

Finite element-based invariant-domain preserving approximation of hyperbolic systems: Beyond second-order accuracy in space[☆]

Jean-Luc Guermond^{a,*}, Murtazo Nazarov^b, Bojan Popov^a

^a Department of Mathematics, Texas A&M University 3368 TAMU, College Station, TX 77843, USA

^b Division of Scientific Computing, Department of Information Technology, Uppsala University, Uppsala, Sweden

ARTICLE INFO

Keywords:

Hyperbolic systems
Riemann problem
Invariant domain
High-order method
Limiting
Finite element method

MSC:

65M60
65M12
35L65
35L45

ABSTRACT

This paper proposes an invariant-domain preserving approximation technique for nonlinear conservation systems that is high-order accurate in space and time. The algorithm mixes a high-order finite element method with an invariant-domain preserving low-order method that uses the closest neighbor stencil. The construction of the flux of the low-order method is based on an idea from Abgrall et al. (2017). The mass flux of the low-order and the high-order methods are identical on each finite element cell. This allows for mass preserving and invariant-domain preserving limiting.

1. Introduction

This paper is concerned with the approximation of hyperbolic systems in conservation form using finite elements of degree two and higher. In particular, the paper proposes answers to some questions raised over the years 2014 to 2016 in the PhD theses of [1, §3.3.2.1] and [2, §4.2.2], where it was observed that the low-order invariant-domain preserving method proposed in [3] was not robust with respect to the polynomial degree. Building on ideas developed in [4] (see Propositions 3.1 and 3.2 therein), we propose here a variation of the invariant-domain preserving method from [3] that behaves better as the polynomial degree increases. In particular, the low-order method is based on the closest neighbor stencil.

To avoid distracting details regarding boundary conditions, we consider the Cauchy problem posed over a space domain $D \subset \mathbb{R}^d$ and a time interval $[0, T]$ with $T > 0$:

[☆] This material is based upon work supported in part by the National Science Foundation, USA grants DMS2110868, by the Air Force Office of Scientific Research, USAF, USA, under grant/contract number FA9550-23-1-0007, and by the Army Research Office, USA under grant/contract number W911NF-19-1-0431. The second author is funded by the Swedish Research Council (VR), Sweden under grant number 2021-04620. The third author is supported by the European Union-NextGenerationEU, through the National Recovery and Resilience Plan of the Republic of Bulgaria, project No BG-RRP-2.004-0008.

* Corresponding author.

E-mail address: guermond@tamu.edu (J.-L. Guermond).

<https://doi.org/10.1016/j.cma.2023.116470>

Received 27 April 2023; Received in revised form 26 August 2023; Accepted 19 September 2023

Available online 16 October 2023

0045-7825/© 2023 Elsevier B.V. All rights reserved.

$$\begin{cases} \partial_t \mathbf{u} + \nabla \cdot \mathbb{f}(\mathbf{u}) = \mathbf{0}, & \text{for } (\mathbf{x}, t) \in D \times (0, T), \\ \mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0(\mathbf{x}), & \text{for } \mathbf{x} \in D. \end{cases} \quad (1.1)$$

The dependent variable takes values in \mathbb{R}^m , $m \geq 1$. We assume that (1.1) has an invariant domain $\mathcal{A} \subset \mathbb{R}^m$. This means that if u_0 takes values in \mathcal{A} almost everywhere in D (in the absence of perturbations due to the boundary conditions), then any admissible solution to the Cauchy problem also takes values in \mathcal{A} almost everywhere in $D \times (0, T)$. We assume that $\mathbb{f} : \mathcal{A} \rightarrow \mathbb{R}^{m \times d}$ is Lipschitz. We say that a numerical approximation of (1.1) is invariant-domain preserving if it leaves \mathcal{A} globally invariant.

There is a vast literature on finite volume and discontinuous Galerkin methods describing techniques that are third- and higher-order accurate in space and are invariant-domain preserving (see e.g., [5–7] for the finite volume literature, and [8–10] for the discontinuous Galerkin literature). By comparison, the continuous finite element literature on this topic is sparse. The objective of this paper is to propose some solutions to fill this gap. Some of the arguments presented in this paper find their root in [4] and have some similarities with the residual distribution method developed in [11–13], [14].

Our starting point is the technique described in the following series of papers [3,15], [16–18]. The main idea behind [3,15–17] consists of combining two methods in the spirit of the flux transport corrected methodology of [19,20] (see also [21,22], and the literature therein for other finite element extensions on this idea): A low-order method that is invariant-domain preserving serves as gate keeper to limit a high-order method which may not be invariant-domain preserving but is somewhat entropy consistent. One important property of most of the methods mentioned above is that they do not have theoretical upper limits on the polynomial degree of the space approximation to be invariant-domain preserving. In principle, the method from [3,15–17] can be implemented with any polynomial degree. This is indeed true, but as observed in the PhD theses of [1, §3.3.2.1] and [2, §4.2.2], the low-order invariant-domain preserving method, which is the gate keeper of the technique, is not robust with respect to the polynomial degree; more precisely, the CFL number that is required to maintain the invariant-domain property decreases very fast as the polynomial degree increases, and the method becomes more and more diffusive as the polynomial degree increases. This phenomenon is also reported in [23, §3.3]. It is also shown in Quezada De Luna [2, §4.2.2] and Anderson et al. [23, §3.3] that Bernstein finite elements behave far better than Lagrange elements in this respect. We have worked on this problem since the observations made in Alrashed [1] and Quezada De Luna [2]. It was clear from the start that a hierarchical decomposition of the space approximation had to be done, but the main roadblock on the way that made progresses slow was to have the low-order method and the high-order method to carry exactly the same mass. A solution to this problem for fourth-order finite differences is proposed in [24]. It is shown therein that the high-order fluxes can be recombined in a conservative manner on the low-order stencil. Building on an original idea from [4], we have found a reasonable solution to the conservation problem for finite elements in 2019, and it is the objective of this paper to expose this solution. In addition to have the stencil of the low-order method only depend on the next neighbors, the key idea is to slightly modify the fluxes of the low-order method so that it carries exactly the same mass as the high-order method while still being conservative and consistent (in the spirit of the so-called residual distribution technique by [4]). The proposed technique is robust with respect to the polynomial degree and is exactly mass preserving on patches.

The paper is organized as follows. We describe the low-order and the high-order space approximation setting in Section 2. The method that is high-order accurate in space is described in Section 3. The low-order method is described in Section 4. The limiting operation ensuring that the method combining the high-order approximation and the low-order one is invariant-domain preserving is described in Section 5. The proposed approach involves a node-based limiting and a cell-based limiting. The method is numerically illustrated in Section 6. Technicalities are collected in the Appendices A to C.

2. Space approximation

The goal of this section is to describe the setting for the space approximation. We restrict ourselves to continuous finite elements. We denote by \mathbb{P}_k and \mathbb{Q}_k the (real) vector spaces composed of the d -variate polynomials of degree at most k and of partial degree at most k , respectively.

2.1. Motivation

One important property of the finite-element-based invariant-domain preserving low-order technique introduced in [3,15,16] is that nowhere in the theory there is a theoretical upper limit on the polynomial degree. In principle, the method can be implemented with any polynomial degree. But as observed in the theses of [1, §3.3.2.1] and [2, §4.2.2], the low-order invariant-domain preserving method is not robust with respect to the polynomial degree. The key reason is that the size of the stencil of the method grows with the polynomial degree. This in turn makes the CFL number that is required to maintain the invariant-domain property decrease as the polynomial degree increases, and the method becomes more and more diffusive. This phenomenon is numerically illustrated in Fig. 4 in Section 4.1. The purpose of the paper is to introduce a hierarchical decomposition of the space approximation to address this problem akin to what is done in [4,8,13,25]. We focus in this paper on continuous Lagrange finite elements.

2.2. High-order finite element setting

Using Ciarlet’s notation, we consider a (high-order) reference Lagrange or Bernstein finite element $(\widehat{K}, \widehat{P}^H, \widehat{\Sigma}^H)$. The superscript H is meant to remind us that the vector space \widehat{P}^H is composed of high-order polynomials. The shape functions of the reference element $(\widehat{K},$

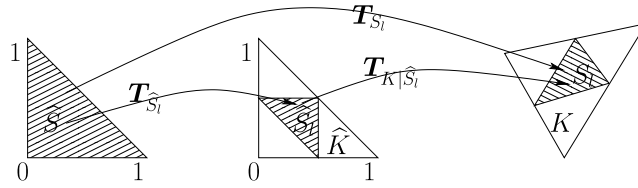


Fig. 1. Let $K \in \mathcal{T}_h$ and $S_l \in \mathcal{T}_K^s$. T_{S_l} maps the reference element \widehat{S} to $\widehat{S}_l \subset \widehat{K}$. $T_{K|_{\widehat{S}_l}}$ (restriction of T_K to \widehat{S}_l) maps \widehat{S}_l to $S_l \subset K$. Then $T_{S_l} := T_{K|_{\widehat{S}_l}} \circ T_{\widehat{S}_l}$.

$\widehat{P}^H, \widehat{\Sigma}^H$) are denoted $\{\widehat{\theta}_n^H\}_{n \in \widehat{\mathcal{J}}}$. When $\{\widehat{\theta}_n^H\}_{n \in \widehat{\mathcal{J}}}$ are Lagrange shape functions, the corresponding reference Lagrange nodes are denoted $\{\widehat{a}_n\}_{n \in \widehat{\mathcal{J}}}$ (here $\widehat{\mathcal{N}}$ is the index set enumerating the reference shape functions). When $\{\widehat{\theta}_n^H\}_{n \in \widehat{\mathcal{J}}}$ are Bernstein shape functions, the domain nodes are denoted $\{\widehat{a}_n\}_{n \in \widehat{\mathcal{J}}}$. Recall that $\dim \widehat{P}^H = \text{card}(\widehat{\mathcal{N}})$.

Let $(\mathcal{T}_h)_{h \in \mathcal{H}}$ be a shape-regular sequence of matching meshes, where \mathcal{H} is a countable set with 0 as unique accumulation point; the index h refers to the typical meshsize of the cells in \mathcal{T}_h . We then introduce the continuous high-order finite element space

$$P^H(\mathcal{T}_h) := \{v \in C^0(D; \mathbb{R}) \mid v|_K \circ T_K \in \widehat{P}^H, \quad \forall K \in \mathcal{T}_h\}, \tag{2.1}$$

where $T_K : \widehat{K} \rightarrow K$ is the bijective geometric transformation that maps the reference element \widehat{K} to the current element K . For the sake of simplicity, we assume in the entire paper that the geometric transformations are affine. The high-order approximation in space of the solution of (1.1) will be done with the space $P^H(\mathcal{T}_h) := (P^H(\mathcal{T}_h))^m$. The global shape functions in $P^H(\mathcal{T}_h)$ are denoted by $\{\varphi_i^H\}_{i \in \mathcal{V}}$. Recall that these functions form a basis of $P^H(\mathcal{T}_h)$, i.e., $\dim(P^H(\mathcal{T}_h)) = \text{card}(\mathcal{V})$. We denote by $j : \mathcal{T}_h \times \widehat{\mathcal{N}} \rightarrow \mathcal{V}$ the connectivity array, which we recall is defined such that

$$\varphi_{iK}^H := \begin{cases} \widehat{\theta}_n^H \circ T_K^{-1} & \text{if there exists } (n, K) \in \widehat{\mathcal{N}} \times \mathcal{T}_h \text{ s.t. } i = j(n, K) \\ 0 & \text{otherwise.} \end{cases} \tag{2.2}$$

For all $i \in \mathcal{V}$ and all $K \in \mathcal{T}_h$, we set

$$\mathcal{S}(i) := \{j \in \mathcal{V} \mid \varphi_j^H \varphi_i^H \not\equiv 0\}, \quad \mathcal{S}(K) := \{i \in \mathcal{V} \mid \varphi_{iK}^H \not\equiv 0\}, \tag{2.3a}$$

$$\mathcal{T}(i) := \{K \in \mathcal{T}_h \mid \varphi_{iK}^H \not\equiv 0\}. \tag{2.3b}$$

We refer to $\mathcal{S}(i)$ as the stencil of the shape function φ_i^H . We introduce the nodes $\{a_i\}_{i \in \mathcal{V}}$ in D such that $a_i = T_K(\widehat{a}_n)$ for all $K \in \mathcal{T}(i)$, where $n \in \widehat{\mathcal{N}}$ is such that $i = j(n, K)$. Recall that $\varphi_j^H(a_i) = \delta_{ij}$.

Finally, for all $i \in \mathcal{V}$ and all $K \in \mathcal{T}(i)$, we define the following quantities which play an important role in the rest of the paper:

$$m_{ij} := \int_D \varphi_i^H(x) \varphi_j^H(x) \, dx, \quad c_{ij}^H := \int_K \varphi_j^H \nabla \varphi_i^H \, dx, \tag{2.4a}$$

$$m_j := \sum_{i \in \mathcal{V}} m_{ij} = \int_D \varphi_j^H(x) \, dx, \quad m_j^K := \int_K \varphi_j^H(x) \, dx. \tag{2.4b}$$

Assumption 2.1. For all $i \in \mathcal{V}$ we have $m_i \geq 0$.

Notice that Assumption 2.1 holds if $m_n^K := \int_K \widehat{\theta}_n^H \, dx \geq 0$ for all $n \in \widehat{\mathcal{N}}$. In particular, this property holds for Lagrange finite elements on triangles up to degree 3 in two dimensions. It holds for all the Lagrange finite elements on quadrangles and hexahedrons. It holds for all Bernstein finite elements.

2.3. Multiscale structure

We now assume for simplicity that the reference element can be subdivided into sub-cells by connecting the reference Lagrange (or domain) points $\{\widehat{a}_n\}_{n \in \widehat{\mathcal{J}}}$ so that each sub-cell contains exactly the same number of Lagrange (or domain) points. Let $\{\widehat{S}_l\}_{l \in \widehat{\mathcal{L}}}$ be the enumerated collection of sub-cells in question, and let us denote $\mathcal{T}_K^s := \{\widehat{S}_l\}_{l \in \widehat{\mathcal{L}}}$ the sub-mesh of \widehat{K} thus formed (here $\widehat{\mathcal{L}}$ is the index set corresponding to the enumeration in question). The superscript s is meant to remind us that we are dealing with subdivided entities (cells or degrees of freedom). For all $l \in \widehat{\mathcal{L}}$, we introduce the index set $\widehat{\mathcal{N}}_l \subsetneq \widehat{\mathcal{N}}$ so that $\{\widehat{a}_n\}_{n \in \widehat{\mathcal{N}}_l}$ is the set of the reference Lagrange (or domain) points that belong to \widehat{S}_l .

To make the above construction more precise, we assume that there exists a reference cell \widehat{S} so that for each sub-cell $\widehat{S}_l \in \mathcal{T}_K^s$ there

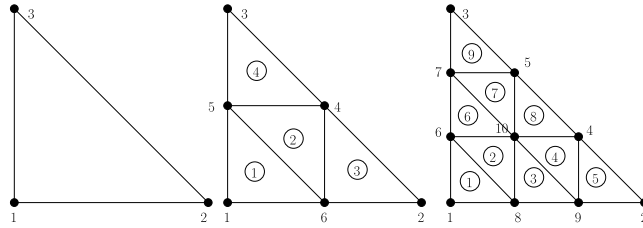


Fig. 2. Degrees of freedom (black dots) and sub-triangle enumeration (circled numbers) for two-dimensional \mathbb{P}_1 , \mathbb{P}_2 , and \mathbb{P}_3 Lagrange elements.

exists an affine geometric transformations $T_{\hat{S}_l}$ so that $\hat{S}_l = T_{\hat{S}_l}(\hat{S})$; see Fig. 1. (It happens in general that $\hat{S} = \hat{K}$, but this is not necessary; for instance, one could subdivide the unit square into triangles. See [14] where mixed subdivisions are considered.) We finally assume that there exists a set of points in \hat{S} , say $\{\hat{z}_{n^s}\}_{n^s \in \hat{\mathcal{N}}^s}$, so that for all $n \in \hat{\mathcal{N}}$ there exist $l \in \hat{\mathcal{L}}$ and $n^s \in \hat{\mathcal{N}}^s$ so that $\hat{a}_n = T_{\hat{S}_l}(\hat{z}_{n^s})$. We formalize this property by introducing the reference hierarchical connectivity array $\hat{j} : \hat{\mathcal{N}}^s \times \hat{\mathcal{L}} \rightarrow \hat{\mathcal{N}}$ such that $T_{\hat{S}_l}(\hat{z}_{n^s}) = \hat{a}_{\hat{j}(n^s, l)}$ for all $n^s \in \hat{\mathcal{N}}^s$ and all $l \in \hat{\mathcal{L}}$. In this paper we adopt the increasing vertex-index enumeration; that is, $(n_1 \leq n_2) \Leftrightarrow (\hat{j}(n_1, l) \leq \hat{j}(n_2, l))$ for all $n_1, n_2 \in \hat{\mathcal{N}}^s$ and all $l \in \hat{\mathcal{L}}$. We finally assume that the above construction satisfies the following property:

Assumption 2.2. There exists a polynomial space \hat{P}^L and a set of linear forms $\hat{\Sigma}^L$ such that $\mathbb{P}_1 \subset \hat{P}^L$ and $(\hat{S}, \hat{P}^L, \hat{\Sigma}^L)$ is a Lagrange finite element based on the Lagrange nodes $\{\hat{z}_{n^s}\}_{n^s \in \hat{\mathcal{N}}^s}$, i.e., $\hat{\sigma}_{n^s}(\hat{p}) = \hat{p}(\hat{z}_{n^s})$, for all $\hat{\sigma}_{n^s} \in \hat{\Sigma}^L$. We denote by $\{\hat{\theta}_{n^s}^L\}_{n^s \in \hat{\mathcal{N}}^s}$ the corresponding reference shape functions.

This assumption means that the cells obtained after subdivision allow for piecewise linear interpolation.

Examples of subdivisions and enumerations for the unit simplex in two dimensions are shown in Fig. 2. For instance, for the triangular \mathbb{P}_3 finite element shown in the rightmost panel of Fig. 2 we have $\hat{j}(1, 1) = 1, \hat{j}(2, 1) = 6, \hat{j}(3, 1) = 8$, and $\hat{j}(1, 3) = 8, \hat{j}(2, 3) = 9, \hat{j}(3, 3) = 10$. The arrays \hat{j} for \mathbb{P}_2 and \mathbb{P}_3 finite elements in two dimensions are given in Table 8. (Notice that we do not impose any restriction on the sign of the determinant of the geometric transformations $T_{\hat{S}_l}, l \in \hat{\mathcal{L}}$.)

2.4. Low-order finite element space

All the mesh cells K in \mathcal{T}_h are subdivided as explained Section 2.3; i.e., all the sub-cells of K are images by T_K of the sub-cells $\{\hat{S}_l\}_{l \in \hat{\mathcal{L}}}$. The collection of the sub-cells of K (see Fig. 3) is denoted

$$\mathcal{T}_K^s := \{T_K(\hat{S}_l)\}_{l \in \hat{\mathcal{L}}}. \tag{2.5}$$

We denote by \mathcal{T}_h^s the mesh obtained by subdividing all the cells. The subdivision process guarantees that the sequence $(\mathcal{T}_h^s)_{h \in \mathcal{H}}$ is shape-regular. Referring to Fig. 1 for an illustration, we introduce the following notation for all $S_l \in \mathcal{T}_K^s$:

$$T_{S_l} := T_{K|\hat{S}_l} \circ T_{\hat{S}_l}. \tag{2.6}$$

Recalling that $(\hat{S}, \hat{P}^L, \hat{\Sigma}^L)$ is a Lagrange finite element (see Assumption 2.2), we define the corresponding low-order Lagrange finite element space

$$P^L(\mathcal{T}_h^s) := \{v \in C^0(D; \mathbb{R}) | v|_{S_l} \circ T_{S_l} \in \hat{P}^L, \quad \forall S_l \in \mathcal{T}_h^s\}. \tag{2.7}$$

We introduce the connectivity array $j^L : \hat{\mathcal{N}}^s \times \mathcal{T}_h^s \rightarrow \mathcal{V}$ defined by setting

$$j^L(n^s, S_l) := j(\hat{j}(n^s, l), K), \tag{2.8}$$

for all $n^s \in \hat{\mathcal{N}}^s$, all $K \in \mathcal{T}_h$, and all $S_l \in \mathcal{T}_K^s$. We denote by $\{\varphi_i^L\}_{i \in \mathcal{V}}$ the global Lagrange shape functions of $P^L(\mathcal{T}_h^s)$. The enumeration is done so that

$$\varphi_{i|S}^L := \begin{cases} \hat{\theta}_{n^s}^L \circ T_S^{-1} & \text{if there exists } (n^s, S) \in \hat{\mathcal{N}}^s \times \mathcal{T}_h^s \text{ s.t. } i = j^L(n^s, S) \\ 0 & \text{otherwise.} \end{cases} \tag{2.9}$$

Lemma 2.3. For all $i, j \in \mathcal{V}$, we have $\varphi_i^L(a_j) = \varphi_i^H(a_j) = \delta_{ij}$.

Proof. This is a consequence of (2.2), (2.8), and (2.9). \square

We are going to make use of the following notation for all $T \in \mathcal{T}_h^s$, all $i \in \mathcal{V}$ and all $K \in \mathcal{T}(i)$:

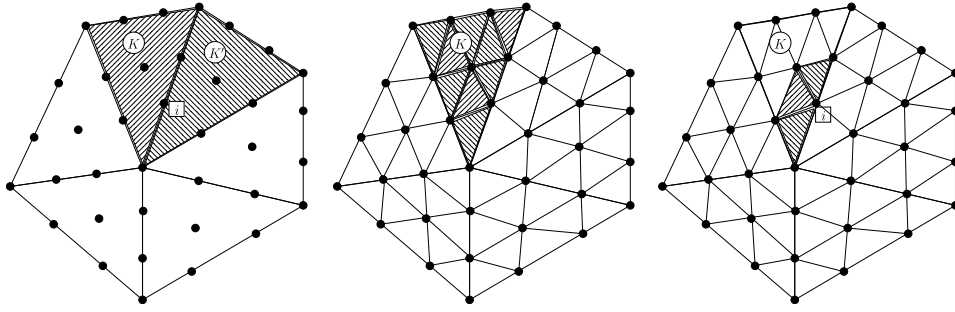


Fig. 3. Left panel: Definition of $\mathcal{S}(i)$ (shaded cells). Center panel: Definition of \mathcal{S}_K^s (shaded sub-cells). Right panel: Definition of $\mathcal{S}_{K,i}^s$ (shaded sub-cells).

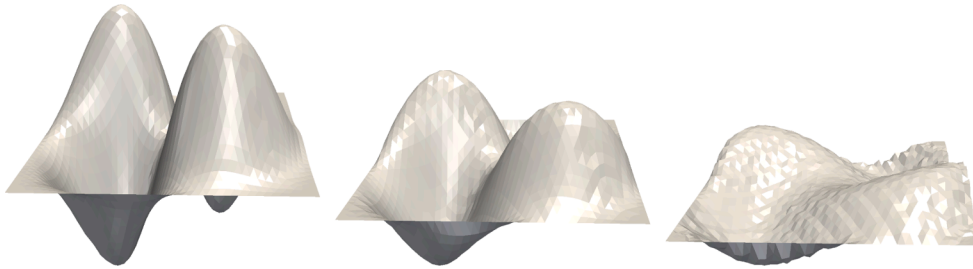


Fig. 4. Loss of robustness of the low-order solution w.r.t. the polynomial degree when using the full high-order stencil. Left: \mathbb{P}_1 , 1927 nodes; center: \mathbb{P}_2 , 1945 nodes; Right: \mathbb{P}_3 , 1888 nodes.

$$\mathcal{S}^L(i) := \{j \in \mathcal{V} \mid \varphi_j^L \varphi_j^L \equiv 0\}, \quad \mathcal{S}^L(S) := \{i \in \mathcal{V} \mid \varphi_{iS}^L \equiv 0\}, \quad (2.10a)$$

$$\mathcal{S}_{K,i}^s := \{S \in \mathcal{S}_K^s \mid \varphi_{iS}^L \equiv 0\}. \quad (2.10b)$$

Notice that $\mathcal{S}^L(S) = j^L(\widehat{\mathcal{N}}^s, S)$ and $\mathcal{S}_{K,i}^s = \{S \in \mathcal{S}_K^s \mid \exists n^s \in \widehat{\mathcal{N}}^s, j^L(n^s, S) = i\}$. The notation introduced above is illustrated in Fig. 3.

For all $K \in \mathcal{T}_h$ and $S \in \mathcal{S}_K^s$, we set $P_K^H := \{\widehat{p} \circ T_K^{-1} \mid \widehat{p} \in \widehat{P}^H\}$ and $P_S^L := \{\widehat{p} \circ T_S^{-1} \mid \widehat{p} \in \widehat{P}^L\}$. We then define the local low-order Lagrange interpolation operator $\Pi_S^L : P_K^H \rightarrow P_S^L$ as follows:

$$\Pi_S^L(u_{h|K}) = \sum_{n^s \in \widehat{\mathcal{N}}^s} U_{j^L(n^s, S)} \varphi_{j^L(n^s, S)}^L, \quad \forall u_h := \sum_{i \in \mathcal{V}} U_i \varphi_i^H \in P^H(\mathcal{T}_h). \quad (2.11)$$

Then we define the corresponding global low-order Lagrange interpolation operator $\Pi_h^L : P^H(\mathcal{T}_h) \rightarrow P^L(\mathcal{T}_h^s)$ by $\Pi_h^L(u_h)|_S := \Pi_S^L(u_{h|K})$ for all $K \in \mathcal{T}_h$ and all $S \in \mathcal{S}_K^s$. Notice that (2.9) implies that

$$\Pi_h^L(u_h) = \sum_{i \in \mathcal{V}} U_i \varphi_i^L, \quad \forall u_h := \sum_{i \in \mathcal{V}} U_i \varphi_i^H \in P^H(\mathcal{T}_h). \quad (2.12)$$

3. High-order method

For completeness, we introduce in this section the high-order approximation of (1.1). No originality is claimed here for we essentially paraphrase [3, §3.2]. The low-order method is introduced in Section 4.

3.1. High-order update

To properly describe the forward Euler time stepping technique, we denote by t^n the current time level, $n \in \mathbb{N}$, and let τ^n be the current time step; i.e., $t^{n+1} = t^n + \tau^n$. The time step may vary at each time level, but to simplify the (already heavy) notation we are going to drop the super-index n and use the symbol τ for the time step. We denote by $u_h^n := \sum_{i \in \mathcal{V}} U_i^n \varphi_i^H \in P^H(\mathcal{T}_h)$ the high-order approximation of (1.1) at t^n , and by induction, we assume $U_i^n \in \mathcal{A}$ for all $i \in \mathcal{V}$. We now briefly describe a way to create a high-order update $u_h^{H,n+1} := \sum_{i \in \mathcal{V}} U_i^{H,n+1} \varphi_i^H \in P^H(\mathcal{T}_h)$ at the time level t^{n+1} using the forward Euler method.

Our first task is to construct a high-order approximation of $\nabla \cdot \mathbb{f}(u)$. When using finite elements it is natural to consider the Galerkin

approximation thereof; that is, for all $i \in \mathcal{V}$ one computes $\int_D \nabla \cdot (\mathbb{f}(\mathbf{u}_h^n)) \varphi_i^H dx$ using an appropriate quadrature. The problem with this expression is that although the values $\{\mathbf{U}_i^n\}_{i \in \mathcal{V}}$ are in \mathcal{A} (by the induction assumption), there is no guarantee that $\mathbf{u}_h^n(x)$ is in \mathcal{A} for all x in D (or at least for all the quadrature points of the quadrature approximating the integral. Recall that it is imperative that $\mathbf{u}_h^n(x) \in \mathcal{A}$ for $\mathbb{f}(\mathbf{u}_h^n(x))$ to make sense.) Notice that $\mathbf{u}_h^n(x) \in \mathcal{A}$ for all $x \in D$ for $\mathbb{P}_1, \mathbb{Q}_1$, and Bernstein finite elements, but this is not the case when using Lagrange elements of degree two and higher. We solve this difficulty as in [3, Eq. (3.7)] by using the Lagrange interpolant of the flux. Let $\Pi_h^H : C^0(D; \mathbb{R}) \rightarrow \mathbf{P}^H(\mathcal{T}_h)$ be the Lagrange interpolation operator. With an obvious abuse of notation, we have $\Pi_h^H(\mathbb{f}(\mathbf{u}_h^n)) := \sum_{i \in \mathcal{V}} \mathbb{f}(\mathbf{u}_h^n(\mathbf{a}_i)) \varphi_i^H$. For Lagrange elements we have $\mathbf{u}_h^n(\mathbf{a}_i) = \mathbf{U}_i^n$, and in this case

$$\Pi_h^H(\mathbb{f}(\mathbf{u}_h^n)) := \sum_{i \in \mathcal{V}} \mathbb{f}(\mathbf{U}_i^n) \varphi_i^H. \tag{3.1}$$

For Bernstein elements we have $\mathbf{u}_h^n(\mathbf{a}_i) = \mathbf{U}_i^n + \mathcal{O}(h^2)$. In this case (3.1) is no longer an identity by a second order approximation. The method has to be modified for higher-order Bernstein polynomials by invoking a change of basis. We omit the details to keep the presentation simple. The simulations reported in the paper are done only for \mathbb{P}_2 Bernstein elements. The term $\int_D \nabla \cdot (\mathbb{f}(\mathbf{u})) \varphi_i^H dx$ is then approximated by $\int_D \nabla \cdot (\Pi_h^H(\mathbb{f}(\mathbf{u}_h^n))) \varphi_i^H dx$. Recalling the definition of c_{ij}^H in (2.4a), the approximation takes the following form: $\sum_{j \in \mathcal{V}} \mathbb{f}(\mathbf{U}_j^n) c_{ij}^H$ (recall that $\mathbb{f}(\mathbf{U}_j^n)$ is a $m \times d$ matrix and c_{ij}^H is $d \times 1$ column vector).

Recalling that $\mathbf{u}_h^{H,n+1} = \sum_{i \in \mathcal{V}} \mathbf{U}_i^{H,n+1} \varphi_i^H$ and using the forward Euler technique for the time approximation, we define the high-order approximation of (1.1) by

$$\sum_{j \in \mathcal{V}} m_{ij} \frac{\mathbf{U}_j^{H,n+1} - \mathbf{U}_j^n}{\tau} = - \sum_{j \in \mathcal{J}(i)} \mathbb{f}(\mathbf{U}_j^n) c_{ij}^H + \sum_{j \in \mathcal{J}^L(i)} d_{ij}^{H,n} (\mathbf{U}_j^n - \mathbf{U}_i^n). \tag{3.2}$$

The term $d_{ij}^{H,n}$ is some high-order graph viscosity that can be defined in many ways. It could be based on a smoothness indicator like in [26, Eq. (12)], [27, p. 6], [28, Thm. 4.1], [29, §4.3], [18, §6.2], or it could be based on an entropy viscosity commutator like in [17, §6.4]. The exact definition of $d_{ij}^{H,n}$ does not really matter at the moment provided it induces a high-order perturbation of (1.1) and satisfies $d_{ij}^{H,n} = d_{ji}^{H,n} \geq 0$. That is to say, we assume that the term $\sum_{j \in \mathcal{J}^L(i)} d_{ij}^{H,n} (\mathbf{U}_j^n - \mathbf{U}_i^n)$ is of the same order as the consistency error of the approximation of the flux when the solution is smooth. Notice that in (3.2) we insist on the connectivity of $d_{ij}^{H,n}$ to be that of the low-order approximation, i.e., the summation is done over the low-order stencil $\mathcal{J}^L(i)$ instead of the high-order stencil $\mathcal{J}(i)$. The way $d_{ij}^{H,n}$ is computed in the numerical simulations reported at the end of the paper is explained in Section 6.1. The following result clarifies the conservation properties of the high-order method.

Lemma 3.1. *The scheme (3.2) has the following conservation property:*

$$\int_D \mathbf{u}_h^{H,n+1} dx = \int_D \mathbf{u}_h^n dx - \tau \int_D \nabla \cdot (\Pi_h^H(\mathbb{f}(\mathbf{u}_h^n))) dx. \tag{3.3}$$

Proof. Notice that $\sum_{i,j \in \mathcal{V}} m_{ij} \mathbf{U}_j^{H,n+1} = \int_D \sum_{j \in \mathcal{V}} \mathbf{U}_j^{H,n+1} \varphi_j^H dx = \int_D \mathbf{u}_h^{H,n+1} dx$ and $\sum_{i,j \in \mathcal{V}} m_{ij} \mathbf{U}_j^n = \int_D \mathbf{u}_h^n dx$. Moreover using $\sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{J}^L(i)} = \sum_{j \in \mathcal{V}} \sum_{i \in \mathcal{J}^L(j)}$ and $d_{ij}^{H,n} = d_{ji}^{H,n}$, we have

$$\sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{J}^L(i)} d_{ij}^{H,n} (\mathbf{U}_j^n - \mathbf{U}_i^n) = \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{J}^L(i)} d_{ij}^{H,n} \mathbf{U}_j^n - \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{J}^L(i)} d_{ji}^{H,n} \mathbf{U}_i^n = 0.$$

Summing (3.2) over $i \in \mathcal{V}$ gives the assertion. \square

4. Low-order method

We introduce in this section the low-order method that will be used as reference for limiting the high-order method. The low-order solution is approximated in space using the Lagrange finite element space $\mathbf{P}^L(\mathcal{T}_h^s)$.

4.1. Motivation

Since the technique described in [3, Eq (3.9)] is a priori independent of the polynomial degree of the approximation space $\mathbf{P}^H(\mathcal{T}_h)$, we can in principle define the low-order update in $\mathbf{P}^H(\mathcal{T}_h)$ as follows:

$$m_i \frac{\mathbf{U}_i^{H,n+1} - \mathbf{U}_i^n}{\tau} = - \sum_{j \in \mathcal{J}(i)} \mathbb{f}(\mathbf{U}_j^n) c_{ij}^H + \sum_{j \in \mathcal{J}(i)} d_{ij}^{L,n} (\mathbf{U}_j^n - \mathbf{U}_i^n), \tag{4.1}$$

where the stencil for the low-order viscosity $d_{ij}^{L,n}$ in (4.1) is $\mathcal{J}(i)$, i.e., the summation is done over $j \in \mathcal{J}(i)$ instead of $j \in \mathcal{J}^L(i)$. A definition of $d_{ij}^{L,n}$ that makes the method invariant-domain preserving is

$$d_{ij}^{L^n} := \max\left(\lambda_{\max}(\mathbf{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_j^n) \left\| \mathbf{c}_{ij}^H \right\|_{\ell^2}, \lambda_{\max}(\mathbf{n}_{ji}, \mathbf{U}_j^n, \mathbf{U}_i^n) \left\| \mathbf{c}_{ji}^H \right\|_{\ell^2}\right), \tag{4.2}$$

where $\lambda_{\max}(\mathbf{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_j^n)$ is any upper bound on the maximum wave speed in the Riemann problem $\partial_t v + \partial_x(\mathbb{f}(v)\mathbf{n}_{ij}) = 0$, where $\mathbf{n}_{ij} := \mathbf{c}_{ij}^H / \left\| \mathbf{c}_{ij}^H \right\|$, with left and right data \mathbf{U}_i^n and \mathbf{U}_j^n respectively (see e.g., [3, Eq (3.16)]). Using ideas, now classical, from [30, p. 163], [31], [32, p. 375], and [33, §5], it is established in [3, Th. 4.1] and [18, Thm. 3.6] that using (4.2) yields a low-order method that is invariant-domain preserving (see Theorem 4.4 for a more precise statement).

As observed in the PhD theses of [1, §3.3.2.1] and [2, §4.2.2], the above scheme is not robust with the polynomial degree. This phenomenon is also reported [25, Tab. 1] and [8, Fig. 1]. To illustrate this observation we show in Fig. 4 the results of three two-dimensional simulations using the low-order method (4.1) with \mathbb{P}_1 , \mathbb{P}_2 and \mathbb{P}_3 Lagrange finite elements on nonuniform triangular meshes. The problem solved is the linear transport equation $\partial_t u + \beta \cdot \nabla u = 0$ with the divergence-free velocity $\beta(x, t) = \cos(\pi t)(- \sin(2\pi x_2)\sin^2(\pi x_1)e_1 + \sin(2\pi x_1)\sin^2(\pi x_2)e_2)$. The initial data is $u_0(x) = 2^8(x_1(1-x_1)x_2(1-x_2))^2 \sin(2\pi x_1)\sin(2\pi x_2)$. The solution is periodic in time with period 1. Fig. 4 shows that graph of the low-order approximation at $t = 1$. The total number of degrees of freedom is approximately the same for the three cases to make the comparison fair (1927 for \mathbb{P}_1 , 1945 for \mathbb{P}_2 , and 1888 for \mathbb{P}_3). We observe that the quality of the approximation deteriorates as the polynomial degree of the approximation increases. (One can verify though that the asymptotic convergence rates are identical, see Section 6.2.) The reason for this behavior is that the cardinality of the high-order stencil $\mathcal{S}(i)$ used in (4.1) increases with the polynomial degree; actually, $\text{card}(\mathcal{S}(i)) \sim k^d$ where k is the polynomial degree and d is the space dimension.

The purpose of the rest of this section is to introduce a low-order method that relies only on the next neighbors to be invariant-domain preserving; that is, we are going to construct a low-order flux based on the stencil $\mathcal{S}^L(i)$ instead of $\mathcal{S}(i)$. The main difficulty in this exercise is to make sure that the low-order solution somehow carries the same mass as the high-order one.

4.2. Reducing the stencil to next neighbors

We are going to use an idea introduced in [4, Prop. 3.1]. Recall the low-order Lagrange interpolation operator $\Pi_S^L : P_K^H \rightarrow P_S^L$ introduced in (2.11). Let us define $\widehat{P}^H := (\widehat{P}^H)^d$. We assume in the rest of the paper that the finite element setting described in Section 2 satisfies the following assumption.

Assumption 4.1. (i) For all $K \in \mathcal{T}_h$, there exists a collection of real numbers $\{\beta^S\}_{S \in \mathcal{S}_K^*}$ such that the following holds true for every $p \in P_K^H := \{\widehat{p} \circ T_K^{-1} | \widehat{p} \in \widehat{P}^H\}$:

$$\int_K \nabla \cdot p \, dx = \sum_{S \in \mathcal{S}_K^*} \beta^S \int_S \nabla \cdot (\Pi_S^L(p)) \, dx. \tag{4.3}$$

(ii) For all $K \in \mathcal{T}_h$ and all $S \in \mathcal{S}_K^s$, there exists a set of real numbers $\{\alpha_j^S\}_{j \in \mathcal{J}^L(S)}$ such that the following identities holds:

$$\forall j \in \mathcal{J}(K), \quad \sum_{S \in \mathcal{S}_{Kj}^s} \alpha_j^S = 1, \quad \forall S \in \mathcal{S}_K^s, \quad \sum_{j \in \mathcal{J}^L(S)} \alpha_j^S m_j^K = \beta^S |S|. \tag{4.4}$$

The existence of the coefficients $\{\beta^S\}_{S \in \mathcal{S}_K^*}$ so that (4.3) holds true in dimensions 2 and 3 is established in [4, Prop. 3.1] and [4, Prop. 3.2]. The existence of the coefficients $\{\alpha_j^S\}_{j \in \mathcal{J}^L(S)}$ so that identities (4.4) holds true is established in the Appendices B, C for various finite elements. It is also shown in Lemmas A.1, A.2 that it suffices that (4.3) and (4.4) hold on the reference element for these properties to hold for every mesh cell K in \mathcal{T}_h if the mesh is affine.

We now want to approximate $\int_D \nabla \cdot (\mathbb{f}(u_h)) \varphi_i^H \, dx$ with the restriction that the approximation in question involves only the next neighbors of the Lagrange point a_i , i.e., we want to involve only the sub-cells in $\bigcup_{K \in \mathcal{T}(i)} \mathcal{S}_{K,i}^s$. Recalling that Π_h^L is the low-order Lagrange interpolation operator introduced in Section 3.1, we now consider $\Pi_h^L(\mathbb{f}(u_h)) := \sum_{i \in \mathcal{I}} \mathbb{f}(u_h)(a_i) \varphi_i^L$. (Notice in passing that $\Pi_h^L(\mathbb{f}(u_h)) = \sum_{i \in \mathcal{I}} \mathbb{f}(u_h)(a_i) \varphi_i^L$ because $\mathbb{f}(u_h)(a_i) = \mathbb{f}(u_h)(a_i) = \mathbb{f}(u_h)(a_i)$.) Since $\Pi_h^L(\mathbb{f}(u_h))$ is a second-order approximation of $\mathbb{f}(u_h)$, $\nabla \cdot (\Pi_h^L(\mathbb{f}(u_h)))$ is a first-order approximation of $\nabla \cdot (\mathbb{f}(u_h))$. Since for every cell K in $\mathcal{T}(i)$ and every sub-cell S in $\mathcal{S}_{K,i}^s$, the quantity $\frac{1}{|S|} \int_S \nabla \cdot (\mathbb{f}(u_h)) \, dx$ is an acceptable first-order approximation of $\nabla \cdot (\mathbb{f}(u_h))(a_i)$, it is also the case of $\frac{1}{|S|} \int_S \nabla \cdot (\Pi_h^L(\mathbb{f}(u_h))) \, dx$. Then recalling that $\sum_{S \in \mathcal{S}_{K,i}^s} \alpha_i^S = 1$, we also have

$$\nabla \cdot (\mathbb{f}(u_h))(a_i) = \sum_{S \in \mathcal{S}_{K,i}^s} \alpha_i^S \frac{1}{|S|} \int_S \nabla \cdot (\Pi_h^L(\mathbb{f}(u_h))) \, dx + \mathcal{O}(h). \tag{4.5}$$

Since the object of interest is $\int_D \nabla \cdot (\mathbb{f}(u)) \varphi_i^H \, dx$, and up to a quadrature based on the Lagrange points we have $\int_D \nabla \cdot (\mathbb{f}(u)) \varphi_i^H \, dx \approx \left(\sum_{K \in \mathcal{T}(i)} m_K^K \right) \nabla \cdot (\mathbb{f}(u))(a_i)$, we finally infer that

$$\int_D \nabla \cdot (\mathbb{f}(\mathbf{u}_h)) \varphi_i^H \, dx = \sum_{K \in \mathcal{T}(i)} m_i^K \sum_{S \in \mathcal{F}_{K,i}^s} \alpha_i^S \frac{1}{|S|} \int_S \nabla \cdot (\Pi_h^L(\mathbb{f}(\mathbf{u}_h))) \, dx + \mathcal{O}(h). \tag{4.6}$$

Proposition 4.2. For every $i \in \mathcal{V}$, the approximations (4.5) and (4.6) are exact if \mathbb{f} is linear and \mathbf{u}_h is linear over the patch composed of the cells in $\mathcal{T}(i)$.

Proof. Let $i \in \mathcal{V}$ and $S \in \mathcal{F}_{K,i}^s$. Then $\Pi_h^L(\mathbb{f}(\mathbf{u}_h))|_S = \sum_{j \in \mathcal{J}^L(S)} \mathbb{f}(\mathbf{u}_h(\mathbf{a}_j)) \varphi_j^L|_S = \mathbb{f}(\sum_{j \in \mathcal{J}^L(S)} \mathbf{u}_h(\mathbf{a}_j) \varphi_j^L|_S)$ by linearity of \mathbb{f} . Since \mathbf{u}_h is linear over K , \mathbf{u}_h is also linear over SK ; hence, $\sum_{j \in \mathcal{J}^L(S)} \mathbf{u}_h(\mathbf{a}_j) \varphi_j^H|_S = \mathbf{u}_h|_S$. This means that $\Pi_h^L(\mathbb{f}(\mathbf{u}_h))|_S = \mathbb{f}(\mathbf{u}_h)|_S$. Then $\nabla \cdot (\Pi_h^L(\mathbb{f}(\mathbf{u}_h)))|_S = \nabla \cdot (\mathbb{f}(\mathbf{u}_h))|_S$ and $\nabla \cdot (\Pi_h^L(\mathbb{f}(\mathbf{u}_h)))|_K = \nabla \cdot (\mathbb{f}(\mathbf{u}_h))|_K$ because \mathbb{f} is linear and $\mathbf{u}_h|_K$ is linear. Hence,

$$\begin{aligned} \sum_{K \in \mathcal{T}(i)} m_i^K \sum_{S \in \mathcal{F}_{K,i}^s} \alpha_i^S \frac{1}{|S|} \int_S \nabla \cdot (\Pi_h^L(\mathbb{f}(\mathbf{u}_h))) \, dx &= \sum_{K \in \mathcal{T}(i)} m_i^K (\nabla \cdot (\mathbb{f}(\mathbf{u}_h)))|_K \sum_{S \in \mathcal{F}_{K,i}^s} \alpha_i^S \\ &= \sum_{K \in \mathcal{T}(i)} (\nabla \cdot (\mathbb{f}(\mathbf{u}_h)))|_K \int_K \varphi_i^L \, dx = \int_D \nabla \cdot (\mathbb{f}(\mathbf{u}_h)) \varphi_i^L \, dx \end{aligned}$$

This concludes the proof. \square **Remark 4.3.** (Second-order Accuracy) Proposition 4.2 shows that the approximation (4.6) is actually formally second-order accurate in space since it is exact when \mathbf{u}_h is piecewise linear (and continuous over D) and \mathbb{f} is linear.

4.3. Low-order update

Recalling that $\{\varphi_j^L\}_{j \in \mathcal{J}}$ are the low-order Lagrange shape functions, we introduce the following quantity c_{ij}^L to simplify the notation:

$$c_{ij}^L := \sum_{K \in \mathcal{T}(i)} m_i^K \sum_{S \in \mathcal{F}_{K,i}^s} \alpha_i^S \frac{1}{|S|} \int_S \nabla \varphi_j^L \, dx. \tag{4.7}$$

This definition implies that

$$\sum_{K \in \mathcal{T}(i)} m_i^K \sum_{S \in \mathcal{F}_{K,i}^s} \alpha_i^S \frac{1}{|S|} \int_S \nabla \cdot (\Pi_h^L(\mathbb{f}(\mathbf{u}_h))) \, dx = \sum_{j \in \mathcal{J}^L(i)} \mathbb{f}(\mathbf{U}_j^n) c_{ij}^L. \tag{4.8}$$

The low-order update is constructed in the high-order space $\mathbf{P}_h^H(\mathcal{T}_h)$, i.e., we set $\mathbf{u}_h^{L,n+1} = \sum_{i \in \mathcal{V}} \mathbf{U}_i^{L,n+1} \varphi_i^H$. Using the heuristics (4.6), we compute the low-order coefficients $(\mathbf{U}_i)_{i \in \mathcal{V}}$ as follows:

$$m_i \frac{\mathbf{U}_i^{L,n+1} - \mathbf{U}_i^n}{\tau} = - \sum_{j \in \mathcal{J}^L(i)} \mathbb{f}(\mathbf{U}_j^n) c_{ij}^L + \sum_{j \in \mathcal{J}^L(i)} d_{ij}^{L,n} (\mathbf{U}_j^n - \mathbf{U}_i^n). \tag{4.9}$$

This expression has the same form as in [3, Eq (3.9)]. Here $d_{ij}^{L,n}$ is the graph viscosity coefficient; the purpose of the graph viscosity is to make the method invariant-domain preserving. Let $\lambda_{\max}(\mathbf{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_j^n)$ be any upper bound on the maximum wave speed in the Riemann problem $\partial_t v + \partial_x(\mathbb{f}(v)\mathbf{n}_{ij}) = 0$ with left and right data \mathbf{U}_i^n and \mathbf{U}_j^n respectively, where $\mathbf{n}_{ij} := c_{ij}^L / \|c_{ij}^L\|$.

Theorem 4.4. (Local invariance) Let $n \geq 0$ and let $i \in \mathcal{V}$. Assume that

$$d_{ij}^{L,n} := \max(\lambda_{\max}(\mathbf{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_j^n) \|c_{ij}^L\|_{\ell^2}, \lambda_{\max}(\mathbf{n}_{ji}, \mathbf{U}_j^n, \mathbf{U}_i^n) \|c_{ji}^L\|_{\ell^2}). \tag{4.10}$$

Let \mathcal{B} be any convex subset of \mathcal{A} that is invariant for (1.1). Assume that τ satisfies the CFL condition $1 \geq 2\tau \sum_{j \in \mathcal{J}^L(i) \setminus \{i\}} \frac{d_{ij}^{L,n}}{m_i}$ and $\mathbf{U}_j^n \subset \mathcal{B}$ for all $j \in \mathcal{J}^L(i)$. Then $\mathbf{U}_i^{n+1} \in \mathcal{B}$. **Proof.** See [3, Th., 4.1] and [18, Thm. 3.6]. \square

4.4. Conservation

One important property of the approximation (4.6) is that it mimics the identity $\sum_{i \in \mathcal{J}(K)} \int_K \nabla \cdot (\mathbb{f}(\mathbf{u})) \varphi_i^H \, dx = \int_K \nabla \cdot (\mathbb{f}(\mathbf{u})) \, dx$, which we recall is a consequence of the partition of unity. More precisely, we have the following result.

Lemma 4.5. (Local conservation) Let Assumption 4.1 be met. The following holds true for all $\mathbf{u}_h \in \mathbf{P}^H(\mathcal{T}_h)$:

$$\sum_{i \in \mathcal{J}(K)} m_i^K \sum_{S \in \mathcal{F}_{K,i}^s} \alpha_i^S \frac{1}{|S|} \int_S \nabla \cdot (\Pi_h^L(\mathbb{f}(\mathbf{u}_h))) \, dx = \int_K \nabla \cdot (\Pi_h^H \mathbb{f}(\mathbf{u}_h)) \, dx. \tag{4.11}$$

Proof. By rearranging the summations and using (4.3) and (4.4), we obtain

$$\begin{aligned} & \sum_{i \in \mathcal{I}(K)} m_i^K \sum_{S \in \mathcal{T}_{K,i}^s} \alpha_i^S \frac{1}{|S|} \int_S \nabla \cdot (\Pi_h^L(\mathbb{f}(\mathbf{u}_h))) \, dx \\ &= \sum_{S \in \mathcal{T}_K^s} \frac{1}{|S|} \int_S \nabla \cdot (\Pi_h^L(\mathbb{f}(\mathbf{u}_h))) \, dx \sum_{i \in \mathcal{I}(S)} m_i^K \alpha_i^S \\ &= \sum_{S \in \mathcal{T}_K^s} \beta^S \int_S \nabla \cdot (\Pi_h^L(\mathbb{f}(\mathbf{u}_h))) \, dx = \int_K \nabla \cdot (\Pi_h^H(\mathbb{f}(\mathbf{u}_h))) \, dx. \quad \square \end{aligned}$$

Lemma 4.5 shows that the low-order flux and the high-order flux produce the same change of mass over each mesh cell in \mathcal{T}_h , just like for the residual distribution scheme described in [4] (see comment after Thm. 2.1 therein).

Contrary to the method proposed in [3], conservation in (4.9) does not arise from skew-symmetry properties. The graph viscosity still has the skew-symmetry property: $d_{ij}^{L,n}(\mathbf{U}_j^n - \mathbf{U}_i^n) = -d_{ji}^{L,n}(\mathbf{U}_i^n - \mathbf{U}_j^n)$ because $d_{ij}^{L,n} = d_{ji}^{L,n}$, but this is no longer the case for the flux terms $(\mathbb{f}(\mathbf{U}_j^n) + \mathbb{f}(\mathbf{U}_i^n))c_{ij}^L$ because $c_{ij}^L \neq c_{ji}^L$. The key reason for the loss of skew-symmetry is that the flux in (4.9) is approximated by using the low-order Lagrange basis functions, $\{\varphi_j^L\}_{j \in \mathcal{V}}$, (i.e., $\Pi_h^L(\mathbb{f}(\mathbf{u}_h)) = \sum_{i \in \mathcal{V}} \mathbb{f}(\mathbf{U}_i) \varphi_i^L$).

Theorem 4.6. (Conservation) Let Assumption 4.1 be met. The following conservation property holds for the low-order scheme (4.9):

$$\int_D \mathbf{u}_h^{L,n+1} \, dx = \int_D \mathbf{u}_h^n \, dx - \tau \int_D \nabla \cdot (\Pi_h^H(\mathbb{f}(\mathbf{u}_h^n))) \, dx. \tag{4.12}$$

Proof. Summing (4.9) over $i \in \mathcal{V}$ and using that $d_{ij}^{L,n} = d_{ji}^{L,n}$, we obtain

$$\sum_{i \in \mathcal{V}} \frac{m_i}{\tau} \mathbf{u}_i^{L,n+1} = \sum_{i \in \mathcal{V}} \frac{m_i}{\tau} \mathbf{U}_i^n - \sum_{K \in \mathcal{T}_h} m_i^K \sum_{S \in \mathcal{T}_{K,i}^s} \alpha_i^S \frac{1}{|S|} \int_S \nabla \cdot (\Pi_h^L(\mathbb{f}(\mathbf{u}_h))) \, dx + \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{V} \setminus \{i\}} d_{ij}^{L,n} \mathbf{U}_j^n - \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{V} \setminus \{i\}} d_{ji}^{L,n} \mathbf{U}_i^n.$$

Then invoking (4.11) from Lemma 4.5 (since Assumption 4.1 holds), we obtain

$$\sum_{i \in \mathcal{V}} \frac{m_i}{\tau} \mathbf{u}_i^{L,n+1} = \sum_{i \in \mathcal{V}} \frac{m_i}{\tau} \mathbf{U}_i^n - \sum_{K \in \mathcal{T}_h} \int_K \nabla \cdot (\Pi_h^H(\mathbb{f}(\mathbf{u}_h))) \, dx.$$

After observing that $\sum_{i \in \mathcal{V}} m_i \mathbf{U}_i^{L,n+1} = \sum_{i \in \mathcal{V}} m_i \int_D \varphi_i^H \mathbf{U}_i^{L,n+1} \, dx = \int_D \mathbf{u}_h^{n+1} \, dx$ and $\sum_{i \in \mathcal{V}} m_i \mathbf{U}_i^n = \int_D \mathbf{u}_h^n \, dx$, the assertion follows readily.

□ **Remark 4.7. (Local Conservation)** A local form of conservation making (4.12) more precise can be established by using the residual distribution technique discussed in [14,34].

5. Limiting

We describe in this section a limiting technique that guarantees that, once limited, the high-order update is invariant-domain preserving.

5.1. Maintaining conservation

Let $\mathcal{B} \subset \mathcal{V}$ be any convex set in the phase space. We recall that a function $\Psi : \mathcal{B} \rightarrow \mathbb{R}$ is said to be quasiconcave if for all finite sets $\{\theta_j\}_{j \in \mathcal{J}}, \{\mathbf{U}_j\}_{j \in \mathcal{J}}$ with $\theta_j \in [0, 1], \sum_{j \in \mathcal{J}} \theta_j = 1$ and $\mathbf{U}_j \in \mathcal{B}$ for all $j \in \mathcal{J}$, the following holds: $\Psi(\sum_{j \in \mathcal{J}} \theta_j \mathbf{U}_j) \geq \min_{j \in \mathcal{J}} \Psi(\mathbf{U}_j)$. Let $i \in \mathcal{V}$ and let $\Psi_i : \mathcal{B} \rightarrow \mathbb{R}$ be a quasiconcave continuous function that is such that $\Psi_i(\mathbf{U}_i^{L,n+1}) \geq 0$. Ways to construct quasiconcave functions satisfying this property are explained in [17, § 4] and [18, § 7.2] (including bound relaxation). Our goal is to correct the high-order update $\mathbf{U}_i^{H,n+1}$ by applying a limiting technique so that once limited the update \mathbf{U}_i^{n+1} satisfies $\Psi_i(\mathbf{U}_i^{n+1}) \geq 0$ as well. We explain in this section how this can be done while maintaining conservation.

In the spirit of the flux transport corrected literature, we proceed as in [20, §III] (see also [21,22], and [17,18]) by subtracting the low-order update (4.9) from the high-order update (3.2). We obtain

$$m_i \mathbf{U}_i^{H,n+1} = m_i \mathbf{U}_i^{L,n+1} - \sum_{j \in \mathcal{I}(i)} \tau \mathbb{f}(\mathbf{U}_j^n) c_{ij}^H + \sum_{j \in \mathcal{L}(i)} \tau \mathbb{f}(\mathbf{U}_j^n) c_{ij}^L + \sum_{j \in \mathcal{L}(i)} \mathbf{A}_{ij}^n + \sum_{j \in \mathcal{I}(i)} \mathbf{C}_{ij}^n, \tag{5.1a}$$

$$\mathbf{A}_{ij}^n := \tau (d_{ij}^{H,n} - d_{ij}^{L,n}) (\mathbf{U}_j^n - \mathbf{U}_i^n), \tag{5.1b}$$

$$\mathbf{C}_{ij}^n := (m_i \delta_{ij} - m_{ij}) (\mathbf{U}_j^{H,n+1} - \mathbf{U}_i^{H,n+1} - (\mathbf{U}_j^n - \mathbf{U}_i^n)). \tag{5.1c}$$

We observe that $\mathbf{A}_{ij}^n = -\mathbf{A}_{ji}^n$ and $\mathbf{C}_{ij}^n = -\mathbf{C}_{ji}^n$ which are important property to maintain mass conservation. From now on we extend the definition of \mathbf{A}_{ij}^n to all $j \in \mathcal{I}(i)$ by setting $\mathbf{A}_{ij}^n = \mathbf{0}$ if $j \in \mathcal{L}(i)$, and we define

$$\mathbf{D}_{ij}^n := \mathbf{A}_{ij}^n + \mathbf{C}_{ij}^n, \quad \forall j \in \mathcal{I}(i). \tag{5.2}$$

The key difference between the present setting and that in [17,18] is that the high-order flux $\sum_{j \in \mathcal{I}(i)} \mathbb{f}(\mathbf{U}_j^n) \mathbf{c}_{ij}^H$ and the low-order flux $\sum_{j \in \mathcal{I}^L(i)} \mathbb{f}(\mathbf{U}_j^n) \mathbf{c}_{ij}^L$ do not cancel each other (unless $P^H(\mathcal{T}_h) = P^L(\mathcal{T}_h^S)$). Actually, we have

$$-\tau \sum_{j \in \mathcal{I}(i)} \mathbb{f}(\mathbf{U}_j^n) \mathbf{c}_{ij}^H + \tau \sum_{j \in \mathcal{I}^L(i)} \mathbb{f}(\mathbf{U}_j^n) \mathbf{c}_{ij}^L = \sum_{K \in \mathcal{T}(i)} \mathbf{B}_{i,K}^n. \tag{5.3}$$

with

$$\mathbf{B}_{i,K}^n := -\tau \int_K \nabla \cdot (\Pi_h^H(\mathbb{f}(\mathbf{u}_h^n))) \varphi_i^H \, dx + \tau m_i^K \sum_{T \in \mathcal{T}_{K,i}^S} \frac{\alpha_i^S}{|S|} \int_S \nabla \cdot (\Pi_h^L(\mathbb{f}(\mathbf{u}_h^n))) \, dx. \tag{5.4}$$

The following result tells us how limiting should be done to keep the method cell-wise conservative.

Lemma 5.1. *Let $K \in \mathcal{T}_h$ and let $\ell_K \in [0, 1]$. Let $i \in \mathcal{V}$ and $j \in \mathcal{I}(i)$, and let $\ell_{ij} \in [0, 1]$ with the assumption that $\ell_{ij} = \ell_{ji}$. Let \mathbf{U}_i^{n+1} be defined by*

$$m_i \mathbf{U}_i^{n+1} = m_i \mathbf{U}_i^{L,n+1} + \sum_{j \in \mathcal{I}(i) \setminus \{i\}} \ell_{ij} \mathbf{D}_{ij}^n + \sum_{K \in \mathcal{T}(i)} \ell_K \mathbf{B}_{i,K}^n. \tag{5.5}$$

Then \mathbf{u}_h^{n+1} and $\mathbf{u}_h^{L,n+1}$ carry the same mass, i.e., $\int_D \mathbf{u}_h^{n+1} \, dx = \int_D \mathbf{u}_h^{L,n+1} \, dx$. **Proof.** Summing (5.5) over $i \in \mathcal{V}$, we obtain

$$\int_D \mathbf{u}_h^{n+1} \, dx = \int_D \mathbf{u}_h^{L,n+1} \, dx + \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{I}(i)} \ell_{ij} \mathbf{D}_{ij}^n + \sum_{i \in \mathcal{V}} \sum_{K \in \mathcal{T}(i)} \ell_K \mathbf{B}_{i,K}^n.$$

Since $\mathbf{D}_{ij}^n = -\mathbf{D}_{ji}^n$ and $\ell_{ij} = \ell_{ji}$, we have $\sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{I}(i)} \ell_{ij} \mathbf{D}_{ij}^n = \mathbf{0}$. For the second term on the right-hand side $\mathbf{R} := \sum_{i \in \mathcal{V}} \sum_{K \in \mathcal{T}(i)} \ell_K \mathbf{B}_{i,K}^n$ we have

$$\begin{aligned} -\tau^{-1} \mathbf{R} &= \sum_{K \in \mathcal{T}_h} \ell_K \sum_{i \in \mathcal{I}(K)} -\tau^{-1} \mathbf{B}_{i,K}^n \\ &= \sum_{K \in \mathcal{T}_h} \ell_K \left[\int_K \nabla \cdot (\Pi_h^H(\mathbb{f}(\mathbf{u}_h^n))) \, dx - \sum_{i \in \mathcal{I}(K)} m_i^K \sum_{S \in \mathcal{T}_{K,i}^S} \frac{\alpha_i^S}{|S|} \int_S \nabla \cdot (\Pi_h^L(\mathbb{f}(\mathbf{u}_h^n))) \, dx \right]. \end{aligned}$$

Then using (4.11) we obtain $\mathbf{R} = \mathbf{0}$. The conclusion follows readily. \square We now show how the degree of freedom limiters ℓ_{ij} and cell limiters ℓ_K can be estimated so that the update defined in (5.5) satisfies $\Psi_i(\mathbf{U}^{n+1}) \geq 0$. We consider two cases. We assume that \mathcal{V} is affine in the first one, then we address the general situation in the second case.

5.2. Affine functionals

Let us consider the simple case when Ψ_i is affine; that is, let us assume that there exist $\mathbf{J}_i \in \mathbb{R}^m$ and $\mathbf{b}_i \in \mathbb{R}$, s.t. $\Psi_i(\mathbf{V}) = \mathbf{J}_i \cdot \mathbf{V} + \mathbf{b}_i$. In this situation, one can apply the limiting technique inspired by [20] (see Eq. (10)–(13) therein and [19]). Let us give the details. Let $K \in \mathcal{T}(i)$ and

$$\mathcal{I}^-(i) := \left\{ j \in \mathcal{I}(i) \mid \mathbf{J}_i \cdot \mathbf{D}_{ij}^n < 0 \right\}, \quad \mathcal{I}^-(K) := \left\{ K \in \mathcal{T}(K) \mid \mathbf{J}_i \cdot \mathbf{B}_{i,K}^n < 0 \right\}. \tag{5.6a}$$

$$\mathbf{P}_i^- := \frac{1}{m_i} \left[\sum_{j \in \mathcal{I}^-(i)} \mathbf{D}_{ij}^n + \sum_{K \in \mathcal{I}^-(i)} \mathbf{B}_{i,K}^n \right], \quad \ell^i := \frac{\min(\Psi_i(\mathbf{U}_i^{L,n+1}), -\mathbf{P}_i^-)}{-\mathbf{P}_i^-}, \tag{5.6b}$$

$$\ell_j^i := \begin{cases} \ell^i & \text{if } \mathbf{J}_i \cdot \mathbf{D}_{ij}^n < 0 \\ 1 & \text{otherwise,} \end{cases} \quad \ell_K^i := \begin{cases} \ell^i & \text{if } \mathbf{J}_i \cdot \mathbf{B}_{i,K}^n < 0 \\ 1 & \text{otherwise,} \end{cases} \tag{5.6c}$$

$$\ell_{ij} := \min \left(\ell_j^i, \ell_i^j \right), \quad \ell_K := \min_{i \in \mathcal{I}(K)} \ell_K^i. \tag{5.6d}$$

with the convention that $\ell^i = 1$ if $\mathcal{I}^-(i) = \emptyset$ and $\mathcal{I}^-(i) = \emptyset$.

Lemma 5.2. (Affine limiting) *Let \mathbf{U}^{n+1} be defined in (5.5). Assume that (5.6) holds, then $\Psi_i(\mathbf{U}^{n+1}) \geq 0$.*

Proof. Using (5.5) and (5.6), and since Ψ_i is affine, we infer that

$$\begin{aligned} \Psi_i(\mathbf{U}^{n+1}) &= \Psi_i(\mathbf{U}^{L,n+1}) + \frac{1}{m_i} \sum_{j \in \mathcal{J}(i)} \ell_{ij} \mathbf{J}_i \cdot \mathbf{D}_{ij}^n + \frac{1}{m_i} \sum_{K \in \mathcal{K}(i)} \ell_K \mathbf{J}_i \cdot \mathbf{B}_{i,K}^n \\ &\geq \Psi_i(\mathbf{U}^{L,n+1}) + \frac{1}{m_i} \sum_{j \in \mathcal{J}(i)} \ell_j^i \mathbf{J}_i \cdot \mathbf{D}_{ij}^n + \frac{1}{m_i} \sum_{K \in \mathcal{K}(i)} \ell_K^i \mathbf{J}_i \cdot \mathbf{B}_{i,K}^n \\ &= \Psi_i(\mathbf{U}^{L,n+1}) + \ell^i \mathbf{P}_i^- \geq 0. \quad \square \end{aligned}$$

Remark 5.3. (Stability with Respect to Roundoff Errors) We have observed that the expression $\ell^i := (-\mathbf{P}_i^-)^{-1} \min(\Psi_i(\mathbf{U}_i^{L,n+1}), -\mathbf{P}_i^-)$ is more stable with respect to roundoff errors than the identity $\ell^i := \min(\Psi_i(\mathbf{U}_i^{L,n+1}), (-\mathbf{P}_i^-)^{-1}, 1)$ often found in the “flux corrected transport” literature.

5.3. Non-affine functionals: Convex limiting

Now we do not make the assumption that Ψ_i is affine. One cannot use Zalezak’s technique consisting of grouping terms in positive and negative contributions. One possible way of dealing with this situation consists of adopting the convex limiting method introduced in [17, § 4.2] and [18, § 7], which is essentially a divide and conquer strategy. Let us give the details.

Lemma 5.4. Let \mathbf{U}^{n+1} be defined in (5.5). Let $\theta \in (0, 1)$ and let $(\lambda_j)_{j \in \mathcal{J}(i)}$ and $(\mu_K)_{K \in \mathcal{K}(i)}$ be real numbers in the open interval $(0, 1)$ such that $\sum_{j \in \mathcal{J}(i) \setminus \{i\}} \lambda_j = 1$ and $\sum_{K \in \mathcal{K}(i)} \mu_K = 1$. Then,

$$\Psi_i(\mathbf{U}^{n+1}) \geq \min \left(\min_{j \in \mathcal{J}(i) \setminus \{i\}} \Psi_i \left(\mathbf{U}_i^{L,n+1} + \frac{\ell_{ij}}{\theta \lambda_j m_i} \mathbf{D}_{ij}^n \right), \min_{K \in \mathcal{K}(i)} \Psi_i \left(\mathbf{U}_i^{L,n+1} + \frac{\ell_K}{(1-\theta) \mu_K m_i} \mathbf{B}_{i,K}^n \right) \right). \tag{5.7}$$

Proof. We have

$$\begin{aligned} \mathbf{U}^{n+1} &:= \mathbf{U}_i^{L,n+1} + \frac{1}{m_i} \sum_{j \in \mathcal{J}(i) \setminus \{i\}} \ell_{ij} \mathbf{D}_{ij}^n + \frac{1}{m_i} \sum_{K \in \mathcal{K}(i)} \ell_K \mathbf{B}_{i,K}^n \\ &= \theta \left(\sum_{j \in \mathcal{J}(i) \setminus \{i\}} \lambda_j \mathbf{U}_i^{L,n+1} + \frac{\ell_{ij}}{\theta m_i} \mathbf{D}_{ij}^n \right) + (1-\theta) \left(\sum_{K \in \mathcal{K}(i)} \mu_K \mathbf{U}_i^{L,n+1} + \frac{\ell_K}{(1-\theta) m_i} \mathbf{B}_{i,K}^n \right) \\ &= \sum_{j \in \mathcal{J}(i) \setminus \{i\}} \theta \lambda_j \left(\mathbf{U}_i^{L,n+1} + \frac{\ell_{ij}}{\theta \lambda_j m_i} \mathbf{D}_{ij}^n \right) + \sum_{K \in \mathcal{K}(i)} (1-\theta) \mu_K \left(\mathbf{U}_i^{L,n+1} + \frac{\ell_K}{(1-\theta) \mu_K m_i} \mathbf{B}_{i,K}^n \right) \end{aligned}$$

Then the assertion is a consequence of Ψ_i being quasiconcave. \square **Lemma 5.5.** Let $i \in \mathcal{V}$ and $\mathbf{P} \in \mathbb{R}^m$. Assume that $\mathbf{U}_i^{L,n+1} + \mathbf{P}$ is in the domain of Ψ_i . Let $\ell(i, \mathbf{P}) \in [0, 1]$ be such that

$$\ell(i, \mathbf{P}) = \begin{cases} 1 & \text{if } \Psi_i(\mathbf{U}_i^{L,n+1} + \mathbf{P}) \geq 0, \\ \max\{\ell \in [0, 1] \mid \Psi_i(\mathbf{U}_i^{L,n+1} + \ell \mathbf{P}) \geq 0\} & \text{otherwise.} \end{cases} \tag{5.8}$$

Then $\Psi_i(\mathbf{U}_i^{L,n+1} + \ell \mathbf{P}) \geq 0$ for every $\ell \in [0, \ell(i, \mathbf{P})]$. **Proof.** Let us verify that the definition of $\ell(i, \mathbf{P})$ makes sense. We first observe that the set $\{\ell \in [0, 1] \mid \Psi_i(\mathbf{U}_i^{L,n+1} + \ell \mathbf{P}) \geq 0\}$ is not empty because $\Psi_i(\mathbf{U}_i^{L,n+1}) \geq 0$. Then this non-empty bounded set must have a maximal element because Ψ_i is continuous. This established that $\ell(i, \mathbf{P})$ is well defined if $\Psi_i(\mathbf{U}_i^{L,n+1} + \mathbf{P}) < 0$. The definition of $\ell(i, \mathbf{P})$ is unambiguous if $\Psi_i(\mathbf{U}_i^{L,n+1} + \mathbf{P}) \geq 0$.

Let $L_0(\Psi_i) := \{\mathbf{U} \in \mathcal{D} \mid \Psi_i(\mathbf{U}) \geq 0\}$. Notice that the set $L_0(\Psi_i)$ is not empty because $\mathbf{U}_i^{L,n+1} \in L_0(\Psi_i)$. It is also convex because Ψ_i is quasiconcave. Then for all $\ell \in [0, \ell_K^i]$ we have $\Psi_i(\mathbf{U}_i^{L,n+1} + \ell \mathbf{P}) \geq 0$ because $\mathbf{U}_i^{L,n+1} \in L_0(\Psi_i)$, $\mathbf{U}_i^{L,n+1} + \ell_K^i \mathbf{P} \in L_0(\Psi_i)$ and $L_0(\Psi_i)$ is convex. \square

Remark 5.6. (Line Search) Computing the maximal element in the set $\{\ell \in [0, 1] \mid \Psi_i(\mathbf{U}_i^{L,n+1} + \ell \mathbf{P}) \geq 0\}$ can be done by a line search when $\Psi_i(\mathbf{U}_i^{L,n+1} + \mathbf{P}) < 0$. The line search problem has a unique solution when Ψ_i is strictly quasiconcave since in this case the function $[0, 1] \ni \ell \rightarrow \Psi_i(\mathbf{U}_i^{L,n+1} + \ell \mathbf{P})$ is strictly monotone decreasing and therefore the equation $\Psi_i(\mathbf{U}_i^{L,n+1} + \ell \mathbf{P}) = 0$ has a unique solution in $(0, 1)$.

Theorem 5.7. Let $\theta, (\lambda_j)_{j \in \mathcal{J}(i)}$, and $(\mu_K)_{K \in \mathcal{K}(i)}$ be real number in $(0, 1)$ as in Lemma 5.4. (i.e., $\sum_{j \in \mathcal{J}(i) \setminus \{i\}} \lambda_j = 1$ and $\sum_{K \in \mathcal{K}(i)} \mu_K = 1$). Let

$$\ell_{ij} := \min \left(\ell \left(i, \frac{1}{\theta \lambda_j m_i} \mathbf{D}_{ij}^n \right), \ell \left(j, \frac{1}{\theta \lambda_j m_i} \mathbf{D}_{ij}^n \right) \right), \quad \ell_K^i := \ell \left(i, \frac{1}{(1-\theta) \mu_K m_i} \mathbf{B}_{i,K}^n \right). \tag{5.9}$$

Then $\Psi_i(\mathbf{U}^{n+1}) \geq 0$. **Proof.** Apply Lemma 5.4 and Lemma 5.5. \square

Remark 5.8. (Parameter θ) The purpose of the coefficients $\theta \in (0, 1)$ is to transform the identity $\mathbf{U}^{n+1} := \mathbf{U}_i^{L,n+1} + \frac{1}{m_i} \sum_{j \in \mathcal{J}(i) \setminus \{i\}} \ell_{ij} \mathbf{D}_{ij}^n + \frac{1}{m_i} \sum_{K \in \mathcal{T}(i)} \ell_K \mathbf{B}_{i,K}^n$ into the following convex combination $\mathbf{U}^{n+1} := \theta \left(\mathbf{U}_i^{L,n+1} + \frac{1}{\theta m_i} \sum_{j \in \mathcal{J}(i) \setminus \{i\}} \ell_{ij} \mathbf{D}_{ij}^n \right) + (1 - \theta) \left(\mathbf{U}_i^{L,n+1} + \frac{1}{(1-\theta)m_i} \sum_{K \in \mathcal{T}(i)} \ell_K \mathbf{B}_{i,K}^n \right)$. The free parameter θ can be used to balance the terms $\frac{1}{m_i} \sum_{j \in \mathcal{J}(i) \setminus \{i\}} \ell_{ij} \mathbf{D}_{ij}^n$ and $\frac{1}{m_i} \sum_{K \in \mathcal{T}(i)} \ell_K \mathbf{B}_{i,K}^n$. Assuming that the terms \mathbf{D}_{ij}^n and $\mathbf{B}_{i,K}^n$ have all the same magnitude, then a good choice consists of setting θ so that $\frac{\text{card}(\mathcal{J}(i))}{\theta} = \frac{\text{card}(\mathcal{T}(i))}{1-\theta}$, that is, $\theta = \frac{\text{card}(\mathcal{J}(i))}{\text{card}(\mathcal{J}(i)) + \text{card}(\mathcal{T}(i))}$. But it is actually interesting to bias θ a little bit according to the magnitude of the coefficients \mathbf{D}_{ij}^n , $\mathbf{B}_{i,K}^n$. More precisely, let us denote by h_i the local meshsize at the node \mathbf{a}_i , and let us set $(\Delta \mathbf{U})_i = \max_{j \in \mathcal{J}(i)} \|\mathbf{U}_j^n - \mathbf{U}_i^n\|_*$, where $\|\mathbf{V}\|_*$ is a norm in \mathbb{R}^m that is dimensionally consistent (i.e., combines the various components of \mathbf{V} in a way that is dimensionally coherent). Then one expects \mathbf{D}_{ij}^n to scale like $\tau m_i h_i^{-1} f'_{i,\max} \times (\Delta \mathbf{U})_i$ where $f'_{i,\max}$ is an upper bound on the norm induced by $\|\cdot\|_*$ of the Jacobian of \mathbb{f} , whereas, owing to Proposition 4.2, one expects $\mathbf{B}_{i,K}^n$ to scale like $\tau m_i h_i^{-1} f''_{i,\max} \times (\Delta \mathbf{U})_i^2$, where $f''_{i,\max}$ is an upper bound on the induced norm of the Hessian of \mathbb{f} . Hence, it is reasonable to expect $\mathbf{B}_{i,K}^n$ to be far smaller than \mathbf{D}_{ij}^n in regions where the solution is smooth. One can then use $\theta = \frac{\sum_{K \in \mathcal{T}(i)} \|\mathbf{B}_{i,K}^n\|_*}{\sum_{j \in \mathcal{J}(i)} \|\mathbf{D}_{ij}^n\|_* + \sum_{K \in \mathcal{T}(i)} \|\mathbf{B}_{i,K}^n\|_*}$ or variations of this idea. Another way to proceed consists of replacing $\theta \lambda_j$ in (5.7) by $\frac{\|\mathbf{D}_{ij}^n\|_*}{\sum_{j \in \mathcal{J}(i) \setminus \{i\}} \|\mathbf{D}_{ij}^n\|_* + \sum_{K \in \mathcal{T}(i)} \|\mathbf{B}_{i,K}^n\|_*}$ and $(1 - \theta) \mu_K$ by $\frac{\|\mathbf{B}_{i,K}^n\|_*}{\sum_{j \in \mathcal{J}(i) \setminus \{i\}} \|\mathbf{D}_{ij}^n\|_* + \sum_{K \in \mathcal{T}(i)} \|\mathbf{B}_{i,K}^n\|_*}$.

6. Numerical illustrations

In this section we briefly illustrate the performance of the method described in this paper.

6.1. Technical details

Two independent codes have been written to verify reproducibility. The first one, `Code 1`, does not use any particular software and is written in Fortran 2003. It is based on Lagrange elements on simplices and is dimension-independent. The second code, `Code 2`, uses the open-source finite element library FEniCS and is written in C++ and python, see e.g., [35]. The implementation in FEniCS is independent of the space dimension and the polynomial degree of the approximation. The numerical quadratures in both codes are chosen to be exact for the mass matrix. All the computations are done on simplices. The computations with \mathbb{P}_1 and \mathbb{P}_3 finite elements are done with the Lagrange bases. The computations with \mathbb{P}_2 finite elements are done with the Bernstein basis as the lumped mass matrix for continuous \mathbb{P}_2 Lagrange elements is singular. We only report tests for scalar conservation equations. Tests for the compressible Euler equations using the present method combined the time stepping techniques developed in [24] will be reported in the second part of this work.

The time stepping is done with the fourth-order five stages strong stability preserving explicit Runge–Kutta method with five stages (see [36, p. 522]). The time step is defined by the expression

$$\tau = s \times \text{CFL} \times \min_{i \in \mathcal{I}} \frac{m_i}{\sum_{j \in \mathcal{J}^L(i) \setminus \{i\}} a_{ij}^{L,n}}, \tag{6.1}$$

with $a_{ij}^{L,n}$ defined in (4.2) and s is the number of stages of the Runge–Kutta method; here, we use $s = 5$. The high-order viscosity for scalar conservation equations is computed by estimating the entropy commutator $\int_D (\eta'(v) \nabla \cdot \mathbb{f}(v) - \mathbb{f}(v) \nabla \eta(v)) \varphi_i^H dx = 0$ with $\eta(v) = \text{exp}(v)$. This is done as follows:

$$X_i(\mathbf{u}_h^n) := m_i^{-1} \int_D \Pi_h^H(\eta'(\mathbf{u}_h^n)) \nabla \cdot (\Pi_h^H(\mathbb{f}(\mathbf{u}_h^n))) \varphi_i^H dx, \tag{6.2a}$$

$$Y_i(\mathbf{u}_h^n) := m_i^{-1} \int_D \Pi_h^H(\mathbb{f}(\mathbf{u}_h^n)) \cdot \nabla \Pi_h^H(\eta(\mathbf{u}_h^n)) \varphi_i^H dx, \tag{6.2b}$$

$$\alpha_i^n := \frac{|X_i(\mathbf{u}_h^n) - Y_i(\mathbf{u}_h^n)|}{|X_i(\mathbf{u}_h^n)| + |Y_i(\mathbf{u}_h^n)| + \epsilon \max_{i \in \mathcal{I}} (|X_i(\mathbf{u}_h^n)| + |Y_i(\mathbf{u}_h^n)|)}, \quad \epsilon = 10^{-2} \tag{6.2c}$$

$$a_{ij}^{H,n} := a_{ij}^{L,n} \max\left(\psi\left(\frac{n}{\alpha_i}\right), \psi\left(\alpha_j^n\right)\right), \tag{6.2d}$$

By definition, the normalized residual α_i^n takes values in $[0, 1]$. The numerator in the definition of α_i^n behaves like the truncation error of the Galerkin method. This quantity approximately behaves like $\mathcal{O}(h^k)$ where k is the polynomial degree of the approximation and h is the mesh size. Finally we set

$$d_{ij}^{H,n} := d_{ij}^{L,n} \max \left(\psi \left(\frac{\alpha_i^n}{\alpha_0} \right), \psi \left(\frac{\alpha_j^n}{\alpha_0} \right) \right), \tag{6.3}$$

where $\alpha_0 \in (0, 1]$, the activation function ψ is defined as follows:

$$\psi(t) := 1 + (t - 1)^2 (6t^2 + 3t + 1) (t - 1 - (t - 1)_+), \tag{6.4}$$

where $(x)_+ = (x + |x|)/2$. Notice that $\psi(0) = 0$, $\psi(\frac{1}{2}) = \frac{1}{2}$ and $\psi(t) = 1$ for all $t \in [\alpha_0, 1]$, and $\psi(t) \sim 10t^3 + \mathcal{O}(t^4)$. As a result, one recovers $d_{ij}^{H,n} = d_{ij}^{L,n}$ if the entropy commutator is larger than α_0 , and $d_{ij}^{H,n} = \mathcal{O}(h^{3k}) \times d_{ij}^{L,n}$ otherwise. The numerical tests reported in this paper are done with $\alpha_0 = 0.4$. An activation function with the same purpose is used in [37, Eq. (8)]. Up to re-scaling, the activation function therein behaves like $\frac{1}{2} (1 + \sin(\frac{\pi}{2}(t - 1)))$.

The limiting is done at each grid point at each time stage of the Runge–Kutta method as explained in Section 5.2 and the process is iterated four times. The local bounds are relaxed using the process explained in [17, §4.7] and we use the minmod version therein to average the estimate of the local second variation (see (4.11) and (4.13) in [17]). We now give the details for completeness and reproducibility. First we estimate the second variation of u_i^n at every grid point $i \in \mathcal{I}$ by setting

$$\Delta^2 U_i^n := \frac{\sum_{j \in \mathcal{J}(i) \setminus \{i\}} \beta_{ij} (U_j^n - U_i^n)}{\sum_{j \in \mathcal{J}(i) \setminus \{i\}} \beta_{ij}}, \tag{6.5}$$

where the coefficients $\beta_{ij} := \int_D \nabla \varphi_i^H \cdot \nabla \varphi_j^H \, dx$ are the entries of the high-order stiffness matrix associated with the weak form of the Laplace operator. This definition of β_{ij} guarantees that $\Delta^2 U_i^n = 0$ if u_i^n is locally linear on the support of φ_i . It is our experience that using the high-order stiffness matrix gives a better estimation of the local second variation than using the low-order one. Then we set

$$\widetilde{\Delta^2} U_i^n := \text{minmod} \left\{ \Delta^2 U_j^n \mid j \in \mathcal{J}^L(i) \right\}, \tag{6.6}$$

where minmod of a finite set of real numbers is zero if there are two numbers of different sign in this set and is equal to the number whose absolute value is the smallest otherwise. Finally, using again the low-order stencil, we define $U_i^{\text{min},n} := \min_{i \in \mathcal{J}^L(i)} U_i^n$ and $U_i^{\text{max},n} := \max_{i \in \mathcal{J}^L(i)} U_i^n$. Notice that we use the low-order stencil to define $\widetilde{\Delta^2} U_i^n$, $U_i^{\text{min},n}$ and $U_i^{\text{max},n}$. The rationale for using the low-order stencil in these definitions is based on the observation that the time step restriction induced by the CFL condition implies that the domain of influence of the degree of freedom i is the support of the low-order shape function φ_i^L . To be certain that the global bounds are not violated, the relaxed bounds are not allowed to exceed the minimum and the maximum values of the initial data. More precisely denoting $U^\flat = \text{ess inf } u_0$ and $U^\sharp = \text{ess sup } u_0$, we set

$$U_i^{\text{min},n} := \max \left(U_i^{\text{min},n} - |\widetilde{\Delta^2} U_i^n|, U^\flat \right) \tag{6.7a}$$

$$U_i^{\text{max},n} := \min \left(U_i^{\text{max},n} + |\widetilde{\Delta^2} U_i^n|, U^\sharp \right). \tag{6.7b}$$

Strictly enforcing the global bounds is a stringent constraint. We impose it to demonstrate that the proposed method is genuinely high-order accurate (as this may not always be the case for other high-order methods proposed in the literature). This relaxation process is applied to all the simulations reported in the paper without exception to demonstrate the robustness of the method. The (affine) functionals corresponding to the above bounds are $\Psi_i^{\text{min}}(V) := V - U_i^{\text{min},n}$ and $\Psi_i^{\text{max}}(V) := U_i^{\text{max},n} - V$.

For all the tests reported below, the L^1 – norm of the difference $u_h - u$ is measured by using high-order quadratures. We use the exact solution at the quadrature points instead of the Lagrange interpolant of the solution to avoid extraneous super-convergence effects that may bias the tests and lead us to draw over-optimistic conclusions. The error are all relative, that is we report $\|u_h - u\|_{L^1(D)} / \|u\|_{L^1(D)}$.

Remark 6.1. (Stencil) One can also use the high-order stencil to define $U_i^{\text{min},n} := \min_{i \in \mathcal{J}(i)} U_i^n$ and $U_i^{\text{max},n} := \max_{i \in \mathcal{J}(i)} U_i^n$. This increases the range of CFL numbers for which optimal convergence rate is observed in the L^∞ -norm for smooth solutions, but this also slightly deteriorates the L^1 -norm convergence rate on problems developing shocks. We have observed that using the high-order stencil in the definition of the second variation $\widetilde{\Delta^2} U_i^n$ significantly deteriorates the L^1 -norm convergence rate on problems developing shocks.

6.2. Swirling flow

We start by demonstrating that progresses have indeed been made on the low-order method since the observations made in [1, §3.3.2.1] and [2, §4.2.2]. We consider the same swirling model problem as in Section 4.1. We solve this problem with continuous \mathbb{P}_2

Table 1
 Problem (6.8), \mathbb{P}_1 , structured meshes (left), unstructured meshes (right), $T = 0.3$, CFL = 0.1.

I	Present \mathbb{P}_2		Standard \mathbb{P}_2		I	Present \mathbb{P}_3		Standard \mathbb{P}_3	
	L^1	Rate	L^1	Rate		L^1	Rate	L^1	Rate
1681	6.13E-01	–	9.19E-01	–	2116	5.59E-01	–	9.74E-01	–
3721	4.96E-01	0.53	8.48E-01	0.20	3721	4.72E-01	0.60	9.66E-01	0.03
14641	3.11E-01	0.68	6.66E-01	0.35	14641	2.89E-01	0.72	8.97E-01	0.11
58081	1.77E-01	0.81	4.56E-01	0.55	32761	2.08E-01	0.82	8.15E-01	0.24
130321	1.24E-01	0.88	3.46E-01	0.69	130321	1.13E-01	0.89	6.22E-01	0.39
231361	9.56E-02	0.91	2.78E-01	0.76	292681	7.73E-02	0.93	4.98E-01	0.55
361201	7.77E-02	0.93	2.33E-01	0.80	519841	5.89E-02	0.95	4.14E-01	0.64

Table 2

Convergence tests: $\mathbb{P}_1, \mathbb{P}_2, \mathbb{P}_3$, structured meshes, Problem (6.8), $T = 0.3$, CFL = 0.1.

Code 1								
I	\mathbb{P}_1, L^1	Rate	I	\mathbb{P}_2, L^1	Rate	I	\mathbb{P}_3, L^1	Rate
8281	5.18E-03	–	8281	3.66E-03	–	3721	1.89E-02	–
14641	2.91E-03	2.02	14641	2.13E-03	1.90	8281	4.56E-03	3.56
32761	1.28E-03	2.05	32761	9.32E-04	2.05	18496	9.92E-04	3.79
58081	7.16E-04	2.02	58081	5.22E-04	2.02	32761	3.12E-04	4.05
90601	4.59E-04	2.01	90601	3.34E-04	2.02	73441	7.38E-05	3.57
160801	2.58E-04	2.01	160801	1.90E-04	1.96	130321	2.00E-05	4.55
251001	1.65E-04	2.00	251001	1.24E-04	1.92	292681	2.60E-06	5.04
Code 2								
I	\mathbb{P}_1, L^1	Rate	I	\mathbb{P}_3, L^1	Rate			
8281	5.39E-03	–	8281	4.19E-03	–			
18496	2.36E-03	2.06	18496	9.81E-04	3.61			
41209	1.03E-03	2.08	40804	2.13E-04	3.86			
92416	4.53E-04	2.02	90601	4.58E-05	3.85			
207025	2.02E-04	2.01	203401	5.11E-06	5.43			
465124	8.97E-05	2.01	456976	5.27E-07	5.61			

Table 3

Convergence tests: $\mathbb{P}_1, \mathbb{P}_2, \mathbb{P}_3$, unstructured meshes, problem (6.8), $T = 0.3$, CFL = 0.1.

Code 1								
I	\mathbb{P}_1, L^1	Rate	I	\mathbb{P}_2, L^1	Rate	I	\mathbb{P}_3, L^1	Rate
335	2.12E-01	–	349	1.62E-01	–	760	3.22E-01	–
1249	3.77E-02	2.63	1273	4.00E-02	2.16	2815	3.13E-02	3.56
4852	7.84E-03	2.32	4865	7.10E-03	2.58	10849	2.68E-03	3.64
11742	3.23E-03	2.01	11749	2.67E-03	2.22	42892	1.54E-04	4.16
19167	1.96E-03	2.02	19149	1.55E-03	2.23	104470	1.99E-05	4.59
46568	8.03E-04	2.02	46565	6.20E-04	2.05	170959	4.89E-06	5.69
185518	2.02E-04	1.99	185469	1.64E-04	1.92	416704	2.09E-06	1.91
Code 2								
I	\mathbb{P}_1, L^1	Rate	I	\mathbb{P}_3, L^1	Rate			
3414	1.55E-02	–	2638	1.01E-01	–			
5754	9.96E-03	1.69	6313	2.00E-02	3.72			
13295	4.46E-03	1.92	13084	3.47E-03	4.80			
16474	3.37E-03	2.60	30106	8.95E-04	3.25			
28860	2.45E-03	1.15	64330	1.34E-04	5.01			
37412	1.82E-03	2.28	146890	4.00E-05	2.92			
65923	1.09E-03	1.81	334648	4.72E-06	5.19			

and \mathbb{P}_3 Lagrange finite elements with final time $t = 1$. For each polynomial degree we first solve the problem with the low-order method presented in the paper (see (4.9)), then we solve it again with the low-order method using the graph viscosity based on the high-order stencil (see (4.1)). This test is done on various meshes with increasing resolution. The relative L^1 -norm of the error is computed at the end of each simulation. The results are reported in Table 1. We clearly observe that the present method has almost no pre-asymptotic range (green columns for the \mathbb{P}_2 approximation and the \mathbb{P}_3 approximation), whereas the standard method has a long pre-asymptotic range (orange columns). For instance, on the finest mesh composed of 519841 degrees of freedom, the error with the new low-order \mathbb{P}_3 approximation is 5.89×10^{-2} whereas the error with the standard method is 4.14×10^{-1} . There is almost one order of magnitude difference between the two approximations. This series of tests confirms the claims made in the paper regarding the low-order viscosity.

6.3. 2D smooth problem

We now consider a two-dimensional linear transport problem $\partial_t u + \nabla \cdot (\beta u) = 0$ with $\beta = (1, 0)^T$ and initial condition

Table 4

Convergence tests: Limited technique vs. Galerkin for \mathbb{P}_2 polynomials on structured and unstructured meshes, problem (6.8), Code 1, $T = 0.3$, CFL = 0.1.

Structured meshes			Unstructured meshes						
I	Lim., L^1	Rate	Gal., L^1	Rate	I	Lim., L^1	Rate	Gal., L^1	Rate
8281	3.66E-03	–	5.84E-03	–	349	1.62E-01	–	2.74E-01	–
14641	2.13E-03	1.90	3.18E-03	2.13	1273	4.00E-02	2.16	6.99E-02	2.11
32761	9.32E-04	2.05	1.37E-03	2.09	4865	7.10E-03	2.58	1.43E-02	2.37
58081	5.22E-04	2.02	7.60E-04	2.05	11749	2.67E-03	2.22	5.84E-03	2.03
90601	3.34E-04	2.02	4.82E-04	2.05	19149	1.55E-03	2.23	3.31E-03	2.32
160801	1.90E-04	1.96	2.69E-04	2.03	46565	6.20E-04	2.05	1.28E-03	2.13
251001	1.24E-04	1.92	1.72E-04	2.02	185469	1.64E-04	1.92	3.14E-04	2.04

Table 5

Problem (6.9), \mathbb{P}_2 and \mathbb{P}_3 , unstructured meshes, $T = 1$, CFL = 0.3. Left: results from Code 1. Right: results from Code 2.

Code 1			Code 1			Code 2		
I	\mathbb{P}_2, L^1	Rate	I	\mathbb{P}_3, L^1	Rate	I	\mathbb{P}_3, L^1	rate
3743	3.26E-01	–	3880	4.63E-01	–	6211	3.91E-01	–
14528	2.02E-01	0.71	15058	2.72E-01	0.78	10981	3.09E-01	0.83
56510	1.17E-01	0.80	57712	1.55E-01	0.84	24571	2.21E-01	0.83
99690	1.05E-01	0.37	99859	1.22E-01	0.86	55081	1.54E-01	0.89
224699	8.43E-02	0.55	223924	8.82E-02	0.81	121807	1.11E-01	0.82
394000	7.18E-02	0.58	394594	7.22E-02	0.70	270901	8.13E-02	0.79

$$u_0(\mathbf{x}) = \begin{cases} e^{\frac{r(x)^2 + 2r_0^2}{r(x)^2 - r_0^2}} & \text{if } r(\mathbf{x}) < r_0, \\ 0 & \text{otherwise,} \end{cases} \tag{6.8}$$

where $r(\mathbf{x}) := \|\mathbf{x} - \mathbf{x}_0\|_{\rho^2}$, $\mathbf{x}_0 = (\frac{7}{20}, \frac{1}{2})$, and $r_0 = 0.2$. We solve the problem in the unit square $D = \{\mathbf{x} := (x_1, x_2) \in \mathbb{R}^2 \mid 0 \leq x_1 \leq 1 \text{ and } 0 \leq x_2 \leq 1\}$ up to $T = 0.3$. The simulations are done on several uniform and nonuniform grids with CFL=0.1.

We do two series of computations with \mathbb{P}_1 , \mathbb{P}_2 and \mathbb{P}_3 continuous finite elements on meshes with various meshsizes. The first series is done on structured meshes and the other one is done on unstructured Delaunay meshes. The convergence results for structured meshes are shown in Table 2. The results for unstructured meshes are shown in Table 3. We observe second-order accuracy for \mathbb{P}_1 elements and fourth-order accuracy for \mathbb{P}_3 elements, both with structured and unstructured meshes. We observe sub-optimal second-order for \mathbb{P}_2 elements. This behavior is well-known and is independent of the limiting technique described in the paper. It is well established in the literature that (unless super-convergence occurs) one recovers optimal rate for \mathbb{P}_2 elements only by adding linear stabilization (e.g., Galerkin Least-Squares, edge stabilization, subgrid viscosity, etc.). To illustrate this point we compare in Table 4 the results obtained with the present method with those obtained with the unlimited Galerkin method with \mathbb{P}_2 elements. We observe that the Galerkin method is only second-order accurate. We also notice also that the limited \mathbb{P}_2 approximation is slightly more precise than the Galerkin method whether the meshes are structured or not.

Remark 6.2. (Relaxation) In tests not reported here for brevity, we have observed that the convergence rate for the \mathbb{P}_3 approximation deteriorates to second-order in the L^1 -norm and, irrespective of the polynomial degree the approximation, reduces to first-order in the L^∞ -norm if one does not relax the local bounds as described in (6.5)–(6.6)–(6.7). This loss of accuracy in the neighborhood of extrema due to limiting is a phenomenon that is well documented in the literature (see, e.g., [38, §3.3], [9, p. 2753]). A typical way to address this issue in the finite volume literature consists of relaxing the slope reconstructions; see, e.g., [39, Eq. (5.7)], [40]. Similar techniques can be used with discontinuous Galerkin approximations as in [9,41]. The relaxation technique (6.5)–(6.6)–(6.7) used in the paper has been introduced in [17, §4.7] (see (4.11) and (4.13) therein).

6.4. Three body rotation

In this example, we consider the linear transport equation with the flux $\beta = 2\pi(-x_2, x_1)^T$ and the following initial condition

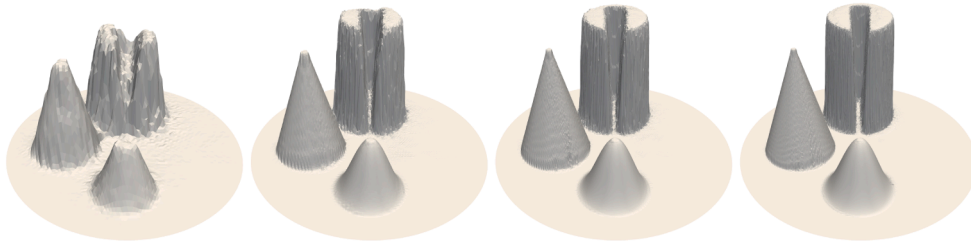


Fig. 5. Three body rotation. From the left to right: approximation with 6211, 55081, 121807 AND 258673 \mathbb{P}_3 nodes with CFL=0.3.

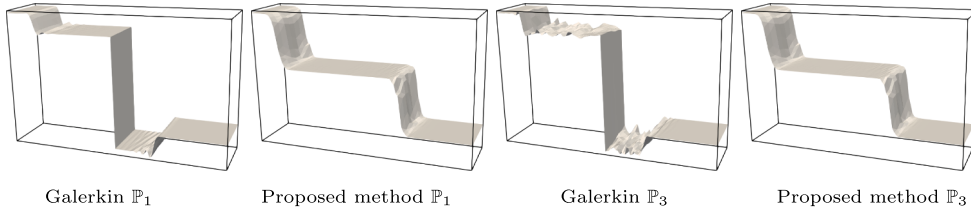


Fig. 6. Non-strictly convex flux (6.10); Solution at $t = 0.75$ with CFL=0.5; 847 grid points.

$$u_0(\mathbf{x}) = \begin{cases} 1 & \text{if } r_1(\mathbf{x}) \leq r_0 \text{ and } (|x_1| \geq \frac{1}{20} \text{ or } x_2 \geq \frac{7}{10}), \\ 1 - \frac{r_2(\mathbf{x})}{r_0} & \text{if } r_2(\mathbf{x}) \leq r_0, \\ \frac{1}{4} \left(1 + \cos\left(\frac{r_3(\mathbf{x})}{r_0} \pi\right) \right) & \text{if } r_3(\mathbf{x}) \leq r_0, \\ 0 & \text{otherwise,} \end{cases} \tag{6.9}$$

where $r_0 := 0.3$, $r_1(\mathbf{x}) := \|\mathbf{x} - \mathbf{x}_1\|_{\ell^2}$, $r_2(\mathbf{x}) := \|\mathbf{x} - \mathbf{x}_2\|_{\ell^2}$, $r_3(\mathbf{x}) := \|\mathbf{x} - \mathbf{x}_3\|_{\ell^2}$, with $\mathbf{x}_1 := (0, \frac{1}{2})$, $\mathbf{x}_2 := (-\frac{1}{2}, 0)$, $\mathbf{x}_3 := (0, -\frac{1}{2})$.

The computational domain $D = \{\mathbf{x} \in \mathbb{R}^2 \mid \|\mathbf{x}\|_{\ell^2} \leq 1\}$ is triangulated using non-nested meshes. The simulations are done with \mathbb{P}_2 and \mathbb{P}_3 continuous finite elements on nonuniform meshes up to $T = 1$ with CFL=0.3. The convergence results are shown in Table 5. We observe that the method converges with a rate that is similar to what is reported in [29, Tab. 3] with \mathbb{P}_1 finite elements. This is normal as the solution is not smooth.

We show in Fig. 5 the graph of the approximate solution obtained on four different meshes using Code 2. These results are comparable to those in [29, Fig. 2].

6.5. Non-strictly convex flux

In this section we illustrate the performance of the method with a flux that is not strictly convex. We consider the two-dimensional scalar conservation equation $\partial_t u + \nabla \cdot \mathbb{f}(u) = 0$ using the flux $\mathbb{f}(v) = (2 - v, 0)^T$ if $v \leq 2$ and $\mathbb{f}(v) = (2v - 4, 0)^T$ otherwise. The initial data is $u_0(x) = 1$ if $x \leq 0$ and $u_0(x) = 3$ otherwise. This flux is convex and Lipschitz, but it is not strictly convex: the velocity is piecewise constant and discontinuous. The entropy solution is

$$u(x, t) = \begin{cases} 1 & \text{if } x \leq -t, \\ 2 & \text{if } -t < x \leq 2t, \\ 3 & \text{if } 2t < x. \end{cases} \tag{6.10}$$

The solution is composed of two contact waves (i.e., the characteristics do not cross) separated by an expansion wave. One contact wave moves to the left at speed -1 , the other moves to the right at speed 2. This test is meant to eliminate high-order methods that are not dissipative enough. The lack of dissipation makes these methods to converge to weak solutions that violate entropy inequalities.

We solve the problem over the domain $D = (-2, 2) \times (0, 1)$ using uniform meshes. The solution is computed with the final time $t = 0.75$ using CFL=0.5. We show in Fig. 6 the graph of the solution obtained with \mathbb{P}_1 polynomials on a mesh composed of $120 \times (2 \times 6)$ elements and with \mathbb{P}_3 polynomials on a mesh composed of $40 \times (2 \times 2)$ elements. In both cases the number of grid points is 847. The solutions in the leftmost panel and in the next to last panel have been obtained with the Galerkin method using \mathbb{P}_1 and \mathbb{P}_3 finite elements, respectively. (We use exactly the same codes as the one used in the accuracy tests reported in Sections 6.3 and 6.4.) We observe that they both converge to a wrong weak solution with a stationary shock at $x = 0$. The reason for this behavior is that the viscosity coefficients d_{ij}^H are set to zero for the Galerkin approximation. The solutions in the second panel and in the last panel have

Table 6

Burgers; unstructured meshes; \mathbb{P}_2 and \mathbb{P}_3 finite elements; $T = 0.5$, CFL = 0.3. Left: results from Code 1. Right: results from Code 2.

Code 1			Code 1			Code 2		
I	\mathbb{P}_2, L^1	Rate	I	\mathbb{P}_3, L^1	Rate	I	\mathbb{P}_3, L^1	Rate
4561	5.41E-02	–	4315	4.67E-02	–	4702	5.59E-02	–
7549	4.10E-02	1.10	16864	2.58E-02	0.87	10315	4.17E-02	0.75
29849	2.27E-02	0.86	66919	1.35E-02	0.93	23338	2.79E-02	0.99
118785	1.33E-02	0.77	119560	1.09E-02	0.75	50962	1.84E-02	1.06
210317	9.37E-03	1.23	266785	6.92E-03	1.13	118477	1.24E-02	0.93
474137	6.77E-03	0.80	472573	4.96E-03	1.17	258673	8.82E-03	0.88
843717	4.98E-03	1.06	1065847	3.54E-03	0.83	469786	6.64E-03	0.95

been obtained with the proposed method using \mathbb{P}_1 and \mathbb{P}_3 finite elements, respectively. We observe that they both converge to the correct solution. Notice also that the graph of the solution in the expansion wave is free of oscillation (which would not be the case if we had only enforced global bounds as it sometimes done in the literature). We have verified that the method converges with the expected rate both with \mathbb{P}_1 and \mathbb{P}_3 finite elements (tables not shown here for brevity).

6.6. Burgers

We solve now a two-dimensional version of Burgers’ equation introduced in [29, §6.1]. The initial data is chosen so that there are sonic points over a one-dimensional manifold across the graph of the solution. This test is meant to weed out methods that are accurate but not dissipative enough to be able to select the entropy solution in expansion waves with sonic points. Here we recall that we use the same codes as the one used in the accuracy tests reported in Section 6.3 and Section 6.4; nothing is changed in the codes. We consider the domain $D = (-0.25, 1.75)^2$ and the following two-dimensional version of Burgers’ equation:

$$\partial_t u + \nabla \cdot (f(u)) = 0, \quad f(u) = \frac{1}{2} (u^2, u^2)^T, \quad u(x, 0) = u_0(x) \text{ a.e. } x \in D, \tag{6.11a}$$

$$\text{where } u_0(x) = \begin{cases} 1 & \text{if } \left| x_1 - \frac{1}{2} \right| \leq 1 \text{ and } \left| x_2 - \frac{1}{2} \right| \leq 1 \\ -a & \text{otherwise.} \end{cases} \tag{6.11b}$$

The exact solution is fully described in [29, §6.1]. The tests are done with $a = 0.75$. The errors are computed at $T = 0.5$. The results of the convergence tests for \mathbb{P}_2 and \mathbb{P}_3 continuous finite elements are reported in Table 6. We observe that the convergence rate in the L^1 -norm is close to 1, which is optimal.

Remark 6.3. (Convergence Rate) The convergence rate for Burgers’ equation is better than that observed for the linear transport equation with discontinuous data because the solution of Burgers’ equation maps $L^\infty(\mathbb{R})$ to the space of functions with bounded variations over \mathbb{R} (say $BV(\mathbb{R})$), see [42, Thm. 11.2.2]; hence, compactness is available. Notice that one can invoke [42, Thm. 11.2.2] here because the flux in (6.11a) is isotropic. We refer the reader to [43, Thm. 2] for further clarifications on this topic.

Remark 6.4. (Relaxation) Recall that relaxation of local bounds is only necessary to reach optimal accuracy for problems with smooth solutions in the L^∞ -norm. We have observed that the convergence rate in the L^1 -norm for Burgers’ equation remains close to 1 if the bounds are not relaxed. Tests not reported here for brevity show that optimal convergence is lost on Burgers’ equation if one uses the averaged version of the relaxation. A key conclusion of this paper is that using the minmod version of the relaxation produces optimal convergence rates in the L^1 -norm for all the cases (smooth linear transport problem, non-smooth linear transport problem, Burgers’ equation). In other words, the minmod relaxation method described in (6.5)–(6.6)–(6.7) is robust.

6.7. Non-convex flux

We finish by illustrating the performance of the method on a two-dimensional scalar conservation equation with a non-convex flux originally proposed in [44]:

$$\partial_t u + \nabla \cdot f(u) = 0, \quad u(x, 0) = u_0(x) = \begin{cases} \frac{14\pi}{4} & \text{if } \sqrt{x^2 + y^2} \leq 1 \\ \frac{\pi}{4} & \text{otherwise,} \end{cases} \tag{6.12}$$

with $f(u) := (\sin u, \cos u)^T$. High-order methods that are too greedy usually fail to converge to the entropy solution as they have the tendency to produce shocks where one should have expansions; see e.g., [44]. The computational domain is $D = [-2, 2] \times [-2.5, 1.5]$. The final time is $t = 1$. The simulation are done with Lagrange \mathbb{P}_1 and Lagrange \mathbb{P}_3 finite elements and with \mathbb{P}_2 Bernstein finite

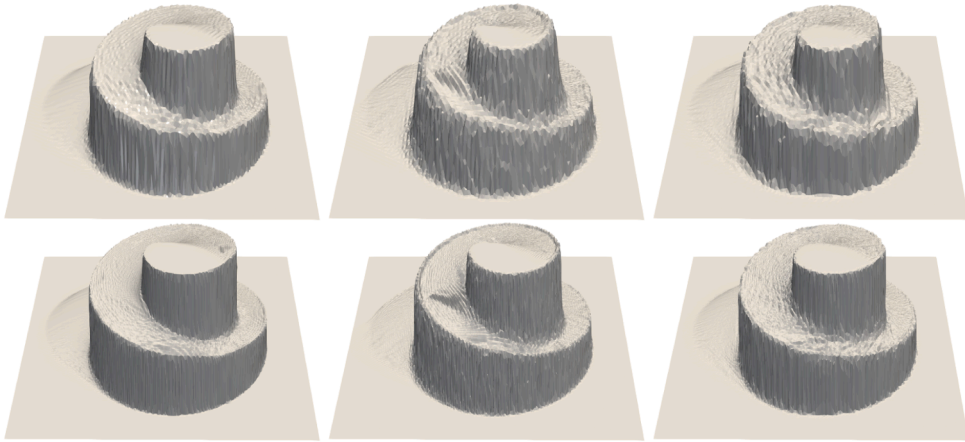


Fig. 7. Problem (6.12) at $t = 1$. From left to right: Lagrange \mathbb{P}_1 , Bernstein \mathbb{P}_2 , Lagrange \mathbb{P}_3 . Top row: approximately 16870 grid points. Bottom row: approximately 67000 grid points.

elements. Two computations are done on unstructured meshes for each polynomial degree. For the \mathbb{P}_1 approximation the first mesh is composed of 33240 triangles and 16861 degrees of freedom and the second mesh is composed of 134066 triangles and 67516 degrees of freedom. For the \mathbb{P}_2 approximation the first mesh is composed of 8312 triangles and 16865 degrees of freedom and the second mesh is composed of 33240 triangles and 66961 degrees of freedom. For the \mathbb{P}_3 approximation the first mesh is composed of 3698 triangles and 16882 degrees of freedom and the second mesh is composed of 14782 triangles and 67000 degrees of freedom.

The results are shown in Fig. 7. The method behaves correctly for each polynomial degree; the spiraling composite wave has the expected shape. These results are identical to what is reported in the literature (see e.g., [44, Fig. 5.11] or [14, Fig. 34]).

7. Conclusions

We have proposed in this paper answers to questions raised in the PhD theses of [1, §3.3.2.1] and [2, §4.2.2], where it was observed that the low-order invariant-domain preserving method proposed in [3] was not robust with respect to the polynomial degree. Building on ideas developed in [4] (see Propositions 3.1 and 3.2 therein), we have proposed a variation of the invariant-domain preserving method from [3] that behaves better as the polynomial degree increases. In particular, the low-order method is based on the closest neighbor stencil. The method has been implemented and tested with continuous \mathbb{P}_2 and \mathbb{P}_3 finite elements. The numerical tests demonstrate that the method behaves as advertised in the theory.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

No data was used for the research described in the article.

Appendix A. Sufficient criterion for Assumption 4.1

We give here a sufficient criterion for the condition (i) of Assumption 4.1 to hold.

Lemma A.1. Let $\widehat{S}_l \in \mathcal{T}_K^s$. Let $\Pi_{\widehat{S}_l}^L : \widehat{\mathcal{P}}^H \rightarrow \widehat{\mathcal{P}}^L$ be the Lagrange interpolation operator defined by $\Pi_{\widehat{S}_l}^L(\widehat{p}) = \sum_{n^s \in \widehat{\mathcal{I}}^s} \widehat{p}(\widehat{a}_{j^s(n^s, l)}^L) \widehat{\theta}_{n^s}^L$. Assume that there exists a collection of numbers $\{\beta^{\widehat{S}_l}\}_{\widehat{S}_l \in \mathcal{T}_K^s}$ such that the following identity holds true for every $\widehat{p} \in \widehat{\mathcal{P}}^H$ and every $k \in \{1 : d\}$:

$$\int_{\widehat{K}} \partial_{\widehat{x}_k} \widehat{p}(\widehat{x}) \, d\widehat{x} = \sum_{l \in \widehat{\mathcal{I}}^s} \beta^{\widehat{S}_l} \int_{\widehat{S}_l} \partial_{\widehat{x}_k} \Pi_{\widehat{S}_l}^L(\widehat{p}) \, d\widehat{x}. \tag{A.1}$$

Assume also that the mapping $T_K : \widehat{K} \rightarrow K$ is affine for all $K \in \mathcal{T}_h$. Then the condition (i) in Assumption 4.1 is met with $\beta^{\widehat{S}_l} := \beta^{\widehat{S}_l}$ where the index l is such that $S = T_K(\widehat{S}_l)$. **Proof.** Let $K \in \mathcal{T}_h$. Let $S_l \in \mathcal{T}_K^s$ be a subs-cell of K with $l \in \widehat{\mathcal{I}}^s$. Recall that $S_l = T_K(\widehat{S}_l)$. Let $T_{K|S_l} : \widehat{S}_l \rightarrow S_l$ be the restriction of T_K to \widehat{S}_l . Let \mathbb{J}_K and \mathbb{J}_{S_l} be the Jacobian matrices of T_K and $T_{K|S_l}$, respectively. Since the mapping T_K is affine, we have

Table 7

\widehat{j} , $\widehat{\beta}$ and $\widehat{\alpha}$ arrays for \mathbb{P}_2 and \mathbb{P}_3 polynomials in one dimension (uniform lattice). The arrays $\widehat{\alpha}$ are only given for the Lagrange bases.

	\mathbb{P}_2		\mathbb{P}_3				\mathbb{P}_2		\mathbb{P}_3		
$l \in \widehat{\mathcal{L}}$	1	2	1	2	3	$l \in \widehat{\mathcal{L}}$	1	2	3		
$\widehat{j}(1, l)$	1	2	1	3	2	$\widehat{\alpha}(1, l)$	1	$\frac{1}{2}$	$\frac{4}{9}$	$\frac{4}{9}$	1
$\widehat{j}(2, l)$	2	3	3	4	4	$\widehat{\alpha}(2, l)$	$\frac{1}{2}$	1	$\frac{5}{9}$	$\frac{4}{9}$	$\frac{5}{9}$
$\widehat{\beta}(l)$	1	1	1	1	1						

Table 8

\widehat{j} for \mathbb{P}_2 and \mathbb{P}_3 polynomials in two dimensions.

	\mathbb{P}_2				\mathbb{P}_3								
$l \in \widehat{\mathcal{L}}$	1	2	3	4	1	2	3	4	5	6	7	8	9
$\widehat{j}(1, l)$	1	4	2	3	1	6	8	4	2	6	5	4	3
$\widehat{j}(2, l)$	5	5	4	4	6	8	9	9	4	7	7	5	5
$\widehat{j}(3, l)$	6	6	6	5	8	10	10	10	9	10	10	10	7
$\widehat{\beta}(l)$	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{2}{3}$	2	$\frac{3}{4}$	$\frac{3}{2}$	$\frac{3}{4}$	$\frac{3}{2}$	$\frac{3}{4}$	$\frac{3}{4}$	$\frac{3}{2}$	$\frac{3}{4}$	$\frac{3}{4}$

$\mathbb{J}_{K|S_l} = \mathbb{J}_{S_l}$ and the determinants $\det(\mathbb{J}_{K|S_l}) = \det(\mathbb{J}_K)$ are constant. Let $\mathbf{p} \in \mathbf{P}_K^H$ with $\mathbf{p} = \sum_{i \in \mathcal{I}(K)} \mathbf{p}(\mathbf{a}_i) \varphi_i^H$. By definition of $\widehat{\mathbf{P}}^H$, there is $\widehat{\mathbf{p}} \in \widehat{\mathbf{P}}$ such that $\widehat{\mathbf{p}} = \mathbf{p} \circ \mathbf{T}_K$. Recall that the chain rule gives $(\nabla \cdot \mathbf{p})(\mathbf{T}_K(\widehat{\mathbf{x}})) \det(\mathbb{J}_K) = \nabla \cdot (\det(\mathbb{J}_K) \mathbb{J}_K^{-1} \widehat{\mathbf{p}})(\widehat{\mathbf{x}})$. Then owing to (A.1) and $\nabla \cdot (\mathbb{J}_K^{-1} \widehat{\mathbf{p}})(\widehat{\mathbf{x}}) = \mathbb{J}_K^{-1} : (\nabla \widehat{\mathbf{p}})(\widehat{\mathbf{x}})$ (since the mapping \mathbf{T}_K is affine), we infer that

$$\begin{aligned} \int_K \nabla \cdot \mathbf{p}(\mathbf{x}) \, dx &= \int_K (\nabla \cdot \mathbf{p})(\mathbf{T}_K(\widehat{\mathbf{x}})) \det(\mathbb{J}_K) \, d\widehat{\mathbf{x}} = \int_K \nabla \cdot (\det(\mathbb{J}_K) \mathbb{J}_K^{-1} \widehat{\mathbf{p}})(\widehat{\mathbf{x}}) \, d\widehat{\mathbf{x}} = \det(\mathbb{J}_K) \mathbb{J}_K^{-1} : \int_K (\nabla \widehat{\mathbf{p}})(\widehat{\mathbf{x}}) \, d\widehat{\mathbf{x}} = \det(\mathbb{J}_K) \mathbb{J}_K^{-1} \\ &: \sum_{l \in \widehat{\mathcal{L}}} \beta^{S_l} \int_{S_l} (\nabla \Pi_{S_l}^L \widehat{\mathbf{p}})(\widehat{\mathbf{x}}) \, d\widehat{\mathbf{x}} = \sum_{l \in \widehat{\mathcal{L}}} \beta^{S_l} \int_{S_l} \nabla \cdot (\det(\mathbb{J}_{S_l}) \mathbb{J}_{S_l}^{-1} \Pi_{S_l}^L \widehat{\mathbf{p}})(\widehat{\mathbf{x}}) \, d\widehat{\mathbf{x}} = \sum_{l \in \widehat{\mathcal{L}}} \beta^{S_l} \int_{S_l} \nabla \cdot (\Pi_{S_l}^L \widehat{\mathbf{p}} \circ \mathbf{T}_K^{-1})(\mathbf{x}) \, dx. \end{aligned}$$

Recalling that $\mathbf{p} := \sum_{i \in \mathcal{I}(K)} \mathbf{p}(\mathbf{a}_i) \varphi_i^H$ and the definitions (2.9) and (2.11), we have

$$\begin{aligned} \left(\Pi_{S_l}^L \widehat{\mathbf{p}} \circ \mathbf{T}_K^{-1} \right)_{|S_l} &= \sum_{n^s \in \mathcal{I}^s} \widehat{\mathbf{p}}(\widehat{\mathbf{a}}_{\widehat{j}(n^s, l)}) \widehat{\theta}_{n^s}^L \circ \mathbf{T}_{K|S_l}^{-1} = \sum_{n^s \in \mathcal{I}^s} \mathbf{p}(\mathbf{a}_{j(n^s, S_l)}) \varphi_{j(n^s, S_l)}^L \\ &=: \Pi_{S_l}^L(\mathbf{p}). \end{aligned}$$

The assertion readily follows. \square The immediate consequence of Lemma A.1 is that, provided the mesh is affine, one just need to verify that (i) holds on the reference element to ascertain that the condition (i) in Assumption 4.1 is met. This type of computation is done for finite elements of arbitrary polynomial degree in two and three dimensions in [4].

Lemma A.2. Assume that there exist two mappings $\widehat{\alpha} : \widehat{\mathcal{N}} \times \widehat{\mathcal{L}} \rightarrow \mathbb{R}$ and $\widehat{\beta} : \widehat{\mathcal{L}} \rightarrow \mathbb{R}$ such that the following holds for all $n \in \widehat{\mathcal{N}}$ and all $l \in \widehat{\mathcal{L}}$:

$$\sum_{\{(n^s, l) | n = \widehat{j}(n^s, l)\}} \widehat{\alpha}(n, l) = 1, \quad \sum_{n^s \in \mathcal{I}^s} \widehat{\alpha}(\widehat{j}(n^s, l), l) \widehat{m}_{j(n^s, l)}^K = \widehat{\beta}(l) |\widehat{S}_l|. \tag{A.2}$$

Assume that the mapping \mathbf{T}_K is affine for all $K \in \mathcal{T}_h$. For all $K \in \mathcal{T}_h$ and all $S_l \in \mathcal{F}_K^s$ with $\mathbf{T}_K(\widehat{S}_l) = S_l$, we set $\beta^{S_l} = \widehat{\beta}(l)$. Similarly we set $\alpha_{j(n^s, l)}^{S_l} := \widehat{\alpha}(n, l)$. Then the conditions (ii) in Assumption 4.1 are met. **Proof.** This is a simple consequence of the identity $\frac{|\widehat{K}|}{|\widehat{K}|} = \frac{|\widehat{S}_l|}{|S_l|}$ since \mathbf{T}_K is affine. \square

Examples for the coefficients $\widehat{\alpha}$ and $\widehat{\beta}$ are given in the sections below.

Appendix B. \widehat{j} , $\widehat{\alpha}$ And $\widehat{\beta}$ for \mathbb{P}_2 and \mathbb{P}_3 in one dimension

The arrays \widehat{j} , $\widehat{\alpha}$ and $\widehat{\beta}$ for \mathbb{P}_2 and \mathbb{P}_3 polynomials in one dimension are shown in Table 7. The arrays $\widehat{\alpha}$ are only given for the Lagrange bases. The nodes are enumerated using the increasing vertex index enumeration convention (see e.g., Definition 10.6 and §21.3.2 in [45]). Recall that the enumeration of the sub-cells is done with the index set \mathcal{L} .

Table 9

$\hat{\alpha}$ for the \mathbb{P}_2 and \mathbb{P}_3 Lagrange bases in two dimensions.

	\mathbb{P}_2				\mathbb{P}_3								
$l \in \widehat{\mathcal{L}}$	1	2	3	4	1	2	3	4	5	6	7	8	9
$\hat{\alpha}(1, l)$	1	$\frac{1}{2}$	1	1	1	y	x	y	1	x	y	x	1
$\hat{\alpha}(2, l)$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{3}$	y	x	y	$\frac{1}{3}$	x	y	x	$\frac{1}{3}$
$\hat{\alpha}(3, l)$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{3}$	z	$\frac{1}{3} - z$	z	$\frac{1}{3}$	$\frac{1}{3} - z$	z	$\frac{1}{3} - z$	$\frac{1}{3}$

Table 10

$\hat{\alpha}$ for the \mathbb{P}_2 and \mathbb{P}_3 Bernstein bases in two dimensions.

	\mathbb{P}_2				\mathbb{P}_3								
$l \in \widehat{\mathcal{L}}$	1	2	3	4	1	2	3	4	5	6	7	8	9
$\hat{\alpha}(1, l)$	1	1	1	1	1	y	x	y	1	x	y	x	1
$\hat{\alpha}(2, l)$	0	1	0	0	$-\frac{1}{12}$	y	x	y	$-\frac{1}{12}$	x	y	x	$-\frac{1}{12}$
$\hat{\alpha}(3, l)$	0	1	0	0	$-\frac{1}{12}$	z	$\frac{1}{3} - z$	z	$-\frac{1}{12}$	$\frac{1}{3} - z$	z	$\frac{1}{3} - z$	$-\frac{1}{12}$

Appendix C. \hat{j} , $\hat{\alpha}$ And $\hat{\beta}$ for \mathbb{P}_2 and \mathbb{P}_3 in two dimensions

The arrays \hat{j} and $\hat{\beta}$ for \mathbb{P}_2 and \mathbb{P}_3 polynomials in two dimensions are shown in Table 8 . The nodes are enumerated using the increasing vertex index enumeration convention.

The arrays $\hat{\alpha}$ for the \mathbb{P}_2 and \mathbb{P}_3 Lagrange bases in two dimensions are shown in Table 9. For \mathbb{P}_3 polynomials the array $\hat{\alpha}$ is defined up to a free parameter y , and we use the definitions $x := \frac{2}{3} - y$ and $z := \frac{1}{3}(\frac{10}{9} - y)$ to save space in the table. In the computation reported in the paper we take $y = \frac{7}{36}$.

The arrays $\hat{\alpha}$ for the \mathbb{P}_2 and \mathbb{P}_3 Bernstein bases in two dimensions are shown in Table 10. For \mathbb{P}_3 polynomials the array $\hat{\alpha}$ is defined up to a free parameter y and we use the definitions $x := \frac{13}{12} - y$ and $z := \frac{5}{3} - 2y$ to save space in the table. In the computation reported in the paper we take $y = \frac{2}{3}$.

References

- [1] F.S.A. Alrashed, Parallel Multiphase Navier-Stokes Solver (Ph.D. thesis –Texas A & M University), ProQuest LLC, Ann Arbor, MI, 2015. URL <https://oaktrust.library.tamu.edu/bitstream/handle/1969.1/155339/ALRASHED-DISSERTATION-2015.pdf>.
- [2] M. Quezada De Luna, High-Order Maximum Principle Preserving (MPP) Techniques for Solving Conservation Laws with Applications on Multiphase Flow (Ph.D. thesis –Texas A & M University), ProQuest LLC, Ann Arbor, MI, 2016. URL <https://oaktrust.library.tamu.edu/bitstream/handle/1969.1/158007/QUEZADADELUNA-DISSERTATION-2016.pdf>.
- [3] J.-L. Guermond, B. Popov, Invariant domains and first-order continuous finite element approximation for hyperbolic systems, SIAM J. Numer. Anal. 54 (4) (2016) 2466–2489.
- [4] R. Abgrall, Q. Viville, H. Beaugendre, C. Dobrzynski, Construction of a p -adaptive continuous residual distribution scheme, J. Sci. Comput. 72 (3) (2017) 1232–1268.
- [5] R. Sanders, A third-order accurate variation nonexpansive difference scheme for single nonlinear conservation laws, Math. Comp. 51 (184) (1988) 535–558.
- [6] X.-D. Liu, E. Tadmor, Third order nonoscillatory central scheme for hyperbolic conservation laws, Numer. Math. 79 (3) (1998) 397–425.
- [7] A. Kurganov, G. Petrova, A third-order semi-discrete genuinely multidimensional central scheme for hyperbolic conservation laws and related problems, Numer. Math. 88 (4) (2001) 683–729.
- [8] W. Pazner, Sparse invariant domain preserving discontinuous Galerkin methods with subcell convex limiting, Comput. Methods Appl. Mech. Engrg. 382 (28) (2021), 113876.
- [9] X. Zhang, C.-W. Shu, Positivity-preserving high order discontinuous Galerkin schemes for compressible Euler equations with source terms, J. Comput. Phys. 230 (4) (2011) 1238–1248.
- [10] X. Zhang, Y. Xia, C.-W. Shu, Maximum-principle-satisfying and positivity-preserving high order discontinuous Galerkin schemes for conservation laws on triangular meshes, J. Sci. Comput. 50 (1) (2012) 29–62.
- [11] R. Abgrall, Residual distribution schemes: current status and future trends, Comput. Fluids 35 (7) (2006) 641–669.
- [12] M. Ricchiuto, R. Abgrall, Explicit Runge-Kutta residual distribution schemes for time dependent problems: second order case, J. Comput. Phys. 229 (16) (2010) 5653–5691.
- [13] R. Abgrall, P. Bacigaluppi, S. Tokareva, High-order residual distribution scheme for the time-dependent Euler equations of fluid dynamics, Comput. Math. Appl. 78 (2) (2019) 274–297.
- [14] F. Vilar, R. Abgrall, A posteriori local subcell correction of high-order discontinuous galerkin scheme for conservation laws on two-dimensional unstructured grids, 2022.
- [15] J.-L. Guermond, M. Nazarov, A maximum-principle preserving C^0 finite element method for scalar conservation equations, Comput. Methods Appl. Mech. Engrg. 272 (2014) 198–213.
- [16] J.-L. Guermond, M. Nazarov, B. Popov, Y. Yang, A second-order maximum principle preserving Lagrange finite element technique for nonlinear scalar conservation equations, SIAM J. Numer. Anal. 52 (4) (2014) 2163–2182.

- [17] J.-L. Guermond, M. Nazarov, B. Popov, I. Tomas, Second-order invariant domain preserving approximation of the Euler equations using convex limiting, *SIAM J. Sci. Comput.* 40 (5) (2018) A3211–A3239.
- [18] J.-L. Guermond, B. Popov, I. Tomas, Invariant domain preserving discretization-independent schemes and convex limiting for hyperbolic systems, *Comput. Methods Appl. Mech. Engrg.* 347 (2019) 143–175.
- [19] J.P. Boris, D.L. Book, Flux-corrected transport. I. SHASTA, a fluid transport algorithm that works, *J. Comput. Phys.* 135 (2) (1997) 170–186, *J. Comput. Phys.* 11 (1) (1973) 38–69.
- [20] S.T. Zalesak, Fully multidimensional flux-corrected transport algorithms for fluids, *J. Comput. Phys.* 31 (3) (1979) 335–362.
- [21] D. Kuzmin, S. Turek, Flux correction tools for finite elements, *J. Comput. Phys.* 175 (2) (2002) 525–558.
- [22] D. Kuzmin, R. Löhner, S. Turek, Flux-corrected transport: Principles, algorithms, and applications, in: *Scientific Computation*, Springer, 2012.
- [23] R. Anderson, V. Dobrev, T. Kolev, D. Kuzmin, M. Quezada de Luna, R. Rieben, V. Tomov, High-order local maximum principle preserving (MPP) discontinuous Galerkin finite element method for the transport equation, *J. Comput. Phys.* 334 (2017) 102–124.
- [24] A. Ern, J.-L. Guermond, Invariant-domain-preserving high-order time stepping: I Explicit Runge-Kutta schemes, *SIAM J. Sci. Comput.* 44 (5) (2022) A3366–A3392.
- [25] D. Kuzmin, M. Quezada de Luna, Subcell flux limiting for high-order Bernstein finite element discretizations of scalar hyperbolic conservation laws, *J. Comput. Phys.* 411 (19) (2020), 109411.
- [26] A. Jameson, W. Schmidt, E. Turkel, Numerical solution of the Euler equations by finite volume methods using Runge-Kutta time-stepping schemes, in: *14th AIAA Fluid and Plasma Dynamics Conference*, 1981, AIAA Paper 1981–1259.
- [27] A. Jameson, Origins and further development of the Jameson-Schmidt-Turkel scheme, *AIAA J.* 55 (5) (2017).
- [28] E. Burman, On nonlinear artificial viscosity, discrete maximum principle and hyperbolic conservation laws, *BIT* 47 (4) (2007) 715–733.
- [29] J.-L. Guermond, B. Popov, Invariant domains and second-order continuous finite element approximation for scalar conservation equations, *SIAM J. Numer. Anal.* 55 (6) (2017) 3120–3146.
- [30] P.D. Lax, Weak solutions of nonlinear hyperbolic equations and their numerical computation, *Comm. Pure Appl. Math.* 7 (1954) 159–193.
- [31] A. Harten, P.D. Lax, B. van Leer, On upstream differencing and Godunov-type schemes for hyperbolic conservation laws, *SIAM Rev.* 25 (1) (1983) 35–61.
- [32] E. Tadmor, Numerical viscosity and the entropy condition for conservative difference schemes, *Math. Comp.* 43 (168) (1984) 369–381.
- [33] B. Perthame, C.-W. Shu, On positivity preserving finite volume schemes for Euler equations, *Numer. Math.* 73 (1) (1996) 119–130.
- [34] R. Abgrall, Some remarks about conservation for residual distribution schemes, *Comput. Methods Appl. Math.* ISSN: 1609-4840 18 (3) (2018) 327–351.
- [35] A. Logg, K.-A. Mardal, G.N. Wells, et al., *Automated Solution of Differential Equations By the Finite Element Method*, Springer, 2012.
- [36] J.F.B.M. Kraaijevanger, Contractivity of Runge-Kutta methods, *BIT* 31 (3) (1991) 482–528.
- [37] P.-O. Persson, J. Peraire, Sub-cell shock capturing for discontinuous galerkin methods, in: *44th AIAA Aerospace Sciences Meeting and Exhibit*, Number AIAA Paper No. 2015-2006-112 in *Aerospace Sciences Meetings*, 2006.
- [38] B. Khobalatte, B. Perthame, Maximum principle on the entropy and second-order kinetic schemes, *Math. Comp.* 62 (205) (1994) 119–131.
- [39] A. Harten, On the symmetric form of systems of conservation laws with entropy, *J. Comput. Phys.* 49 (1) (1983) 151–164.
- [40] A. Harten, S. Osher, Uniformly high-order accurate nonoscillatory schemes. I, *SIAM J. Numer. Anal.* 24 (2) (1987) 279–309.
- [41] X. Zhang, C.-W. Shu, On maximum-principle-satisfying high order schemes for scalar conservation laws, *J. Comput. Phys.* 229 (9) (2010) 3091–3120.
- [42] C.M. Dafermos, in: *Hyperbolic Conservation Laws in Continuum Physics*, third ed. *Grundlehren der mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*, vol. 325, Springer-Verlag, Berlin, 2010.
- [43] D. Hoff, The sharp form of Oleĭnik’s entropy condition in several space variables, *Trans. Amer. Math. Soc.* 276 (2) (1983) 707–714.
- [44] A. Kurganov, G. Petrova, B. Popov, Adaptive semidiscrete central-upwind schemes for nonconvex hyperbolic conservation laws, *SIAM J. Sci. Comput.* 29 (6) (2007) 2381–2401.
- [45] A. Ern, J.-L. Guermond, in: *Finite Elements I—Approximation and Interpolation*, *Texts in Applied Mathematics*, vol. 72, Springer, Cham, 2021.