

Weierstraß-Institut
für Angewandte Analysis und Stochastik
Leibniz-Institut im Forschungsverbund Berlin e. V.

Preprint

ISSN 2198-5855

**On the consistency of Runge–Kutta methods up to order three
applied to the optimal control of scalar conservation laws**

Michael Hintermüller^{1,2}, Nikolai Strogies¹

submitted: October 27, 2017

¹ Weierstrass Institute
Mohrenstr. 39
10117 Berlin
Germany
E-Mail: michael.hintermueller@wias-berlin.de
nikolai.strogies@wias-berlin.de

² Humboldt-Universität zu Berlin
Unter den Linden 6
10099 Berlin
Germany
E-Mail: hint@math.hu-berlin.de

No. 2442
Berlin 2017



2010 *Mathematics Subject Classification.* 49J15, 49J20, 35L65, 65L06, 65M06, 65M12.

Key words and phrases. Optimal control, conservation laws, discretization methods, RK methods, TVD-RK.

This work is supported by the German Research Foundation (DFG) within project B02 of CRC TRR 154.

Edited by
Weierstraß-Institut für Angewandte Analysis und Stochastik (WIAS)
Leibniz-Institut im Forschungsverbund Berlin e. V.
Mohrenstraße 39
10117 Berlin
Germany

Fax: +49 30 20372-303
E-Mail: preprint@wias-berlin.de
World Wide Web: <http://www.wias-berlin.de/>

On the consistency of Runge–Kutta methods up to order three applied to the optimal control of scalar conservation laws

Michael Hintermüller, Nikolai Strogies

Abstract

Higher-order Runge-Kutta (RK) time discretization methods for the optimal control of scalar conservation laws are analyzed and numerically tested. The hyperbolic nature of the state system introduces specific requirements on discretization schemes such that the discrete adjoint states associated with the control problem converge as well. Moreover, conditions on the RK-coefficients are derived that coincide with those characterizing strong stability preserving Runge-Kutta methods. As a consequence, the optimal order for the adjoint state is limited, e.g., to two even in the case where the conservation law is discretized by a third-order method. Finally, numerical tests for controlling Burgers equation validate the theoretical results.

1 Introduction

We investigate discretization techniques for problems of optimal control subject to scalar conservation laws in one space dimension which, in conservative form, are given as

$$\begin{aligned} y_t + [f(y)]_x &= 0 & \text{in } \mathcal{Q} := (0, T] \times \mathbb{R}, \\ y(0, x) &= u(x) & \text{in } \mathbb{R}. \end{aligned} \quad (1)$$

Here, $f \in C^2(\mathbb{R})$ is a nonlinear convex flux function that is uniformly convex with $f'' \geq c > 0$. Partial differential equations like (1) might, even for smooth initial data, develop shocks (see, e.g., [8]) and thus require the consideration of weak solutions that satisfy additional conditions, guaranteeing uniqueness of solutions. In case of conservation laws, usually modeling physical processes, the relevant solution is called entropy solution. It is known, that the map of control to entropy solution, $u(\cdot) \mapsto y(t, \cdot)$, is usually not differentiable in $L^1(\mathbb{R})$ but *shift-differentiable* in $BV(\mathbb{R})$. This notion of directional differentiability was introduced and discussed for balance laws, inhomogeneous conservation laws with an additional source term, in [9, 33] and extended to strictly hyperbolic systems of balance laws in [4, 11].

In this paper, we consider the model problem

$$\begin{aligned} \text{minimum} \quad & \frac{1}{2} \int_0^1 (y(T, x) - y^d(x))^2 dx + \mathfrak{R}(u) =: \mathcal{J}(y, u) \\ \text{over} \quad & (y, u) \in (L^\infty(\mathcal{Q}) \cap C([0, T]; L^1_{loc}(\mathbb{R}))) \times L^\infty(\mathbb{R}) \\ \text{subject to} \quad & y \text{ solves (1) for } u, \end{aligned} \quad (P)$$

for a desired state $y^d \in PC^1(\mathbb{R})$, the control entering as initial data for (1) and a suitable convex cost functional $\mathfrak{R}(u)$ with an effective domain $dom_{\text{eff}}(\mathfrak{R})$ embedding compactly into $L^1(\mathbb{R})$ and being

coercive with respect to $\|u\|_{L^\infty(\mathbb{R})}$. In case of scalar conservation laws, theoretical results for optimality conditions of (P) have been discussed for example in [33] in case of an unbounded domain and in [30] for bounded domains with a switching control at the boundary. The main focus of these papers lies on obtaining a representation of the reduced gradient of

$$\frac{1}{2} \int_0^1 (y(T, x) - y^d(x))^2 dx \quad (2)$$

with respect to perturbations in the control. If the conservation law satisfies the weakened one-sided Lipschitz condition (OSL), the reduced gradient is given by the solution to the adjoint equations (5) below, a linear conservation law with discontinuous coefficients. It is known that such equations, even for Lipschitz continuous terminal conditions, do not admit unique solutions (see, e.g., [6]) and the relevant *reversible solution* is identified utilizing solutions to the sensitivity equation (6) below.

In case of systems of balance laws, theoretical results concerning optimality conditions for a distributed control have been formulated in [10] for an unbounded domain and in [12] for bounded domains in terms of generalized tangent vectors, the first order variations of the solution.

The numerical treatment of problems of optimal control subject to scalar conservation laws has been studied for example in [3]. In case of systems, a first step for function space consistent numerical methods is considered in [23], where a numerical method for the computation of generalized tangent vectors for a system of conservation laws has been introduced.

Addressing (P) numerically is a delicate task since the non uniqueness of solutions to the adjoint equation requires suitable discretization techniques for both, state and adjoint equation, respectively. In [18] numerical results for Burgers' equation have been obtained where the discrete solutions to the adjoint equation can converge to an incorrect solution if the discretization scheme is not chosen properly. In [16, 17] the convergence behavior of a discretization of the primal equation based on the Lax-Friedrichs flux and a mesh dependent artificial viscosity has been studied. The consistency of the discretized problems with (P) has been proven in [33] for monotone discretization schemes of (1) satisfying certain assumptions. Such discretizations can be interpreted as explicit Euler time discretizations of a system of ordinary differential equations representing a semi discretization of the conservation law. We will investigate higher order Runge-Kutta time discretization methods applied to this semi discretization and derive conditions for the coefficients such that the resulting full discretization of (P) is still consistent.

In the context of optimal control subject to ordinary differential equations, the application of Runge-Kutta (RK) time discretization schemes has been investigated for example in [5, 21].

The application of RK-schemes to semi-discretizations of conservation laws has been a subject of investigation for a long time (see, e.g., [20, 25]). Given the discretization of a conservation law that is total variation diminishing (TVD) for the basic explicit Euler time stepping, these RK methods are used to construct higher order approximations of the solutions with respect to time while preserving the overall TVD property of the full discretization scheme.

Notation. Throughout this paper, Δx denotes the width of the spatial discretization, Δt the size of the time step and $\lambda = \Delta t / \Delta x$. The terminal time is denoted by T and $N = T / \Delta t$ represents the number of time steps. Moreover, we utilize several function spaces and refer to [1, 15] for details. Besides the standard Lebesgue and Sobolev spaces L_{loc}^p and $W^{k,p}(\Omega)$ with corresponding norms,

we consider the space of piecewise continuously differentiable functions with possible discontinuities at finitely many points, the space of all bounded functions that can be approximated pointwise almost everywhere by a sequence of Lipschitz continuous functions $\{w_n\}$ bounded in $C(\mathbb{R}) \cap W_{loc}^{1,1}(\mathbb{R})$ and the space of bounded functions, $PC^1(\Omega)$, $B_{Lip}(\mathbb{R})$ and $B(\mathbb{R})$ respectively, all equipped with the sup norm. Finally, $\mathcal{S}_{\mathcal{M}} := C([0, T]; \mathcal{M}_{loc}(\mathbb{R}) - w(\mathcal{M}_{loc}(\mathbb{R}), C_c(\mathbb{R})))$ denotes the space of continuous functions on $[0, T]$ with values in the local Borel measures on \mathbb{R} , $\mathcal{M}(\mathbb{R})$, with the weak topology induced by continuous functions with compact support.

The paper is organized as follows. In section 2 we provide existence results for solutions of the primal and adjoint equations and for (P) . Moreover, we introduce the semi discretization of (1) and recall conditions for the consistency of the fully discretized problems with (P) for time discretizations based on Euler's method. In section 3 we analyze RK schemes and the resulting discretizations of (1) and present conditions for consistency of the resulting discretizations. In section 4 we briefly discuss the fully discretized problems and convergence order of the time discretization and in section 5 we present numerical results validating our theory of section 3.

2 Preliminaries

In this section we study appropriate solution concepts for the state equation (1), associated adjoint equations, and we discuss the use of the adjoint when computing the reduced gradient of $\mathcal{J}(y, u)$.

2.1 The state equation and its adjoints

As outlined above, even for smooth initial data solutions to (1) might develop shocks and require to study weak solutions $y \in L^\infty(\mathbb{R}_+ \times \mathbb{R})$ that satisfy the identity

$$\int_{\mathbb{R} \times \mathbb{R}_+} y \phi_t + f(y) \phi_x dx dt + \int_{\mathbb{R}} u(\cdot) \phi(0, \cdot) dx = 0 \text{ for all } \phi \in C_c^\infty(\mathbb{R}_+ \times \mathbb{R}).$$

In general, weak solutions are not unique and the physically relevant solution, referred to as *entropy solution* (see, e.g., [27]), is characterized as follows.

Definition 1. Consider $\eta(y) := |y - k|$ and $q(y) := \text{sign}(y - k)(f(y) - f(k))$. A weak solution of (1) is an entropy solution if it satisfies

$$\int_{\mathbb{R} \times \mathbb{R}_+} \eta(y) \phi_t + q(y) \phi_x dx dt + \int_{\mathbb{R}} \eta(u(\cdot)) \phi(0, \cdot) dx \geq 0$$

for all $\phi \in C_c^\infty(\mathbb{R} \times \mathbb{R}_+)$, $\phi \geq 0$, and $k \in \mathbb{R}$.

The following result provides the existence of such solutions, provides some properties and goes back to [33, Proposition 3.4.1].

Proposition 1. For every $u \in L^\infty(\mathbb{R})$ there exists a unique entropy solution

$$y \in L^\infty(\mathcal{Q}) \cap C([0, T]; L^1_{loc}(\mathbb{R}))$$

that satisfies $\|y(t)\|_{L^\infty(\mathbb{R})} \leq \|u\|_{L^\infty(\mathbb{R})}$ for all $t \in [0, T]$,

$$\|y^1(t) - y^2(t)\|_{L^1(\mathbb{R})} \leq \|u_0^1 - u_0^2\|_{L^1(\mathbb{R})} \text{ for all } t \in (0, T]$$

with $y^i(t)$ denoting the solution of (1) for initial data $u_0^i \in L^\infty(\mathbb{R})$, $i \in \{1, 2\}$ and, for all controls with $\|u\|_{L^\infty(\mathbb{R})} \leq M_u$, $M_u > 0$, there exists some $C = C(M_u, u, t) > 0$ such that

$$y_x(t) \leq C(M_u, u, t) \quad (3)$$

holds in the sense of distributions.

Consequently, (1) is well defined which allows us to consider (P) next. Note here that the y depends on the control u entering as initial data in the state equation. Let the desired state y^d be an element of $PC^1(\mathbb{R})$. If the effective domain of $\mathfrak{R}(\cdot)$, $dom_{\text{eff}}(\mathfrak{R})$, embeds compactly into $L^1(\mathbb{R})$, and $\mathfrak{R}(\cdot)$ is coercive with respect to $\|u\|_{L^\infty(\mathbb{R})}$, then (P) admits a solution (see, e.g., [33, Theorem 2.4.2]).

Proposition 1 allows for considering $y = y(u)$, i.e., the entropy solution of (1) depends on the control. In order to establish a gradient related descent algorithm for solving (P) iteratively, a gradient representation of the reduced objective in (2) is necessary. In fact, let $u \in PC^1(\mathbb{R})$ and consider fixed positions of discontinuities in the initial data. Then the reduced objective associated with (P) is given by $\hat{\mathcal{J}}(u) := \mathcal{J}(y(u), u)$ and the application of its gradient to a direction $\delta u \in PC^1(\mathbb{R})$ is given by

$$d_u \mathcal{J}(y(u), u) \cdot \delta u = (p(0, \cdot), \delta u)_{L^2(\mathbb{R})} + d_u \mathfrak{R}(u) \cdot \delta u \quad (4)$$

with $L^2(\mathbb{R})$ -scalar product, derivative of the cost term $d_u \mathfrak{R}(u)$ and p denoting the solution to the adjoint equation

$$\begin{aligned} p_t + f'(y(t, x))p_x &= 0 & \text{in } \mathcal{Q}, \\ p(T, x) &= p^T(x) & \text{in } \mathbb{R}, \end{aligned} \quad (5)$$

with final time data

$$p^T(x) = \int_0^1 (y(T, x+) + \tau[y(T, x)] - y^d(x)) d\tau.$$

Here, $[y(T, x)] := y(T, x-) - y(T, x+)$ denotes a possible jump of the entropy solution to (1) at final time T at x . In general, linear conservation laws with discontinuous coefficients as in (5), even for continuous end data p^T , do not admit unique solutions. While entropy solutions are the physically relevant ones for the nonlinear state equation (1), the proper concept of solutions for (5) is called *reversible solution* and relies on the *one-sided Lipschitz continuity condition*, i.e., that there exists an $\alpha \in L^1(0, T)$ with

$$\partial_x f'(y)(t, \cdot) \leq \alpha(t). \quad (\text{OSL})$$

Next we recall the definition of reversible solutions from [6] and [33] for regular and more general final time data p^T respectively.

Definition 2. In case of regular final time data $p^T \in C_{loc}^{0,1}(\mathbb{R})$, a Lipschitz continuous solution p to (5) is a reversible solution if and only if there exist Lipschitz continuous solutions p_1, p_2 to (5) with $\partial_x p_1 \geq 0$ and $\partial_x p_2 \leq 0$ such that $p = p_1 - p_2$.

In case of $p^T \in B_{Lip}(\mathbb{R})$, they are defined as broad solutions along the generalized backward characteristics.

Broad solutions are solutions to first-order partial differential equations that are, in case of linear conservation laws, constant along the characteristic lines of the problem; see [13] for the corresponding definitions. In case of nonlinear conservation laws, generalized characteristics have to be considered; see [33] for a discussion. The following result is proven in [33, Corollary 4.2.11].

Proposition 2. Let f satisfy (OSL). Then, for end data $p^T \in B_{Lip}(\mathbb{R})$ there exists a unique reversible solution $p \in B(\mathcal{Q}) \cap C^{0,1}([0, T]; L_{loc}^1(\mathbb{R})) \cap B([0, T]; BV_{loc}(\mathbb{R})) \cap BV_{loc}(\mathcal{Q}^{cl})$ of (5) fulfilling the maximum principle

$$\|p(t)\|_{B(I)} \leq \|p^T\|_{B(J)}$$

with $I = [z_1, z_2]$ and $J = [z_1 - \|f'(y)\|_{L^\infty}(T-t), z_2 + \|f'(y)\|_{L^\infty}(T-t)]$ for all $z_1, z_2 \in \mathbb{R}$, $z_1 < z_2$ and $t \in [0, T]$.

In case of entropy solutions to (1), (OSL) of the flux term f is ensured by estimate (3). Moreover, final time data $p^T \in B_{Lip}(\mathbb{R})$ are, for example, elements of $PC^1(\mathbb{R})$.

Reversible solutions are closely related to *duality solutions* (compare Definition 3 below) that characterize directional derivatives of solutions y to the state equation (1) with respect to perturbations in the control u as shown in [7, Theorem 3.1]. In particular, for the perturbed position of a shock in the control u , such sensitivities are measures and defined as weak solutions to

$$\begin{aligned} \mu_t + (f'(y)\mu)_x &= 0 & \text{in } \mathcal{Q}, \\ \mu(0) &= \delta u, \end{aligned} \tag{6}$$

for $\delta u \in \mathcal{M}_{loc}(\mathbb{R})$ and y denoting the solution to (1), given the control u (see, e.g., [33, Example 3.1.1.]).

Definition 3. Let $f'(y)$ satisfy (OSL) and consider $\delta u \in \mathcal{M}_{loc}(\mathbb{R})$. Solutions $\mu \in \mathcal{S}_{\mathcal{M}}$ to (6) are called *duality solution* if for any $\tau \in (0, T]$, any $p^\tau \in B_{Lip}(\mathbb{R})$ with compact support and any reversible solution p of

$$\begin{aligned} p_t + f'(y(t, x))p_x &= 0 & \text{in } (0, \tau) \times \mathbb{R}, \\ p(\tau, x) &= p^\tau(x) & \text{in } \mathbb{R}, \end{aligned} \tag{7}$$

we have

$$\int_{\mathbb{R}} p^\tau \mu(\tau, dx) = \int_{\mathbb{R}} p(0, x) \delta u(dx)$$

Their existence result follows from [33, Theorem 4.3.7].

Proposition 3. *Let $f'(y)$ satisfy (OSL). Then for any $\delta u \in \mathcal{M}_{loc}(\mathbb{R})$ there exists a unique duality solution $\mu \in \mathcal{S}_{\mathcal{M}}$ to (6).*

Moreover, the following relation of reversible and duality solutions has been established in [33, Theorem 4.4.1.]: for every $p^\tau \in B_{Lip}(\mathbb{R})$, $\tau \in (0, T]$, p is a reversible solution to (7) if it is Borel-measurable and satisfies for all $\sigma \in (0, \tau)$ and all $\delta u \in \mathcal{M}(\mathbb{R})$ the duality relation

$$\int_{\mathbb{R}} p^\tau \mu(\tau, dx) = \int_{\mathbb{R}} p(\sigma, x) \delta u(dx), \quad (8)$$

where μ is the duality solution of (6) on $(\sigma, \tau) \times \mathbb{R}$ for the control δu .

2.2 Discrete schemes

Since state, adjoint and sensitivity equations, respectively, in general do not admit unique solutions, discretization schemes have to be chosen that approximate the relevant entropy, reversible and duality solution properly. While convergent schemes for (1) are available in terms of monotone schemes, the respective discretization of (5) and (6) has to be derived from such schemes along with properties of reversible and duality solutions.

Monotone schemes operate on cell averages of the solution by accounting for their evolution over time steps. Thus, given a uniform mesh of width Δx on \mathbb{R} , the discretization of the initial state is obtained by averaging the function in each interval or cell j of the spatial discretization, i.e.,

$$\mathbf{u}_{0j} := T^j(u) = \Delta x^{-1} \int_{x_j}^{x_{j+1}} u(x) dx. \quad (9)$$

with $\Delta x := x_{j+1} - x_j$ for $j \in \mathbb{N}$. A semidiscretization of (1) utilizing the method of lines provides a system of ordinary differential equations (ODE),

$$\dot{\mathbf{y}}_j = F(\mathbf{y})_j = -\Delta x^{-1} [f^\Delta(y_j^n, y_{j+1}^n) - f^\Delta(y_{j-1}^n, y_j^n)], \quad \dot{\mathbf{y}}(0) = \mathbf{u}_0. \quad (10)$$

Here, $f^\Delta : \mathbb{R}^2 \rightarrow \mathbb{R}$ represents a suitable numerical flux function that quantifies the flux between the intervals $j+1$ and j . It is assumed to be at least Lipschitz continuous and thus guarantees the existence of a unique, Lipschitz continuous solution to (10). Depending on the regularity of the chosen numerical flux and the underlying conservation law, the solution of the semidiscrete system of ODE's might attain even higher regularity with respect to t . In fact, for Burgers' equation and the Engquist-Osher flux (see Section 5), the solution to (10) exists and is twice continuously differentiable with respect to time. Using an explicit Euler method with time step size Δt for the time discretization we obtain

$$y_j^{n+1} = y_j^n - \lambda (f^\Delta(y_j^n, y_{j+1}^n) - f^\Delta(y_{j-1}^n, y_j^n)), \quad (11)$$

a so-called centered three-point scheme since the average of the state in cell j only depends on the average in the neighboring cells. In more general schemes, the numerical flux function and consequently the average of the state in cell j , might depend on cell averages of the state in further cells, increasing

the number of arguments of f^Δ . In the discussion below this number, compared to (11), will increase but we restrict ourselves to even numbers of arguments for the numerical flux function, i.e., $f^\Delta : \mathbb{R}^{2K} \rightarrow \mathbb{R}$ with $K = 1, 2, \dots$ and formally, f^Δ depending on $(y_{j-K}^n, \dots, y_{j+K}^n)$. Thus, the integer K defines the domain of determinacy in that in the general, fully discrete scheme, y_j^{n+1} depends on $y_{j-K}^{n+1}, \dots, y_{j+K}^{n+1}$; cf. (11) for the case of $K = 1$.

Based on (11), discretization schemes of (5) and (6) can be derived as demonstrated next. Their convergence depends on properties of certain coefficients that are introduced in this process.

Derivatives of entropy solutions to (1) with respect to variations in the initial data are characterized by duality solutions to (6). Thus, if the numerical flux f^Δ is differentiable, a sensitivity scheme can be established that characterizes derivatives of the discrete approximations y_j^n defined by (11) with respect to the discrete initial data \mathbf{u} for $n = 1, \dots, N$ and $j \in \mathbb{N}$. It is given by

$$\mu_j^{n+1} = \mu_j^n - \lambda [f_2^\Delta(y_j^n, y_{j+1}^n) \mu_{j+1}^n + (f_1^\Delta(y_j^n, y_{j+1}^n) - f_2^\Delta(y_{j-1}^n, y_j^n)) \mu_j^n - f_1^\Delta(y_{j-1}^n, y_j^n) \mu_{j-1}^n] \quad (12)$$

along with the initial data $\mu_j^0 = \delta \mathbf{u}_j$. Here, $f_l^\Delta(v_1, v_2)$ represents the partial derivative of f^Δ with respect to the l -th argument. This scheme can be written explicitly as

$$\mu_j^{n+1} = \sum_{k=-1}^1 D_{j,k}^n \mu_{j+k}^n \quad \text{for} \quad \begin{aligned} D_{j,-1}^n &= \lambda f_1^\Delta(y_{j-1}^n, y_j^n), \\ D_{j,0}^n &= 1 - \lambda (f_1^\Delta(y_j^n, y_{j+1}^n) - f_2^\Delta(y_{j-1}^n, y_j^n)), \\ D_{j,1}^n &= -\lambda f_2^\Delta(y_j^n, y_{j+1}^n). \end{aligned} \quad (13)$$

Note that if these coefficients are non negative, (11) forms a monotone discretization scheme since they represent the partial derivatives in, e.g., [32, Definition 5.1]. Similar to the relation of duality and reversible solutions in the continuous setting, the discretization scheme for the adjoint equation has to satisfy the discrete analogue of (8), obtained for $\sigma = t, \tau = t + \Delta t$ and rescaling the sum representing the integrals, given by

$$\sum_j p_j^{n+1} \mu_j^{n+1} = \sum_j p_j^n \mu_j^n, \quad (14)$$

for any solution $\{\mu_j^n\}$ of (12) with bounded support, i.e., $\mu_j^0 \neq 0$ at finitely many $j \in \mathbb{N}$ only. In case of tracking type objective functionals for desired states with bounded support, the control is assumed to have bounded support as well. Multiplying (13) by p_j^{n+1} , a summation over j and reordering provides

$$\sum_j p_j^{n+1} \mu_j^{n+1} = \sum_j \sum_{k=-1}^1 D_{j-k,k}^n p_{j-k}^{n+1} \mu_j^n.$$

A substitution into (14) yields the discretization scheme for the adjoint equation, given by

$$p_j^n = p_j^{n+1} + \lambda (D_{j+1,-1}^n (p_{j+1}^{n+1} - p_j^{n+1}) - D_{j-1,1}^n (p_j^{n+1} - p_{j-1}^{n+1})), \quad (15)$$

a time discretization for the system of ordinary differential equations

$$\dot{\mathbf{p}}_j = \Delta x^{-1} [f_2(y_{j-1}(t), y_j(t))(p_{j+1}(t) - p_j(t)) - f_1(y_j(t), y_{j+1}(t))(p_j(t) - p_{j-1}(t))] \quad (16)$$

utilizing an implicit Euler method along with discrete final time condition $\mathbf{p}(T) = d_{\mathbf{y}(T)}\mathcal{J}(\mathbf{y}(T), \mathbf{u})$. In case of (P), the latter is given as $\mathbf{p}(T)_j = \Delta x(\mathbf{y}(T)_j - T^j(y^d))$.

The following set of coefficients allows for quantifying the difference of neighboring values of the solution to (15) after one time step in the adjoint scheme.

$$p_{j+1}^{n+1} - p_j^{n+1} = \sum_{k=-1}^1 C_{j,k}^n (p_{j+k+1}^n - p_{j+k}^n), \quad \text{with} \quad \begin{aligned} C_{j,-1}^n &= -\lambda f_2^\Delta(y_{j-1}^n, y_j^n), \\ C_{j,0}^n &= 1 + \lambda(f_2^\Delta(y_j^n, y_{j+1}^n) - f_1^\Delta(y_j^n, y_{j+1}^n)), \\ C_{j,1}^n &= \lambda f_1^\Delta(y_{j+1}^n, y_{j+2}^n). \end{aligned} \quad (17)$$

Finally, the existence of reversible and duality solutions depend on (OSL), respectively, to ensure convergence of their corresponding discrete approximations, the schemes have to be consistent with it. A sufficient condition for (OSL)-consistency that relies on the coefficients

$$l_j^{n+1} := \Delta x^{-1}(y_{j+1}^n - y_j^n), \quad l_{j,K}^{n,+} := \max(0, l_{j-K}^n, \dots, l_{j+K}^n). \quad (18)$$

was studied in [33, Section 6.4.4.] and allows to establish a bound on $\partial_x f'(y)(t, \cdot)$.

By the coefficients introduced above, we are now able to present the following theorem, collecting results from [33, Theorem 6.4.10 and Theorem 6.4.15.]. It characterizes the consistency of nonlinear programs obtained from (P) by discretizing the state system with suitable schemes like (11) and obtaining gradient information of the reduced objective by associated adjoint schemes like (15). The formulation covers centered three-point schemes, i.e., (11) but also holds for larger domains of determinacy (see discussion below (11)).

Theorem 1 ([33]). *For $u \in PC^1(\mathbb{R})$, which provides $y(T, \cdot) \in PC(\mathbb{R})$, and $y^d \in PC^1(\mathbb{R})$, let the flux function $f \in C^2(\mathbb{R})$ satisfy $f'' \geq c > 0$ and consider $K \geq 1, K \in \mathbb{N}$ and $M_y > 0$ sufficiently large. Moreover, let the numerical flux f^Δ fulfill the following conditions:*

- 1 $f^\Delta \in C_{loc}^{1,1}(\mathbb{R}^{2K})$ and the numerical flux is consistent, i.e., $f^\Delta(y, \dots, y) = f(y)$;
- 2 The coefficients $D_{j,l}^n$ defined in (13) are non negative for all $y_j^n \in [-M_y, M_y]$, $j \in \mathbb{N}$;
- 3 The coefficients $C_{j,l}^n$ defined in (17) are non negative for all $y_j^n \in [-M_y, M_y]$, $j \in \mathbb{N}$;
- 4 The discrete state y_j^n is contained in $[-M_y, M_y]$ for $j \in \mathbb{N}, n = 1, \dots, N$ and there exist some $\nu > 0$ such that

$$l_j^{n+1} \leq l_{j,K}^{n,+} - \Delta t \nu (l_{j,K}^{n,+})^2; \quad (19)$$

- 5 The partial derivatives $f_{y_i}^\Delta$ are non decreasing on $[-M_y, M_y]^{2K}$.

Then the solution to the discrete sensitivity equation (12) converges to the duality solution of (6) in $B([0, T]; \mathcal{M}_{loc}(\mathbb{R}) - w(\mathcal{M}_{loc}(\mathbb{R}), C_c(\mathbb{R})))$ and the solution to the discrete adjoint equation (15) converges to the solution of (5) in $L^r(Q \setminus \bigcup_k D_k)$ with $r \in [1, \infty)$ and D_k depending on the shock positions in the terminal condition for the adjoint equation and the height of the corresponding jumps. Moreover, $\mathbf{p}^0 \rightarrow p(0, \cdot)$ as $\Delta x \rightarrow 0$ at least in $L^r(\mathbb{R} \setminus (\bigcup_k D_k \cap \{t = 0\}))$ where \mathbf{p}^0 is considered to be piecewise constant on the corresponding intervals of the discretization.

Convergence of the discrete state y_j^n to y as $\Delta x \rightarrow 0$ follows from the theory for monotone discretization schemes for conservation laws which, as discussed above, is ensured by the non-negativity of the coefficients $D_{j,l}^n$. The monotonicity also determines the constant M_y since, for such schemes, time iterates satisfy the maximum principle (see, e.g., [32, Theorem 13.36]), given as

$$\max_j \{|\mathbf{y}_j^n|\} \leq \max_j \{|\mathbf{u}_j|\}, \text{ for all } n = 1, \dots, N.$$

Consequently, the control introduces bounds to the discrete approximation of the solution, \mathbf{y}_j^n with $n = 1, \dots, N$ and $j \in \mathbb{N}$. By the assumed coercivity of $\mathfrak{R}(\cdot)$, the discrete initial condition has to be bounded in $L^\infty(\mathbb{R})$ as $\Delta x \rightarrow 0$, allowing to establish a bound on M_y . Condition 2 of Theorem 1 poses explicit conditions on the numerical flux to ensure $D_{j,-1}^n \geq 0$ and $D_{j,1}^n \geq 0$ while $D_{j,0}^n \geq 0$ typically can be ensured by a restriction on Δt and thus on λ ; compare (13). In case of $K = 1$, the numerical flux functions (EO) and (LF) (compare Section 5) are known to satisfy all conditions of Theorem 1 under certain conditions on Δt (see [33, Chapter 6]).

After introducing the central convergence result for nonlinear programs obtained from (P) by applying suitable discretization techniques, we will now study the impact of higher order Runge-Kutta (RK) time discretization methods applied to (10). If the solution to the system of ODE's is sufficiently smooth, these schemes are more accurate in that the truncation error depends on some power of the chosen time step Δt . Thus, the term higher order should not be mistaken with the overall truncation error of the discretization strategy for (1) that still heavily relies on the discretization method with respect to the spatial variable. The section closes with introducing RK schemes in the well known form of Butcher arrays and the Shu-Osher representation.

Definition 4. Consider a uniform time step Δt . In vector notation, an s stage RK scheme is given by

$$\begin{aligned} \mathbf{y}^{n,i} &= \mathbf{y}^n + \Delta t \sum_{l=1}^s a_{il} F(\mathbf{y}^{n,l}), \quad i = 1, \dots, s, \\ \mathbf{y}^{n+1} &= \mathbf{y}^n + \Delta t \sum_{i=1}^s b_i F(\mathbf{y}^{n,i}). \end{aligned}$$

The coefficients $A = (a_{ij})$ and weights $b = (b_j)$ with $1 \leq i, j \leq s$ represent the Butcher array and characterize the method.

Here we only consider explicit RK-methods, thus restricting the coefficients to $a_{ij} = 0$ for $j \geq i$. If the coefficients in A and b satisfy conditions formulated, e.g., in [21, Table 2], the truncation error of the approximation of the ODE's is of higher order.

In numerical methods for conservation laws, a different representation of RK methods is used frequently (see, e.g., [20, 25]).

Definition 5. The Shu-Osher representation of a s -stage RK method is a convex combination of forward Euler steps, parametrized by two sets of coefficients $\{\alpha_{ij}\}$ and $\{\beta_{ij}\}$ for $i = 1, \dots, s$ and $j = 0, \dots, s-1$,

defined by

$$\mathbf{y}^{(i)} = \sum_{l=0}^{i-1} \alpha_{il} \mathbf{y}^{(l)} - \beta_{il} h F(\mathbf{y}^{(l)}),$$

$$\mathbf{y}^{(0)} = \mathbf{y}^n, \quad \mathbf{y}^{n+1} = \mathbf{y}^{(s)}$$

and satisfying $\alpha_{ij} \geq 0$, $\sum_{j=0}^{i-1} \alpha_{ij} = 1$.

This format allows for studying high-order total variation diminishing (TVD) discretizations with respect to time for a given spatial discretization that is TVD. The resulting, so called, TVD-RK methods preserve the TVD-property of the original discretization and have a higher order of accuracy with respect to time in terms of the truncation error. In fact, they are strong stability preserving (SSP) methods and ensure

$$\mathcal{C}(\mathbf{y}(t) + \Delta t F(\mathbf{y}(t))) \leq \mathcal{C}(\mathbf{y}(t)) \text{ for all } \Delta t \leq c_{SSP} \Delta t_E, \quad (20)$$

for arbitrary convex functionals $\mathcal{C} : \mathbb{R}^N \rightarrow \mathbb{R}$ including norms and the TV-seminorm, and time steps Δt up to a multiple of the time step Δt_E for the original discretization, employing an explicit Euler method. The scaling factor, $c_{SSP} := \min_{ij} \frac{\alpha_{ij}}{|\beta_{ij}|}$ is called SSP coefficient.

We will derive a close relationship between RK schemes that can be used in the context of Theorem 1 and strong stability preserving RK methods. As a consequence, known restrictions from SSP-RK methods apply in our case as well. In particular, it is known that explicit SPP methods are restricted to order $p \leq 4$ (see [31]). Moreover, we aim for time stepping methods with order matching the number of stages to avoid the computational effort of calculating and storing additional stages. There exist no combination of coefficients in the Butcher tableau of any RK method with number of stages matching the order such all coefficients in the Shu-Osher representation are non-negative (see [19]) in case of convergence order $p > 3$. We will see in the next section, that we are restricted to RK methods that have a Shu-Osher representation with $\alpha_{ij} \geq 0, \beta_{ij} \geq 0$. Consequently, we are restricted to RK-schemes of order at most 3.

Equation (11) forms a centered three-point scheme and Theorem 1 holds with $K = 1$. However, in case of higher order time stepping methods, this holds for every intermediate step. This increases the domain of determinacy for the full time step which is the number of values in a time slice contributing to values in the following time slice.

3 Consistency of the Runge-Kutta time stepping

We are now ready to study optimization RK schemes, i.e., we are interested in RK-discretizations of (1) which yield consistent RK-schemes for the adjoint (5) with high approximation order. Similar to SSP-methods in the context of strong stability preserving methods, we assume that an explicit Euler time stepping for (10) with a properly chosen time step size Δt provides a full discretization of (1) such that

the conditions of Theorem 1 are satisfied and identify RK schemes that preserve this properties, thus allowing for an application of the result. Consequently, the basic assumption for the following is given by

- (A) The numerical flux f^Δ and Δt are chosen such that conditions 1. to 5. of Theorem 1 are met by (11).

Writing $f(v_1, v_2)$ instead of $f^\Delta(v_1, v_2)$ in a slight misuse of notation, we introduce $F_j^{n,s} : \mathbb{R}^3 \rightarrow \mathbb{R}$, its gradient and the vector valued functions $\tilde{F}_j^{n,s} : \mathbb{R}^4 \rightarrow \mathbb{R}^3$ and $\hat{F}_j^{n,s} : \mathbb{R}^4 \rightarrow \mathbb{R}^3$ utilized for validating the assumptions of Theorem 1 in case of a multistage time discretization of (10):

$$\begin{aligned} F_j^{n,s} &:= y_j^{n,s} - \lambda[f(y_j^{n,s}, y_{j+1}^{n,s}) - f(y_{j-1}^{n,s}, y_j^{n,s})], \\ \nabla F_j^{n,s} &:= [-\lambda f_2(y_j^{n,s}, y_{j+1}^{n,s}), 1 - \lambda(f_1(y_j^{n,s}, y_{j+1}^{n,s}) - f_2(y_{j-1}^{n,s}, y_j^{n,s})), \lambda f_1(y_{j-1}^{n,s}, y_j^{n,s})], \\ \tilde{F}_j^{n,s} &:= [\lambda f_1(y_j^{n,s}, y_{j+1}^{n,s}), 1 - \lambda(f_1(y_j^{n,s}, y_{j+1}^{n,s}) - f_2(y_{j-1}^{n,s}, y_j^{n,s})), -\lambda f_2(y_{j-1}^{n,s}, y_j^{n,s})], \\ \hat{F}_j^{n,s} &:= [\lambda f_1(y_{j+1}^{n,s}, y_{j+2}^{n,s}), 1 - \lambda(f_1(y_j^{n,s}, y_{j+1}^{n,s}) - f_2(y_j^{n,s}, y_{j+1}^{n,s})), -\lambda f_2(y_{j-1}^{n,s}, y_j^{n,s})]. \end{aligned}$$

Here, $n = 0, \dots, N - 1$, s and $j \in \mathbb{N}$ denote time step, stage of the RK method and cell number respectively. By this notation, the general RK-scheme with three stages (RK3-scheme) can be rewritten as follows:

$$\begin{aligned} y_j^{n,1} &= y_j^n, \\ y_j^{n,2} &= c_{20}y_j^n + c_{21}F_j^{n,1}, \\ y_j^{n,3} &= c_{30}y_j^n + c_{31}F_j^{n,1} + c_{32}F_j^{n,2}, \\ y_j^{n+1} &= c_{40}y_j^n + c_{41}F_j^{n,1} + c_{42}F_j^{n,2} + c_{43}F_j^{n,3}. \end{aligned} \quad (21)$$

This forms a particular Shu-Osher representation (see Definition 5) of the general Runge-Kutta method in that we assume $\alpha_{ij} = \beta_{ij}$, fixing $c_{SSP} = 1$. In particular we find the structure of the original Euler time step preserved in that $F_j^{n,i}$, $i = 1, \dots, 3$ coincides with (11) evaluated at the intermediate steps of the RK scheme and $\nabla F_j^{n,i}$ and $\hat{F}_j^{n,i}$ represent the coefficients from (13) and (17), evaluated at the intermediate steps i of the RK method for $i = 1, \dots, 3$. The coefficients in (21) are given next.

Definition 6. Given a time discretization of (10) by a RK3-method, the coefficients in (21) are defined by

$$\begin{aligned} c_{21} &= a_{21}, & c_{43} &= b_3, \\ c_{20} &= 1 - a_{21}, & c_{42} &= b_2 - b_3 a_{32}, \\ c_{32} &= a_{32}, & c_{41} &= b_1 - b_3 a_{31} - (b_2 - b_3 a_{32}) a_{21}, \\ c_{31} &= a_{31} - a_{32} a_{21}, & c_{40} &= 1 - \sum_{j=1}^3 c_{4,j}, \\ c_{30} &= 1 - \sum_{j=1}^2 c_{3,j}, \end{aligned}$$

General RK-schemes with two stages (RK2-scheme) can be interpreted as RK3-schemes with $a_{3,1} = a_{3,2} = b_3 = 0$ in the corresponding Butcher array. In this case, representation (21) follows by substituting these values into the coefficients of Definition 6. The differentiability of the numerical flux

(according to Condition 1 of Theorem 1) allows to obtain the following sensitivity scheme, where $\bar{\mu}_j^{n,s} := [\mu_{j+1}^{n,s}, \mu_j^{n,s}, \mu_{j-1}^{n,s}]$:

$$\begin{aligned}\mu_j^{n,1} &= \mu_j^n, \\ \mu_j^{n,2} &= c_{20}\mu_j^n + c_{21}\nabla F_j^{n,1} \cdot \bar{\mu}_j^{n,1}, \\ \mu_j^{n,3} &= c_{30}\mu_j^n + c_{31}\nabla F_j^{n,1} \cdot \bar{\mu}_j^{n,1} + c_{32}\nabla F_j^{n,2} \cdot \bar{\mu}_j^{n,2}, \\ \mu_j^{n+1} &= c_{40}\mu_j^n + c_{41}\nabla F_j^{n,1} \cdot \bar{\mu}_j^{n,1} + c_{42}\nabla F_j^{n,2} \cdot \bar{\mu}_j^{n,2} + c_{43}\nabla F_j^{n,3} \cdot \bar{\mu}_j^{n,3}.\end{aligned}\quad (22)$$

Assuming $b_j \neq 0$ for $j = 1, \dots, s$ and given an RK-method with coefficients from the corresponding Butcher-tableau (see Definition 4), the adjoint scheme is will be derived by substituting the intermediate stages into the equation defining μ_j^{n+1} and employing the discrete duality condition (14) followed by collecting suitable terms in the result. To keep notation short, we restrict ourselves to a RK2-scheme. For more stages, the adjoint scheme is obtained analogously. Any time step is, with $y_j^{n,1} = y_j^n$ and $\mu_j^{n,1} = \mu_j^n$, given as follows:

$$\mu_j^{n,2} = \mu_j^n - a_{21}\lambda(f_2(y_j^n, y_{j+1}^n)\mu_{j+1}^n + (f_1(y_j^n, y_{j+1}^n) - f_2(y_{j-1}^n, y_j^n))\mu_j^n - f_1(y_{j-1}^n, y_j^n)\mu_{j-1}^n), \quad (23)$$

$$\begin{aligned}\mu_j^{n+1} &= \mu_j^n - b_1\lambda(f_2(y_j^n, y_{j+1}^n)\mu_{j+1}^n + (f_1(y_j^n, y_{j+1}^n) - f_2(y_{j-1}^n, y_j^n))\mu_j^n - f_1(y_{j-1}^n, y_j^n)\mu_{j-1}^n) \\ &\quad - b_2\lambda(f_2(y_j^{n,2}, y_{j+1}^{n,2})\mu_{j+1}^{n,2} + (f_1(y_j^{n,2}, y_{j+1}^{n,2}) - f_2(y_{j-1}^{n,2}, y_j^{n,2}))\mu_j^{n,2} - f_1(y_{j-1}^{n,2}, y_j^{n,2})\mu_{j-1}^{n,2}).\end{aligned}\quad (24)$$

Multiplying (24) by p_j^{n+1} , summing over j and utilizing the bounded support of $\{\mu_i^n\}$, we can reorder the sum and obtain

$$\begin{aligned}&\sum_{j \in \mathbb{N}} \mu_j^{n+1} p_j^{n+1} \\ &= \sum_{j \in \mathbb{N}} \mu_j^n (-b_1\lambda f_2(y_{j-1}^n, y_j^n) p_{j-1}^{n+1} + (1 - b_1\lambda(f_1(y_j^n, y_{j+1}^n) - f_2(y_{j-1}^n, y_j^n))) p_j^{n+1} \\ &\quad + b_1\lambda f_1(y_j^n, y_{j+1}^n) p_{j+1}^{n+1}) \\ &\quad + b_2 \sum_{j \in \mathbb{N}} \mu_j^{n,2} (-\lambda f_2(y_{j-1}^{n,2}, y_j^{n,2}) p_{j-1}^{n+1} - \lambda(f_1(y_j^{n,2}, y_{j+1}^{n,2}) - f_2(y_{j-1}^{n,2}, y_j^{n,2})) p_j^{n+1} \\ &\quad + \lambda f_1(y_j^{n,2}, y_{j+1}^{n,2}) p_{j+1}^{n+1}).\end{aligned}$$

Substituting (23) in this expression and further reordering of the sum provides, just as in the case of a single Euler step, an expression for $\sum_{j \in \mathbb{N}} \mu_j^{n+1} p_j^{n+1}$ that depends on μ_i^n and p_i^{n+1} only, providing a

formula for p_j^n . Collecting suitable terms, we get

$$\begin{aligned}
p_j^{n+1,1} &= p_j^{n+1}, \\
p_j^{n+1,2} &= p_j^{n+1} - \lambda a_{12}^\dagger (\lambda f_2(y_{j-1}^{n,2}, y_j^{n,2}) p_{j-1}^{n+1} \\
&\quad + \lambda (f_1(y_j^{n,2}, y_{j+1}^{n,2}) - f_2(y_{j-1}^{n,2}, y_j^{n,2})) p_j^{n+1} - \lambda f_1(y_j^{n,2}, y_{j+1}^{n,2}) p_{j+1}^{n+1}), \\
p_j^n &= p_j^{n+1} \\
&\quad - \lambda b_2 (f_2(y_{j-1}^{n,2}, y_j^{n,2}) p_{j-1}^{n+1,1} \\
&\quad\quad + (f_1(y_j^{n,2}, y_{j+1}^{n,2}) - f_2(y_{j-1}^{n,2}, y_j^{n,2})) p_j^{n+1,1} - f_1(y_j^{n,2}, y_{j+1}^{n,2}) p_{j+1}^{n+1,1}) \\
&\quad - \lambda b_1 (f_2(y_{j-1}^{n,1}, y_j^{n,1}) p_{j-1}^{n+1,2} \\
&\quad\quad + (f_1(y_j^{n,1}, y_{j+1}^{n,1}) - f_2(y_{j-1}^{n,1}, y_j^{n,1})) p_j^{n+1,2} - f_1(y_j^{n,1}, y_{j+1}^{n,1}) p_{j+1}^{n+1,2}).
\end{aligned}$$

Applying this technique to a general RK3-scheme provides the following update rule for the time step $n + 1 \rightarrow n$ of the adjoint equation:

$$\begin{aligned}
p_j^{n+1,1} &= p_j^{n+1}, \\
p_j^{n+1,2} &= p_j^{n+1} \\
&\quad - \lambda a_{12}^\dagger (f_2(y_{j-1}^{n,3}, y_j^{n,3}) p_{j-1}^{n+1,1} \\
&\quad\quad + (f_1(y_j^{n,3}, y_{j+1}^{n,3}) - f_2(y_{j-1}^{n,3}, y_j^{n,3})) p_j^{n+1,1} - f_1(y_j^{n,3}, y_{j+1}^{n,3}) p_{j+1}^{n+1,1}), \\
p_j^{n+1,3} &= p_j^{n+1} \\
&\quad - \lambda a_{13}^\dagger (f_2(y_{j-1}^{n,3}, y_j^{n,3}) p_{j-1}^{n+1,1} \\
&\quad\quad + (f_1(y_j^{n,3}, y_{j+1}^{n,3}) - f_2(y_{j-1}^{n,3}, y_j^{n,3})) p_j^{n+1,1} - f_1(y_j^{n,3}, y_{j+1}^{n,3}) p_{j+1}^{n+1,1}) \\
&\quad - \lambda a_{23}^\dagger (f_2(y_{j-1}^{n,2}, y_j^{n,2}) p_{j-1}^{n+1,2} \\
&\quad\quad + (f_1(y_j^{n,2}, y_{j+1}^{n,2}) - f_2(y_{j-1}^{n,2}, y_j^{n,2})) p_j^{n+1,2} - f_1(y_j^{n,2}, y_{j+1}^{n,2}) p_{j+1}^{n+1,2}), \\
p_j^n &= p_j^{n+1} \tag{25} \\
&\quad - \lambda b_3 (f_2(y_{j-1}^{n,3}, y_j^{n,3}) p_{j-1}^{n+1,1} \\
&\quad\quad + (f_1(y_j^{n,3}, y_{j+1}^{n,3}) - f_2(y_{j-1}^{n,3}, y_j^{n,3})) p_j^{n+1,1} - f_1(y_j^{n,3}, y_{j+1}^{n,3}) p_{j+1}^{n+1,1}) \\
&\quad - \lambda b_2 (f_2(y_{j-1}^{n,2}, y_j^{n,2}) p_{j-1}^{n+1,2} \\
&\quad\quad + (f_1(y_j^{n,2}, y_{j+1}^{n,2}) - f_2(y_{j-1}^{n,2}, y_j^{n,2})) p_j^{n+1,2} - f_1(y_j^{n,2}, y_{j+1}^{n,2}) p_{j+1}^{n+1,2}) \\
&\quad - \lambda b_1 (f_2(y_{j-1}^{n,1}, y_j^{n,1}) p_{j-1}^{n+1,3} \\
&\quad\quad + (f_1(y_j^{n,1}, y_{j+1}^{n,1}) - f_2(y_{j-1}^{n,1}, y_j^{n,1})) p_j^{n+1,3} - f_1(y_j^{n,1}, y_{j+1}^{n,1}) p_{j+1}^{n+1,3}).
\end{aligned}$$

Formally, this corresponds to an implicit multistage method with the following tableau of coefficients:

$$\begin{array}{c}
A_{RK2}^\dagger \\
b_{RK2}^\dagger
\end{array}
= \frac{\begin{array}{cc} 0 & \frac{b_2 a_{21}}{b_1} \\ 0 & 0 \end{array}}{\begin{array}{cc} b_2 & b_1 \end{array}}, \quad
\begin{array}{c}
A_{RK3}^\dagger \\
b_{RK3}^\dagger
\end{array}
= \frac{\begin{array}{ccc} 0 & \frac{b_3 a_{32}}{b_2} & \frac{b_3 a_{31}}{b_1} \\ 0 & 0 & \frac{b_2 a_{21}}{b_1} \\ 0 & 0 & 0 \end{array}}{\begin{array}{ccc} b_3 & b_2 & b_1 \end{array}}$$

We observe that the adjoint RK schemes correspond to the methods established in [21] in case of optimal control of ordinary differential equations. Similar to (21), the adjoint scheme can be rewritten in a form that preserves the structure of the single-step adjoint scheme with coefficients $\tilde{c}_{l,k}$.

Proposition 4. *Let (A) be satisfied. For $s \in \{2, 3\}$ consider an s -stage RK-scheme for the time discretization of (10). Moreover, let $b_j \neq 0$ hold for $j = 1, \dots, s$ and the coefficients from the Butcher tableau of the time stepping scheme satisfy the conditions of Table 1. Then we have $D_{j,l}^n \geq 0$ for $1 \leq n \leq N, j \in \mathbb{N}, -s \leq l \leq s$ and $s \in \{2, 3\}$.*

Stages	Conditions
2	$a_{21} \geq 0, 1 - a_{21} \geq 0, b_2 \geq 0, b_1 - b_2 a_{21} \geq 0$
3	$a_{32} \geq 0, a_{31} - a_{32} a_{21} \geq 0, 1 - a_{31} - a_{32}(1 - a_{21}) \geq 0, b_3 \geq 0, b_2 - b_3 a_{32} \geq 0, b_3 a_{32} + b_3 a_{31} + b_2 a_{21} - b_3 a_{32} a_{21} \geq 0, b_1 - b_2 a_{21} - b_3(a_{31} - a_{32} a_{21}) \geq 0$

Table 1: Order Conditions for RK-Scheme

Proof. The proof establishes conditions in case of RK3 schemes only, but RK2 schemes are treated analogously. Any intermediate stage $i = 1, 2, 3$ and the full time step, $\mu_j^{n+1}, \mu_j^{n,3}, \mu_j^{n,2}, \mu_j^{n,1} = \mu_j^n$, will be expressed in the form

$$\mu_j^{n,i} = \sum_{l=-1}^1 \dot{D}_{j,l}^{n,i} \mu_{j+l}^{n,1} + \tilde{D}_{j,l}^{n,i} \mu_{j+l}^{n,2} + \hat{D}_{j,l}^{n,i} \mu_{j+l}^{n,3} \quad (26)$$

for coefficients $\dot{D}_{j,l}^{n,i}, \tilde{D}_{j,l}^{n,i}, \hat{D}_{j,l}^{n,i} \in \mathbb{R}$ to be determined. For unifying the representation, we utilize the superscript $(n, 4)$ for μ_j^{n+1} and since we consider explicit schemes, it holds that

$$\dot{D}_{j,l}^{n,k} \equiv \tilde{D}_{j,l}^{n,k} \equiv \hat{D}_{j,l}^{n,k} \equiv 0$$

for $k \geq i$.

The proof consists of two parts. Recall that assumption (A) ensures positivity of the coefficients $D_{j,l}^n$ in (13) only if $\mathbf{y}_j^n \in [-M_y, M_y]$ holds for all $j \in \mathbb{N}$. Thus, these coefficients, evaluated at intermediate steps of the RK method, are not necessarily non-negative. Consequently, we first assume $D_{j,l}^n \geq 0$ for $\mathbf{y}_j^n \in \mathbb{R}$ for all $j \in \mathbb{N}$ and derive conditions on the coefficients in the Butcher array of the RK-scheme ensuring positivity of the corresponding coefficients in case of a multistage method. In a second step, we discuss the restriction $\mathbf{y}_j^n \in [-M_y, M_y]$ required in Theorem 1.

As outlined above, (22) and the intermediate steps preserve the original single-step structure which, by assumption (A), ensures positivity of each component in $\nabla F_j^{n,i}, i = 1, \dots, s$. By the assumption on the numerical flux and Δt and considering (22) rewritten in the form of (26) we obtain

$$\dot{D}_{j,l}^{n,4} = c_{40}[0, 1, 0] + c_{41} \nabla F_j^{n,1}, \quad \tilde{D}_{j,l}^{n,4} = c_{42} \nabla F_j^{n,2}, \quad \hat{D}_{j,l}^{n,4} = c_{43} \nabla F_j^{n,3}.$$

Non-negativity of these coefficients follows for $c_{4,k} \geq 0, k = 0, \dots, 3$. Similarly, $\dot{D}_{j,l}^{n,3} \geq 0, \tilde{D}_{j,l}^{n,3} \geq 0$ if $c_{3,k} \geq 0, k = 0, \dots, 2$ and $\dot{D}_{j,l}^{n,2} \geq 0$ if $c_{20} \geq 0, c_{21} \geq 0$. As a consequence, the coefficients of the

reduced time step satisfy $D_{j,l}^n \geq 0$ for $1 \leq n \leq N$, $j \in \mathbb{N}$, $-3 \leq l \leq 3$ by consisting of sums and products of non negative numbers only.

For the second part of the proof we observe, that the conditions derived in Table 1 exactly match the conditions that characterize strong stability preserving Runge-Kutta methods with full Eulerian time step (see, e.g., [25, Theorem 2]). The latter RK-methods ensure the maximum principle, not only for the full time step but for each intermediate stage as well. Consequently, given $y_j^n \in [-M_y, M_y]$ for all $j \in \mathbb{N}$, we have $y_j^{n+1} \in [-M_y, M_y]$ and in particular $y_j^{n,i} \in [-M_y, M_y]$ for $i = 1, \dots, s$. Thus, the arguments from step 1 remain true. \square

Now we analyze the coefficients $C_{j,l}^n$ in the representation of $p_{j+1}^n - p_j^n$ and find the following result.

Proposition 5. *Let the assumptions of Proposition 4 be satisfied. Then the coefficients $C_{j,l}^n$ for $1 \leq n \leq N$, $j \in \mathbb{N}$ and $-3 \leq l \leq 3$ are non-negative if the conditions of Table 1 are satisfied.*

Proof. Similar to the discretization scheme for the conservation law and its sensitivity equation, for $\bar{p}_j^{n+1,s} = [p_{j-1}^{n,s}, p_j^{n+1,s}, p_{j+1}^{n+1,s}]$ we find the following representation of the adjoint discretization scheme (25):

$$\begin{aligned} p_j^{n+1,1} &= p_j^{n+1}, \\ p_j^{n+1,2} &= \tilde{c}_{20} p_j^{n+1} + \tilde{c}_{21} \tilde{F}_j^{n,3} \cdot \bar{p}_j^{n+1,1}, \\ p_j^{n+1,3} &= \tilde{c}_{30} p_j^{n+1} + \tilde{c}_{31} \tilde{F}_j^{n,3} \cdot \bar{p}_j^{n+1,1} + \tilde{c}_{32} \tilde{F}_j^{n,2} \cdot \bar{p}_j^{n+1,2}, \\ p_j^n &= \tilde{c}_{40} p_j^{n+1} + \tilde{c}_{41} \tilde{F}_j^{n,3} \cdot \bar{p}_j^{n+1,1} + \tilde{c}_{42} \tilde{F}_j^{n,2} \cdot \bar{p}_j^{n+1,2} + \tilde{c}_{43} \tilde{F}_j^{n,1} \cdot \bar{p}_j^{n+1,3}. \end{aligned}$$

Abbreviating $\Delta_+ \bar{p}_j^{n+1,s} = [p_{j+2}^{n+1,s} - p_{j+1}^{n+1,s}, p_{j+1}^{n+1,s} - p_j^{n+1,s}, p_j^{n+1,s} - p_{j-1}^{n+1,s}]$ we obtain

$$\begin{aligned} p_{j+1}^{n+1,1} - p_j^{n+1,1} &= p_{j+1}^{n+1} - p_j^{n+1}, \\ p_{j+1}^{n+1,2} - p_j^{n+1,2} &= \tilde{c}_{20}(p_{j+1}^{n+1} - p_j^{n+1}) + \tilde{c}_{21} \hat{F}_j^{n,3} \cdot \Delta_+ \bar{p}_j^{n+1,1}, \\ p_{j+1}^{n+1,3} - p_j^{n+1,3} &= \tilde{c}_{30}(p_{j+1}^{n+1} - p_j^{n+1}) + \tilde{c}_{31} \hat{F}_j^{n,3} \cdot \Delta_+ \bar{p}_j^{n+1,1} + \tilde{c}_{32} \hat{F}_j^{n,2} \cdot \Delta_+ \bar{p}_j^{n+1,2}, \\ p_{j+1}^n - p_j^n &= \tilde{c}_{40}(p_{j+1}^{n+1} - p_j^{n+1}) + \tilde{c}_{41} \hat{F}_j^{n,3} \cdot \Delta_+ \bar{p}_j^{n+1,1} + \tilde{c}_{42} \hat{F}_j^{n,2} \cdot \Delta_+ \bar{p}_j^{n+1,2} \\ &\quad + \tilde{c}_{43} \hat{F}_j^{n,1} \cdot \Delta_+ \bar{p}_j^{n+1,3}. \end{aligned} \quad (27)$$

As in Proposition 4 and by the assumptions on the numerical flux and Δt we find the coefficients $C_{j,l}^n$ for $1 \leq n \leq N$, $j \in \mathbb{N}$ and $-3 \leq l \leq 3$ to be non-negative if $\tilde{c}_{j,k} \geq 0$ holds for all $j = 2, \dots, 4$ and $k = 0, \dots, 3$.

We first consider the RK2-scheme with $\tilde{c}_{20} = 1 - a_{23}^\dagger = 1 - b_1^{-1} b_2 a_{23}$, $\tilde{c}_{21} = a_{23}^\dagger = b_1^{-1} b_2 a_{21}$ and $b_2 - b_1 a_{23}^\dagger = b_2(1 - a_{21})$. Here, non-negativity is ensured by the conditions presented in Table 1.

In case of a RK3-scheme we have $\tilde{c}_{32} = a_{23}^\dagger = b_1^{-2} b_2 a_{21}$, $\tilde{c}_{31} = a_{13}^\dagger - a_{23}^\dagger a_{12}^\dagger = b_1^{-1} b_3 (a_{31} - a_{32} a_{21})$ and $\tilde{c}_{30} = 1 - a_{23}^\dagger - a_{13}^\dagger + a_{23}^\dagger a_{12}^\dagger \geq 0$ which holds if $b_1 - b_2 a_{21} - b_3 a_{31} + b_3 a_{32} a_{21} \geq 0$ is satisfied. Moreover, $\tilde{c}_{40} = b_1 a_{23}^\dagger + b_1 a_{13}^\dagger + b_2 a_{12}^\dagger - b_1 a_{12}^\dagger a_{23}^\dagger \geq 0$ and $\tilde{c}_{41} = b_3 - b_1 a_{13}^\dagger - b_2 a_{23}^\dagger + b_1 a_{12}^\dagger a_{23}^\dagger \geq 0$ if $b_2 a_{21} + b_3 a_{31} + b_3 a_{32} - b_3 a_{32} a_{21} \geq 0$ and $b_3(1 - a_{31} - a_{32} + a_{32} a_{21}) \geq 0$. Finally we have

$\tilde{c}_{42} = b_2 - b_1 a_{13}^\dagger = b_2 - b_3 a_{31}$ and $\tilde{c}_{43} = b_1$. Thus, the conditions of Table 1 imply non negativity of the coefficients $\tilde{c}_{j,k}$ and, by the representation (27), of $C_{j,l}^n$ for $n = 1, \dots, N$, $j \in \mathbb{N}$ and $k = -s, \dots, s$.

□

The next result provides conditions ensuring that Assumption 4 of Theorem 1 is satisfied.

Proposition 6. *Let the assumptions of Proposition 4 be satisfied and let the RK-scheme satisfy the conditions of Table 2. Then there exists a constant $\nu_s > 0$ such that*

$$l_j^{n+1,0} \leq l_{j,s}^{n,0,+} - \Delta t \nu_s (l_{j,s}^{n,0,+})^2$$

is satisfied (cf. (18)). Here, the constant ν_s depends on the number of stages s .

Stages	Conditions
2	$a_{21} \neq 0, \quad b_2 \neq 0$
3	$a_{21} \neq 0, \quad a_{32} \neq 0, \quad b_3 \neq 0$

Table 2: Conditions for OSLC consistency

Proof. We start proving the claim for $s = 2$. To establish the existence of $\nu_2 > 0$ we recall (21):

$$y_j^{n,2} = (1 - a_{21})y_j + a_{21}(y_j - \lambda(f(y_j, y_{j+1}) - f(y_{j-1}, y_j))), \quad (28)$$

$$y_j^{n+1} = b_2 a_{21} y_j + (b_1 - b_2 a_{21})(y_j - \lambda(f(y_j, y_{j+1}) - f(y_{j-1}, y_j))) \\ + b_2 (y_j^{n,2} - \lambda(f(y_j^{n,2}, y_{j+1}^{n,2}) - f(y_{j-1}^{n,2}, y_j^{n,2}))), \quad (29)$$

with non-negative weights by assumption. Since the discretization scheme, utilizing a single Euler step, satisfies (19), substituting (28) into (18) provides

$$l_j^{n,2} \leq (1 - a_{21})l_{j,0}^{n,0} + a_{21}(l_{j,1}^{n,0,+} - \Delta t \nu (l_{j,1}^{n,0,+})^2) \leq l_{j,1}^{n,0,+} - a_{21} \Delta t \nu (l_{j,1}^{n,0,+})^2 \quad (30)$$

where we used $0 < a_{21} \leq 1$ and for the ν form the explicit Euler discretization. Similarly, (29) can be estimated by

$$l_j^{n+1,0} \leq b_2 a_{21} l_{j,0}^{n,0} + (b_1 - b_2 a_{21})(l_{j,1}^{n,0,+} - \nu \Delta t (l_{j,1}^{n,0,+})^2) + b_2 (l_{j,1}^{n,1,+} - \nu \Delta t (l_{j,1}^{n,1,+})^2) \quad (31)$$

with

$$l_{j,1}^{n,2,+} = \max(0, l_{j-1}^{n,2}, l_j^{n,2}, l_{j+1}^{n,2}) \\ \leq \max(0, l_{j+1,1}^{n,0,+} - \Delta t \nu a_{21} (l_{j+1,1}^{n,0,+})^2, l_{j,1}^{n,0,+} - \Delta t \nu a_{21} (l_{j,1}^{n,0,+})^2, l_{j-1,1}^{n,0,+} - \Delta t \nu a_{21} (l_{j-1,1}^{n,0,+})^2)$$

where we utilized (30). For $\hat{\nu} = (4\lambda M_y a_{21} \nu \max\{1, \Delta t\})^{-1}$ we find the function $x - \Delta t \lambda \nu \hat{\nu} a_{21} x^2$ to be non-negative and monotonically increasing for $x \in [0, 2M_y]$. Consequently, we have

$$l_{j,1}^{n,2,+} = \max(l_{j-1}^{n,2}, l_j^{n,2}, l_{j+1}^{n,2}) \leq \max(l_{j+1,1}^{n,0,+}, l_{j,1}^{n,0,+}, l_{j-1,1}^{n,0,+}) - \Delta t \hat{\nu} (\max(l_{j+1,1}^{n,0,+}, l_{j,1}^{n,0,+}, l_{j-1,1}^{n,0,+}))^2.$$

For (31), by using $\max(l_{j+1,1}^{n,0,+}, l_{j,1}^{n,0,+}, l_{j-1,1}^{n,0,+}) = l_{j,2}^{n,0,+}$ and dropping the negative terms, the following upper bound can be established:

$$l_{j,2}^{n,0,+} - \Delta t \hat{\nu} b_2 (l_{j,2}^{n,0,+})^2.$$

Thus, for $\nu_2 = \hat{\nu} b_2$, the RK2 method satisfies the OSL-consistency conditions (19) where positivity is ensured by the conditions of Table 2. In case of $s = 3$ we first note that

$$l_{j,1}^{n,3} \leq l_{j,2}^{n,0,+} - \Delta t \hat{\nu} a_{32} (l_{j,2}^{n,0,+})^2$$

and $a_{32} > 0$ hold true. Now, the same arguments as in the case $s = 2$ provide the existence of $0 < \nu_3 = b_3 (4\lambda M_y a_{32} \hat{\nu} \max\{1, \Delta t\})^{-1}$ with

$$l_j^{n+1,0} \leq l_{j,3}^{n,0,+} - \Delta t \nu_3 (l_{j,3}^{n,0,+})^2,$$

which concludes the proof. \square

Finally we discuss assumption 5 of Theorem 1.

Proposition 7. *Let the assumptions of Proposition 4 be satisfied. Let a RK-scheme of order s with $s = 2, 3$ satisfy the conditions of Table 1. Then the numerical flux*

$$f^\Delta(y_{j-s+1}^n, \dots, y_{j+s}^n) = \sum_{i=1}^s b_i f^\Delta(y_j^{n,i}, y_{j+1}^{n,i})$$

satisfies Assumptions 1 and 5 of Theorem 1, respectively.

Proof. By assumption we have $f^\Delta(y, y) = f(y)$ for the original numerical flux. Consequently, one easily finds $y_j^{n,i} = y_j^n = y$ for $i = 1, \dots, s$ and $y_l = y$ for $l = j - s + 1, \dots, j + s$. A necessary condition for higher-order time stepping schemes is $\sum b_i = 1$ and consequently $f^\Delta(y_{j-s+1}^n, \dots, y_{j+s}^n) = f(y)$.

The second assertion follows from the corresponding properties of the original numerical flux, representation (21) and the conditions from Table 1 as an application of the chain rule. \square

Before summarizing the previous results, we revise the assumption $b_j \neq 0$ for $j = 1, \dots, s$, originally imposed for establishing the discretization scheme of the adjoint. The conditions presented in Table 2 not only enable the results of Proposition 7 but are also necessary for the a higher convergence order of the discretization of the state system. These order-conditions for RK-schemes are summarized, e.g., in [21, Table 2] and given by $b_2 a_{21} = 1/2$ in case of a RK2 scheme and $b_2 a_{21} + b_3 (a_{31} + a_{32}) = 1/2$, $b_3 a_{32} a_{21} = 1/6$ and $b_2 a_{21}^2 + b_3 (a_{31} + a_{32})^2 = 1/3$ for RK3 methods. Consequently, RK2 and RK3 schemes of order two and three automatically satisfy the conditions of Table 2 and further have strictly positive coefficients b_2 and b_3 . Strict positivity of b_1 now follows from the conditions in Table 1.

Theorem 2. *Let a semi-discretization of the conservation law (1) be given such that a time discretization by the explicit Euler method with a suitable time step Δt provides a full discretization (11) satisfying (A). Moreover, suppose that a RK-scheme for the time discretization of (10) is of order $s \leq 3$. Let the RK-scheme satisfy the conditions formulated in Table 1. Then the resulting full discretization of (1) provides a consistent discretization of (P) in the sense of Theorem 1.*

The set of RK methods that satisfy the conditions of Table 1 and 2 is discussed next. In [26], Runge-Kutta methods in terms of the Butcher-Arrays have been studied in the context of strong-stability-preserving time stepping. For $s = 2$, Heun's method

$$\begin{array}{cc} 0 & 0 \\ 1 & 0 \\ \hline 1/2 & 1/2 \end{array}, \quad (32)$$

and in case of $s = 3$ the method defined by the array

$$\begin{array}{ccc} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 1/4 & 1/4 & 0 \\ \hline 1/6 & 1/6 & 2/3 \end{array} \quad (33)$$

are the only time stepping schemes satisfying the conditions of Table 1 and 2. They provide a method that has $c_{SSP} = 1$ in (20) and the associated number of stages matches the order of convergence.

4 Problems of optimal control and order of convergence for the adjoint scheme

The problem of optimal control subject to a system of ordinary differential equations resulting from a semidiscretization of (1) is given by

$$\begin{aligned} & \text{minimize } \frac{\Delta x}{2} |\mathbf{y}(T) - \mathbf{y}^d|_{l^2}^2 + \mathfrak{R}(\mathbf{u}) = \mathcal{J}^\Delta(\mathbf{y}(T), \mathbf{u}), \\ & \text{over } (\mathbf{y}(T), \mathbf{u}) \in \mathbb{R}^N \times \mathbb{R}^{\tilde{N}}, \\ & \text{subject to } \mathbf{y} \text{ satisfies (10),} \end{aligned} \quad (P_\Delta)$$

and can be interpreted as a Mayer problem in the context of optimal control of ODE's where N, \tilde{N} depend on the spatial discretization. First order necessary optimality conditions for (P_Δ) , as given in [14, Theorem 4.2.i], involve the adjoint ODE system

$$\begin{aligned} \dot{\mathbf{p}} &= -[F(\mathbf{y})\mathbf{p}]_{\mathbf{y}} \\ &= \Delta x^{-1} [f_2(y_{j-1}(t), y_j(t))(p_{j+1}(t) - p_j(t)) - f_1(y_j(t), y_{j+1}(t))(p_j(t) - p_{j-1}(t))] \end{aligned}$$

with terminal condition $\mathbf{p}(\mathbf{T}) = \mathcal{J}_{\mathbf{y}(T)}^\Delta(\mathbf{y}(T), \mathbf{u})$ as discussed in [21]. Comparing the adjoint system obtained from the optimal control of ordinary differential equations with (16), obtained via sensitivity scheme and discrete duality relation, we find them to be equal.

This equivalence allows for utilizing the order conditions on RK time stepping methods in the context of optimal control of ODE's from [21] also in the context of optimal control for scalar conservation laws. Consulting the conditions formulated in [21, Table 1], Heun's scheme has convergence order 2 with

respect to time for both, state and adjoint system. In case of the third order method, the additional condition $(b_2 a_{21} + b_3 a_{31})^2/b_1 + (b_3 a_{32})^2/b_2 = 1/3$ is not satisfied by the coefficients in (33) and consequently, the approximation of the adjoint is merely of order 2. This observation coincides with the results of [22], where certain numerical flux functions were investigated and conditions on the coefficients of the associated Shu-Osher representation that guarantee total variation stability of the adjoint discretization were derived. In addition, an upper bound for the order of convergence was established, namely imposing the SSP property on the time discretization of the state system and stability on the discretization of the adjoint system limits the order of convergence for the latter scheme to 2. As outlined before, the consideration of (32) and (33) is sufficient since they represent the only methods with $c_{SSP} = 1$ in (20) and order matching the number of stages.

Let $\mathcal{S}^s : \mathbb{R}^{\tilde{N}} \rightarrow \mathbb{R}^N$ denote the solution operator of the underlying full discretization of (1) with s stages. Then the fully discrete problem is given as

$$\begin{aligned} & \text{minimize } \mathcal{J}^\Delta(\mathbf{y}(T), \mathbf{u}), \\ & \text{over } (\mathbf{y}(T), \mathbf{u}) \in \mathbb{R}^N \times \mathbb{R}^{\tilde{N}}, \\ & \text{subject to } \mathbf{y}(T) = \mathcal{S}^s(\mathbf{u}), \end{aligned} \quad (34)$$

where the discrete desired state is obtained by averaging the given function as discussed, in case of the initial datum, in (9).

By design, \mathcal{S}^s is Lipschitz continuous with a constant depending on the Lipschitz constant of the chosen numerical flux, the number of time steps s , and the coefficients of the RK-scheme even if these coefficients do not meet the conditions of Table 1. Consequently, following the direct method from the calculus of variations, there exists a solution of (34) independently of the chosen mesh width Δx .

5 Numerical Experiments

We end this paper by a report on numerical experiments associated with (P). In fact, utilizing a Tikhonov-type cost for the control (R) and employing a standard gradient descent scheme with Armijo line search, a numerical study when using Burgers' equation as the underlying state system is conducted.

5.1 Regularization Term

A possible choice for the regularization term is given by $\mathfrak{R}(u) = \|u\|_{BV(\mathbb{R})}$ that ensures boundedness of u in $L^1(\mathbb{R})$ and coercivity with respect to $L^\infty(\mathbb{R})$ due to the continuous embedding $BV(\mathbb{R}) \hookrightarrow L^\infty(\mathbb{R})$. To avoid additional problems due to the non-differentiability of this choice, our numerical experiments utilize an $H^1(\mathbb{R})$ -type cost. However, the assumptions on the cost term ensure a bounded support of the optimal control u^* and the associated solution $y^* = y(u^*)$ of (1), provided the desired state has bounded support as well. Since $u \equiv 0$ is a feasible control with corresponding solution $y^0 := y(0)$ of (1), we obtain

$$|u^*|_{L^\infty(\mathbb{R})} \leq \mathfrak{R}(u^*) \leq \mathcal{J}(y^*, u^*) \leq \mathcal{J}(y(0), 0) = \frac{1}{2} |y^0(T) - y^d|_{L^2(\mathbb{R})}^2 =: M_u.$$

By the maximum principle (see Proposition 1), we further have $|y^*(t, \cdot)|_{L^\infty(\mathbb{R})} \leq M_u$ for all $t \in [0, T]$, bounding the overall characteristic speed by

$$\max_{\eta \in [-M_u, M_u]} |f'(\eta)| = M_f.$$

The maximum exists by assumptions on the flux function f and the extreme value theorem. Consequently, if (a, b) with $-\infty < a < b < \infty$ denotes the support of y^d , (\tilde{a}, \tilde{b}) with $\tilde{a} = a - M_f T$ and $\tilde{b} = b + M_f T$ represents the part of the domain of u^* that might influence $y^*(T, \cdot)$ on (a, b) directly. Finally, optimality enforces u^* to tend to zero outside of (\tilde{a}, \tilde{b}) . As a consequence we consider

$$\mathfrak{R}(u) = \frac{\alpha}{2} \|u\|_{H_0^1(\Omega)}^2$$

with $\alpha > 0$ and $(\tilde{a}, \tilde{b}) \subset \Omega$ chosen large enough, to allow u^* to tend to zero. Following the discussion above, this restriction of the domain does not change the problem since optimal control and the corresponding state are zero outside Ω and $[0, T] \times \Omega$ anyway. Although this kind of regularization enforces the initial data to be continuous, shock phenomena may still occur in the solution of (1). For example, Burgers' equation with smooth initial data u and some $x \in \mathbb{R}$ with $u'(x) < 0$ develops shocks in finite time $T = -1/\min\{u'(x)\}$; see [28].

5.2 The Algorithm

Since the adjoint equation allows for a compact gradient representation of (2), we utilize a steepest descent method to solve the discretized problem (34) that is given in Algorithm 1.

Algorithm 1 Solution Algorithm

- 1: Choose $u^{(1)}$, set $k = 1$.
- 2: **while** stopping criterion not satisfied **do**
- 3: Solve discretization of primal equation (10) according to (21) to obtain $y_{(k)}^T$
- 4: Evaluate $\mathcal{J}^\Delta(y_{(k)}, u^{(k)})$
- 5: Solve adjoint equation (16) according to (25) to obtain $p_{(k)}^0$
- 6: Compute the update direction $\delta u^{(k)}$
- 7: Perform line search to obtain $\theta > 0$ with

$$\mathcal{J}^\Delta(y(u_\theta^{(k)}), u_\theta^{(k)}) - \mathcal{J}^\Delta(y(u^{(k)}), u^{(k)}) \leq -\theta \sigma \|\delta u^{(k)}\|_{H_0^1(\Omega)}^2 \text{ for } u_\theta^{(k)} := u^{(k)} + \theta \delta u^{(k)} \quad (35)$$

- 8: Update control, i.e. $u^{(k+)} := u_\theta^{(k)}$, set $k := k + 1$
 - 9: **end while**
-

Recalling (4), for a control $u^{(k)} \in H_0^1(\Omega)$ the reduced gradient of (P) , i.e., $p_{(k)}^0 - \alpha \Delta u^{(k)}$, is an element of $H^{-1}(\Omega)$ when considering the function space setting. Thus, the update direction $\delta u^{(k)}$ in $H_0^1(\Omega)$ is given by the Riesz representative (see, e.g., [24]) $\delta u^{(k)} = v - \alpha u^{(k)}$ with $v \in H_0^1(\Omega)$ solving

$$\Delta v = p_{(k)}^0 \text{ in } \Omega.$$

The line search is realized by a backtracking method halving the trial step in each iteration until the Armijo condition (35) is satisfied with $\sigma = 10^{-4}$ as suggested in [29].

The algorithm terminates when the relative stopping criterion

$$\|\delta u^{(k)}\|_{H_0^1(\Omega)} \leq \varepsilon_G + \varepsilon_D \|\delta u^{(1)}\|_{H_0^1(\Omega)}$$

is satisfied with $\varepsilon_G = \varepsilon_D = 10^{-6}$.

In all numerical tests, the weight for the regularization term is $\alpha = 10^{-5}$.

5.3 Examples

The numerical tests consider Burgers' equation

$$y_t + \frac{1}{2}[y^2]_x = 0 \quad (36)$$

on the computational domain $(0, 3)$ with homogenous Dirichlet boundary data. The chosen computational domain is large enough such that the support of the primal and adjoint equations, respectively, are included as discussed above. The control is discretized by standard P_1 conforming finite elements on a mesh of width Δx and its integral average on the cells is computed exactly for the approximation.

We utilize the following numerical flux functions. First, we use the Engquist-Oscher scheme employing the numerical flux

$$f^{EO} : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad f^{EO}(v_1, v_2) = \int_{\bar{v}}^{v_1} f'(\xi)^+ d\xi + \int_{\bar{v}}^{v_2} f'(\xi)^- d\xi - f(\bar{v}) \quad (EO)$$

with $\bar{v} \in \mathbb{R}$ arbitrary and $f'(\xi)^+ := \max\{0, f'(\xi)\}$, $f'(\xi)^- := \min\{0, f'(\xi)\}$ denoting positive and negative part of $f'(\xi)$, respectively. If it exists, \bar{v} is chosen to be the sonic point of the flux function with $f'(\bar{v}) = 0$ to simplify computations. Applied to (36), the solution of the semi discretization (10) admits a solution that is twice continuously differentiable with respect to time. In case of (EO), it was shown in [33] that the assumptions of Theorem 1 are met for time steps Δt such that $\lambda = \Delta t / \Delta x$ satisfies $\lambda \sup_{|y| \leq M_y} |f'(y)| \leq (1 - \rho)2^{-1}$ and some $\rho \in (0, 1)$.

Second, we also utilize the modified Lax-Friedrichs scheme with numerical flux

$$f^{LF} : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad f^{LF}(v_1, v_2) = \frac{1}{2}(f(v_1) + f(v_2)) + \frac{\gamma \Delta x}{2\Delta t}(v_1 - v_2) \quad (LF)$$

for a parameter $\gamma \in (0, 1)$ (the original Lax-Friedrichs numerical flux is obtained for $\gamma = 1$). For this choice, the solution of (10) is smooth. Again, it was shown in [33] that the assumptions of Theorem 1 are satisfied for a time discretization with an explicit Euler method and time steps such that

$$\lambda \sup_{|y| \leq M_y} |f'(y)| \leq \min\{(1 - \rho) \min\{\gamma, 2(1 - \gamma)\}, 1 - \gamma\}$$

holds with some $\rho \in (0, 1)$.

Example 1 presents the performance of the algorithm and demonstrates the necessity of employing RK schemes that satisfy the conditions in Table 1 to obtain a consistent discretization of (P) for the same step length as in case of the basic explicit Euler time stepping. We consider the desired, piecewise constant state

$$y^d(x) = \begin{cases} \frac{1}{3}, & x \in [\frac{6}{5}, \frac{5}{3}], \\ -\frac{1}{10}, & x \in (\frac{5}{3}, \frac{7}{4}], \\ 0, & \text{else,} \end{cases}$$

depicted in red in Figure 2, a spatial discretization of width $\Delta x = 500^{-1}$, $T = 1$ and set

$$\lambda = 1.0.$$

Moreover, we utilize the Engquist-Osher numerical flux and thus, λ is chose sufficiently small to satisfy $\lambda \sup_{|y| \leq M_y} |f'(y)| < 2^{-1}$ for all initial conditions with $|u| < 1/2$ since $f'(y) = y$ and the maximum principle for entropy solutions (see Proposition 1) ensures $\|y\|_{L^\infty(\mathcal{Q})} \leq \|u\|_{L^\infty(\mathbb{R})}$.

For Heun's scheme (32), Algorithm 1 converged within 149.587 iterations for the initialization $u^{(1)} \equiv 0$. The optimal control can be found in Figure 1 while the corresponding state at $t = T$ is presented in Figure 2.

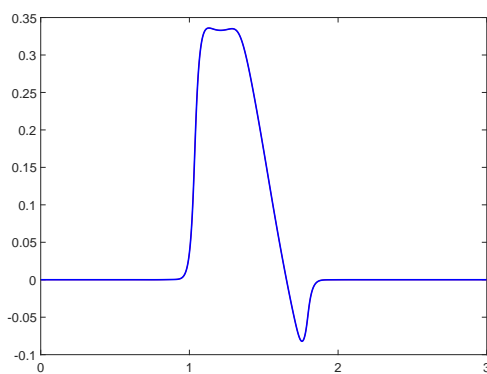


Figure 1: Optimal control

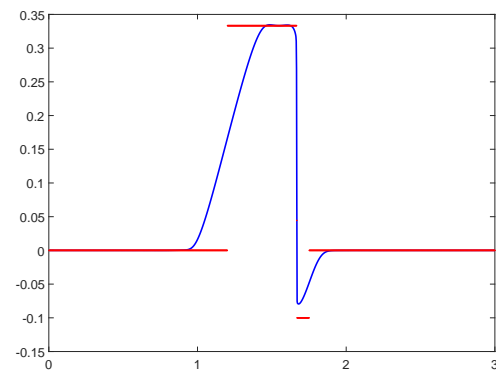


Figure 2: Solution of (1) (blue) vs desired state (red)

Figure 3 displays the final time data for the adjoint state, the difference of desired state and solution to the state system evaluated at $t = T$, clearly containing discontinuities. This, in particular, requires the adjoint scheme to be TVD-stable to avoid spurious oscillations in the solution. Figure 4 displays this solution to the discretized adjoint equation evaluated at $t = 0$, not showing oscillations.

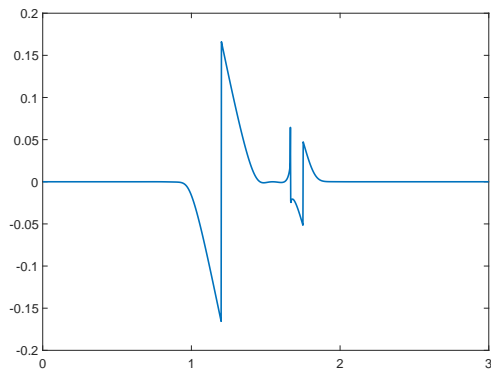
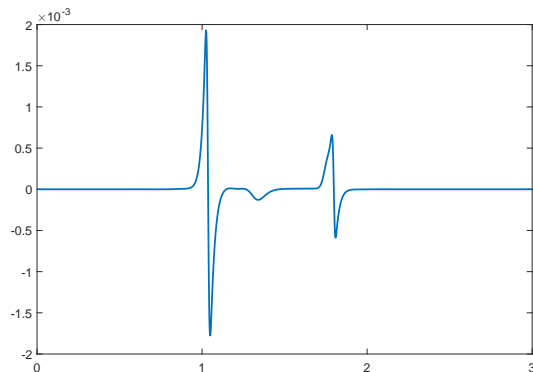


Figure 3: Final time data for the adjoint equation.

Figure 4: Solution to the adjoint equation in $t = 0$.

To report on the convergence behavior of the algorithm, we show the behavior of the objective value and the $H_0^1(\Omega)$ -norm of the update direction in Figure 5 and 6, respectively.

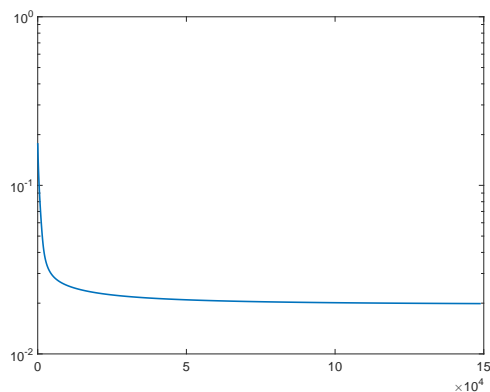


Figure 5: Evolution of the objective value.

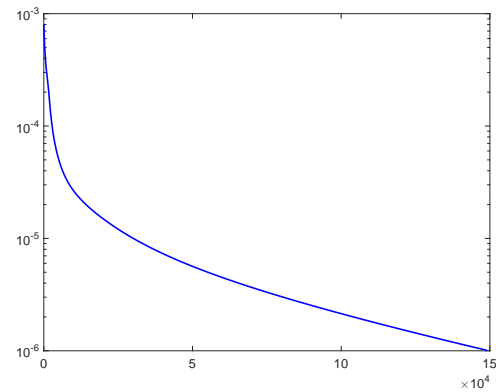


Figure 6: Evolution of the norm of the update direction.

For the second purpose of this example we change the initialization of Algorithm 1 to

$$u^{(1)}(x) = \max\{0, \min\{ax, a - ax\}\}, \quad u^{(1)} \in H_0^1(\Omega)$$

with $a = 2.075$ and consider $\Delta x = 1000^{-1}$. Moreover, we verify in every iteration and each time step of step 3 and 5 of Algorithm 1, that the total variation of state and adjoint at this time is bounded by the total variation if the initial and final time data respectively, i.e., whether

$$|\mathbf{y}^i|_{TV} \leq |\mathbf{y}^0|_{TV} \text{ and } |\mathbf{p}^{i-1}|_{TV} \leq |\mathbf{p}^{N_T}|_{TV}$$

holds for $i = 1, \dots, N_T$ with N_T denoting the number of timesteps in the discretization scheme of state

and adjoint equation and the total variation seminorm $|v|_{TF} = \sum_{j=1}^{N-1} |v_{j+1} - v_j|$. If this is not the case the algorithm terminates since the corresponding iterate does not satisfy the requirements of Theorem 1. The choice of the initialization $u^{(1)}$ violates $\lambda \sup_{|y| \leq M_y} |f'(y)| \leq (1 - \rho)2^{-1}$ for some $\rho \in (0, 1)$ but on the given mesh, Algorithm 1 still converges in case of the Heun scheme (32). Figure 8 displays the corresponding optimal control while Figure 7 displays the norm of the update direction along the iterations of the algorithm.

In addition we consider the RK2 scheme

$$\begin{array}{cc} 0 & 0 \\ \frac{1}{2(1-10^{-4})} & 0 \\ 10^{-4} & (1-10^{-4}) \end{array}, \quad (37)$$

clearly violating the conditions of Table 1 as

$$0 \leq b_1 - b_2 a_{21} = 10^{-4} - 1/2.$$

In the same setting, Algorithm 1 terminates in the very first iteration as the computed discrete state fails to satisfy $|y^i|_{TV} \leq |y^0|_{TV}$.

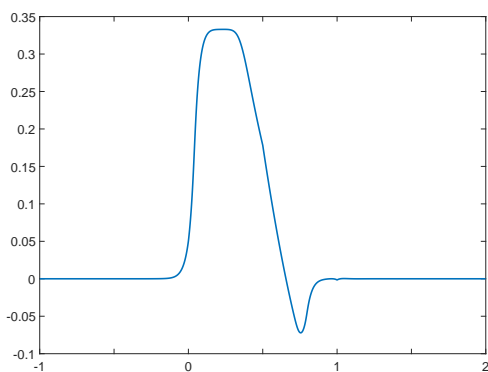


Figure 7: Optimal Control.

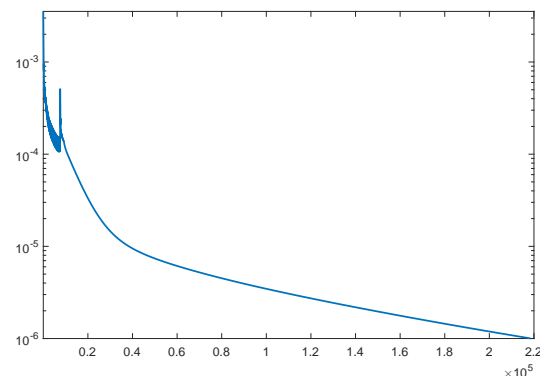


Figure 8: Evolution of the norm of the update direction.

In this example we reported on the convergence behavior of Algorithm 1 and demonstrated the importance of choosing RK-schemes that satisfy the conditions of Table 1. Although the problem with the RK2 scheme (37) can be circumvented by choosing a smaller time step, this is not recommended in the context of optimal control since this increases the numerical effort to solve the problem with respect to computing time along with storage requirements.

Example 2 considers a setting, where shocks appear neither in the solution to the state equation nor in the final time data due to the chosen desired state that is taken from [2]. In the latter work, long time behavior of optimization algorithms for Burgers' equation with respect to several methodologies was

studied. We only transform the domain such that the target function is contained in $(0, 3)$. The desired state is given by

$$y^d(x) = \frac{3}{2000} \left(-e^{-(5\sqrt{20}-\phi(x))^2} + e^{-(2\sqrt{20}+\phi(x))^2} + \sqrt{\pi}\phi(x)(\operatorname{erf}(5\sqrt{20}-\phi(x)) + \operatorname{erf}(2\sqrt{20}+\phi(x))) \right)$$

and it is depicted in Figure 9. Here, φ represents a linear transformation of arguments from $[0, 3]$ to the interval $[-50, 100]$, a domain in the scale of the original desired state from [2]. To compensate for gradient scaling in the initial state because of the transformed domain, we refrain from long time behavior analysis and consider $T = 1$ as well as $\Delta x = 5 \cdot 10^{-2}$. This setting allows for a discussion concerning the approximation order of the solutions to the state and adjoint equations as we can expect them to be regular enough for the Taylor-expansion that provides the order of convergence of the time stepping. For both numerical flux functions, we will obtain a reference solution based on the RK3 scheme (33) and compare it with solutions obtained by the explicit Euler and Heun's method. Moreover, we will analyze the order of convergence of the RK methods at the example of the Engquist-Osher numerical flux function. To this end, a further reference solution based on the RK3 scheme and for a time step $\Delta \tilde{t} = 2^{-4} \Delta t$ with $\Delta t = \lambda \Delta x$ is generated. Then, solutions to the problem of optimal control are generated for the RK1, RK2 and RK3 scheme with time step sizes $\Delta \tilde{t} = 2^{-j} \Delta t$, $j = 0, \dots, 3$ are obtained and compared to the reference solution by providing

$$E_y = \max_{j,n} |y_j^n - \tilde{y}_j^n|, \quad E_p = \max_{j,n} |p_j^n - \tilde{p}_j^n|.$$

First we use the numerical flux of the Lax-Friedrichs scheme (LF) with $\gamma = 1/2$ and $\lambda = 2.75$ since in this case the regularity of the solution to (10) with respect to time is sufficient for the application of a RK3-scheme. Concerning the numerical flux functions, this choice of λ satisfies the conditions for the explicit Euler time stepping as formulated in [33]. The solution obtained by (33) provides the reference solution for this test and the corresponding optimal control is presented in Figure 10.

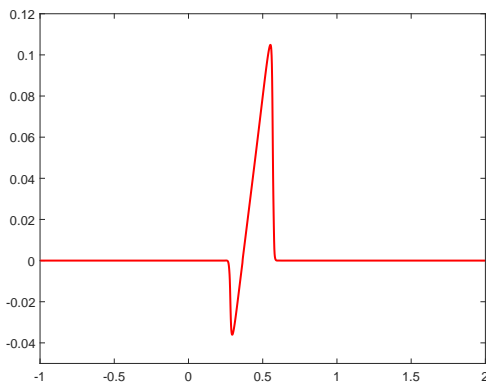


Figure 9: Desired State

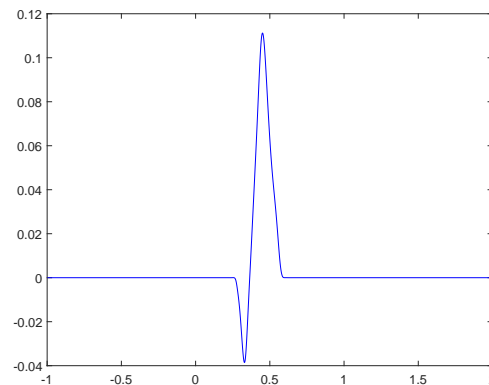


Figure 10: Opt. control for the RK3 scheme.

Further, we computed the solution for the same parameters utilizing an explicit Euler time discretization and Heun's method (32). In Figure 11 (a), (b) and (c) we present the solutions to the semi-discrete

forward problems (10), with the controls obtained by Algorithm 1, evaluated at $t = T$ and the difference to the desired state. In Figure 11 (d), the differences of solutions to the lower order time discretization schemes and the reference solution are plotted at $t = T$.

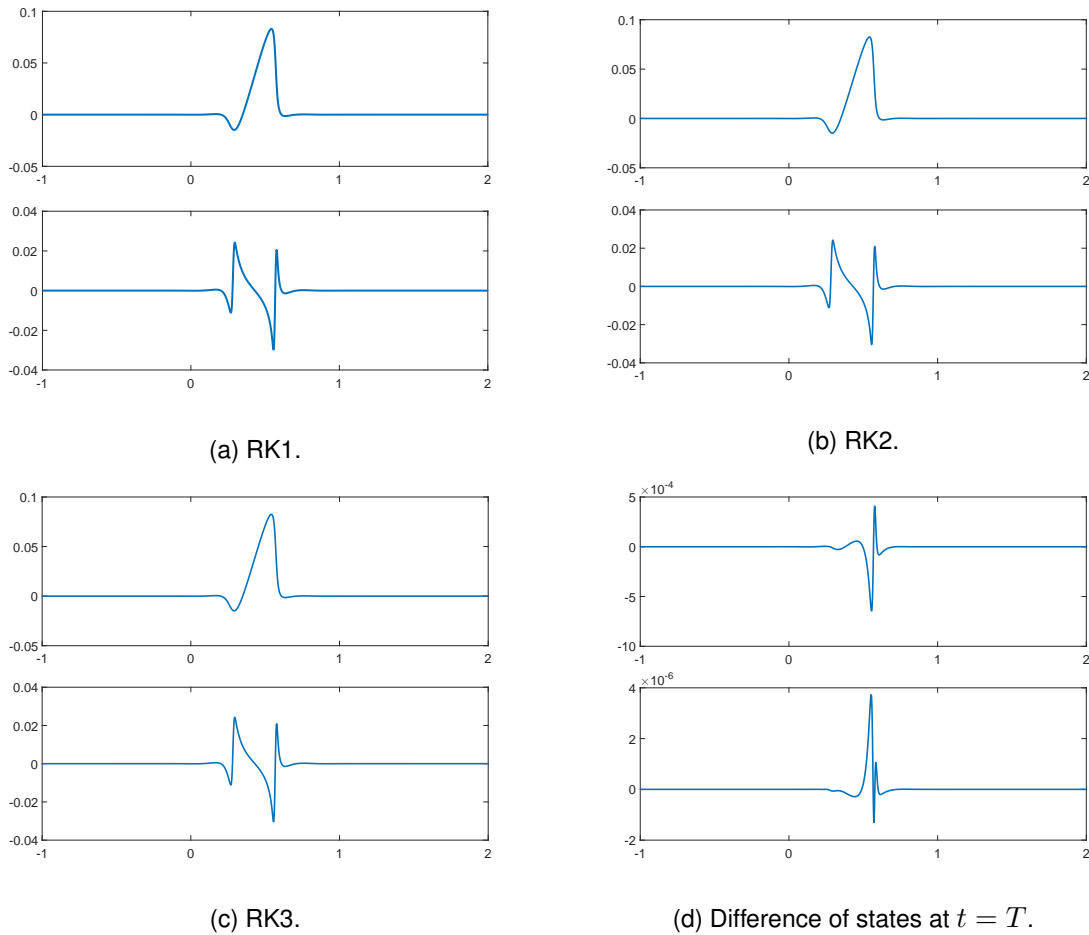


Figure 11: States at $t = T$ and differences to desired state.

As it can be seen in Figure 11 (d), the RK2-scheme approximates the solution obtained from applying the RK3-scheme more accurately than the explicit Euler scheme. In Figure 12 we have depicted the differences in the optimal control obtained by Algorithm 1 for the lower order schemes to the reference solutions. Finally, the following tables provides E_y and E_p .

RK1	RK2
$6.4 \cdot 10^{-4}$	$1.3 \cdot 10^{-6}$

RK1	RK2
$1.2 \cdot 10^{-5}$	$2.2 \cdot 10^{-7}$

Table 3: L^∞ -discrepancy of states to reference state E_y .

Table 4: L^∞ -discrepancy of adjoint to reference adjoint E_p .

Again, we observe that the second-order scheme approximates the solution of the RK3-scheme better

than the first-order scheme by two orders of magnitude, which, by $\Delta t = 5 \cdot 10^{-2}/\lambda \approx 2 \cdot 10^{-2}$ corresponds to the behavior expected for the approximation order for RK-schemes.

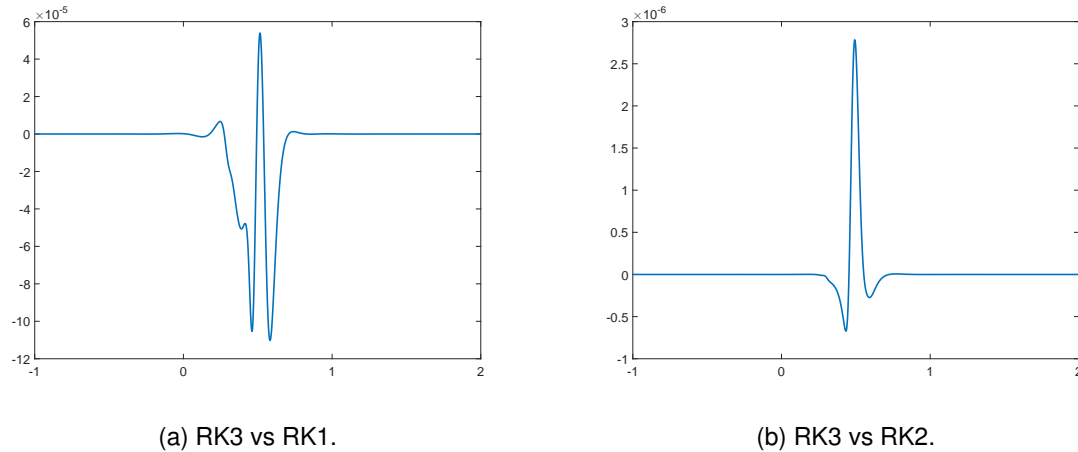


Figure 12: Difference of the optimal control.

Next we test the problem of optimal control with a spatial discretization based on the Engquist-Osher scheme (EO) with $\lambda = 1.0$. The scheme is less diffusive but only provides a solution of (10) which is twice continuously differentiable with respect to time. Again we obtained the reference solution by the RK3-scheme defined in (33) and depict the corresponding optimal control in Figure 13.

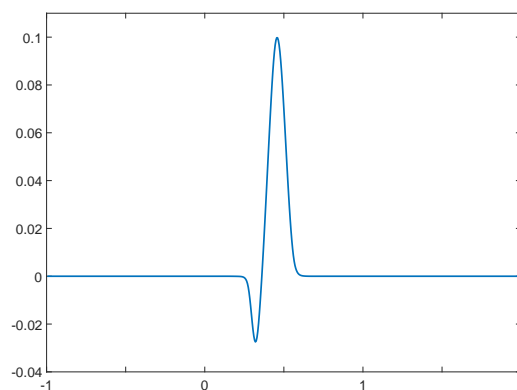


Figure 13: Optimal control for RK3 scheme.

Further, we computed the solution for the same parameters utilizing an explicit Euler time discretization and Heun's method (32). Again, we depict the optimal states evaluated at $t = T$ and the difference to the desired state in Figure 14 (a), (b) and (c) while we find the differences of the optimal states for the RK1- and RK2-scheme at $t = T$ and the reference solution in Figure 14 (d).

In Figure 15 we have depicted the differences in the controls obtained by Algorithm 1 for the lower order time discretization methods to the reference solutions. Again we observe a better approximation of the

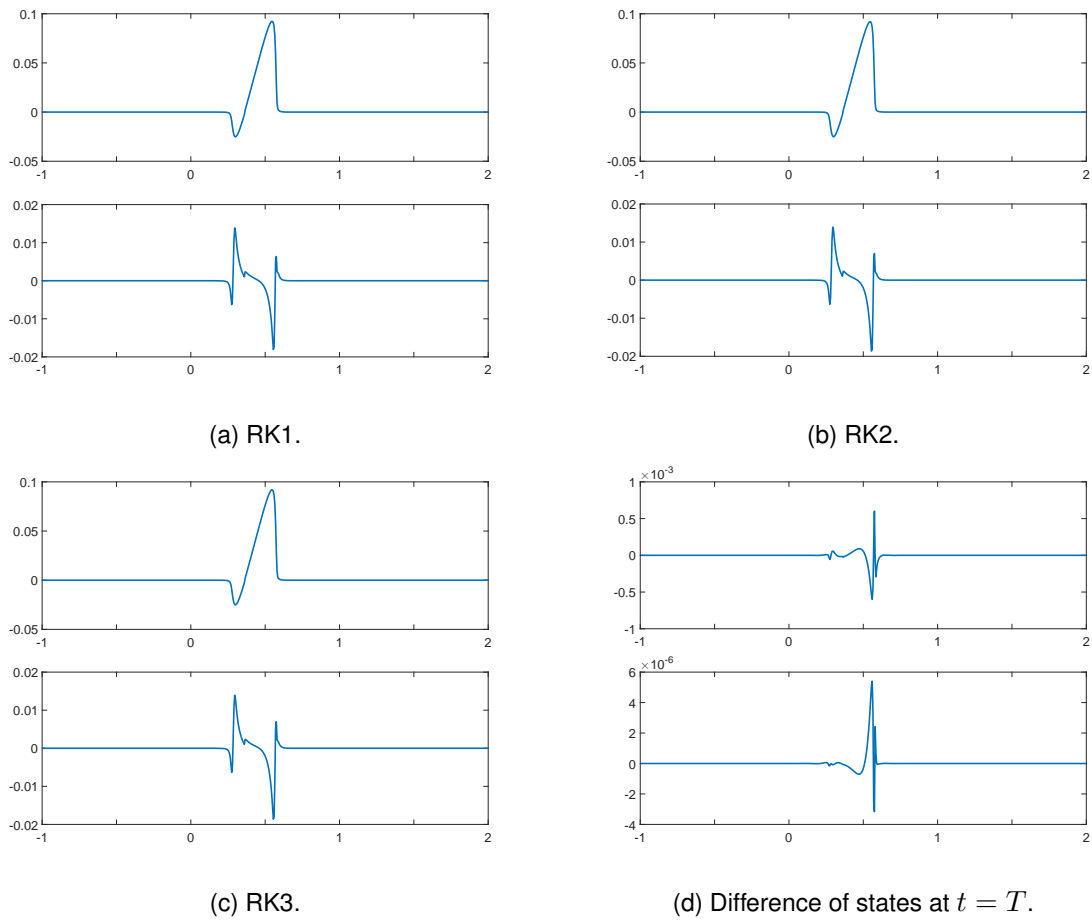


Figure 14: States at $t = T$ and differences to desired state.

reference solution in case of the RK2 scheme.

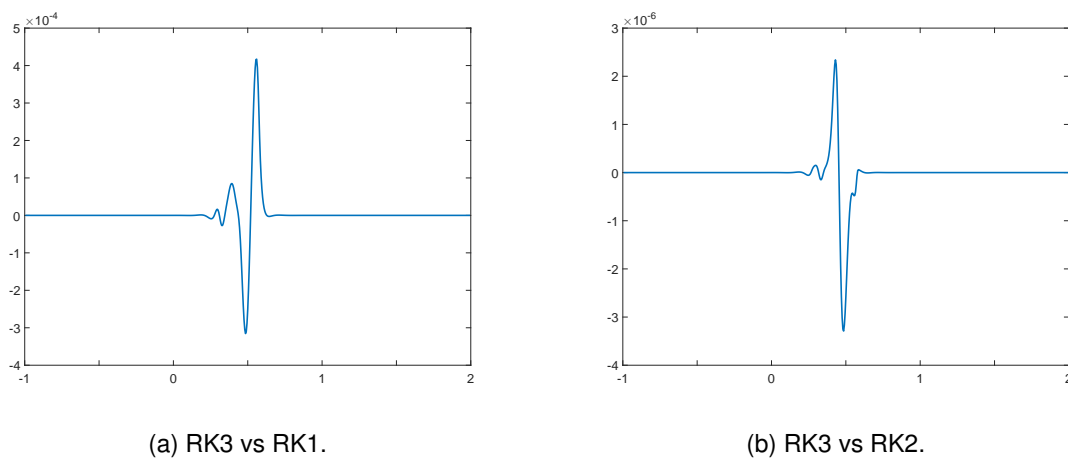


Figure 15: Difference of the optimal control.

In order to quantify the convergence order, we next generate a new reference solution by utilizing the RK3 scheme and a time step $\Delta\tilde{t} = 2^{-4}\Delta t$. Then we compare solutions to the problem of optimal control for the RK1, RK2 and RK3 schemes, respectively, and time steps $\Delta\tilde{t} = 2^{-j}\Delta t$ for $j = 0, \dots, 3$, i.e., we compare the states to the optimal control and the corresponding adjoints. The following tables provide E_y and E_p in case of the Engquist-Osher numerical flux function.

2^{-j}	RK1	RK2	RK3
1	$6.4 \cdot 10^{-4}$	$5.4 \cdot 10^{-6}$	$5.1 \cdot 10^{-8}$
1/2	$3.2 \cdot 10^{-4}$	$1.3 \cdot 10^{-6}$	$6.2 \cdot 10^{-9}$
1/4	$1.6 \cdot 10^{-4}$	$3.3 \cdot 10^{-7}$	$7.6 \cdot 10^{-10}$
1/8	$8.0 \cdot 10^{-5}$	$8.3 \cdot 10^{-8}$	$8.4 \cdot 10^{-11}$

Table 5: L^∞ -discrepancy of states to reference state E_y .

2^{-j}	RK1	RK2	RK3
1	$7.0 \cdot 10^{-5}$	$4.7 \cdot 10^{-5}$	$4.7 \cdot 10^{-5}$
1/2	$3.4 \cdot 10^{-5}$	$2.2 \cdot 10^{-5}$	$2.2 \cdot 10^{-5}$
1/4	$1.5 \cdot 10^{-5}$	$9.7 \cdot 10^{-6}$	$9.7 \cdot 10^{-6}$
1/8	$6.2 \cdot 10^{-6}$	$3.3 \cdot 10^{-6}$	$3.3 \cdot 10^{-6}$

Table 6: L^∞ -discrepancy of adjoints to reference adjoint E_p .

This example demonstrated the convergence order of the corresponding discretization schemes. In particular we observe the order bound for the adjoint in Table 6.

References

- [1] R. A. Adams and J. J. F. Fournier. *Sobolev Spaces*, volume 140 of *Pure and Applied Mathematics (Amsterdam)*. Elsevier/Academic Press, Amsterdam, second edition, 2003.
- [2] N. Allahverdi, A. Pozo, and E. Zuazua. Numerical aspects of large-time optimal control of Burgers equation. *ESAIM Math. Model. Numer. Anal.*, 50(5):1371–1401, 2016.
- [3] M. K. Banda and M. Herty. Adjoint IMEX-based schemes for control problems governed by hyperbolic conservation laws. *Comput. Optim. Appl.*, 51(2):909–930, 2012.
- [4] S. Bianchini. On the shift differentiability of the flow generated by a hyperbolic system of conservation laws. *Discrete Contin. Dynam. Systems*, 6(2):329–350, 2000.
- [5] J. F. Bonnans and J. Laurent-Varin. Computation of order conditions for symplectic partitioned Runge-Kutta schemes with application to optimal control. *Numer. Math.*, 103(1):1–10, 2006.

- [6] F. Bouchut and F. James. One-dimensional transport equations with discontinuous coefficients. *Nonlinear Anal.*, 32(7):891–933, 1998.
- [7] F. Bouchut and F. James. Differentiability with respect to initial data for a scalar conservation law. In *Hyperbolic problems: theory, numerics, applications, Vol. I (Zürich, 1998)*, volume 129 of *Internat. Ser. Numer. Math.*, pages 113–118. Birkhäuser, Basel, 1999.
- [8] A. Bressan. *Hyperbolic systems of conservation laws*, volume 20 of *Oxford Lecture Series in Mathematics and its Applications*. Oxford University Press, Oxford, 2000. The one-dimensional Cauchy problem.
- [9] A. Bressan and G. Guerra. Shift-differentiability of the flow generated by a conservation law. *Discrete Contin. Dynam. Systems*, 3(1):35–58, 1997.
- [10] A. Bressan and A. Marson. A maximum principle for optimally controlled systems of conservation laws. *Rend. Sem. Mat. Univ. Padova*, 94:79–94, 1995.
- [11] A. Bressan and A. Marson. A variational calculus for discontinuous solutions of systems of conservation laws. *Comm. Partial Differential Equations*, 20(9-10):1491–1552, 1995.
- [12] A. Bressan and W. Shen. Optimality conditions for solutions to hyperbolic balance laws. In *Control methods in PDE-dynamical systems*, volume 426 of *Contemp. Math.*, pages 129–152. Amer. Math. Soc., Providence, RI, 2007.
- [13] A. Bressan and W. Shen. Optimality conditions for solutions to hyperbolic balance laws. In *Control methods in PDE-dynamical systems*, volume 426 of *Contemp. Math.*, pages 129–152. Amer. Math. Soc., Providence, RI, 2007.
- [14] L. Cesari. *Optimization - Theory and Applications*, volume 17 of *Applications of Mathematics (New York)*. Springer-Verlag, New York, 1983. Problems with ordinary differential equations.
- [15] N. Dunford and J. T. Schwartz. *Linear Operators. Part I*. Wiley Classics Library. John Wiley & Sons, Inc., New York, 1988. General theory, With the assistance of William G. Bade and Robert G. Bartle, Reprint of the 1958 original, A Wiley-Interscience Publication.
- [16] M. Giles and S. Ulbrich. 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(3):882–904, 2010.
- [17] M. Giles and S. Ulbrich. Convergence of linearized and adjoint approximations for discontinuous solutions of conservation laws. Part 2: Adjoint approximations and extensions. *SIAM J. Numer. Anal.*, 48(3):905–921, 2010.
- [18] M. B. Giles. Discrete adjoint approximations with shocks. In *Hyperbolic problems: theory, numerics, applications*, pages 185–194. Springer, Berlin, 2003.

- [19] S. Gottlieb and C.-W. Shu. Total variation diminishing Runge-Kutta schemes. *Math. Comp.*, 67(221):73–85, 1998.
- [20] S. Gottlieb, C.-W. Shu, and E. Tadmor. Strong stability-preserving high-order time discretization methods. *SIAM Rev.*, 43(1):89–112 (electronic), 2001.
- [21] W. W. Hager. Runge-Kutta methods in optimal control and the transformed adjoint system. *Numer. Math.*, 87(2):247–282, 2000.
- [22] S. Hajian, M. Hintermüller, and S. Ulbrich. Total variation diminishing schemes in optimal control of scalar conservation laws, 2017. WIAS-Preprint, DOI 10.20347/WIAS.PREPRINT.2383, Submitted.
- [23] M. Herty and B. Piccoli. A numerical method for the computation of tangent vectors to 2×2 hyperbolic systems of conservation laws. *Commun. Math. Sci.*, 14(3):683–704, 2016.
- [24] M. Hinze, R. Pinnau, M. Ulbrich, and S. Ulbrich. *Optimization with PDE constraints*, volume 23 of *Mathematical Modelling: Theory and Applications*. Springer, New York, 2009.
- [25] D. I. Ketcheson. Highly efficient strong stability-preserving Runge-Kutta methods with low-storage implementations. *SIAM J. Sci. Comput.*, 30(4):2113–2136, 2008.
- [26] J. F. B. M. Kraaijevanger. Contractivity of Runge-Kutta methods. *BIT*, 31(3):482–528, 1991.
- [27] S. N. Kružkov. First order quasilinear equations with several independent variables. *Mat. Sb. (N.S.)*, 81 (123):228–255, 1970.
- [28] R. J. LeVeque. *Numerical Methods for Conservation Laws*. Lectures in Mathematics ETH Zürich. Birkhäuser Verlag, Basel, second edition, 1992.
- [29] J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer, New York, second edition, 2006.
- [30] S. Pfaff and S. Ulbrich. Optimal boundary control of nonlinear hyperbolic conservation laws with switched boundary data. *SIAM J. Control Optim.*, 53(3):1250–1277, 2015.
- [31] S. J. Ruuth and R. J. Spiteri. Two barriers on strong-stability-preserving time discretization methods. In *Proceedings of the Fifth International Conference on Spectral and High Order Methods (ICOSAHOM-01) (Uppsala)*, volume 17, pages 211–220, 2002.
- [32] E. F. Toro. *Riemann Solvers and Numerical Methods for Fluid Dynamics*. Springer-Verlag, Berlin, third edition, 2009. A Practical Introduction.
- [33] S. Ulbrich. Optimal control of nonlinear hyperbolic conservation laws with source terms. *Technische Universität München*, 2001.