.10 \times 10^{-1}$	$3.60 \times 10^{-1}$	30	$2.41 \times 10^{-1}$	$3.45 \times 10^{-1}$
$2\pi/100$	70	$2.10 \times 10^{-1}$	$2.00 \times 10^{-1}$	90	$1.17 \times 10^{-1}$	$1.53 \times 10^{-1}$
$2\pi/200$	313	$1.10 \times 10^{-1}$	$1.00 \times 10^{-1}$	271	$4.63 \times 10^{-2}$	$5.19 \times 10^{-2}$
$2\pi/400$	1033	$5.00 \times 10^{-2}$	$4.88 \times 10^{-2}$	5226	$3.47 \times 10^{-2}$	$4.11 \times 10^{-2}$
$2\pi/800$	5845	$2.47 \times 10^{-2}$	$2.44 \times 10^{-2}$	1953	$1.17 \times 10^{-2}$	$1.41 \times 10^{-2}$

**Table 6** Example 2: Maximum values of the adjoint solution  $p$ , computed at time  $t = 0$  outside and inside the region of the extremal backwards characteristics

$\Delta x$	No. of iterations	$\max_{j:x_j \in I \cup I_r}  p_j(0) $	$\max_{j:x_j \in [x_l, x_r]}  p_j(0) $
$2\pi/50$	30	$7.46 \times 10^{-2}$	$6.29 \times 10^{-1}$
$2\pi/100$	90	$3.38 \times 10^{-2}$	$4.11 \times 10^{-1}$
$2\pi/200$	271	$1.05 \times 10^{-2}$	$4.26 \times 10^{-1}$
$2\pi/400$	5226	$8.61 \times 10^{-3}$	$1.12 \times 10^{-1}$
$2\pi/800$	1953	$3.08 \times 10^{-3}$	$1.36 \times 10^{-1}$

holds. The desired profile of  $A$  is then obtained by solving a minimization problem for a given state  $u_d(x)$  obtained as a solution of (4.9), (4.10) for a prescribed area  $A_d(x)$  and the following boundary data  $u(0) = u_\ell = 1.299$  and  $u(1) = u_r = 0.506$ . The function  $A_d$  is defined as the following cubic polynomial:

$$A_d(x) = -1.19x^3 + 1.785x^2 + 0.1x + 1.05.$$

In [44, Section 7.1], a time dependent version of the original duct design problem has been considered subject to the nonlinear hyperbolic balance law constraint:

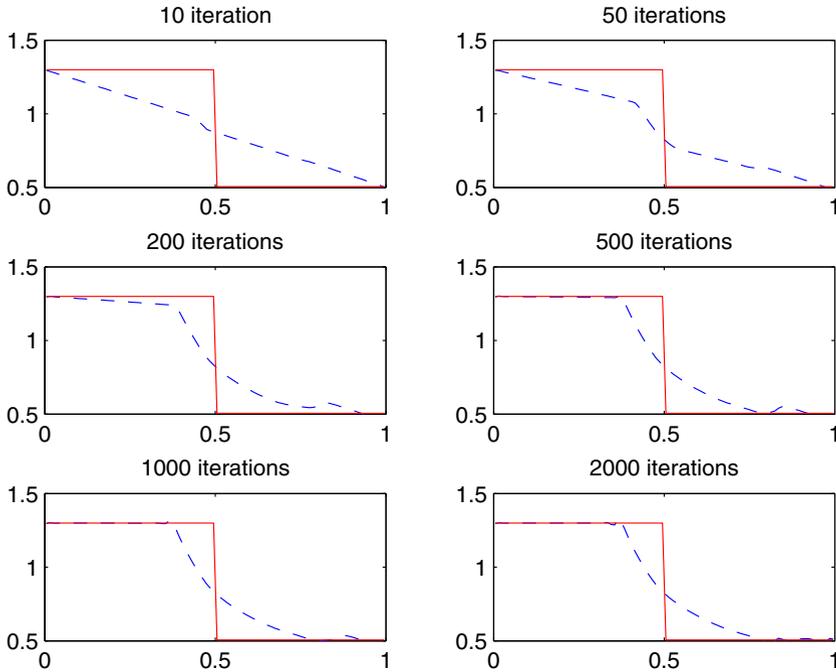
$$\begin{aligned} u_t + f(u)_x &= h(u, A), & x \in (0, 1), \quad t > 0, \\ u(x, 0) &= u_0(x), & x \in (0, 1), \\ u(0, t) &= u_\ell(1 + 0.15 \sin(4\pi t)), \quad u(1, t) = u_r, & t > 0, \end{aligned} \tag{4.11}$$

with  $f$  and  $h$  given by (4.10). The terminal velocity profile  $u_d$  is given by the solution of (4.11) with  $A = A_d$ . The optimization for  $A$  is then performed using a tracking-type functional on the full spatial and temporal grid and a regularization of  $A$ . In this paper, we consider a slightly different optimization problem: For a given  $A = A_d$  and a desired flow profile  $u_d$  we identify the initial flow condition  $u_0$ . Hence, we study an optimization problem (4.1a) subject to the constraint (4.11). In our numerical experiments, we choose the parameters  $\gamma = \frac{1}{6}$  and  $H = 1.2$  as in [16]. The terminal state  $u_d$  is obtained as the numerical solution of (4.11) (computed by the 1-D version of the second-order semi-discrete central-upwind scheme from Sect. 3.1) at time  $T = 0.15$  from Riemann initial data with

$$u(x, 0) = \begin{cases} u_\ell, & x < 0.5, \\ u_r, & x > 0.5. \end{cases} \tag{4.12}$$

The recovered initial condition  $u_0^{(m)}(x)$  is shown in Fig. 5. As in the previous nonsmooth example, the recovered initial condition is not unique since the terminal state is discontinuous. In Fig. 6, we plot the obtained solution of the optimal control problem (4.1a), (4.11). As one can see, the convergence in this example is much slower than on the previous ones, but after about 2000 iterations we recover  $u^{(2000)}(x, T = 0.15)$ , which almost coincides with  $u_d$ .

The presented results were obtained using the initial guess  $u_0^{(0)}(x) = u_\ell + (u_r - u_\ell)x$ , which satisfies the prescribed boundary condition at time  $t = 0$ . We have used a uniform finite-volume grid with  $\Delta x = 1/100$  and taken 400 characteristics points, which were placed uniformly in the interval  $(0, 1)$  at the terminal time  $T = 0.15$ . Notice that some of the traced characteristics curves leave the computational domain at times  $t \in (0, T)$ . This may lead to appearing and widening gaps at/next the shock area. We fill these gaps using the linear interpolation technique discussed in Remark 3.8.



**Fig. 5** Example 3: Recovered initial data of the optimal control problem,  $u_0^{(m)}(x, T = 0.15)$  with  $m = 10, 50, 200, 500, 1000$  and  $2000$  (dashed line) and  $u_0(x)$  given by (4.12) (solid line)

### 4.2 The two-dimensional case

We finally turn to the 2-D example.

#### Example 4: 2-D Burgers equation

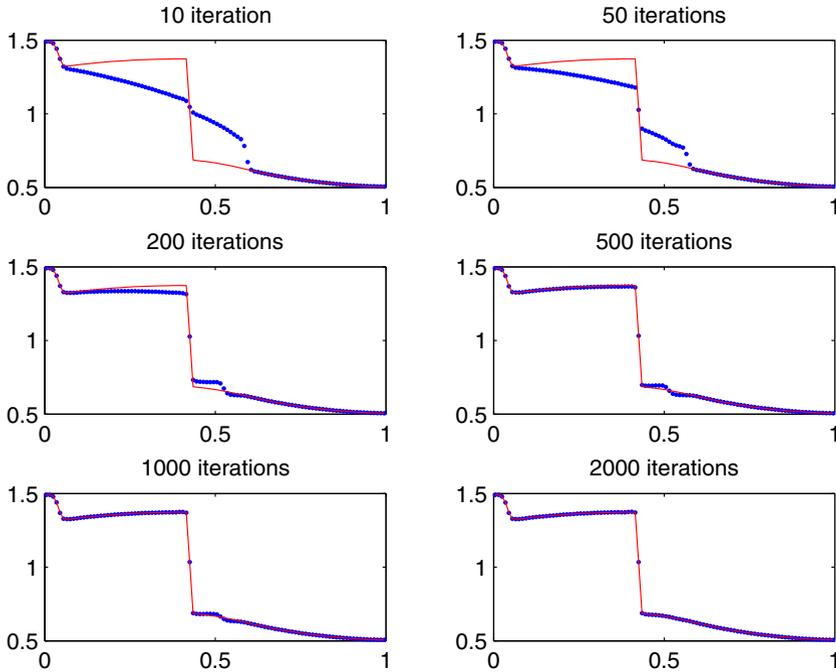
We numerically solve the 2-D optimization problem (1.1a)–(1.1c) subject to the inviscid Burgers equation:

$$u_t + \left(\frac{u^2}{2}\right)_x + \left(\frac{u^2}{2}\right)_y = 0. \tag{4.13}$$

The optimal control problem is solved in the domain  $[0, 2\pi] \times [0, 2\pi]$  with the period boundary conditions and the terminal state  $u_d$  obtained by a numerically solving (using the second-order semi-discrete central-upwind scheme described in Sect. 3.1) Eq. (4.13) subject to the following initial data:

$$u(x, y, 0) = \frac{1}{2} + \sin^2\left(\frac{1}{2}x\right) \sin^2\left(\frac{1}{2}y\right). \tag{4.14}$$

The solution was computed on a uniform finite-volume grid with  $\Delta x = \Delta y = 2\pi/100$ . We have started the Lagrangian method of characteristics with  $10^4$  points uniformly



**Fig. 6** Example 3: Recovered solution of the optimal control problem,  $u^{(m)}(x, T = 0.15)$  with  $m = 10, 50, 200, 500, 1000$  and  $2000$  (plotted with *points*) and the terminal state  $u_d$  (*solid line*)

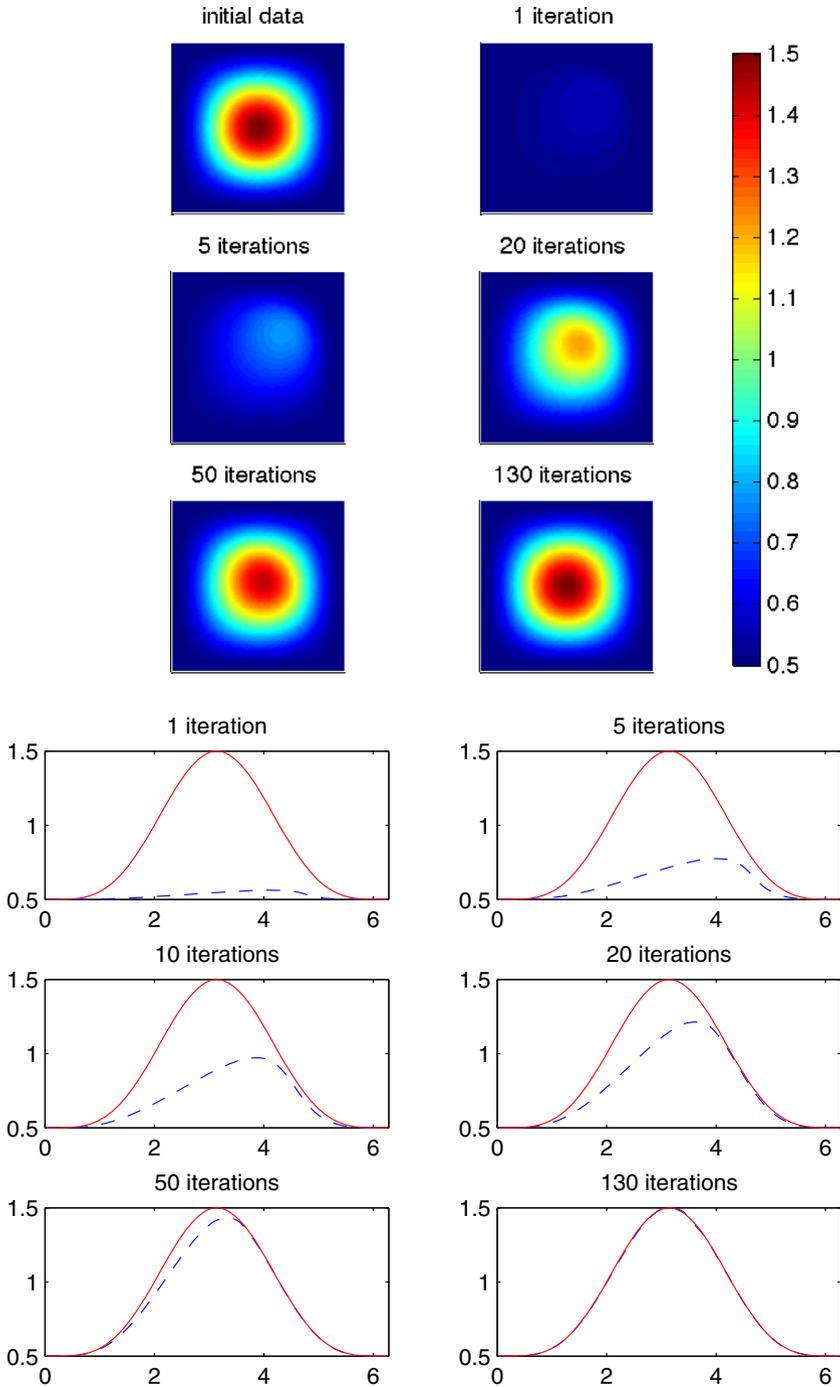
distributed throughout the computational domain at time  $t = T$ . We show 2-D plots of the optimal controls and the corresponding optimal states at times  $T = 1$  and  $T = 3$  (both the top view and the 1-D diagonal slice along the line  $y = x$ ). The results for  $T = 1$  are shown in Figs. 7 and 8, while the results for  $T = 3$  are presented in Figs. 9 and 10. In the former cases, when the terminal state is smooth, the solution of the optimization problem exhibits quite fast convergence in recovering both the initial data and the terminal state. In the latter case of a nonsmooth terminal state, the convergence is slower, and the control  $u_0$  given by (4.14) is not fully recovered due to lack of uniqueness. Nevertheless, the optimal state  $u^{(200)}(x, y, T = 3)$  in Fig. 10 almost coincides with  $u_d$ .

### 5 A convergence analysis

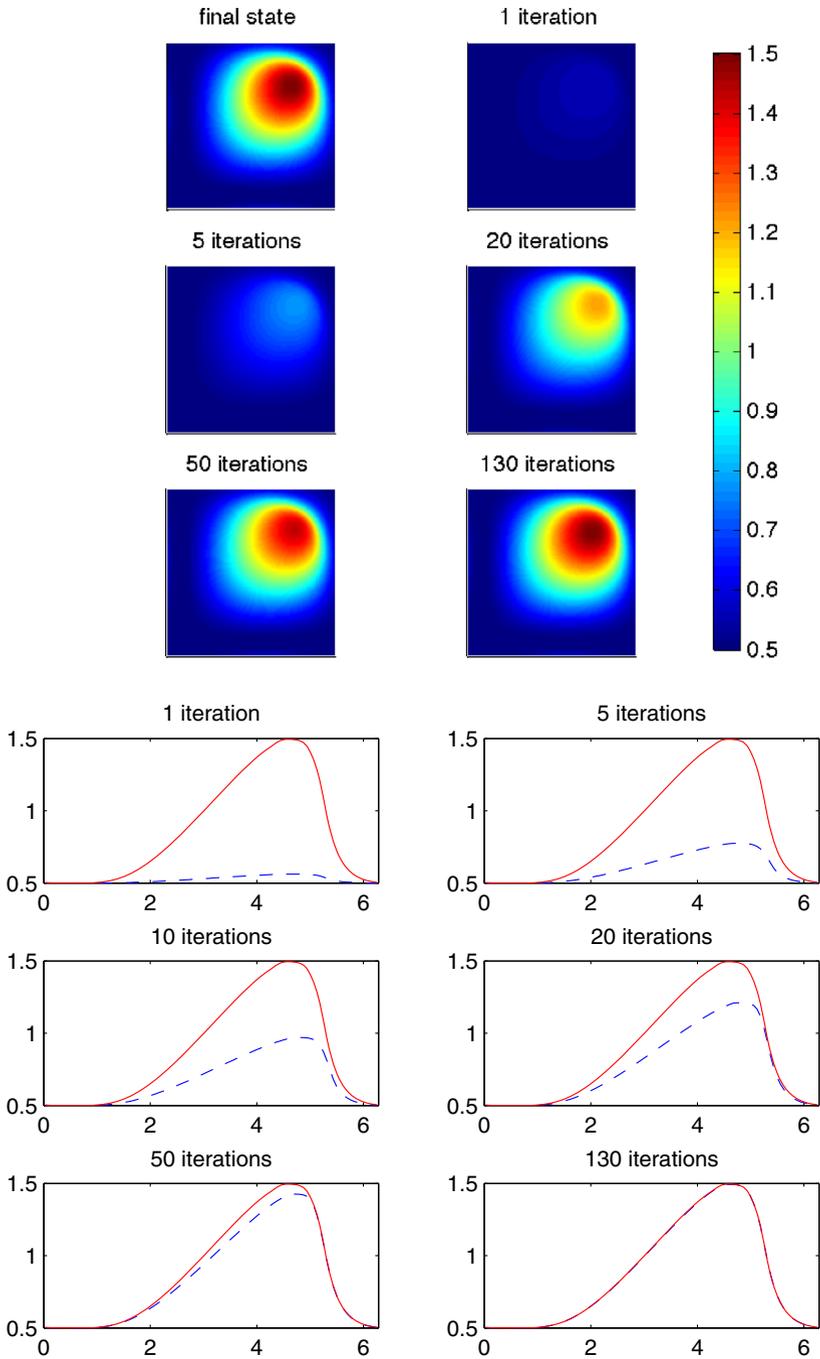
In this section, we discuss convergence properties of the proposed method in the 1-D case. Here, the derivation closely follows [44] and we apply the results from [44] to the presented scheme in order to proof its convergence.

To this end, we consider the problem (4.1b) in  $\mathbb{R}$  with no source term ( $h(u, x, t \equiv 0)$ ):

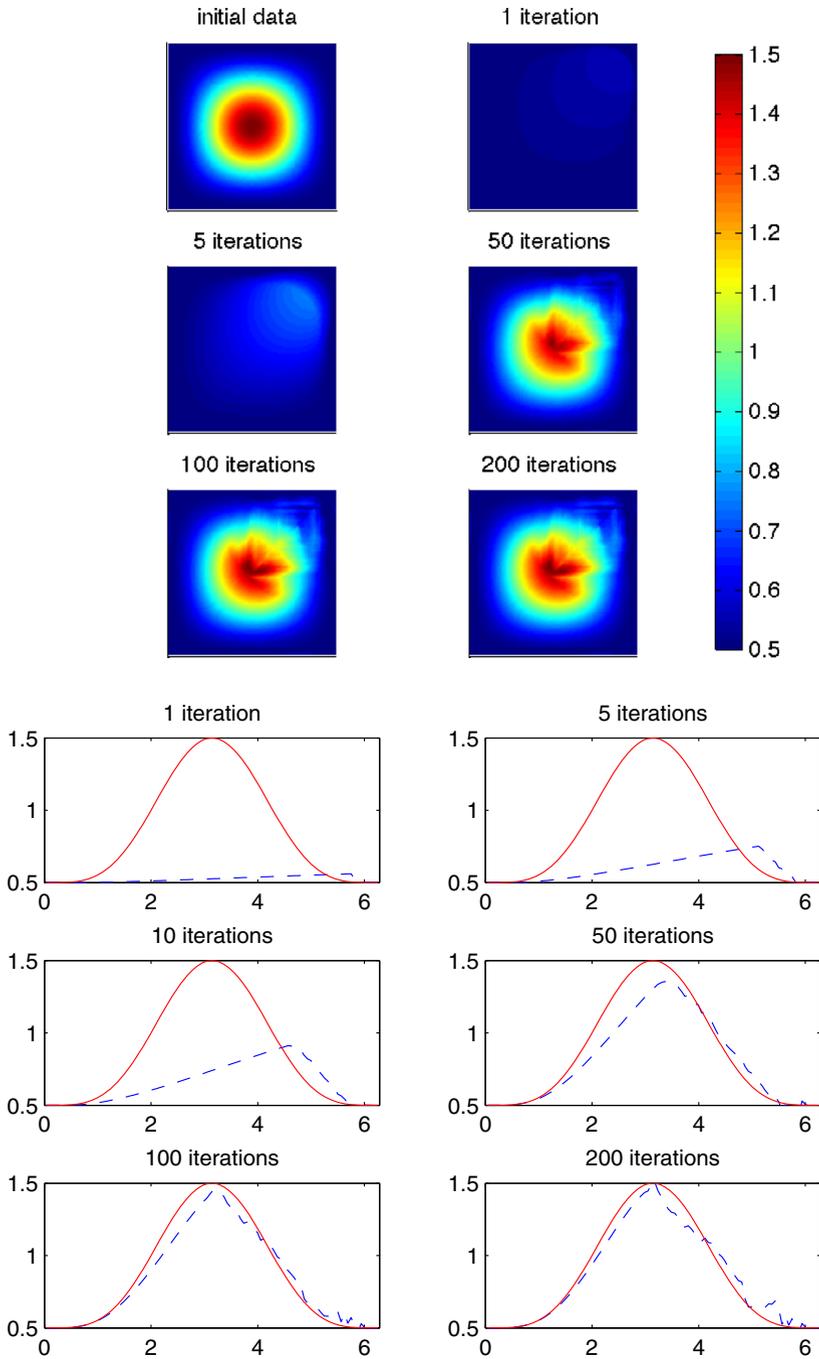
$$\begin{aligned}
 u_t + f(u)_x &= 0, & x \in \mathbb{R}, \quad t > 0, \\
 u(x, 0) &= u_0(x), & x \in \mathbb{R},
 \end{aligned}
 \tag{5.1}$$



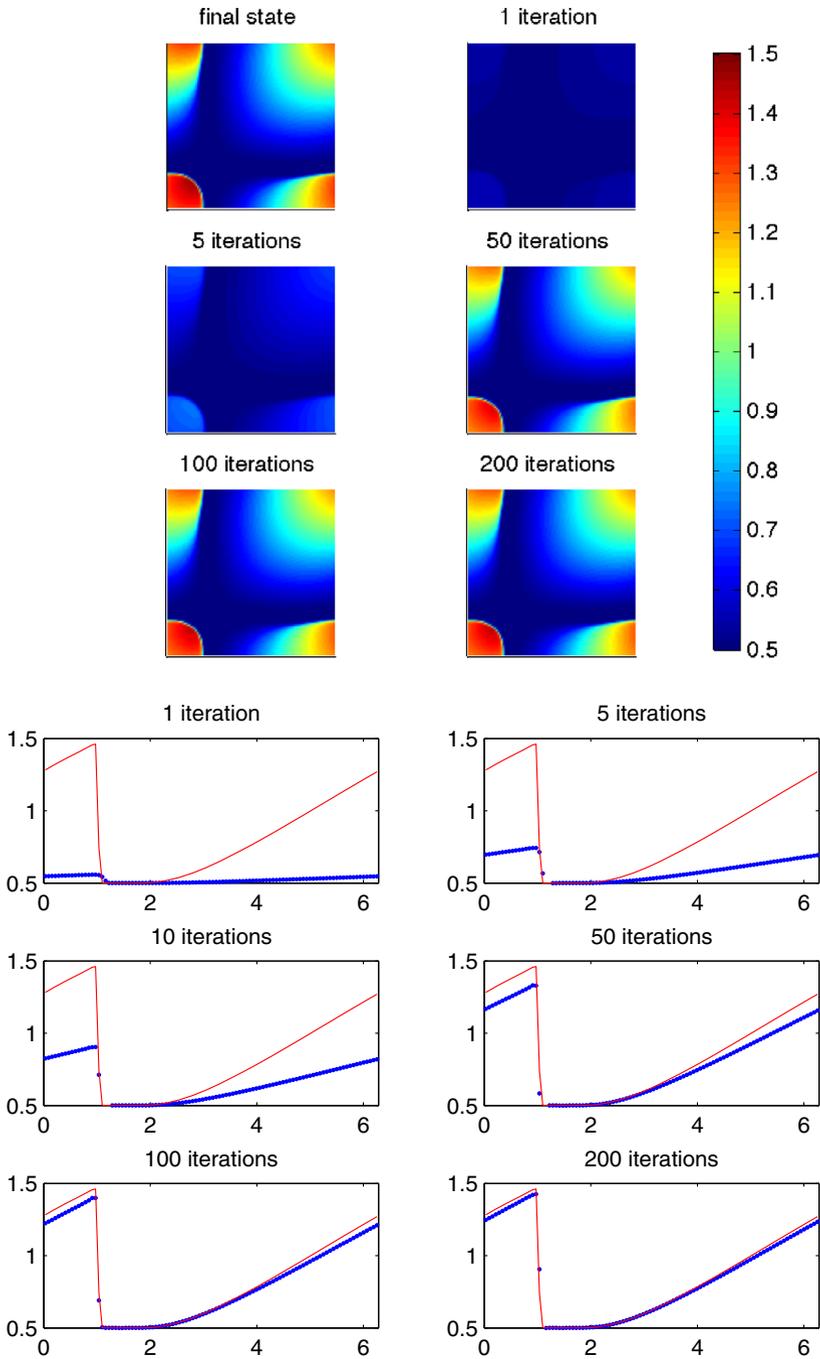
**Fig. 7** Example 4 ( $T = 1$ ): Recovered initial data  $u_0^{(m)}(x, y)$  with  $m = 1, 5, 20, 50$  and  $130$  and the exact initial data (4.14) (top left). The bottom part contains the corresponding plots along the diagonal  $y = x$



**Fig. 8** Example 4: Recovered solution of the optimal control problem,  $u^{(m)}(x, y, T = 1)$ , with  $m = 1, 5, 10, 20, 50$  and  $130$  and the terminal state  $u_d$  (top left). The bottom part contains the corresponding plots along the diagonal  $y = x$



**Fig. 9** Example 4 ( $T = 3$ ): Recovered initial data  $u_0^{(m)}(x, y)$  with  $m = 1, 5, 10, 50, 100$  and  $200$  and the exact initial data (4.14) (top left). The bottom part contains the corresponding plots along the diagonal  $y = x$



**Fig. 10** Example 4: Recovered solution of the optimal control problem,  $u^{(m)}(x, y, T = 3)$ , with  $m = 1, 5, 10, 50, 100$  and 200 and the terminal state  $u_d$  (top left). The bottom part contains the corresponding plots along the diagonal  $y = x$

and a strictly convex flux function  $f \in C^2(\mathbb{R})$ , that is,

$$f'' \geq c > 0 \quad \text{for some } c > 0. \tag{5.2}$$

The adjoint problem is then of the type

$$\begin{aligned} p_t + v(x, t)p_x &= 0, & x \in \mathbb{R}, \quad t \in (0, T), \\ p(x, T) &= p_T(x) & x \in \mathbb{R}, \end{aligned} \tag{5.3}$$

where  $v(x, t) = f'(u(x, t))$  in the case under consideration. It can be shown that if  $p_T$  is Lipschitz continuous and  $v \in L^\infty((0, T) \times \mathbb{R})$  is OSLC, that is, it satisfies,

$$v_x(\cdot, t) \leq \alpha(t), \quad \alpha \in L^1(0, T), \tag{5.4}$$

then there exists a reversible solution of (5.3). We refer, for instance, to [25, Definition 2.1, Theorem 2.2] and [44, Definition 3.7.2, Theorem 3.7.3] for the notion of reversible solutions. Further results on (5.3) can be found in [4].

*Remark 5.1* The adjoint equation (5.3) as well as the assumption (5.4) is well-defined for  $v(x, t) = f'(u(x, t))$  only if  $u$  does not contain any shocks. Therefore, the following convergence analysis is only applicable where the characteristics are *outside* the region entering a shock.

Under the above assumptions convergence properties within a general theory were established in [25,44] for some first-order numerical methods. The assumptions are restrictive but to the best of our knowledge no results for multidimensional problems, general flux functions or higher-order methods are available. Even though our numerical experiments clearly demonstrate the convergence of the proposed second-order Eulerian–Lagrangian method, we have been unable to formally prove this. Instead, we significantly simplify our method by adding a substantial amount of numerical diffusion to both the Eulerian and Lagrangian parts as described below, and prove the convergence for the simplified method only.

We first reduce the order of the method to the first one. The central-upwind scheme then becomes the HLL scheme [27], which can be written in the fully discrete form as

$$\bar{u}_j^{n+1} = \bar{u}_j^n - \lambda \left( F_{j+\frac{1}{2}}^n - F_{j-\frac{1}{2}}^n \right), \tag{5.5}$$

where  $\bar{u}_j^n \approx \frac{1}{\Delta x} \int_{x_{j-\frac{1}{2}}}^{x_{j+\frac{1}{2}}} u(x, t^n) dx$  are the computed cell averages,  $\lambda := \Delta t / \Delta x$  and the numerical fluxes are given by

$$F_{j+\frac{1}{2}}^n = \frac{a_{j+\frac{1}{2}}^+ f(\bar{u}_j^n) - a_{j+\frac{1}{2}}^- f(\bar{u}_{j+1}^n)}{a_{j+\frac{1}{2}}^+ - a_{j+\frac{1}{2}}^-} + \frac{a_{j+\frac{1}{2}}^+ a_{j+\frac{1}{2}}^-}{a_{j+\frac{1}{2}}^+ - a_{j+\frac{1}{2}}^-} (\bar{u}_{j+1}^n - \bar{u}_j^n), \tag{5.6}$$

with the following one-sided local speeds:

$$a_{j+\frac{1}{2}}^+ = \max \left\{ f'(\bar{u}_j^n), f'(\bar{u}_{j+1}^n), 0 \right\}, \quad a_{j+\frac{1}{2}}^- = \min \left\{ f'(\bar{u}_j^n), f'(\bar{u}_{j+1}^n), 0 \right\}. \tag{5.7}$$

The amount of numerical viscosity present in the HLL scheme (5.5)–(5.7) can be reduced by setting

$$a_{j+\frac{1}{2}}^+ = -a_{j+\frac{1}{2}}^- = a_{j+\frac{1}{2}} := \max \left\{ |f'(\bar{u}_j^n)|, |f'(\bar{u}_{j+1}^n)| \right\}, \tag{5.8}$$

which leads to the Rusanov scheme [41] with the numerical flux

$$F_{j+\frac{1}{2}}^n = \frac{1}{2} \left( f(\bar{u}_j^n) + f(\bar{u}_{j+1}^n) \right) - \frac{a_{j+\frac{1}{2}}}{2} (\bar{u}_{j+1}^n - \bar{u}_j^n). \tag{5.9}$$

However, the numerical diffusion in the Rusanov scheme (5.5), (5.8), (5.9) is still not large enough and we thus further reduce it by considering the modified Lax–Friedrichs scheme studied, for instance, in [44, Section 6.5.2]: We replace the local speeds  $a_{j+\frac{1}{2}}$  by the global ones, which results in the following numerical flux:

$$F_{j+\frac{1}{2}}^n = \frac{1}{2} \left( f(\bar{u}_j^n) + f(\bar{u}_{j+1}^n) \right) - \frac{\gamma}{2\lambda} (\bar{u}_{j+1}^n - \bar{u}_j^n), \quad \gamma = \lambda \max_j a_{j+\frac{1}{2}}. \tag{5.10}$$

Assuming that at time  $T$  the set of the characteristics points coincides with the finite-volume grid, the 1-D method of characteristics for the adjoint problem (5.3) is

$$\begin{cases} \frac{dx_i^c(t)}{dt} = f'(u(x_i^c(t), t)), & x_i^c(T) = x_i, \\ \frac{dp_i^c(t)}{dt} = 0, & p_i^c(T) = p_T(x_i^c(T)) = p_i^{N_T}. \end{cases} \tag{5.11}$$

Here, the value of  $u(x_i^c(t), t)$  is obtained from the piecewise constant reconstruction

$$\tilde{u}(x, \cdot) = \sum_j \chi_j(x) \bar{u}_j(\cdot), \tag{5.12}$$

that is,

$$u(x_i^c(t^n), t^n) = \tilde{u}(x_i^c(t^n), t^n) = \bar{u}_j^n,$$

provided the characteristic point  $x_i^c$  is located in cell  $j$  at time  $t^n$ . The ODE system (5.11) is then to be integrated backward in time starting from  $t = T$ .

In order to establish a convergence result, we modify the method of characteristics as follows. At the end of each time step, the obtained solution is first projected to the finite-volume cell centers  $\{x_j\}$  (this procedure adds a substantial amount of a numerical diffusion to the nondiffusive method of characteristics) and then the method

is restarted. This way, at time  $t = t^{n+1}$  each characteristics will start from the cell center and provided the CFL number is smaller than  $1/2$ , it will stay in the same cell at the end of the time step. Therefore, the new location of the characteristics can be obtained using the Euler method:

$$x_j^c(t^n) = x_j - \Delta t f_j^m, \text{ where } f_j^m := f'(\bar{u}_j^n). \tag{5.13}$$

The resulting Lagrangian method then reduces to a first-order Eulerian (finite-difference) one, which can be described as follows:

- If  $f_j^m = 0$ , then

$$p_j^n = p_j^{n+1}; \tag{5.14}$$

- If  $f_j^m > 0$ , then

$$p_j^n = (1 - \beta_j^n) p_j^{n+1} + \beta_j^n p_{j+1}^{n+1}, \tag{5.15}$$

where

$$\beta_j^n = \frac{\lambda f_j^m}{1 - \lambda(f_{j+1}^m - f_j^m)^+} = \frac{\lambda f_j^m}{1 - \lambda(\Delta^+ f_j^m)^+}; \tag{5.16}$$

- If  $f_j^m < 0$ , then

$$p_j^n = (1 - \gamma_j^n) p_j^{n+1} + \gamma_j^n p_{j-1}^{n+1}, \tag{5.17}$$

where

$$\gamma_j^n = \frac{-\lambda f_j^m}{1 - \lambda(f_j^m - f_{j-1}^m)^+} = \frac{-\lambda f_j^m}{1 - \lambda(\Delta^+ f_{j-1}^m)^+}. \tag{5.18}$$

In (5.16) and (5.18), we have used the standard notations  $\Delta^+ \omega_j := \omega_{j+1} - \omega_j$ ,  $\xi^+ := \max(\xi, 0)$ , and  $\xi^- := \min(\xi, 0)$ .

It should be observed that  $\forall j, n$ ,

$$\beta_j^n > 0, \quad \gamma_j^n > 0, \quad \frac{1}{2} \leq 1 - \lambda(\Delta^+ f_j^m)^+ \leq \frac{3}{2}, \quad \lambda(\Delta^+ f_j^m)^+ \leq \frac{1}{2}, \tag{5.19}$$

as long as the following CFL condition is imposed:

$$\lambda \leq \frac{1}{4 \max_{j,n} |f_j^m|}. \tag{5.20}$$

*Remark 5.2* Notice that if  $f_{j+1}^m \geq f_j^m$ , then (5.15), (5.16) is a result of the linear interpolation between the point values of  $p$  at the characteristic points  $x_j^c(t^n)$  and  $x_{j+1}^c(t^n)$ , while if  $f_{j+1}^m < f_j^m$ , then (5.15), (5.16) reduces to a simple first-order upwinding:

$$p_j^n = p_j^{n+1} + \lambda f_j^m (p_j^{n+1} - p_{j+1}^{n+1}).$$

A similar argument holds for (5.17), (5.18).

Next, we reformulate the above cases and rewrite equations (5.14)–(5.18) in the following compact form (see also [44, (6.45)] or [25, (3.9)]):

$$\begin{aligned}
 p_j^n &= \beta_j^n H(f_j^m) p_{j+1}^{n+1} + \gamma_j^n H(-f_j^m) p_{j-1}^{n+1} + \left(1 - \beta_j^n H(f_j^m) - \gamma_j^n H(-f_j^m)\right) p_j^n \\
 &= \sum_{l=-1}^1 B_j^{n,l} p_{j-l}^{n+1},
 \end{aligned}
 \tag{5.21}$$

where

$$\begin{aligned}
 B_j^{n,-1} &:= \lambda v_j^{n,0} = \beta_j^n H(f_j^m) = \lambda \frac{(f_j^m)^+}{1 - \lambda(\Delta^+ f_{j-1}^m)^+}, \\
 B_j^{n,1} &:= -\lambda v_{j-1}^{n,1} = \gamma_j^n H(-f_j^m) = -\lambda \frac{(f_j^m)^-}{1 - \lambda(\Delta^+ f_{j-1}^m)^+}, \\
 B_j^{n,0} &:= 1 + \lambda(v_{j-1}^{n,1} - v_j^{n,0}) = 1 - \gamma_j^n H(-f_j^m) - \beta_j^n H(f_j^m),
 \end{aligned}
 \tag{5.22}$$

and  $H$  is the Heaviside function. We also note that the difference  $\Delta^+ p_j^n$  can be written as

$$\Delta^+ p_j^n = \sum_{l=-1}^1 D_j^{n,l} \Delta^+ p_{j-l}^{n+1},
 \tag{5.23}$$

where

$$D_j^{n,-1} := \lambda v_{j+1}^{n,0}, \quad D_j^{n,1} := -\lambda v_{j-1}^{n,1}, \quad D_j^{n,0} := 1 + \lambda(v_{j-1}^{n,1} - v_j^{n,0}).
 \tag{5.24}$$

Taking into account (5.19) and the CFL condition (5.20), it is easy to check that

$$B_j^{n,l} \geq 0, \quad D_j^{n,l} \geq 0, \quad \forall j, n \text{ and } l \in \{-1, 0, 1\}.
 \tag{5.25}$$

In what follows, we prove the convergence of the discrete scheme (5.21), (5.22) towards the solution of (5.3) with  $v = f'(u)$ . We follow the lines of [25,44] and assume that:

(A1) There exist constants  $\Delta_0, M_u > 0$  such that for all  $\Delta t = \lambda \Delta x \leq \Delta_0$ , we have

$$\|\tilde{u}\|_\infty \leq M_u, \quad \tilde{u}(\cdot, t) \rightarrow u(\cdot, t) \text{ in } L^1_{loc}(\mathbb{R}) \quad \forall t \in [0, T],
 \tag{5.26}$$

where  $\tilde{u}$  is given by (5.12) and  $u$  is the entropy solution of (5.1), (5.2);

(A2) There exists a function  $\mu \in L^1(0, T)$  such that for all  $\Delta t = \lambda \Delta x \leq \Delta_0$ , the discrete OSLC holds, that is,

$$\Delta^+ \tilde{u}_j^n \leq \lambda \int_{t^n}^{t^{n+1}} \mu(s) ds \quad \forall j, n.
 \tag{5.27}$$

It was proved in [44] that solutions obtained by the scheme (5.5), (5.10) satisfy the assumptions (A1) and (A2). The OSLC property and rigorous convergence rate estimates were established for several first-order [38, 39] and second-order [30] schemes. Numerous numerical experiments reported in the literature suggest that solutions of both the HLL (5.5)–(5.7) and Rusanov (5.5), (5.9) schemes also satisfy the assumptions (A1) and (A2). However, to the best of our knowledge, no rigorous proof of this is available.

Equipped with the above assumptions, we apply [44, Theorem 6.3.4] to our scheme. To this end, we verify that the discrete functions  $v_j^{n,l}$  in (5.22) fulfill the following conditions which are similar to the assumptions (A1) and (A2) (see [25, Theorem 2.4]):

(A3) There exist a constant  $M_a > 0$  such that for all  $\Delta t = \lambda \Delta x \leq \Delta_0$

$$\|\tilde{v}^l\|_\infty \leq M_a, \quad \tilde{v} = \sum_{l=0}^1 \tilde{v}^l \rightarrow v = f'(u) \text{ in } L^1_{\text{loc}}((0, T) \times \mathbb{R}) \text{ as } \Delta x \rightarrow 0. \tag{5.28}$$

Here,

$$\tilde{v}^l(x, t) := \sum_{j,n} v_j^{n,l} \chi_{Q_j^n}(x, t),$$

where  $\chi_{Q_j^n}$  is a characteristic function of the space-time volume  $Q_j^n := [x_{j-\frac{1}{2}}, x_{j+\frac{1}{2}}) \times [t^n, t^{n+1})$ ;

(A4) There exists a function  $\alpha \in L^1(0, T)$  such that for all  $\Delta t = \lambda \Delta x \leq \Delta_0$

$$\sum_{l=0}^1 \Delta^+ v_j^{n,l} \leq \lambda \int_{t^n}^{t^{n+1}} \alpha(s) ds \quad \forall j, n. \tag{5.29}$$

We start by proving (5.28) and note that

$$\|\tilde{v}^l\|_\infty \leq 2 \max_{j,n} |f''_j| =: M_a < \infty,$$

since  $\|\tilde{u}\|_\infty$  is bounded, and

$$\tilde{v}(x, t) = \sum_{l=0}^1 \tilde{v}^l(x, t) = \sum_{j,n} \frac{(f''_j)^+ + (f''_{j+1})^-}{1 - \lambda(\Delta^+ f''_j)^+} \chi_{Q_j^n}(x, t).$$

Taking into account (5.26), (5.27), we have that for all  $R > 0$  and all  $\Delta x = \lambda \Delta x \leq \lambda$

$$\begin{aligned}
 & \| \tilde{u}(\cdot + \Delta x, \cdot) - u \|_{L^1((-R,R) \times (0,T))} \\
 &= \| \tilde{u}(\cdot + \Delta x, \cdot) - u(\cdot + \Delta x, \cdot) + u(\cdot + \Delta x, \cdot) - u \|_{L^1((-R,R) \times (0,T))} \\
 &\leq \| \tilde{u} - u \|_{L^1((-R-1,R+1) \times (0,T))} + \| u(\cdot + \Delta x, \cdot) \\
 &\quad - u \|_{L^1((-R,R) \times (0,T))} \rightarrow 0 \text{ as } \Delta x \rightarrow 0.
 \end{aligned}
 \tag{5.30}$$

We also have

$$\begin{aligned}
 \left| f'(u) - \sum_{l=0}^1 \tilde{v}^l \right| &\leq \sum_{j,n} \left[ \left| (f'(u))^+ - v_j^{n,0} \right| + \left| (f'(u))^- - v_j^{n,1} \right| \right] \chi_{Q_j^n}(x, t) \\
 &\leq 2 \sum_{j,n} \left[ \left| (f'(u))^+ - (f''_j)^+ \right| + \lambda \left| (f'(u))^+ (\Delta^+ f''_j)^+ \right| \right. \\
 &\quad \left. + \left| (f'(u))^- - (f''_{j+1})^- \right| + \lambda \left| (f'(u))^- (\Delta^+ f''_j)^+ \right| \right] \chi_{Q_j^n}(x, t)
 \end{aligned}$$

and

$$\begin{aligned}
 \left| (f'(u))^+ (\Delta^+ f''_j)^+ \right| + \left| (f'(u))^- (\Delta^+ f''_j)^+ \right| &\leq 2 \| f'(u) \|_\infty |f''_{j+1} - f''_j| \\
 &\leq 2 \| f'(u) \|_\infty \max_{\tilde{u}} \{ |f''(\tilde{u})| \} |\bar{u}^n_{j+1} - \bar{u}^n_j|
 \end{aligned}$$

provided  $u \in L^\infty(\mathbb{R} \times (0, T))$  and  $\tilde{u}$  satisfies (5.26) and (5.27). We then denote by  $M_0 := \|u\|_\infty$ ,  $M_1 := \max_{|u| \leq \max\{M_0, M_u\}} |f'(u)|$ , and  $M_2 := \max_{|u| \leq \max\{M_0, M_u\}} |f''(u)|$  to obtain

$$\left| f'(u) - \sum_{l=0}^1 \tilde{v}^l \right| \leq 4M_2 \sum_{j,n} \left[ |u - \bar{u}^n_j| + \lambda M_1 |\bar{u}^n_{j+1} - \bar{u}^n_j| \right] \chi_{Q_j^n}(x, t).$$

The last inequality together with (5.30) lead to (5.28) since  $\tilde{u}(x + \Delta x, t) = \sum_{j,n} \bar{u}^n_{j+1} \chi_{Q_j^n}(x, t)$ .

It remains to prove (5.29). Let  $\alpha \in L^1(0, T)$ ,  $\Delta_0 > 0$  such that (5.26) and (5.27) hold with  $\mu(t) = \alpha(t)$ , and assume without lost of generality that  $\alpha(t) \geq 0, \forall t$ . We also denote by

$$\begin{aligned}
 \Omega &:= \Delta^+ v_j^{n,0} + \Delta^+ v_{j-1}^{n,1} \\
 &= \frac{(f''_{j+1})^+}{1 - \lambda(\Delta^+ f''_{j+1})^+} - \frac{(f''_j)^+}{1 - \lambda(\Delta^+ f''_j)^+} + \frac{(f''_{j+1})^-}{1 - \lambda(\Delta^+ f''_j)^+} \\
 &\quad - \frac{(f''_j)^-}{1 - \lambda(\Delta^+ f''_{j-1})^+}.
 \end{aligned}$$

To establish the required estimate, we distinguish between several cases, which are treated using (5.19) and (5.20) to obtain:

- If  $f'_{j+1} \geq 0$  and  $f'_j \leq 0$ , then

$$\begin{aligned}\Omega &= \frac{f'_{j+1}}{1 - \lambda(\Delta^+ f'_{j+1})^+} - \frac{f'_j}{1 - \lambda(\Delta^+ f'_{j-1})^+} \\ &\leq 2 \left( f'_{j+1} - f'_j \right) \leq 2M_2(\Delta^+ \bar{u}_j^+)^+;\end{aligned}$$

- If  $f'_{j+1} \geq 0$ ,  $f'_j > 0$  and  $f'_{j+1} > f'_j$ , then

$$\begin{aligned}\Omega &= \frac{f'_{j+1}}{1 - \lambda(\Delta^+ f'_{j+1})^+} - \frac{f'_j}{1 - \lambda\Delta^+ f'_j} \\ &= \frac{(1 - \lambda f'_{j+1})\Delta^+ f'_j + \lambda f'_j(\Delta^+ f'_{j+1})^+}{(1 - \lambda\Delta^+(f'_{j+1})^+)(1 - \lambda\Delta^+ f'_j)} \\ &\leq M_2 \left[ 4(\Delta^+ \bar{u}_j^+)^+ + (\Delta^+ \bar{u}_{j+1}^+)^+ \right];\end{aligned}$$

- If  $f'_{j+1} \geq 0$ ,  $f'_j > 0$  and  $f'_{j+1} \leq f'_j$ , then

$$\Omega = \frac{f'_{j+1}}{1 - \lambda(\Delta^+ f'_{j+1})^+} - f'_j \leq \frac{f'_j}{1 - \lambda(\Delta^+ f'_{j+1})^+} - f'_j \leq \frac{1}{2}M_2(\Delta^+ \bar{u}_{j+1}^+)^+;$$

- If  $f'_{j+1} < 0$ ,  $f'_j > 0$ , then

$$\Omega = \frac{f'_{j+1}}{1 - \lambda(\Delta^+ f'_{j+1})^+} - \frac{f'_j}{1 - \lambda(\Delta^+ f'_{j-1})^+} < 0;$$

- If  $f'_{j+1} < 0$ ,  $f'_j \leq 0$  and  $f'_{j+1} \leq f'_j$ , then

$$\begin{aligned}\Omega &= f'_{j+1} - \frac{f'_j}{1 - \lambda(\Delta^+ f'_{j-1})^+} \leq f'_j - \frac{f'_j}{1 - \lambda(\Delta^+ f'_{j-1})^+} \\ &\leq \frac{1}{2}M_2(\Delta^+ \bar{u}_{j-1}^+)^+;\end{aligned}$$

- If  $f''_{j+1} < 0$ ,  $f''_j \leq 0$  and  $f''_{j+1} > f''_j$ , then

$$\begin{aligned} \Omega &= \frac{f''_{j+1}}{1 - \lambda \Delta^+ f''_j} - \frac{f''_j}{1 - \lambda(\Delta^+ f''_{j-1})^+} \\ &= \frac{(1 + \lambda f''_j) \Delta^+ f''_j - \lambda f''_{j+1} (\Delta^+ f''_{j-1})^+}{(1 - \lambda \Delta^+ f''_j)(1 - \lambda \Delta^+ (f''_{j-1})^+)} \\ &\leq M_2 \left[ 4(\Delta^+ \bar{u}^n_j)^+ + (\Delta^+ \bar{u}^n_{j-1})^+ \right]; \end{aligned}$$

Assuming as before that  $u$  and  $\tilde{u}$  satisfy (5.26), (5.27), we conclude that in all of the cases  $\Omega$  can be estimated by

$$\Omega \leq c\lambda \int_{t^n}^{t^{n+1}} \alpha(s) ds$$

with  $c = c(M_2)$ .

The established estimates together with [44, Theorem 6.3.4] lead to the following convergence result.

**Theorem 5.1** *Assume that  $f \in C^2(\mathbb{R})$  satisfies the assumption (5.2). Let  $p_T \in Lip(\mathbb{R})$  and  $u \in L^\infty(\mathbb{R} \times (0, T))$  and satisfy (5.26), (5.27). Assume that the discretization of  $p_T$  is consistent, that is, there exist constants  $M_T > 0$  and  $L_T > 0$  such that*

$$\|\tilde{p}_T\|_\infty \leq M_T, \quad \sup_{x \in \mathbb{R}} \left| \frac{\tilde{p}_T(x + \Delta x) - \tilde{p}_T(x)}{\Delta x} \right| \leq L_T, \tag{5.31}$$

$\tilde{p}_T \rightarrow p_T$  in  $[-R, R]$  for all  $R > 0$  as  $\Delta x \rightarrow 0$ .

Assume also that the CFL condition (5.20) holds. Then, the solution to the adjoint scheme (5.21) converges locally uniformly to the unique reversible solution  $p \in Lip(\mathbb{R} \times [0, T])$  of (5.3) with  $v = f'(u)$ , that is,

$$\tilde{p} \rightarrow p \text{ in } [-R, R] \times [0, T] \text{ for all } R > 0 \text{ as } \Delta t = \lambda \Delta x \rightarrow 0.$$

Here,  $\tilde{p}_T$  and  $\tilde{p}$  are piecewise constant approximations of  $p_T(x)$  and the computed solution  $\{p^n_j\}$ , respectively.

*Remark 5.3* While equations (5.26), (5.27) and (5.31) mimic the assumptions (D2), (D3) and (C1) in [44], it should be observed that (5.27) can be weakened the same way it was done in [44]. Similarly, one could use [25, Theorem 3.7] to establish the convergence result for the adjoint equation.

The previous simplifications allow to state a convergence result for the gradients for a smooth version of the cost functional (1.1c). For a given nonnegative function  $\phi_\delta \in Lip_0(\mathbb{R})$  with the support in  $[-\frac{\delta}{2}, \frac{\delta}{2}]$  and  $\int_{\mathbb{R}} \phi_\delta(x) dx = 1$  and  $\psi \in C^1_{loc}(\mathbb{R}^2)$ ,

let  $J_\delta$  be given by

$$J_\delta(u_0) := \int_0^1 \psi((\phi_\delta * u)(x, T), (\phi_\delta * u_d)(x)) dx. \tag{5.32}$$

Here,  $u$  is the entropy solution of the initial value problem (5.1) and  $*$  denotes a convolution in  $x$ . For  $J_\delta$  to be well-posed we assume  $u_d \in L^\infty(\mathbb{R})$ . We discretize  $J_\delta$  by

$$\tilde{J}_\delta(\tilde{u}_0) = \sum_k \psi((\phi_\delta * \tilde{u})(x_k, T), (\phi_\delta * \tilde{u}_d)(x_k)) \Delta x, \tag{5.33}$$

where  $(\tilde{\cdot})$  denotes a corresponding piecewise polynomial (piecewise constant for first-order methods) approximation.

The gradient of  $J_\delta$  exists in the sense of Frechet differentials, see [44, Theorem 5.3.1]. Using Theorem 5.1 and [44, Theorem 6.4.8], one may obtain the following convergence result.

**Theorem 5.2** *Assume that  $f \in C^2(\mathbb{R})$  satisfies the assumption (5.2). Let  $J_\delta$  be given by (5.32) and assume that  $u_0 \in L^\infty(\mathbb{R} \times (0, T))$  such that  $(u_0)_x \leq K$ . Let*

$$\tilde{p}(x_j, T) = \sum_k \phi_\delta(x_j - x_k) \partial_1 \psi((\phi_\delta * \tilde{u})(x_k, T), (\phi_\delta * \tilde{u}_d)(x_k)) \Delta x,$$

where  $\partial_1 \psi$  denotes a partial derivative of  $\psi$  with respect to its first component.

Let  $\tilde{u}(x, \cdot)$  be an approximate solution of (5.1) obtained by (5.5), (5.10), (5.12) and  $\tilde{p}$  be a piecewise constant approximation of the solution computed by (5.21). Let the CFL conditions (5.20) and

$$\lambda \max_{|u| \leq \|u_0\|_\infty} |f'(u)| \leq \min\{(1 - \rho) \min(\gamma, 2 - 2\gamma), 1 - \gamma\}$$

hold for some fixed value of  $\rho \in (0, 1)$ .

Then,  $\tilde{p}(\cdot, 0)$  is an approximation to the Frechet derivative of  $J_\delta$  with respect to  $u_0$  in the following sense:

$$\tilde{p}(\cdot, 0) \rightarrow p(\cdot, 0) = \nabla J_\delta(u_0) \text{ in } L^r(\mathbb{R}) \text{ as } \Delta t = \lambda \Delta x \rightarrow 0,$$

for all  $r \geq 1$ . Herein,  $p$  is the reversible solution of (5.3) with the terminal data

$$p_T(x) = \int_0^1 \phi_\delta(x - z) \partial_1 \psi((\phi_\delta * u)(z, T), (\phi_\delta * u_d)(z)) dz.$$

*Remark 5.4* A similar result holds under the assumption that  $(u_0)_x|_{\mathbb{R} \setminus E} \leq K$  for some closed set  $E$ , see [44, Chapter 6] for more details.

**Acknowledgments** The work of A. Chertock was supported in part by the NSF Grants DMS-0712898 and DMS-1115682. The work of M. Herty was supported by the DAAD 54365630, 55866082, EXC128. The work of A. Kurganov was supported in part by the NSF Grant DMS-1115718 and the German Research Foundation DFG under the Grant No. INST 247/609-1. The authors also acknowledge the support by NSF RNMS Grant DMS-1107444.

## References

1. Baines, M., Cullen, M., Farmer, C., Giles, M., Rabbitt, M. (eds.): 8th ICFD Conference on Numerical Methods for Fluid Dynamics. Part 2. Wiley, Chichester (2005). Papers from the Conference held in Oxford, 2004, *Int. J. Numer. Methods Fluids* **47**(10–11) (2005).
2. Banda, M., Herty, M.: Adjoint IMEX-based schemes for control problems governed by hyperbolic conservation laws. *Comput. Optim. Appl.* **51**(2), 909–930 (2012)
3. Bianchini, S.: On the shift differentiability of the flow generated by a hyperbolic system of conservation laws. *Discret. Contin. Dyn. Syst.* **6**, 329–350 (2000)
4. Bouchut, F., James, F.: One-dimensional transport equations with discontinuous coefficients. *Nonlinear Anal.* **32**, 891–933 (1998)
5. Bouchut, F., James, F.: Differentiability with respect to initial data for a scalar conservation law. In: *Hyperbolic Problems: Theory, Numerics, Applications*, Vol. I (Zürich, 1998), vol. 129. *Int. Ser. Numer. Math.* Birkhäuser, Basel, pp. 113–118 (1999)
6. Bouchut, F., James, F.: Duality solutions for pressureless gases, monotone scalar conservation laws, and uniqueness. *Commun. Partial Differ. Equ.* **24**, 2173–2189 (1999)
7. Bouchut, F., James, F., Mancini, S.: Uniqueness and weak stability for multi-dimensional transport equations with one-sided Lipschitz coefficient. *Ann. Sc. Norm. Super. Pisa Cl. Sci* **5**(4), 1–25 (2005)
8. Bressan, A., Guerra, G.: Shift-differentiability of the flow generated by a conservation law. *Discret. Contin. Dyn. Syst.* **3**, 35–58 (1997)
9. Bressan, A., Lewicka, M.: Shift differentials of maps in BV spaces. In: *Nonlinear Theory of Generalized Functions* (Vienna, 1997), vol. 401. *Chapman & Hall/CRC Res. Notes Math.* Chapman & Hall/CRC, Boca Raton, FL, pp. 47–61 (1999)
10. Bressan, A., Marson, A.: A variational calculus for discontinuous solutions to conservation laws. *Commun. Partial Differ. Equ.* **20**, 1491–1552 (1995)
11. Bressan, A., Shen, W.: Optimality conditions for solutions to hyperbolic balance laws, control methods in PDE-dynamical systems. *Contemp. Math.* **426**, 129–152 (2007)
12. Calamai, P., Moré, J.: Projected gradient methods for linearly constrained problems. *Math. Program.* **39**, 93–116 (1987)
13. Castro, C., Palacios, F., Zuazua, E.: An alternating descent method for the optimal control of the inviscid Burgers equation in the presence of shocks. *Math. Models Methods Appl. Sci.* **18**, 369–416 (2008)
14. Chertock, A., Kurganov, A.: On a hybrid finite-volume particle method, *M2AN Math. Model. Numer. Anal.* **38**, 1071–1091 (2004)
15. Chertock, A., Kurganov, A.: On a practical implementation of particle methods. *Appl. Numer. Math.* **56**, 1418–1431 (2006)
16. Cliff, E., Heinkenschloss, M., Shenoy, A.: An optimal control problem for flows with discontinuities. *J. Optim. Theory Appl.* **94**, 273–309 (1997)
17. Frank, P., Subin, G.: A comparison of optimization-based approaches for a model computational aerodynamics design problem. *J. Comput. Phys.* **98**, 74–89 (1992)
18. Giles, M., Ulbrich, S.: Convergence of linearized and adjoint approximations for discontinuous solutions of conservation laws. Part 1: Linearized approximations and linearized output functionals. *SIAM J. Numer. Anal.* **48**, 882–904 (2010)
19. Giles, M., Ulbrich, S.: Convergence of linearized and adjoint approximations for discontinuous solutions of conservation laws. Part 2: Adjoint approximations and extensions. *SIAM J. Numer. Anal.* **48**, 905–921 (2010)
20. Giles, M.B.: Analysis of the accuracy of shock-capturing in the steady quasi 1d-euler equations. *Int. J. Comput. Fluid Dynam.* **5**, 247–258 (1996)
21. Giles, M.B.: Discrete adjoint approximations with shocks. In: *Hyperbolic Problems: Theory, Numerics, Applications*, pp. 185–194. Springer, Berlin (2003)

22. Giles, M.B., Pierce, N.A.: Analytic adjoint solutions for the quasi-one-dimensional Euler equations. *J. Fluid Mech.* **426**, 327–345 (2001)
23. Giles, M.B., Pierce, N.A.: Adjoint error correction for integral outputs. In: *Error Estimation and Adaptive Discretization Methods in Computational Fluid Dynamics*, vol. 25, Lect. Notes Comput. Sci. Eng., Springer, Berlin, pp. 47–95 (2003)
24. Giles, M.B., Süli, E.: Adjoint methods for PDEs: a posteriori error analysis and postprocessing by duality. *Acta Numer.* **11**, 145–236 (2002)
25. Gosse, L., James, F.: Numerical approximations of one-dimensional linear conservation equations with discontinuous coefficients. *Math. Comput.* **69**, 987–1015 (2000)
26. Gottlieb, S., Shu, C.-W., Tadmor, E.: Strong stability-preserving high-order time discretization methods. *SIAM Rev.* **43**, 89–112 (2001)
27. Harten, A., Lax, P., van Leer, B.: On upstream differencing and Godunov-type schemes for hyperbolic conservation laws. *SIAM Rev.* **25**, 35–61 (1983)
28. James, F., Sepúlveda, M.: Convergence results for the flux identification in a scalar conservation law. *SIAM J. Control Optim.* **37**, 869–891 (1999) (electronic)
29. Kelley, C.: *Iterative methods for optimization*. Frontiers in Applied Mathematics. Philadelphia, PA: Society for Industrial and Applied Mathematics. xv 180 p (1999)
30. Kurganov, A.: *Conservation Laws: Stability of Numerical Approximations and Nonlinear Regularization*, PhD Dissertation, Tel-Aviv University, School of Mathematical Sciences (1998)
31. Kurganov, A., Lin, C.-T.: On the reduction of numerical dissipation in central-upwind schemes. *Commun. Comput. Phys.* **2**, 141–163 (2007)
32. Kurganov, A., Noelle, S., Petrova, G.: Semi-discrete central-upwind scheme for hyperbolic conservation laws and Hamilton–Jacobi equations. *SIAM J. Sci. Comput.* **23**, 707–740 (2001)
33. Kurganov, A., Tadmor, E.: New high resolution central schemes for nonlinear conservation laws and convection-diffusion equations. *J. Comput. Phys.* **160**, 241–282 (2000)
34. Kurganov, A., Tadmor, E.: Solution of two-dimensional riemann problems for gas dynamics without riemann problem solvers. *Numer. Methods Partial Differ. Equ.* **18**, 584–608 (2002)
35. Lie, K.-A., Noelle, S.: On the artificial compression method for second-order nonoscillatory central difference schemes for systems of conservation laws. *SIAM J. Sci. Comput.* **24**, 1157–1174 (2003)
36. Liu, Z., Sandu, A.: On the properties of discrete adjoints of numerical methods for the advection equation. *Int. J. Numer. Methods Fluids* **56**, 769–803 (2008)
37. Nessyahu, H., Tadmor, E.: Nonoscillatory central differencing for hyperbolic conservation laws. *J. Comput. Phys.* **87**, 408–463 (1990)
38. Nessyahu, H., Tadmor, E.: The convergence rate of approximate solutions for nonlinear scalar conservation laws. *SIAM J. Numer. Anal.* **29**, 1505–1519 (1992)
39. Nessyahu, H., Tadmor, E., Tassa, T.: The convergence rate of Godunov type schemes. *SIAM J. Numer. Anal.* **31**, 1–16 (1994)
40. Pierce, N.A., Giles, M.B.: Adjoint and defect error bounding and correction for functional estimates. *J. Comput. Phys.* **200**, 769–794 (2004)
41. Rusanov, V.: The calculation of the interaction of non-stationary shock waves with barriers. *Ž. Vyčisl. Mat. i Mat. Fiz.* **1**, 267–279 (1961)
42. Spellucci, P.: *Numerical Methods of Nonlinear Optimization (Numerische Verfahren der nichtlinearen Optimierung)*. ISNM Lehrbuch. Basel: Birkhäuser. 576 S (1993)
43. Sweby, P.: High resolution schemes using flux limiters for hyperbolic conservation laws. *SIAM J. Numer. Anal.* **21**, 995–1011 (1984)
44. Ulbrich, S.: *Optimal control of nonlinear hyperbolic conservation laws with source terms*, Habilitation thesis, Fakultät für Mathematik, Technische Universität München, <http://www3.mathematik.tu-darmstadt.de/hp/optimierung/ulbrich-stefan/> (2001)
45. Ulbrich, S.: Adjoint-based derivative computations for the optimal control of discontinuous solutions of hyperbolic conservation laws. *Syst. Control Lett.* **48**, 313–328 (2003)
46. Ulbrich, S.: On the superlinear local convergence of a filter-sqp method. *Math. Program. Ser. B* **100**, 217–245 (2004)
47. van Leer, B.: Towards the ultimate conservative difference scheme. V. A second-order sequel to Godunov’s method. *J. Comput. Phys.* **32**, 101–136 (1979)