

Mass conservative limiting and applications to the approximation of the steady-state radiation transport equations

Jean-Luc Guermond^{a,*,1}, Zuodong Wang^{b,c}

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

^b CERMICS, Ecole des Ponts, 77455 Marne-la-Vallee Cedex 2, France

^c SERENA, Centre Inria de Paris, 48 rue Barrault, 75647 Paris, France

ARTICLE INFO

MSC:

35L65

65M12

65M60

76V05

Keywords:

Limiting

Advection equation

Radiation transport equation

Stiff sources

Conservation equations

Asymptotic preserving

Invariant domains

ABSTRACT

A limiting technique for scalar transport equations is presented. The originality of the method is that it does not require solving nonlinear optimization problems nor does it rely on the construction of a low-order approximation. The method has minimal complexity and is numerically demonstrated to maintain high-order accuracy. The performance of the method is illustrated on the radiation transport equation.

1. Introduction

The objective of the paper is to present two simple limiting techniques for scalar-valued partial differential equations with a structure like the radiation transport equation. The first limiting method is iterative, locally mass conservative, and does not involve solving any nonlinear optimization problem. At convergence, this method enforces the local minimum/maximum principle (assuming that local bounds are known). Two iterations of the method are in general sufficient, but although the method is observed to converge quickly, there is no guarantee that the bounds are enforced everywhere after a number of iterations that is independent of the meshsize. We then propose a second limiting method that is globally mass conservative and that also does not involve solving any nonlinear optimization problem. The second limiting is applied after the first one. The purpose of this second post-processing is simply to certify that some global bound like positivity is exactly and unconditionally enforced while preserving the total mass of the solution. The combination of these two limiting techniques is illustrated using continuous finite elements stabilized with the continuous interior penalty (CIP) technique (a.k.a. edge stabilization) from Douglas and Dupont [7] and Burman and Hansbo [5]. The method presented

* Corresponding author.

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

¹ This material is based upon work supported in part by the National Science Foundation grant DMS2110868, the Air Force Office of Scientific Research, USAF, under grant/contract number FA9550-18-1-0397, the Army Research Office, under grant number W911NF-19-1-0431, and the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contracts B640889. The support of INRIA through the International Chair program is acknowledged.

<https://doi.org/10.1016/j.jcp.2024.113531>

Received 29 April 2024; Received in revised form 4 October 2024; Accepted 22 October 2024

Available online 28 October 2024

0021-9991/© 2024 Elsevier Inc. All rights are reserved, including those for text and data mining, AI training, and similar technologies.

in the paper is not restricted to continuous elements and CIP though. It can be used with other spatial discretization and other types of stabilization as well. For instance, one can use discontinuous elements (of degree $p \geq 1$) stabilized with the upwind numerical flux. One can also use continuous elements stabilized with other methods like Galerkin Least-Squares, Local Projection Stabilization, Orthogonal Subscale Stabilization, and Subgrid Viscosity.

The original motivation for the work presented here is the solution to the radiation transport equation. Based on our experience on nonlinear hyperbolic systems, we have tried for many years to use upwinding, artificial viscosity, and nonlinear variations thereof to enforce positivity of the angular intensity. But it is well established in the literature that one needs to be careful with this type of method in the context of the radiation transport equations. For instance, it is shown in Adams [1] (see also Larsen [14], Larsen et al. [15], and [9, Rmk. 5.2], [11, §3.2]), that upwinding and artificial viscosity may lock in the diffusion limit if the artificial viscosity is strong enough so as to guarantee that the local minimum/maximum principle holds. We have long tried to modify the artificial viscosity techniques to overcome this difficulty. For instance in [11] we have adopted a technique inspired by Gosse and Toscani [8] consisting of scaling the transport term, the scattering term, and the artificial diffusion by $\frac{1}{1+\sigma^s h}$ where h is the mesh size and σ^s is the scattering cross section. This gives a method that is indeed positive, convergent, and asymptotic preserving, but it is not locally conservative on nonuniform meshes. Moreover, it does not behave properly with grazing incidences in the diffusion regime unless the grazing boundary data is replaced by its angular average weighted with Chandrasekhar's H-function, thereby seriously diminishing the usefulness of the method as estimating the H-function is nontrivial. In the paper we propose instead to rely on the Galerkin approximation augmented with a traditional linear stabilization technique and to enforce the local minimum/maximum principle by using a mass conservative post-processing technique.

An important aspect of the limiting technique we propose is the estimation of local bounds. The approach we use for this purpose consists of using the method of characteristics as in Yee et al. [28, Eq. (16)]. Note that here we do not propose to use the method of characteristics as a solution method but just as a means to estimate guaranteed local upper and lower bounds. We refer to Lathrop [16], Sanchez and McCormick [25, §II.C.2] for early reviews on solution methods based on the method of characteristics. Once local bounds are found, the next problem consists of enforcing these bounds. This can be done as in Maginot et al. [19] where the authors develop a nonlinear characteristics-based methods that is positivity-preserving. Another idea is to invoke strategies like the so-called flux corrected transport method of Zalesak combining in a nonlinear fashion a low-order and a high-order approximation [29]. The problem with this approach is that it relies on a low-order approximation which itself relies in one way or another on artificial viscosity, which as explained above leads to locking. Another possibility more aligned with what we propose in the paper is to enforce local bounds by using nonlinear optimization techniques. The problem is then to find a maximum-principle-satisfying object that minimizes some distance to the maximum-principle-violating solution while maintaining constant the total mass (either locally or globally). This type of method is developed in e.g., Bochev et al. [3,4], Yee et al. [28] Peterson et al. [24], and references therein. A fast convergent algorithm for finding such a minimizer is proposed in Liu et al. [18], Ancellin et al. [2]. Here we propose a slightly different approach consisting of just computing a quasi-minimizer with a set of two methods whose algorithmic complexity is significantly lower than that of computing a genuine minimizer. We finally mention that if positivity is the only property one is interested in, one can use “fix-up” techniques like that presented in Lewis and Miller [17, Chap. 4]. This iterative method sets negative fluxes to zero and continues the iterations with the other degree of freedom to restore balance. This idea is generalized to the discontinuous Galerkin setting in Hamilton et al. [12]. Other improvements of the fix-up technique are presented in Maginot et al. [20].

The paper is organized as follows. We introduce in §2 the two limiting techniques mentioned above. We choose to present them without invoking the radiation transport equation as these two methods are quite general and can be used for other purposes. We recall in §3 how local bounds can be extracted for the scalar-valued linear transport equations using the method of characteristics. We also introduce in this section a second-order relaxation technique that is essential to achieve accuracy beyond second-order. We present in §4 the Galerkin approximation of the radiation transport equation stabilized with the continuous interior penalty technique. The full solution method is explained in Algorithms A.1 and A.2. The proposed method is illustrated on the scalar linear transport equation in §5 and on the radiation transport equation in §6.

2. Iterative limiting algorithm

We present in this section two simple conservative iterative limiting algorithms. The first one preserves mass locally. The second one only preserves mass globally, and is only used as theoretical safeguard. These two algorithms are quite general. They are not specific to the radiation transport equations and can be used in many different contexts.

2.1. The setting

Let $\{x_i\}_{i \in \mathcal{V}}$ be degrees of freedom that have to be limited, $x_i \in \mathbb{R}$ for all i in the index set \mathcal{V} . Let $\mathcal{M} := \sum_{i \in \mathcal{V}} m_i x_i$ be the total mass, where the coefficients m_i are nonnegative and are called mass at i . For all $i \in \mathcal{V}$, let u_i^{\min} and u_i^{\max} be the local minimum and maximum we want to enforce on the i -th degree of freedom. It is henceforth assumed that these local bounds u_i^{\min} and u_i^{\max} satisfy the following reasonable estimate:

$$\sum_{i \in \mathcal{V}} m_i u_i^{\min} \leq \mathcal{M} \leq \sum_{i \in \mathcal{V}} m_i u_i^{\max}. \tag{2.1}$$

Our objective is to post-process the degrees of freedom $\{x_i\}_{i \in \mathcal{V}}$ to enforce the local bounds while maintaining mass conservation, either locally or globally.

We assume that the lumped mass is non-negative, i.e., $m_i \geq 0$ for all $i \in \mathcal{V}$. This assumption holds, for example, when using finite elements with Bernstein basis of any polynomial degree, and for Lagrange bases of degree 1, 2 and 3 with equally distributed interpolation nodes. We introduce a concept of locality by introducing a notion of stencil. For all $i \in \mathcal{V}$ we assume that there exists a collection of indices $\mathcal{I}(i) \subsetneq \mathcal{V}$ that can exchange mass with i . We call this index set stencil. For instance, for finite elements we have $j \in \mathcal{I}(i)$ if $\varphi_j \varphi_i \neq 0$, where φ_i and φ_j are the global shape function associated with the degrees of freedom i and j . We also define $\mathcal{I}(i)^* := \mathcal{I}(i) \setminus \{i\}$. For all the finite element tests reported in §5 and §6, the set $\mathcal{I}(i)$ is the standard stencil, i.e., $\mathcal{I}(i)$ collects the indices of all the global shape functions φ_j that are such that $\varphi_j \varphi_i \neq 0$.

2.2. Local conservative limiting

We describe here an iterative algorithm that is locally conservative. The algorithm consists of looping over all the indices i in \mathcal{V} and proceeds as follows. Let $i \in \mathcal{V}$. If $m_i = 0$, then the i -th degree of freedom does not contribute to the global mass; we simply set $y_i = \min(u_i^{\max}, \max(u_i^{\min}, x_i))$ and continue to the next degree of freedom in the list \mathcal{V} . This possible modification of the value at i does not change the mass, either locally or globally. Let us assume now that $m_i > 0$. Then either $x_i \leq u_i^{\max}$ (i.e., the maximum principle is satisfied) or $u_i^{\max} < x_i$ (i.e., the maximum principle is violated). If $x_i \leq u_i^{\max}$, we do nothing and continue to the next index in the loop. If $u_i^{\max} < x_i$, then we compute

$$a_i^+ := \sum_{j \in \mathcal{I}(i)^*} m_j \max(0, u_j^{\max} - x_j), \quad b_i^+ := \max(x_i - \frac{a_i^+}{m_i}, u_i^{\max}), \tag{2.2a}$$

$$\ell_i^+ := \begin{cases} 0 & \text{if } 0 = a_i^+ \\ m_i \frac{x_i - b_i^+}{a_i^+} & \text{otherwise,} \end{cases} \tag{2.2b}$$

and the actual limiting with respect to u_i^{\max} is done by setting

$$y_j := x_j + \ell_i^+ \max(0, u_j^{\max} - x_j), \quad \forall j \in \mathcal{I}(i)^* \tag{2.3a}$$

$$y_i := b_i^+. \tag{2.3b}$$

Lemma 2.1. *Let $i \in \mathcal{V}$. If $u_i^{\max} < x_i$, then the following holds for all $\{y_j\}_{j \in \mathcal{I}(i)}$ given by (2.3):*

- (i) *There is local mass conservation, i.e., $m_i(y_i - x_i) + \sum_{j \in \mathcal{I}(i)^*} m_j(y_j - x_j) = 0$.*
- (ii) *$x_j \leq y_j \leq \max(x_j, u_j^{\max})$ for all $j \in \mathcal{I}(i)^*$.*
- (iii) *$u_i^{\max} \leq y_i \leq x_i$. Moreover, $y_i < x_i$ if $0 < a_i^+$.*

Proof. (i) Assume first that $a_i^+ \neq 0$. The equation (2.3a) implies that the mass transferred to the node $j \in \mathcal{I}(i)^*$ is $m_j \ell_i^+ \max(0, u_j^{\max} - x_j)$, while the equation (2.3b) implies that the mass transferred to the node i is $m_i(b_i^+ - x_i)$. The total mass exchange is then

$$\begin{aligned} \Delta_i &:= m_i(y_i - x_i) + \sum_{j \in \mathcal{I}(i)^*} m_j(y_j - x_j) \\ &= m_i(b_i^+ - x_i) + \ell_i^+ \sum_{j \in \mathcal{I}(i)^*} m_j \max(0, u_j^{\max} - x_j) = m_i(b_i^+ - x_i) + \ell_i^+ a_i^+. \end{aligned}$$

But since $a_i^+ \neq 0$, (2.2b) gives $\ell_i^+ = m_i \frac{x_i - b_i^+}{a_i^+}$; hence

$$\Delta_i = m_i(b_i^+ - x_i) + m_i \frac{x_i - b_i^+}{a_i^+} a_i^+ = 0.$$

Otherwise, if $a_i^+ = 0$, then $b_i^+ = \max(x_i, u_i^{\max}) = x_i$ and $\ell_i^+ := 0$. Then again $\Delta_i = 0$.

(ii) The assertion is trivial when $a_i^+ = 0$. Let assume now that $a_i^+ > 0$. As $x_i - \frac{a_i^+}{m_i} \leq b_i^+$, we have $\ell_i^+ := m_i \frac{x_i - b_i^+}{a_i^+} \leq 1$. Moreover, the assumptions $u_i^{\max} < x_i$ and $b_i^+ := \max(x_i - \frac{a_i^+}{m_i}, u_i^{\max})$, imply that $b_i^+ \leq \max(x_i, u_i^{\max}) < x_i$; hence, $0 \leq \ell_i^+ < 1$. As $0 \leq \ell_i^+ \leq 1$, we infer that $x_j \leq x_j + \ell_i^+ \max(0, u_j^{\max} - x_j) := y_j \leq x_j + \max(0, u_j^{\max} - x_j) = \max(x_j, u_j^{\max})$.

(iii) By (2.2a), we have $u_i^{\max} \leq b_i^+$ and $b_i^+ = \max(x_i - \frac{a_i^+}{m_i}, u_i^{\max}) \leq \max(x_i, u_i^{\max}) = x_i$; hence, $u_i^{\max} \leq y_i := b_i^+ \leq x_i$. Let us now assume that $a_i^+ \neq 0$. Notice that if $b_i^+ = u_i^{\max}$ then $y_i = b_i^+ < x_i$ (whether a_i^+ is equal to zero or not does not matter here). On the other hand, if $b_i^+ = x_i - \frac{a_i^+}{m_i}$, then the assumption $a_i^+ \neq 0$ implies that $y_i = b_i^+ < x_i$. \square

The same idea as above can be used to enforce the minimum principle. Let $i \in \mathcal{V}$. Either $u_i^{\min} \leq x_i$ (i.e., the local minimum principle is satisfied) or $x_i < u_i^{\min}$ (i.e., the local minimum principle is violated). If $u_i^{\min} \leq x_i$, do nothing and continue to the next index in the list \mathcal{V} . If $x_i < u_i^{\min}$, then compute

$$a_i^- := \sum_{j \in \mathcal{I}(i)^*} m_j \max(0, x_j - u_j^{\min}), \quad b_i^- := \min(x_i + \frac{a_i^-}{m_i}, u_i^{\min}), \tag{2.4a}$$

$$\ell_i^- := \begin{cases} 0 & \text{if } 0 = a_i^- \\ m_i \frac{x_i - b_i^-}{a_i^-} & \text{otherwise,} \end{cases} \tag{2.4b}$$

and set

$$y_j := x_j + \ell_i^- \max(0, x_j - u_j^{\min}), \quad \forall j \in \mathcal{I}(i) \setminus \{i\} \tag{2.5a}$$

$$y_i := b_i^-. \tag{2.5b}$$

Lemma 2.2. *Let $i \in \mathcal{V}$. If $x_i < u_i^{\min}$, then the following holds for all $\{y_j\}_{j \in \mathcal{I}(i)}$ given by (2.5):*

- (i) $m_i(y_i - x_i) + \sum_{j \in \mathcal{I}(i)^*} m_j(y_j - x_j) = 0$, i.e., (2.5) is locally mass conservative
- (ii) $\min(x_j, u_j^{\min}) \leq y_j \leq x_j$ for all $j \in \mathcal{I}(i)^*$.
- (iii) $x_i \leq y_i \leq u_i^{\min}$. Moreover, $x_i < y_i$ if $0 < a_i^-$.

The algorithm (2.2)–(2.5) guarantees that the local maximum decreases and the local minimum increases until reaching the prescribed values (see Items (iii) and (iii) in Lemmas 2.1 and 2.2). The statements made in Items (ii) and (ii) in Lemmas 2.1 and 2.2 guarantee that by correcting $x_i \rightarrow y_i$, the neighboring values that are already in bounds stay in bounds after limiting. The algorithm is iterative, but there is no guarantee that the prescribed bounds are reached in a finite number of iterations. Extensive numerical tests show though that convergence is quick when one starts from a reasonable solution (e.g., a linearly stabilized Galerkin approximation). In all the tests reported in the paper the above algorithm is only applied two times in a row. The limiting algorithm is summarized in Algorithm 2.1. This algorithm is based on Jacobi iterations. One can also reformulate the algorithm by replacing the Jacobi iteration by a Gauss-Seidel one. We have not noticed significant differences between the two algorithms in our Numerical experiments.

Algorithm 2.1 Local conservative limiting (2.2)–(2.5) using Jacobi iterations.

```

Require: Bounds  $\{u_i^{\max}, u_i^{\min}\}_{i \in \mathcal{V}}$ ,  $\{x_i\}_{i \in \mathcal{V}}$ .
1: for  $i \in \mathcal{V}$  do ▷ Loop over dofs
2:   if  $m_i = 0$  then
3:      $y_i = \min(u_i^{\max}, \max(u_i^{\min}, x_i))$ .
4:   else if  $u_i^{\max} < x_i$  then ▷ Local maximum principle violated
5:     Compute  $a_i^+$ ,  $b_i^+$ ,  $\ell_i^+$  using (2.2).
6:      $y_i = b_i^+$ 
7:     for  $j \in \mathcal{I}(i)^*$  do
8:        $y_j = x_j + \ell_i^+ \max(0, u_j^{\max} - x_j)$ 
9:     end for
10:  else if  $x_i < u_i^{\min}$  then ▷ Local minimum principle violated
11:    Compute  $a_i^-$ ,  $b_i^-$ ,  $\ell_i^-$  using (2.4).
12:     $y_i = b_i^-$ 
13:    for  $j \in \mathcal{I}(i)^*$  do
14:       $y_j = x_j + \ell_i^- \max(0, x_j - u_j^{\min})$ 
15:    end for
16:  else ▷ Do nothing
17:     $y_i = x_i$ 
18:  end if
19: end for
20: return  $\{y_i\}_{i \in \mathcal{V}}$ 

```

2.3. Global conservative limiting

As the iterative algorithm (2.2)–(2.5) does not guarantee that the prescribed bounds are achieved in a finite number of iterations, we now propose a final conservative post-processing that can be used to make sure that global bounds that are essential to the physics are strictly enforced.

Consider the set of degrees of freedom $\{x_i\}_{i \in \mathcal{V}}$, the associated mass $\mathcal{M} := \sum_{i \in \mathcal{V}} m_i x_i$, and the local bounds $\{u_i^{\min}, u_i^{\max}\}_{i \in \mathcal{V}}$, which we recall are assumed to satisfy (2.1). We compute a new set of limited values $\{z_i\}_{i \in \mathcal{V}}$ as follows:

For all $i \in \mathcal{V}$, compute,
$$y_i := \min((\max(x_i, u_i^{\min}), u_i^{\max}). \tag{2.6a}$$

Then compute
$$\alpha^+ := \max\left(0, \frac{\mathcal{M} - \sum_{j \in \mathcal{V}} m_j y_j}{\sum_{j \in \mathcal{V}} m_j (u_j^{\max} - y_j)}\right), \tag{2.6b}$$

and
$$\alpha^- := \max\left(0, \frac{\mathcal{M} - \sum_{j \in \mathcal{V}} m_j y_j}{\sum_{j \in \mathcal{V}} m_j (u_j^{\min} - y_j)}\right). \tag{2.6c}$$

For all $i \in \mathcal{V}$, set,
$$z_i := y_i + \alpha^+(u_i^{\max} - y_i) + \alpha^-(u_i^{\min} - y_i). \tag{2.6d}$$

Lemma 2.3. *The following holds for all $i \in \mathcal{V}$:*

$$u_i^{\min} \leq z_i \leq u_i^{\max}, \tag{2.7}$$

$$\sum_{i \in \mathcal{V}} m_i z_i = \sum_{i \in \mathcal{V}} m_i x_i. \tag{2.8}$$

Proof. (1) By definition $0 \leq \alpha^+$. If $\mathcal{M} - \sum_{i \in \mathcal{V}} m_i y_i \leq 0$ then $\alpha^+ = 0$. If $\mathcal{M} - \sum_{i \in \mathcal{V}} m_i y_i \geq 0$, then using the assumption (2.1) together with $\sum_{i \in \mathcal{V}} m_i (u_i^{\max} - y_i) \geq 0$ (which holds owing to (2.6a)), we infer that

$$\alpha^+ \leq \frac{\mathcal{M} - \sum_{i \in \mathcal{V}} m_i y_i}{\sum_{i \in \mathcal{V}} m_i (u_i^{\max} - y_i)} \leq \frac{\sum_{i \in \mathcal{V}} m_i u_i^{\max} - \sum_{i \in \mathcal{V}} m_i y_i}{\sum_{i \in \mathcal{V}} m_i (u_i^{\max} - y_i)} = 1.$$

Hence $0 \leq \alpha^+ \leq 1$. We proceed similarly to show that $0 \leq \alpha^- \leq 1$.

(2) Observing that $u_i^{\min} - y_i \leq 0, 0 \leq u_i^{\max} - y_i$ (which holds owing to (2.6a)), $0 \leq \alpha^+ \leq 1$, and $0 \leq \alpha^-$, we obtain

$$\begin{aligned} z_i &:= y_i + \alpha^+(u_i^{\max} - y_i) + \alpha^-(u_i^{\min} - y_i) \leq y_i + \alpha^+(u_i^{\max} - y_i) \\ &\leq y_i + (u_i^{\max} - y_i) = u_i^{\max}. \end{aligned}$$

We proceed similarly to prove that $u_i^{\min} \leq z_i$.

(3) Since $0 \leq \sum_{i \in \mathcal{V}} m_i (u_i^{\max} - y_i)$ and $\sum_{i \in \mathcal{V}} m_i (u_i^{\min} - y_i) \leq 0$, we have

$$\begin{aligned} \sum_{i \in \mathcal{V}} m_i z_i &= \sum_{i \in \mathcal{V}} m_i y_i + \alpha^+ \sum_{i \in \mathcal{V}} m_i (u_i^{\max} - y_i) + \alpha^- \sum_{i \in \mathcal{V}} m_i (u_i^{\min} - y_i) \\ &= \sum_{i \in \mathcal{V}} m_i y_i + \max(0, \mathcal{M} - \sum_{i \in \mathcal{V}} m_i y_i) + \min(0, \mathcal{M} - \sum_{i \in \mathcal{V}} m_i y_i) = \mathcal{M}. \end{aligned}$$

This completes the proof. \square

Remark 2.4 (Local vs. global bounds). Although the bounds u_i^{\min} and u_i^{\max} invoked in the algorithm can in principle be local, we insist that the algorithm (2.6) should only be used to make sure that global bounds are enforced because the mass conservation mechanism is not local. Notice that the local and global limiting algorithms can be used to limit the solution of time-dependent nonlinear conservation equations as well. Preliminary tests (not shown here for brevity) reveal that these two algorithms perform very well when combined as described above. Preliminary tests, (not shown here for brevity) also show that the global limiting algorithm should not be used alone. When used alone, the algorithm may make the approximation to converge to a solution that does not satisfy the Rankin-Hugoniot condition. The purpose of the global limiting is purely theoretical. Its purpose is simply to certify that some global bound (like positivity) is exactly and unconditionally enforced while preserving the total mass of the solution, but the bulk of the job is done by the local limiting algorithm.

3. Bounds for the transport equation

A critical aspect of the algorithm (2.2)–(2.5) is the estimation of the local bounds $\{u_i^{\min}, u_i^{\max}\}_{i \in \mathcal{V}}$. We explain in this section how this can be done for scalar transport equations. To simplify the presentation, we assume that the degrees of freedoms enumerated with the list \mathcal{V} are Lagrange finite elements associated with the Lagrange nodes $\{\mathbf{x}_i\}_{i \in \mathcal{V}}$ from a mesh \mathcal{T}_h based on a domain $D \subset \mathbb{R}^d$, $d \in \{1, 2, 3\}$. In the rest of this section, the symbol x is a position vector in one dimension and \mathbf{x} is a position vector in two and higher dimensions.

3.1. Computations of the bounds

We start by explaining the method we use to compute bounds in one space dimension. We then generalize the method to high space dimensions later. No originality is claimed here as most of what is said in this section can be found in the radiation transport literature, see e.g., Lathrop [16].

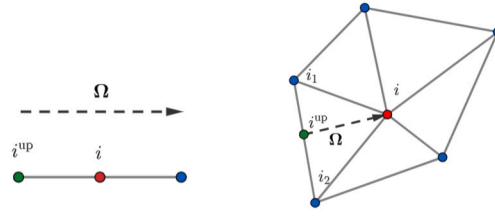


Fig. 3.1. Examples for inflow nodes: left panel: 1D case; right panel: 2D case. (For interpretation of the colors in the figure(s), the reader is referred to the web version of this article.)

3.1.1. One-dimensional case

Let us consider the one-dimensional transport equation

$$\Omega(x)\partial_x u(x) + \sigma(x)u(x) = q(x). \tag{3.1}$$

We assume to simplify the presentation that both Ω and σ are piecewise constant over the mesh cells. We also assume that Ω is not equal to 0; otherwise, (3.1) is trivial and there is nothing to limit. For each $i \in \mathcal{V}$, let x_i be one of the Lagrangian node in the mesh where one wants to estimate an upper and a lower bound. Let K be the (unique) cell in the mesh \mathcal{T}_h such that the segment $\{x_i - s\Omega(s) \mid s \in \mathbb{R} > 0\} \cap K$ is not empty. This cell exists if x_i is not on the inflow boundary (if x_i is on the inflow boundary one can set $u_i^{\min} = u_i^{\max} = \alpha^\partial(x_i)$ where α^∂ is the inflow boundary data of the problem). Then one defines $x_i^{\text{up}} := \partial K \cap \{x_i - s\Omega(s) \mid s \in \mathbb{R} > 0\}$. The point x_i^{up} is the farthest away from x_i in K along the direction $-\Omega$. This point is at the upwind boundary of K (see the left panel in Fig. 3.1 for example). Then, setting $\Omega := \Omega|_K$ and $\sigma := \sigma|_K$, the exact solution to (3.1) is such that

$$u(x_i) = u(x_i^{\text{up}})e^{\frac{\sigma}{\Omega}(x_i^{\text{up}} - x_i)} + \int_{x_i^{\text{up}}}^{x_i} \frac{q(s)}{\Omega} e^{\frac{\sigma}{\Omega}s} ds. \tag{3.2}$$

Setting $q^{\min} := \min_{x \in K} q(x)$, $q^{\max} := \max_{x \in K} q(x)$, this gives $u_i^{\min} \leq u(x_i) \leq u_i^{\max}$ where

$$u_i^{\min} := u(x_i^{\text{up}})e^{\frac{\sigma}{\Omega}(x_i^{\text{up}} - x_i)} + \frac{q^{\min}}{\sigma}(1 - e^{\frac{\sigma}{\Omega}(x_i^{\text{up}} - x_i)}), \tag{3.3a}$$

$$u_i^{\max} := u(x_i^{\text{up}})e^{\frac{\sigma}{\Omega}(x_i^{\text{up}} - x_i)} + \frac{q^{\max}}{\sigma}(1 - e^{\frac{\sigma}{\Omega}(x_i^{\text{up}} - x_i)}). \tag{3.3b}$$

As these expressions are not robust with respect to σ (in particular they do not make sense for $\sigma = 0$), we instead define $\Delta x := \left| \frac{x_i^{\text{up}} - x_i}{\Omega} \right|$, and use the following bounds motivated by Taylor expansion, when $\frac{\sigma}{\Omega}(x_i - x_i^{\text{up}}) \leq 0.005$:

$$u_i^{\min} := u(x_i^{\text{up}})e^{\frac{\sigma}{\Omega}(x_i^{\text{up}} - x_i)} + q^{\min} \left(\Delta x - \sigma \frac{\Delta x^2}{2} + \sigma^2 \frac{\Delta x^3}{6} - \sigma^3 \frac{\Delta x^4}{24} \right), \tag{3.4a}$$

$$u_i^{\max} := u(x_i^{\text{up}})e^{\frac{\sigma}{\Omega}(x_i^{\text{up}} - x_i)} + q^{\max} \left(\Delta x - \sigma \frac{\Delta x^2}{2} + \sigma^2 \frac{\Delta x^3}{6} \right). \tag{3.4b}$$

So assuming that $u(x_i^{\text{up}})$ is known, or one has access to some good approximation thereof by some iteration process, the bounds (3.3) (or (3.4) if $\sigma\Delta x$ is small) give an estimate of the local bounds of the solution at x_i .

3.1.2. Two-dimensional case

The high-dimensional case is similar to the one-dimensional one. We consider $\mathcal{T}(i)$ the collection of the cells containing \mathbf{x}_i and we set $\mathbf{x}_i^{\text{up}} := \bigcap_{K \in \mathcal{T}(i)} \partial K \cap \{\mathbf{x}_i - s\Omega(s) \mid s \in \mathbb{R} > 0\}$. The point \mathbf{x}_i^{up} (green one) is the farthest away from \mathbf{x}_i (red node) in $\mathcal{T}(i)$ along the direction $-\Omega$ (see the right panel in Fig. 3.1 for example). One key difference with the one-dimensional case is that \mathbf{x}_i^{up} is not necessarily a Lagrangian node. Then if one only knows u at the vertices of the triangulation, then one needs to reconstruct u at \mathbf{x}_i^{up} . For instance, in two space dimensions one can proceed as follows:

1. Compute $\mathbf{x}_i^{\text{up}} := \bigcap_{K \in \mathcal{T}(i)} \partial K \cap \{\mathbf{x}_i - s\Omega(s) \mid s \in \mathbb{R} > 0\}$.
2. Define $u(\mathbf{x}_i^{\text{up}}) := u(\mathbf{x}_{i_1}) + \theta(u(\mathbf{x}_{i_2}) - u(\mathbf{x}_{i_1}))$, where $\theta = \frac{\|(\mathbf{x}_i - \mathbf{x}_{i_1}) \times \Omega\|_{\ell^2}}{\|(\mathbf{x}_{i_2} - \mathbf{x}_{i_1}) \times \Omega\|_{\ell^2}}$.

3.2. Relaxation

As for any local limiting technique that we know of, second-order relaxation of the bounds must be applied to avoid order reduction. We refer for instance to [10, §4.7] where this question is discussed at length. Another way this issue is addressed in the finite volume literature consists of relaxing the slope reconstructions; see Harten and Osher [13], Schmidtman et al. [26, §2.1]. Here

we adopt the methodology proposed in [10, §4.7]. Let us assume that we have (approximate) knowledge of u at the Lagrange node $\{\mathbf{x}_i\}_{i \in \mathcal{V}}$, say $u(\mathbf{x}) \approx \sum_{i \in \mathcal{V}} u_i \varphi_i(\mathbf{x})$ where $\{\varphi_i\}_{i \in \mathcal{V}}$ are the global Lagrange shape functions. We estimate the local curvature of u by

$$\Delta_i^2 := \frac{\sum_{j \in \mathcal{I}(i)^*} \beta_{ij} (u_i - u_j)}{\sum_{j \in \mathcal{I}(i)^*} \beta_{ij}}, \quad \text{and set} \quad \overline{\Delta_i^2} := \min \text{mod} \{ \Delta_j^2 \}_{j \in \mathcal{I}(i)}, \quad (3.5)$$

where $\beta_{ij} = \int_D \nabla \varphi_j \cdot \nabla \varphi_i \, dx$ are the stiffness coefficients of the Laplace operator. Observe that $\sum_{j \in \mathcal{I}(i)^*} \beta_{ij} = -\beta_{ii} = -\int_D (\nabla \varphi_i)^2 \, dx \neq 0$. The relaxed local bounds are then defined as follows:

$$\overline{u_i^{\min}} := \max(u_i^{\min} - \overline{\Delta_i^2}, u^{\min}), \quad \overline{u_i^{\max}} := \min(u_i^{\max} + \overline{\Delta_i^2}, u^{\max}), \quad (3.6)$$

where u^{\min} and u^{\max} are global bounds, if known (see Algorithm A.3). Relaxation is essential to recovers optimal convergence rates, see numerical tests in §5.

4. Approximation of the radiation transport equations

We are going to illustrate the local mass conserving limiting algorithm (2.2)–(2.5) and the global mass conserving limiting algorithm (2.6) on the radiation transport equation. We introduce here the radiation transport equation and the associated finite element approximation.

4.1. The model problem

The computational domain D is assumed to be an open, bounded, connected polyhedron in \mathbb{R}^3 . The boundary of D is denoted by ∂D . The symbol \mathbf{n} denotes the outer unit normal on ∂D . The unit sphere in \mathbb{R}^3 is denoted by S . The surface of the unit sphere is denoted $|S|$; recall that $|S| = 4\pi$. We set $\mathcal{O} := \partial D \times S$, and to define the boundary conditions we introduce the inflow boundary $\mathcal{O}_- := \{(\mathbf{x}, \boldsymbol{\Omega}) \in \mathcal{O} \mid \boldsymbol{\Omega} \cdot \mathbf{n}(\mathbf{x}) < 0\}$.

Given a non-negative source term $q : D \times S \rightarrow \mathbb{R}_+$, and a non-negative boundary data $\alpha^\partial : \mathcal{O}_- \rightarrow \mathbb{R}_+$, we look for $\Psi : D \times S \rightarrow \mathbb{R}_+$ so that

$$\boldsymbol{\Omega} \cdot \nabla \Psi(\mathbf{x}, \boldsymbol{\Omega}) + \sigma^t(\mathbf{x}) \Psi(\mathbf{x}, \boldsymbol{\Omega}) = \sigma^s(\mathbf{x}) \overline{\Psi}(\mathbf{x}) + q(\mathbf{x}, \boldsymbol{\Omega}), \quad \text{in } D \times S \quad (4.1a)$$

$$\Psi(\mathbf{x}, \boldsymbol{\Omega}) = \alpha^\partial(\mathbf{x}, \boldsymbol{\Omega}), \quad \text{in } \mathcal{O}_- \quad (4.1b)$$

$$\overline{\Psi}(\mathbf{x}) := \frac{1}{|S|} \int_S \Psi(\mathbf{x}, \boldsymbol{\Omega}) \, d\boldsymbol{\Omega}, \quad \text{in } D, \quad (4.1c)$$

The dependent variable $\Psi(\mathbf{x}, \boldsymbol{\Omega})$ is called angular intensity or angular flux, and $\overline{\Psi}(\mathbf{x})$ is called scalar intensity or scalar flux. The coefficient $\sigma^s : D \rightarrow \mathbb{R}_+$ is the scattering cross section and $\sigma^t : D \rightarrow \mathbb{R}_+$ is the total cross section with $\sigma^t \geq \sigma^s$. At some occasions we are also going to use the absorption cross section $\sigma^a := \sigma^t - \sigma^s$.

Our goal is to construct an approximation of (4.1) that is positivity preserving and asymptotic preserving in the diffusion limit. We also want to make sure that the above properties hold with grazing incidences and inhomogeneous materials.

4.2. Angular discretization

To simplify the presentation of the method we use the discrete ordinate technique to do the discretization with respect to the angles. The resulting angular quadrature is denoted $\{\mu_l, \boldsymbol{\Omega}_l\}_{l \in \mathcal{L}}$ and is assumed to satisfy

$$\sum_{l \in \mathcal{L}} \mu_l = |S|, \quad \sum_{l \in \mathcal{L}} \mu_l \boldsymbol{\Omega}_l = \mathbf{0}, \quad \sum_{l \in \mathcal{L}} \boldsymbol{\Omega}_l |c \cdot \boldsymbol{\Omega}_l| = \mathbf{0}, \quad \sum_{l \in \mathcal{L}} \mu_l \boldsymbol{\Omega}_l \otimes \boldsymbol{\Omega}_l = \frac{|S|}{3} \mathbb{I}, \quad (4.2)$$

for all $c \in \mathbb{R}^3$, where \mathbb{I} is the 3×3 identity matrix. We define $L := \text{card}(\mathcal{L})$. All the simulations reported in the paper are done with the S_N technique (Gauss-Legendre quadrature along the polar axis and equi-distributed angles along the azimuth with $\frac{1}{8}N(N+2)$ angles per octant).

4.3. Space discretization

We are going to use continuous finite elements stabilized with the continuous interior penalty (CIP) technique (a.k.a. edge stabilization) from Douglas and Dupont [7] and Burman and Hansbo [5]. The method presented in the paper is not restricted to continuous elements and CIP though. It can be used with other types of stabilization. For instance, one can use discontinuous elements (of degree $p \geq 1$) stabilized with upwind numerical flux. One can also use continuous elements stabilized with methods like Galerkin Least-Squares (GaLS) (or its SUPG variation), Local Projection Stabilization (LPS), Orthogonal Subscale Stabilization (OSS), and Subgrid Viscosity (SGV).

Let $(\mathcal{T}_h)_{h \in H}$ be a shape-regular family of unstructured matching meshes exactly covering D . For simplicity we assume that all the elements are generated from a reference element denoted \hat{K} . The geometric transformation mapping \hat{K} to an arbitrary element

$K \in \mathcal{T}_h$ is denoted $T_K : \widehat{K} \rightarrow K$. We now introduce a reference finite element $(\widehat{K}, \widehat{P}, \widehat{\Sigma})$, which we assume, for simplicity, to be a Lagrange element. We define the following scalar-valued finite element space:

$$V_h = \{v \in C^0(D; \mathbb{R}) \mid v|_K \circ T_K \in \widehat{P}, \forall K \in \mathcal{T}_h\}. \tag{4.3}$$

The global shape functions are denoted by $\{\varphi_i\}_{i \in \mathcal{V}}$. Recall that $V_h = \text{span}\{\varphi_i\}_{i \in \mathcal{V}}$.

Given any mesh cell K in \mathcal{T}_h , we denote by h_K the diameter of K and \mathbf{n}_K the outward unit normal at the boundary of K . We set $h := \max_{K \in \mathcal{T}_h} h_K$. The collection of the mesh faces is denoted \mathcal{F}_h . The set of interfaces is denoted by \mathcal{F}_h° . The set of boundary faces is denoted by \mathcal{F}_h^∂ . Each interface F in \mathcal{F}_h° is oriented using a unit normal vector \mathbf{n}_F . Letting K_l, K_r be the two cells sharing an interface $F \in \mathcal{F}_h^\circ$, we adopt the convention that \mathbf{n}_F points from K_l to K_r . Every boundary face $F \in \mathcal{F}_h^\partial$ is oriented by the unit normal $\mathbf{n}_F := \mathbf{n}_D$. For all $F \in \mathcal{F}_h$ with $F = \partial K_l \cap \partial K_r$, and all function w smooth enough to have traces on F , we define the jump of w across F as

$$[[w]]_F := \lim_{K_l \ni y \rightarrow x} w(x) - \lim_{K_r \ni y \rightarrow x} w(x). \tag{4.4}$$

4.4. Finite element approximation

First, we define the sesquilinear form associated with the continuous interior penalty. We start by defining $\sigma_\epsilon^t = \sigma^t + \epsilon \sigma_\epsilon$, $\sigma_\epsilon^s = \sigma^s + \epsilon \sigma_\epsilon$, where $\sigma_\epsilon = \epsilon \sup_{x \in D} \sigma^t(x)$ and $\epsilon > 0$. The tests reported in the paper are done with $\epsilon = 10^{-10}$. Then, for all $\psi_h, \phi_h \in V_h$, we set

$$s_h(\psi_h, \phi_h) := \varpi \sum_{F \in \mathcal{F}_h^\circ} h_F^2 (\llbracket \nabla \psi_h \rrbracket_{\theta, F}, \llbracket \nabla \phi_h \rrbracket_{\theta, F})_{L^2(F)}, \tag{4.5a}$$

$$\llbracket \nabla \phi_h \rrbracket_{\theta, F} := (\theta_r \nabla \phi_h|_{K_l} - \theta_l \nabla \phi_h|_{K_r}) \cdot \mathbf{n}_F, \tag{4.5b}$$

$$\theta_r = \frac{\sigma_{\epsilon, l}^t}{\sigma_{\epsilon, l}^t + \sigma_{\epsilon, r}^t}, \quad \theta_l = \frac{\sigma_{\epsilon, r}^t}{\sigma_{\epsilon, l}^t + \sigma_{\epsilon, r}^t}, \tag{4.5c}$$

In all the simulations reported in the paper the parameters, ϖ , and h_F are defined by

$$\varpi := \frac{d^2}{(1+p)^4}, \quad h_F = \frac{\frac{1}{2}(|K_l| + |K_r|)}{|F|}, \tag{4.6}$$

where d is the space dimension and p is the polynomial degree of the approximation.

Next, for all $k \in \mathcal{L}$, we define the bilinear form associated with the operator $\psi \mapsto \Omega_k \cdot \nabla \psi + \sigma^t \psi$ and the bilinear form we use to weakly enforce boundary conditions. For all $k \in \mathcal{L}$, we set

$$t_k(\psi, \phi) = \int_D (\Omega_k \cdot \nabla \psi_{h,k}(x) + \sigma^t(x) \psi_{h,k}(x)) \phi_i(x) dx, \tag{4.7}$$

$$b_k(\psi, \phi) := \sum_{F \in \mathcal{F}_h^\partial} \int_F \frac{1}{2} (|\Omega_k \cdot \mathbf{n}| - \Omega_k \cdot \mathbf{n}) \psi(x) \phi(x) ds. \tag{4.8}$$

The approximation of the angular flux ψ_h is done in $(V_h)^L$. We set $\psi_h := (\psi_{h,1}, \dots, \psi_{h,L}) \in V_h \times \dots \times V_h$, with $\psi_{h,k} := \sum_{j \in \mathcal{V}} \Psi_{ik} \varphi_j \in V_h$ for all the angular direction in the quadrature $k \in \mathcal{L}$. Let α_k^∂ be the value of the boundary incidence along the quadrature angle Ω_k . The discrete ordinate Galerkin approximation of (4.1) consists of seeking $\Psi_h \in (V_h)^L$ so that the following holds for all $k \in \mathcal{L}$ and all $i \in \mathcal{V}$:

$$t_k(\psi_{h,k}, \varphi_i) + s_h(\psi_{h,k}, \varphi_i) + b_k(\psi_{h,k}, \varphi_i) = \int_D \sigma^s(x) \bar{\psi}_h(x) \varphi_i(x) dx + \int_D q(x) \varphi_i(x) dx + b_k(\alpha_k^\partial, \varphi_i). \tag{4.9}$$

We are also going to make use of the diffusion approximation with weakly enforced Dirichlet boundary condition. For this purpose, for all $\phi_h, r_h \in V_h$, we set

$$a(\phi_h, r_h) := \int_D \frac{1}{3\sigma_\epsilon^s(x)} \nabla \phi_h(x) \cdot \nabla r_h(x) dx + \int_D \sigma^a(x) \phi_h(x) r_h(x) dx + \frac{1}{4} \int_{\partial D} \phi_h(x) r_h(x) ds, \tag{4.10}$$

where we recall that $\frac{1}{2|S^1|} \int_{\Omega \in S^1} (|\Omega \cdot \mathbf{n}| + \Omega \cdot \mathbf{n}) ds = \frac{1}{|S^1|} \int_{\Omega \cdot \mathbf{n} < 0} |\Omega \cdot \mathbf{n}| ds = \frac{1}{4}$.

Table 5.1
Problem (5.1). \mathbb{P}_1 , \mathbb{P}_2 , and \mathbb{P}_3 continuous finite elements. Relative error in the L^1 -norm.

\mathbb{P}_1			\mathbb{P}_2			\mathbb{P}_3		
I	L^1 -Err	rate	I	L^1 -Err	rate	I	L^1 -Err	rate
101	1.22E-03	–	101	3.74E-04	–	100	1.07E-04	–
201	2.49E-04	2.31	201	4.13E-05	3.20	202	6.89E-06	3.90
401	5.59E-05	2.16	401	5.07E-06	3.04	400	4.54E-07	3.98
801	1.34E-05	2.07	801	6.31E-07	3.01	799	2.86E-08	4.00
1601	3.28E-06	2.03	1601	7.89E-08	3.00	1600	1.78E-09	4.00

4.5. Solution method

There are many solution methods to solve (4.9). The method that is used does not really matter for the purpose of the paper which we recall is about conservative limiting not involving computing a low-order solution based on artificial viscosity. The key idea is that limiting is done after (4.9) is solved. We explain in Appendix A the method that we use for all the simulations reported in the paper. As the purpose of the paper is just to discuss limiting, we have adopted a simple source iteration technique preconditioned with a diffusion approximation and using a minimum residual technique.

5. Numerical illustrations, scalar transport equation

The objective of this section is to illustrate the limiting technique proposed in the paper. We start by testing the method on the scalar advection equation. Examples involving the radiation transport equation are reported in §6.

5.1. Numerical details

The tests are done with continuous finite elements in one and two space dimensions. Unless specified otherwise, the simulations realized in one dimension are done on uniform meshes and those realized in two dimensions are done on unstructured Delaunay meshes. In all the tests reported below the quadratures are exact for the mass matrix. The index I stands for the number of degrees of freedom (or gridpoints) of the approximation. Given a nonzero function $u \in L^1(D)$ and u_h its finite element approximation, we call relative error in the L^1 -norm the quantity $\|u - u_h\|_{L^1(D)} / \|u\|_{L^1(D)}$.

5.2. 1D smooth solution

We start by solving a transport problem in one space dimension with a smooth solution. We let $D = (0, 8)$ and solve

$$\Omega \partial_x u + \sigma u = q, \quad \text{a.e. } x \in D, \quad u(0) = 0. \tag{5.1}$$

with $\Omega = 1$, $\sigma(x) = 1$, $q(x) = \Omega \pi \sin(\pi x) + \sigma(x)(1 - \cos(\pi x))$. The solution is $u(x) = 1 - \cos(\pi x)$.

We test the method with \mathbb{P}_1 , \mathbb{P}_2 , and \mathbb{P}_3 continuous finite elements on a series of meshes. We compute the relative error in the L^1 -norm. The results are shown in Table 5.1. We observe the expected convergence rates: 2, 3, and 4 for \mathbb{P}_1 , \mathbb{P}_2 , and \mathbb{P}_3 finite elements, respectively.

5.3. 1D non-smooth solution

We continue with one-dimensional transport problem with a non-smooth solution. We consider $D = (0, 1)$ and solve

$$\Omega \partial_x u + \sigma u = q, \quad \text{a.e. } x \in D, \quad u(0) = 0. \tag{5.2}$$

with $\Omega = 1$, and the scalar fields $\sigma(x)$ and $q(x)$ are piecewise constants and given by

$$\sigma(x) = \begin{cases} s_1 & x_0 \leq x \leq x_1 \\ s_2 & x_1 < x \leq x_2 \\ s_3 & x_2 < x \leq 1 \end{cases} \quad q(x) = \begin{cases} q_1 & x_0 \leq x \leq x_1 \\ q_2 & x_1 < x \leq x_2 \\ q_3 & x_2 < x \leq 1. \end{cases} \tag{5.3}$$

The exact solution is given by

$$u(x) = \begin{cases} \frac{q_1}{s_1} (1 - \exp(s_1(x_0 - x))) & x_0 \leq x \leq x_1 \\ u_1 \exp(s_2(x_1 - x)) & x_1 < x \leq x_2 \\ u_2 \exp(s_3(x_2 - x)) + \frac{q_3}{s_3} (1 - \exp(s_3(x_2 - x))) & x_2 < x \leq 1, \end{cases} \tag{5.4}$$

where $u_1 = \frac{q_1}{s_1} (1 - \exp(s_1(x_0 - x_1)))$ and $u_2 = u_1 \exp(s_2(x_1 - x_2))$. In the simulations reported below we use $x_0 = 0$, $x_1 = 0.3$, $x_2 = 0.6$, $s_1 = 1$, $s_2 = 10^3$, $s_3 = 2$, $q_1 = 1$, $q_2 = 0$, $q_3 = 1$.

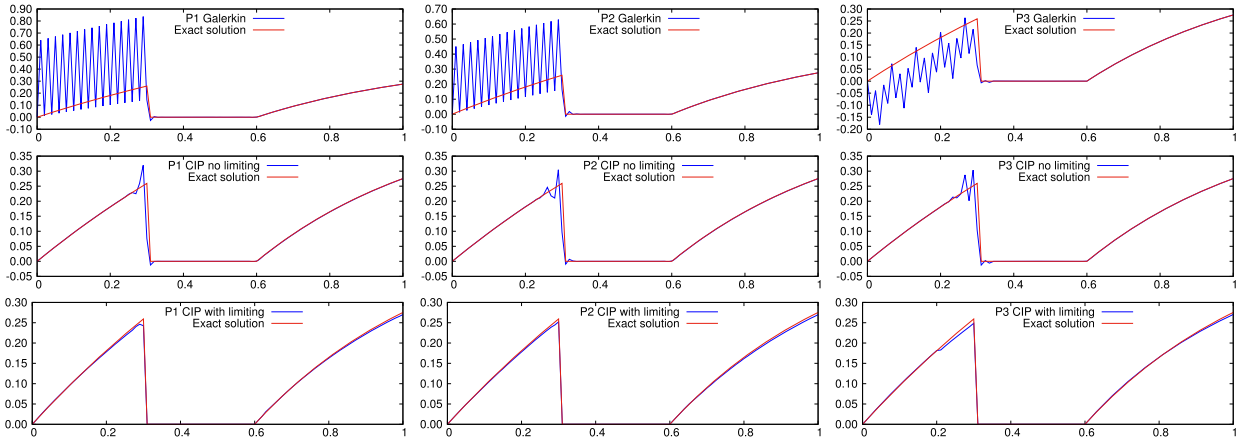


Fig. 5.1. Nonsmooth solution (5.2). Left column: \mathbb{P}_1 , 101 dofs. Center column: \mathbb{P}_2 , 101 dofs. Right column: \mathbb{P}_3 , 91 dofs. Top row: Galerkin. Center row: CIP without limiting. Bottom row: CIP with limiting.

Table 5.2
Non-smooth solution (5.2). \mathbb{P}_1 , \mathbb{P}_2 , and \mathbb{P}_3 continuous finite elements. Relative error in the L^1 -norm.

\mathbb{P}_1			\mathbb{P}_2			\mathbb{P}_3		
I	L^1 -Err	rate	I	L^1 -Err	rate	I	L^1 -Err	rate
101	2.08E-02	–	101	2.20E-02	–	91	1.87E-02	–
201	7.95E-03	1.40	201	7.33E-03	1.60	181	6.31E-03	1.58
401	2.44E-03	1.71	401	1.42E-03	2.38	361	1.34E-03	2.24
801	6.10E-04	2.00	801	2.03E-04	2.81	721	1.78E-04	2.92
1601	1.42E-04	2.10	1601	2.36E-05	3.11	1441	1.59E-05	3.49

We test the method with \mathbb{P}_1 , \mathbb{P}_2 , and \mathbb{P}_3 continuous finite elements. We show in Fig. 5.1 the graph of the \mathbb{P}_1 , \mathbb{P}_2 , and \mathbb{P}_3 approximation on a mesh composed of 101, 101 and 91 grid points, respectively. The panels in the top row of the figure show the unstabilized Galerkin solution. Oscillations are clearly visible in the interval $(0, x_1)$. The panels in the second row of the figure show the solution stabilized with the CIP method. Most of the oscillations are gone at the exception of overshoots localized at the interface located at x_1 . We finally show in the panels of the third row the results obtained with CIP stabilization and the limiting technique proposed in the paper. The maximum principle is satisfied.

The results of the convergence tests are shown in Table 5.2. We observe the optimal convergence rate close to $k + 1$ for \mathbb{P}_k approximation when the mesh is fine enough to capture the boundary layer located at x_1 . The rate is closer is 1.5 and 2 when the mesh is coarse.

5.4. 2D slip line

We now consider the two-dimensional unit square $D = (0, 1)^2$ and solve the problem

$$\Omega \partial_x u = 0, \quad \text{a.e. } x \in D, \quad u_{\Gamma_1} = 0, \quad u_{\Gamma_2} = 1, \tag{5.5}$$

with $\Omega = (0, 1)^2$, and $\Gamma_1 = \{(0, y) \mid 0 < y \leq 1\}$, $\Gamma_2 = \{(x, 0) \mid 0 \leq x \leq 1\}$. The exact solution is discontinuous; it exhibits a slip line align the axis $\{x = y, x > 0\}$. The solution is given by

$$u(x, y) = \begin{cases} 0 & x < y \\ 1 & y \leq x. \end{cases} \tag{5.6}$$

We show in the top three panels of Fig. 5.2 the graph of the limited CIP solution obtained with \mathbb{P}_1 , \mathbb{P}_2 and \mathbb{P}_3 elements on meshes composed of 9535, 9529, and 9541 grid points, respectively. We show in the bottom three panels of the figure the meshes and 11 isolines $\{0.05, 0.1, \dots, 0.9, 0.95\}$ (4 triangles are shown for each \mathbb{P}_2 cell, and 9 triangles are shown for each \mathbb{P}_3 cell)

We report in Table 5.3 the convergence rates for \mathbb{P}_1 , \mathbb{P}_2 , and \mathbb{P}_3 continuous finite elements. We observe a rate close to 0.73 for linear elements. The convergence rate for \mathbb{P}_2 and \mathbb{P}_3 elements is equal to 1 which is optimal as the solution is only in BV . This test confirms the near-optimality of the method for \mathbb{P}_2 and \mathbb{P}_3 elements.

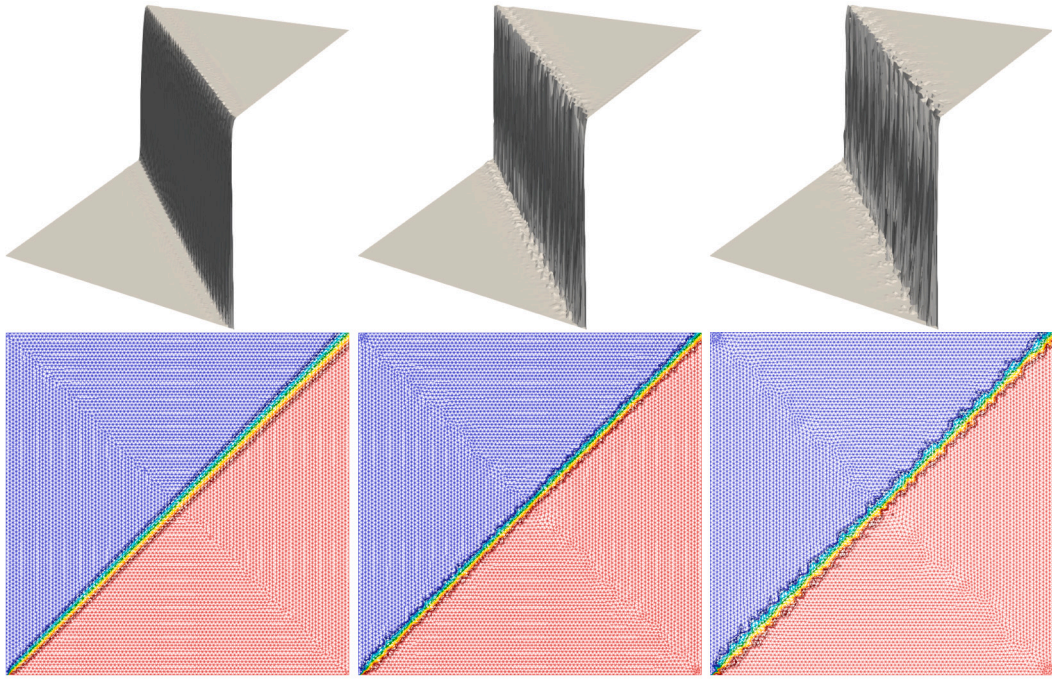


Fig. 5.2. Slip line problem (5.5). Non-uniform Delaunay meshes. Left: \mathbb{P}_1 , 9535 grid points; Center: \mathbb{P}_2 , 9529 grid points; \mathbb{P}_3 , 9541 grid points.

Table 5.3
Problem (5.5). \mathbb{P}_1 , \mathbb{P}_2 , and \mathbb{P}_3 continuous finite elements. Relative error in the L^1 -norm.

\mathbb{P}_1			\mathbb{P}_2			\mathbb{P}_3		
I	L^1 -Err	rate	I	L^1 -Err	rate	I	L^1 -Err	rate
961	5.58E-02	–	1681	4.17E-02	–	961	7.00E-02	–
3721	3.47E-02	0.70	6561	2.17E-02	0.96	3721	3.68E-02	0.95
14641	2.12E-02	0.72	25921	1.11E-02	0.98	14641	1.89E-02	0.98
58081	1.29E-02	0.73	103041	5.59E-03	0.99	58081	9.55E-03	0.99
231361	7.75E-03	0.73	410881	2.81E-03	0.99	231361	4.81E-03	0.99

5.5. 2D non smooth problem

Let $D = (0, 1)^2$. We solve the two-dimensional version of the problem (5.2),

$$\Omega \cdot \nabla u + \sigma u = q, \quad \text{a.e. } x \in D, \quad u(0, y) = 0, \quad 0 \leq y \leq 1. \tag{5.7}$$

with $\Omega = (1, 0)$, and, with the notation $x := (x, y)$, the scalar fields $\sigma(x)$ and $q(x)$ given by

$$\sigma(x, y) = \begin{cases} s_1 & x_0 \leq x \leq x_1 \\ s_2 & x_1 < x \leq x_2 \\ s_3 & x_2 < x \leq 1 \end{cases} \quad q(x, y) = \begin{cases} q_1 & x_0 \leq x \leq x_1 \\ q_2 & x_1 < x \leq x_2 \\ q_3 & x_2 < x \leq 1. \end{cases} \tag{5.8}$$

The exact solution is given by

$$u(x, y) = \begin{cases} \frac{q_1}{s_1}(1 - \exp(s_1(x_0 - x))) & x_0 \leq x \leq x_1 \\ u_1 \exp(s_2(x_1 - x)) & x_1 < x \leq x_2 \\ u_2 \exp(s_3(x_2 - x)) + \frac{q_3}{s_3}(1 - \exp(s_3(x_2 - x))) & x_2 < x \leq 1. \end{cases} \tag{5.9}$$

We use the same set of the coefficients $x_0, x_1, x_2, s_1, s_2, s_3$, and q_1, q_2, q_3 as in §5.3 (i.e., $x_0 = 0, x_1 = 0.3, x_2 = 0.6, s_1 = 1, s_2 = 10^3, s_3 = 2, q_1 = 1, q_2 = 0, q_3 = 1$).

The results of the convergence tests done on uniform meshes are shown in Table 5.4. The convergence rate varies between 1.1 and 2. Recall that the convergence rate 2 is the optimal since the gradient of the solution is only in BV.

Table 5.4
Problem (5.7). \mathbb{P}_1 , \mathbb{P}_2 , and \mathbb{P}_3 continuous finite elements. Relative error in the L^1 -norm.

\mathbb{P}_1			\mathbb{P}_2			\mathbb{P}_3		
I	L^1 -Err	rate	I	L^1 -Err	rate	I	L^1 -Err	rate
961	8.07E-02	–	1681	5.20E-02	–	961	7.01E-02	–
3721	3.83E-02	1.10	6561	2.25E-02	1.23	3721	3.20E-02	1.16
14641	1.65E-02	1.23	25921	8.13E-03	1.48	14641	1.25E-02	1.38
58081	5.96E-03	1.48	103041	2.13E-03	1.94	58081	3.64E-03	1.79
231361	1.73E-03	1.79	410881	4.85E-04	2.14	231361	7.64E-04	2.26

Table 6.1
Data for the one-dimensional test cases.

	#zones	5						#zones	1			#zones	1	
Case 1	Length	2.0	1.0	2.0	1.0	2.0	Case 2	Length	10.0		Case 3	Length	10.0	
	σ_s	0.0	0.0	0.0	0.9	0.9		σ_s	100.0			σ_s	0.09999	
	σ_r	50.0	5.0	0.0	1.0	1.0		σ_r	100.0			σ_r	0.1	
	q	25.	0.0	0.0	.5	0.0		q	0.0			q	0.5	
	#cells	25	25	25	25	25		#cells	100			#cells	100	
B.C.	Vac.					B.C.	$\psi_2(0) = 1$		B.C.	Vac.		B.C.	Vac.	

6. Radiation transport

In this section, we report the tests done on the radiation transport equation using the algorithm described in the paper.

6.1. Numerical details

The positive- and asymptotic-preserving algorithm defined in Algorithms A.1 and A.2 is implemented with continuous finite elements of degree $p \in \{1, 2, 3\}$ on simplices. The meshes in one dimension are uniform. The meshes in two space dimension are non-uniform and composed of triangles. The angular discretization is done with the Gauss-Chebyshev S_N quadrature. In one dimension, the x_1 -component of the angles are the N quadrature points of the Gaussian-Legendre quadrature over $[-1, 1]$ and the weights are the weights of the Gaussian-Legendre quadrature. In two-dimensions we use the standard triangular S_N quadrature. As the x_3 -direction is ignored, there are four quadrants and the total number of angular directions is $\frac{1}{2}N(N + 2)$. Unless specified otherwise, the units are cm for lengths, nb.part./cm²·s·sr for ψ , nb.part./cm²·s for $\bar{\psi}$, nb.part./cm³·s for q , and cm⁻¹ for the cross sections.

6.2. One-dimensional benchmark tests

We start by illustrating the performance of the method in one space dimension. We compare the results given by the proposed the method with those given by the unlimited dG1 approximation stabilized by using the upwind flux. We use the angular quadrature S_8 (8 discrete directions in 1D) for all the cases. The angles, characterized by their x_1 -component, are enumerated in increasing order from 1 to 8. The data for the four cases considered here are reported in Table 6.1. In each case we give the length of the domain and the number of zones composing the domain. For each zone we give the values of σ^l , σ^s , and q (constants), and we also give the number of cells composing the zone. The boundary condition for cases 1, 3, and 4 are $\psi_{h|\partial D_-} = 0$ (this is the so-called vacuum boundary condition). We enforce a grazing incidence for case 2; we set $\psi_{h,k|\partial D_-} = 0$ for $k \neq 5$, $1 \leq k \leq 8$, and $\psi_{h,5}(0) = 1.0$.

The results of the simulations are reported in Fig. 6.1. We show in Panels 6.1a-6.1c the total scalar flux, $2\bar{\psi}_h$, obtained with the dG1 approximation (labeled dG1) and with the proposed method (labeled cG1). We observe a fair agreement between the two methods given the number of grid points (recall that dG1 has two times as many grid points as cG1). Panel 6.1d shows the angular flux $\psi_{h,1}$ for case 4 in the last cell close to the right boundary. For this case the dG1 approximation gives negative values at $x = 100$ on the angular fluxes 1, 2, and 3 (the values are -0.24 , -0.22 , -0.066 , respectively (approximated to 2 digits)). In all the cases the proposed method is always nonnegative.

6.3. Grazing incidences

We now focus on the second test case discussed in §6.2. A grazing incidence is enforced on the left boundary of the domain (only $\psi_{h,5}(0)$ is nonzero). We also have $\sigma^l h = \sigma^s h = 10$; that is, the diffusive regime is dominant. The conjunction of these two conditions produces a boundary layer as shown in Malvagi and Pomraning [21], Chandrasekhar [6]. Moreover, as established in Theorem 5.4 in [9], convergence of the dG approximation only occurs in a weak norm; more precisely, convergence on $\bar{\psi}_h$ occurs in the Sobolev space $H^s(D)$ only for smoothness indices s strictly less than $\frac{1}{2}$. That is to say, the values of $\bar{\psi}_h$ at the boundary do not converge (recall that boundary traces do not exist when $s < \frac{1}{2}$). Let us for a moment denote by $\alpha^\partial(\Omega, \mathbf{x})$ the boundary data of the problem. Then as observed in Adams [1, Eq. 66] (see also Prop. 3.6 in [9]) the leading term of the diffusion expansion of the dG approximation at the boundary is

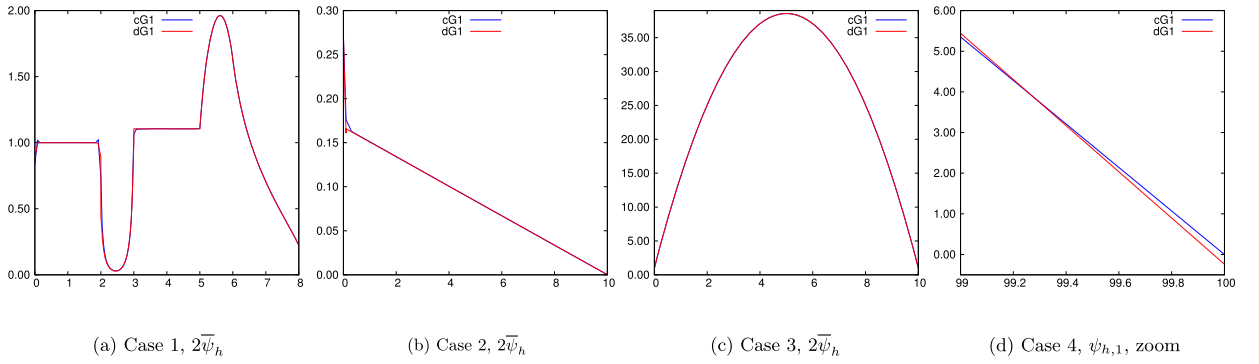


Fig. 6.1. Comparisons between the present method using cG1 and the upwind dG1 method.

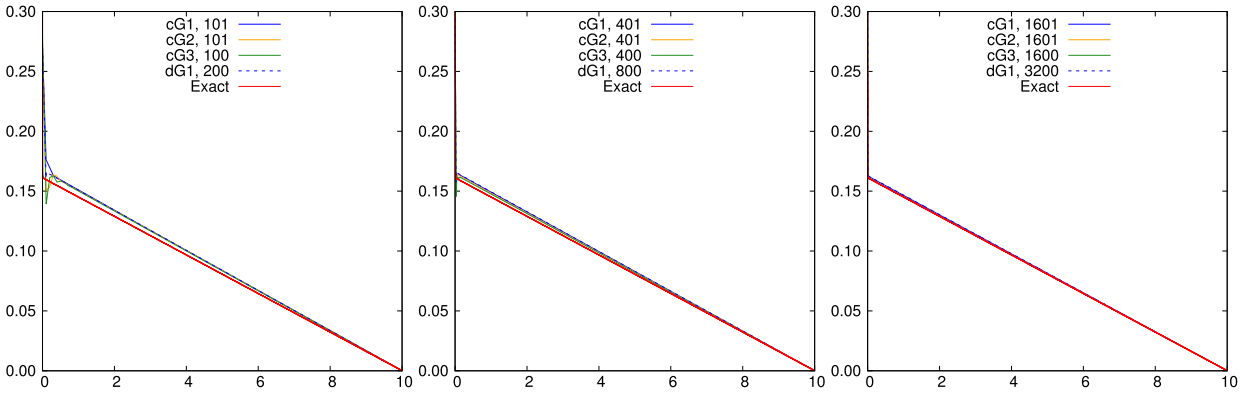


Fig. 6.2. Total scalar flux, $2\bar{\psi}_h$, for the grazing problem with cG1, cG2, cG3, and dG1. Left: 100 grid points. Center: 400 grid points. Right: 1600 grid points.

$$\bar{\psi}_h(\mathbf{x})|_{\partial D} \approx \frac{1}{2\pi} \int_{\Omega \cdot \mathbf{n} \leq 0} (|\Omega \cdot \mathbf{n}| + \frac{3}{2} |\Omega \cdot \mathbf{n}|^2) \alpha^\partial(\Omega, \mathbf{x}) \, d\Omega, \tag{6.1}$$

which is almost equal to $\frac{1}{2\pi} \int_{\Omega \cdot \mathbf{n} \leq 0} \frac{\sqrt{3}}{2} |\Omega \cdot \mathbf{n}| H(|\Omega \cdot \mathbf{n}|) \alpha^\partial(\Omega, \mathbf{x}) \, d\Omega$, (where H is Chandrasekar’s function), because $\frac{\sqrt{3}}{2} \mu H(\mu) \approx 0.91\mu + 1.635\mu^2 \approx \mu + \frac{3}{2}\mu^2$ up to a “a few percents” over the interval $\mu \in [0, 1]$. Hence, although, strictly speaking, the dG approximation is not asymptotic preserving when there are grazing incidences (unless the mesh or the shape functions are designed to resolve boundary layers), due to the above observation, it is a common practice in the radiation transport literature to say that the dG approximation with the upwind numerical flux is asymptotic preserving regardless of whether the boundary data is isotropic or not.

Note that the above comments regarding dG are independent of the polynomial degree, i.e., (6.1) does not depend on p . Note also that the numerical results reported in the second panel of Fig. 6.1 show that the result holds true as well for the stabilized cG1 approximation. The key here is that in both cases the boundary condition is enforced weakly using the bilinear form b_k defined in (4.8). We show in Fig. 6.2 simulations done with the proposed method using cG p , $p \in \{1, 2, 3\}$, on meshes composed on 100 (left panel), 400 (center panel), and 1600 (right panel) grid points. In each case, the number of cells is adjusted so as to maintain the same number of grid points for all $p \in \{1, 2, 3\}$. We also report in this figure the results from the dG1 approximation with 200, 800, and 3200 grid points. The red line labeled “exact” is obtained by using cG3 with 300000 grid points. We observe that, for each number of degrees of freedom all the methods give almost exactly the same results (regardless of whether the method is dG or cG). We notice also that with 100 grid points the numerical solution is indeed “a few percents” away from the exact solution, as claimed in [1, Eq. 66]. The reader is invited to zoom on the right panel where we used 1600 gridpoints: the dG1 solution and the three cG solutions all align on one line that is still slightly away from the exact solution. This confirms the theoretical result established in Theorem 5.4 in [9] where it is shown that when grazing incidences are enforced, convergence occurs in weak norms only and the convergence rate is weak (the convergence rate behaves like $\mathcal{O}(h^{\frac{1}{2}})$ in one space dimension).

Remark 6.1 (Comparison with [11]). The positive and asymptotic preserving method presented in [11] behaves properly in the diffusion limit with grazing incidences only if the asymptotic boundary value (6.1) (or Chandrasekar’s exact value) is enforced weakly (just enforcing $\alpha^\partial(\Omega, \mathbf{x})$ is not asymptotic preserving). This is not the case here. The present method properly works in every regime by just enforcing the boundary conditions using the bilinear form b_k defined in (4.8).

#zones	2	
Length	10.0	10.0
σ_s	0.0	100.0
σ_t	0.0	100.0
q	0.0	0.0
B.C.	$\psi_5(0) = 1$	

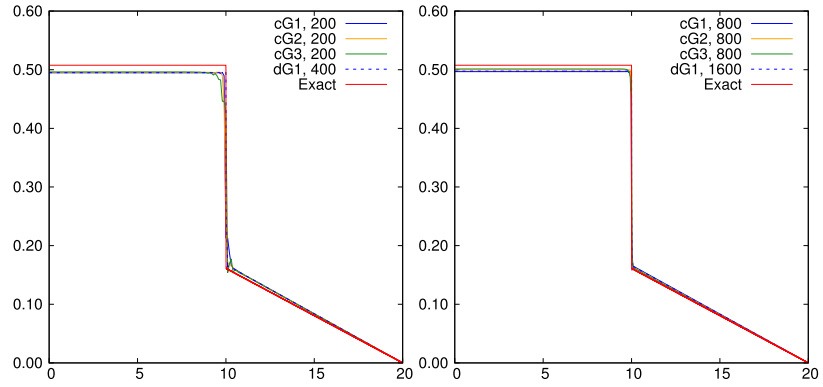


Fig. 6.3. Total scalar flux, $2\bar{\psi}_h$, for the grazing problem with a vacuum interface using cG1, cG2, cG3, and dG1. Left: data. Center: 200 grid points. Right: 800 grid points.

Table 6.2
Diffusion limit. Convergence test with respect to the mesh-size and ϵ .

ϵ	I	rel($\ e\ _{L^2}$)	rate	rel($\ \nabla e\ _{L^2}$)	rate	ϵ	I	rel($\ e\ _{L^2}$)	rate	rel($\ \nabla e\ _{L^2}$)	rate
10^{-3}	140	3.48E-02	-	7.57E-02	-	10^{-5}	140	1.25E-02	-	2.24E-02	-
	507	5.19E-03	2.96	1.91E-02	2.14		507	3.16E-03	2.13	6.98E-03	1.81
	1927	2.35E-03	1.18	5.88E-03	1.77		1927	7.66E-04	2.12	2.65E-03	1.45
	7545	2.91E-03	-31	2.32E-03	1.36		7545	1.70E-04	2.21	7.11E-04	1.93
	29870	3.05E-03	-07	1.44E-03	0.69		29870	2.81E-05	2.62	2.38E-04	1.59
10^{-4}	140	1.52E-02	-	3.92E-02	-	10^{-6}	140	1.23E-02	-	1.95E-02	-
	507	3.39E-03	2.33	1.20E-02	1.84		507	3.14E-03	2.12	5.75E-03	1.90
	1927	6.35E-04	2.51	4.23E-03	1.56		1927	7.84E-04	2.08	1.98E-03	1.60
	7545	1.82E-04	1.83	1.19E-03	1.86		7545	1.92E-04	2.06	7.06E-04	1.51
	29870	2.70E-04	-57	3.14E-04	1.93		29870	4.59E-05	2.08	2.35E-04	1.60

We finally repeat the above test by adding a vacuum region on the left while still enforcing the same grazing incidence boundary condition as above. This test is meant to assess the behavior of the proposed method in the presence of interfaces with vacuum and grazing incidences. We show the results in Fig. 6.3 using 200 grid points in the center panel and 800 grid points in the right panel. We observe that the proposed method properly behaves. The cG and dG results almost coincide in both cases. We observe again, that due to the grazing incidence all the methods are (almost) asymptotic preserving up to “a few percents”.

6.4. Diffusion limit with constant cross sections

To verify that the method does not lock in the diffusion regime, we consider the two-dimensional test reported in [11, §5.2.1]. The computational domain is $D = (0, 1)^2$. The cross sections are constant $\sigma^t = \sigma^s = \frac{1}{\epsilon}$ with $\epsilon > 0$. The source is isotropic and given by $q(\mathbf{x}) = \epsilon^2 \pi^2 \sin(\pi x_1) \sin(\pi x_2)$. When $\epsilon \rightarrow 0$, the regime becomes dominated by diffusion. The asymptotic limit in this case is $\psi^0(\mathbf{x}) = \sin(\pi x_1) \sin(\pi x_2)$.

We solve (4.1) with the method described in the paper using linear elements ($p = 1$). The tests are done using five meshes with 140, 507, 1927, 7545, and 29870 Lagrange nodes, respectively; the corresponding mesh-sizes are approximately $h \approx 0.1, 0.5, 0.025, 0.125$, and 0.00625 . We use the S_6 quadrature (24 angular directions).

We show in Table 6.2 the results of the test for the following four values of the small parameter $\epsilon \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$. We report in this table the relative L^2 -norm and the relative H^1 -semi-norm of the difference between $\bar{\psi}_h$ and the Lagrange interpolant of the asymptotic limit ψ^0 . We observe that, as proved in [9, Th. 5.3] for the upwind dG1 approximation, the scalar flux $\bar{\psi}_h$ converges optimally to ψ^0 when ϵ is significantly smaller than the mesh-size. The convergence order is $\mathcal{O}(h^2)$ in the L^2 -norm. It seems that some super-closeness phenomenon occurs in the H^1 -semi-norm since the rate behaves like $\mathcal{O}(h^{1.5})$. Tests with quadratic and cubic elements demonstrate the same behavior. These tests are not reported for brevity.

6.5. Reflection effects

We now reproduce a test case proposed in [11] using exactly the same finite element grids and angular discretization in order to illustrate that the proposed method is indeed more accurate than that from the reference.

We consider the two-dimensional domain $D = (0, 1)^2$ composed of two regions: one that is optically thick and one without any scattering. The cross sections are distributed as follows: $\sigma^t(\mathbf{x}) = 100, \sigma^s(\mathbf{x}) = 99$ if $x_2 \geq 0.5$ (optically thick and diffusive zone), and $\sigma_t(\mathbf{x}) = \sigma_s(\mathbf{x}) = 0$ if $x_2 \leq 0.5$ (void). The angular approximation is done with the S_6 quadrature (24 angular directions). We enforce a

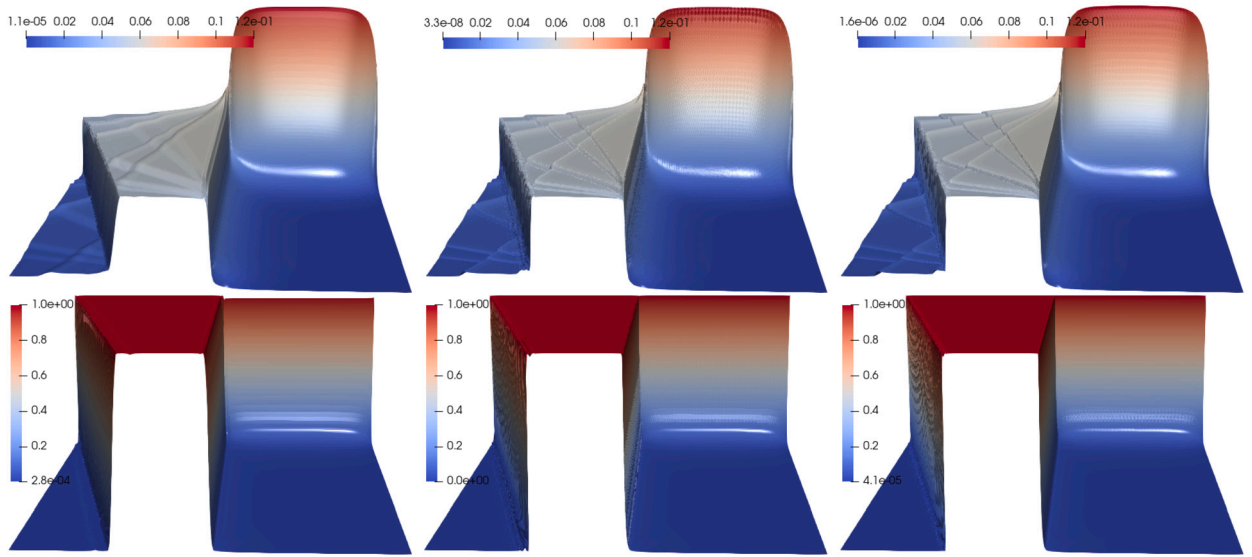


Fig. 6.4. Grazing and reflection effects: Top: scalar intensity $\bar{\psi}_h$. Bottom: first angular intensity, $\psi_{h,1}$. Left, \mathbb{P}_1 , 76230 grid points. Center \mathbb{P}_2 , 303893 grid points. Right: \mathbb{P}_3 , 682990 grid points.

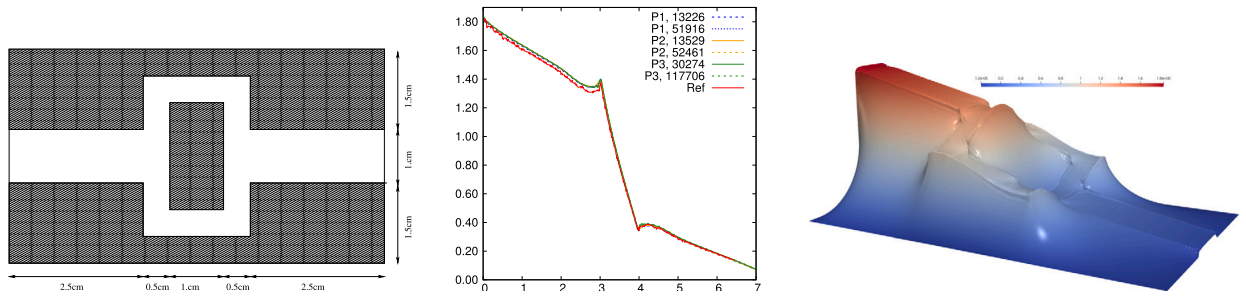


Fig. 6.5. Crooked pipe problem. Left: setting of the problem. Center: profile of the total scalar flux, $|S|\bar{\psi}_h$ for various meshes and polynomial degrees. The symbol “Ref.” stands for the data from [22]. Right: total scalar flux distribution (\mathbb{P}_3).

grazing incidence boundary condition on the leftmost boundary: this boundary is illuminated with intensity 1 along the first direction of the quadrature $\Omega_1 := (0.93802334, 0.25134260, 0.23861919)$ (eight digits truncation). All the other incoming angular fluxes are set to zero on this boundary. All the incoming angular fluxes are set to zero on the other three boundaries. The approximation in space for the asymptotic-preserving method is done on a non-uniform grid composed of 151434 triangles. There are 76230 \mathbb{P}_1 grid points (i.e., 1 829 520 dofs in total), 303893 \mathbb{P}_2 grid points (i.e., 7 239 432 dofs in total), 682990 \mathbb{P}_3 grid points (i.e., 16 391 760 dofs in total).

The results are shown in Fig. 6.4. Comparing these results with what is shown in the two leftmost panels in Fig. 3 in [11], we observe that the present method is significantly more accurate than that in [11] while being positivity-preserving and asymptotic preserving. We also notice the ray effect in the vacuum region $\{x_2 \leq 0.5\}$, which is an artifact of the S_N method (recall that we are just using 24 angular directions on purpose). That the ray effect is so crisply captured clearly demonstrates that the space approximation is very accurate; one cannot discern any smoothing induced by numerical dissipation.

6.6. Crooked pipe problem

We now solve a problem known in the literature as the crooked pipe problem. We adopt here the setting used in Olivier et al. [22, §7.3]. The geometry of the problem is shown in the left panel of Fig. 6.5, $D = (0, 7) \times (-2, 2)$. In the pipe we have $\sigma^s = 0.2$, $\sigma^a = 10^{-3}$, $q = 10^{-7}$. The characteristics of the material composing the walls are $\sigma^s = 200$, $\sigma^a = 10^{-3}$, and $q = 10^{-7}$. The boundary condition is $\Psi(x, y) = 2/|S| = 1/2\pi$ on $\{x = 0, |y| \leq 0.5\}$ and $\Psi(x, y) = 0$ on the rest of the boundary.

We use unstructured triangular Delaunay meshes. The meshes are generated so that the pipe/wall interface is exactly represented. The computations are done on various meshes and with various polynomial degrees ranging from \mathbb{P}_1 to \mathbb{P}_3 : \mathbb{P}_1 with 13226 grid points, \mathbb{P}_1 with 51916 grid points, \mathbb{P}_2 with 13529 grid points, \mathbb{P}_2 with 524621 gridpoints, \mathbb{P}_3 with 30274 grid points, and \mathbb{P}_3 with 117706 grid points.

We show in the center panel of Fig. 6.5, the profile of the total scalar intensity, $|S|\bar{\psi}_h$, along the segment $\{0 \leq x \leq 7, y = 0\}$. There are six different simulations. The angular discretization is done as in [22, §7.3] using the S_{24} angular quadrature ($\frac{1}{2}N(N+2) = 312$

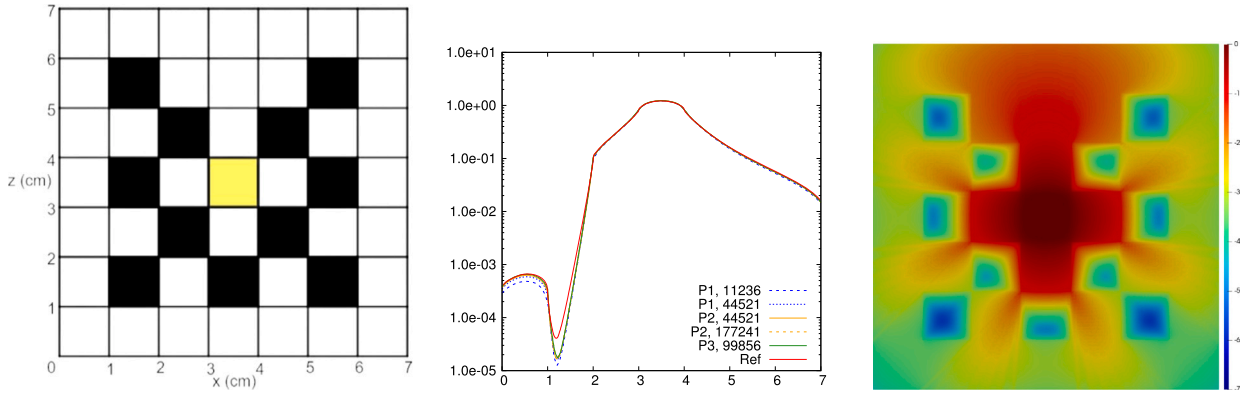


Fig. 6.6. Lattice problem. Left: setting of the problem. Center: profile of the total scalar flux, $|S|\bar{\psi}_h$ for various meshes and polynomial degrees. The symbol “Ref.” stands for the data collected from [23]. Right: total scalar flux distribution in logscale (\mathbb{P}_3 , 99856 grid points).

angles). We also report in this panel the results given in Olivier et al. [22, Fig. 10] (red line). We observe that the results from the six simulations with the proposed method collapse on a single curve, suggesting that all the spatial features are resolved. There are slight discrepancies of a few per cents with the results from [22], but overall the agreement is satisfactory.

6.7. Lattice problem

We continue with a benchmark test from Peng and McClarren [23, §5.2] called “Lattice problem” therein. The computational domain is $D = (0, 7)^2$. The material is organized in a checkerboard fashion. Each elementary region has size 1×1 . The details of the geometry are shown in the left panel of Fig. 6.6. There are 11 purely absorbing regions ($\sigma' = 10, \sigma^s = 0$, black boxes in Fig. 6.6), there are 37 regions that are purely scattering ($\sigma' = 1, \sigma^s = 1$, white boxes in Fig. 6.6), and there is one region with a strong source and scattering material ($q = 1, \sigma' = 1, \sigma^s = 1$, yellow box). The homogeneous Dirichlet boundary condition is enforced over the entire boundary of the domain.

The meshes are generated so that the material interfaces are exactly represented. The computations are done on three meshes with polynomial degrees ranging from \mathbb{P}_1 to \mathbb{P}_3 . Five different simulations done: \mathbb{P}_1 with 11236 grid points, \mathbb{P}_1 with 44521 grid points, \mathbb{P}_2 with 44521 grid points, \mathbb{P}_2 with 177241 grid points, and \mathbb{P}_3 with 99856 grid points. The angular quadrature is done with 312 angular directions (this is the S_{24} quadrature).

We show in the center panel of Fig. 6.6 the profile of the total scalar intensity, $|S|\bar{\psi}_h$, along the segment $\{x = 3.5, 0 \leq y \leq 7\}$. We also report in this figure results from Peng and McClarren [23, Fig. 8(e)]. These results are shown in red and identified with the symbol “Ref.” The representation is done in log scale along the y-axis. We observe that the agreement between the present method and the reference results is overall quite satisfactory.

6.8. Hohlräum

We finish with a test that is purely qualitative and loosely inspired from the hohlraum problem in Southworth et al. [27, IV.B]. The objective is to give some feeling on how the method behaves when solving a somewhat realistic problem. The problem under consideration is a very simplistic representation of a hohlraum device used in the indirect-drive approach of inertial confinement fusion. The dimensions are not given in the reference, but we use a square domain $D = (0, 10)^2$ (i.e., 10cm×10cm). These values are significantly larger than those of an actual hohlraum, but as the problem is linear, everything can be recalled by a length scale. The zero incidence boundary condition is imposed. The thickness of the walls of the cavity is 0.3. The width of the opening at the top and bottom of the cavity is 5.1. The “spherical” capsule inside the cavity is centered at (5., 5.). Its internal radius is 3.3 and its external radius is 3.6. The opening on the right side of the capsule (simulating a filling hole) is a cone with half angle equal to 5° and vertex located at (5., 5.). The sources are meant to simulate the heating of the wall by lasers, hence the sources are only located in the walls of the cavity in the region composed of the points $\mathbf{x} = (x_1, x_2)$ where $(0 \leq x_1 \leq 0.3 \text{ or } 9.7 \leq x_1 \leq 10)$ and $(2.7 \leq x_2 \leq 3.3 \text{ or } 4 \leq x_2 \leq 6 \text{ or } 6.7 \leq x_2 \leq 7.3)$. The constant value of the source is arbitrarily set to $q = 10^8$. The cross sections are distributed as follows. “Gold” wall of the cavity: $\sigma' = 10^2, \sigma^s = 2.5$. “Helium filling” around the spherical capsule: $\sigma' = 10^{-4}, \sigma^s = 10^{-4}$. “Plastic” capsule wall: $\sigma' = 10, \sigma^s = 6$. “Hydrogen fuel” inside the capsule: $\sigma' = 10^{-2}, \sigma^s = 10^{-2}$. Again, these values are inspired from [27, IV.B] and are not to be taken as actual values. Note that the problem is heterogeneous. It involves streaming and diffusive regions.

We show in the center and right panels of Fig. 6.7 the results of two simulations using continuous \mathbb{P}_1 and \mathbb{P}_3 elements. The meshsize is approximately 0.0125 for the \mathbb{P}_1 approximation (682261 grid points) and 0.035 for the \mathbb{P}_3 approximation (864754 grid points). We use the S_{12} angular quadrature as in [27, IV.B]. This makes 84 directions. The figure shows the scalar flux, $\bar{\psi}_h$, in both cases. The results look visually similar to what is reported in Figure 3(b) in [27] although the color scheme is slightly different.

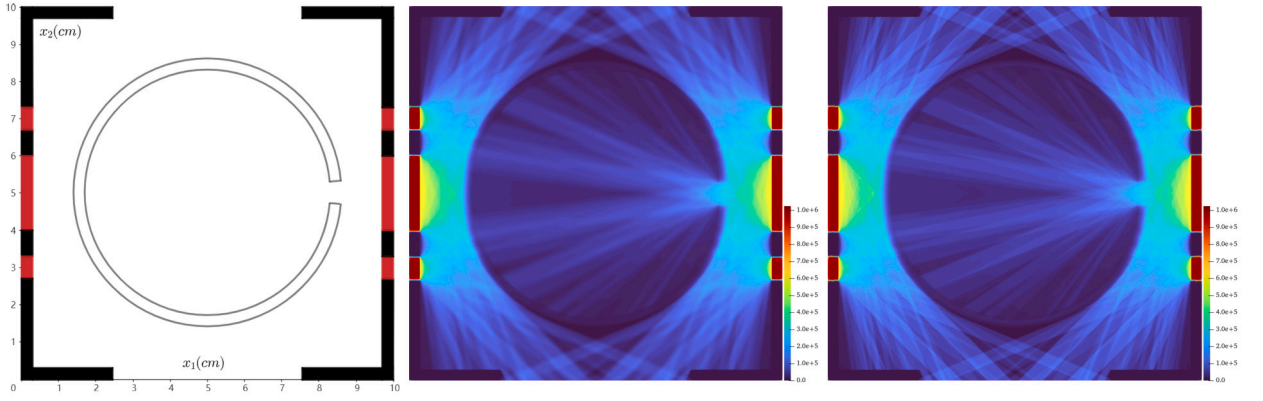


Fig. 6.7. Hohlraum problem. Left: geometry of the problem. Center, scalar flux, $\bar{\psi}_h, \mathbb{P}_1$, 682261 grid points, S_{12} quadrature. Right, scalar flux, $\bar{\psi}_h, \mathbb{P}_3$, 864754 grid points, S_{12} quadrature.

CRedit authorship contribution statement

Jean-Luc Guermond: Writing – review & editing, Writing – original draft, Visualization, Validation, Supervision, Software, Resources, Project administration, Methodology, Investigation, Funding acquisition, Formal analysis, Data curation, Conceptualization.
Zuodong Wang: Writing – review & editing, Visualization, Validation, Investigation, Data curation.

Declaration of competing interest

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Jean-Luc Guermond reports financial support was provided by National Science Foundation. Jean-Luc Guermond reports financial support was provided by AFOSR. Jean-Luc Guermond reports financial support was provided by ARO. Jean-Luc Guermond reports financial support was provided by Lawrence Livermore National Laboratory. If there are other authors, they declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements

The first author heartily thanks Dr. Dominic Caron for many fruitful discussions they had at TAMU from 2020 to 2022 over possible extensions of [11]. These investigations unfortunately did not bear fruits but they eventually led the first author to the strategy presented in the paper (stabilized Galerkin + characteristics-based limiting). The first author also thanks the organizers of the 2024 “Mathematics for neutronics” one-day workshop, Feb 24, 2014 at Lab. J.-L. Lions, Paris, for the stimulating discussions that lead to the completion of this exhausting project. Finally, the first author thanks M. Adams, J. Morel (TAMU) and T. Bailey (LLNL) for their indefectible support for this project over the years. The second author sincerely thanks Dr. Ari Rappaport for his kind help, his patience, and meaningful discussions about programming.

Appendix A. Preconditioned source iteration and limiting

We explain here the preconditioned source iteration method that is used in the radiation transport tests reported in the paper. We start with a few definitions to help us make the algorithm more concise. We define the mapping $\Phi : V_h \rightarrow V_h$ so that for all $p \in V_h$ (not to be confused with the polynomial degree), $\Phi(p)$ solve the diffusion equation $a(\Phi(p), v) = \int_D p v dx$ for all $v \in V$ (notice that the boundary condition is homogeneous). For all angular quadrature index $k \in \mathcal{L}$, we define the mapping $\Psi_k^0 : V_h \rightarrow V_h$ so that for all scalar flux $\phi \in V_h$, $\Psi_k^0(\phi)$ solves (4.9) with homogeneous boundary condition (i.e., zero incidence) and zero source term (i.e., $q \equiv 0$). Likewise we define $\Psi_k^* \in V_h$ so that Ψ_k^* solves (4.9) with the correct boundary condition and the correct source term q . In summary, $\Phi(p)$, $\Psi_k^0(\phi)$ and Ψ_k^* are defined so that the following holds for all $v \in V_h$:

$$a(\Phi(p), v) = \int_D p v dx, \tag{A.1}$$

$$t(\Psi_k^0(\phi), v) + s_h(\Psi_k^0(\phi), v) + b_h(\Psi_k^0(\phi), v) = \int_D \sigma^s \phi v dx, \tag{A.2}$$

$$t(\Psi_k^*, v) + s_h(\Psi_k^*, v) + b_h(\Psi_k^*, v) = \int_D q v dx + b_k(\alpha_k^d, v). \tag{A.3}$$

We set $\Psi^0(\phi) := (\Psi_1^0(\phi), \dots, \Psi_L^0(\phi)) \in V_h^L$ and $\Psi^* := (\Psi_1^*, \dots, \Psi_L^*) \in V_h^L$. Then $\Psi(\phi) := \Psi^0(\phi) + \Psi^*$ solves (4.9) iff $\overline{\Psi(\phi)} = \phi$. Hence, we have to solve the linear system: Find $\phi \in V_h$ so that

$$\phi - \overline{\Psi^0(\phi)} = \overline{\Psi^*}. \tag{A.4}$$

This can be done in a multitude of ways. The simple algorithm we use in the numerical simulations reported in the paper proceeds as follows. We initialize the algorithm by setting $\phi^0 = \Phi(q)$. Using Krylov's method to construct the solution, it is natural to define the first search direction to be the residual of (A.4), i.e., $p^0 = \phi^0 - \overline{\Psi^0(\phi^0)} - \overline{\Psi^*}$. Let ϕ^n (here $n \in \mathbb{N}$ is the loop index) be some estimate of ϕ . We add the direction induced by the diffusion limit $r^n := \Phi(\sigma^s p^n)$. Then we search for $\alpha^n, \beta^n \in \mathbb{R}$ so that the new update

$$\phi^{n+1} = \phi^n + \alpha^n p^n + \beta^n r^n, \tag{A.5}$$

minimizes $\|\phi^{n+1} - \overline{\Psi(\phi^{n+1})}\|$ where $\|\cdot\|$ is some norm induced by some inner product (\cdot, \cdot) that can be chosen by the user. As the next search direction is the residual $p^{n+1} = \phi^{n+1} - \overline{\Psi^0(\phi^{n+1})}$, after some algebraic manipulations we obtain

$$p^{n+1} = p^n + \alpha^n (p^n - \overline{\Psi^0(p^n)}) + \beta^n (r^n - \overline{\Psi^0(r^n)}). \tag{A.6}$$

To avoid stalling, which is a standard for steepest descent methods, we enforce p^{n+1} to be orthogonal to p^n . Setting $dp := p^n - \overline{\Psi^0(p^n)}$ and $dr := r^n - \overline{\Psi^0(r^n)}$, the two constraints on α^n and β^n are then

$$\|p^n\|^2 + \alpha^n (p^n, dp^n) + \beta^n (p^n, dr^n) = 0, \tag{A.7}$$

$$\min \left[\|p^n\|^2 + 2\alpha^n (p^n, dp^n) + 2\beta^n (p^n, dr^n) + \|\alpha^n dp^n + \beta^n dr^n\|^2 \right]. \tag{A.8}$$

The solution to this quadratique system is

$$\beta^n = \|p^n\|^2 \frac{\|dp^n\|^2 (p^n, dr^n) - (dp^n, dr^n)(p^n, dp^n)}{2(d p^n, d r^n)(p^n, d p^n)(p^n, d r^n) - \|d p^n\|^2 (p^n, d r^n)^2 - \|d r^n\|^2 (p^n, d p^n)^2}, \tag{A.9}$$

$$\alpha^n = \frac{-\|p^n\|^2 - \beta^n (p^n, dr^n)}{(p^n, dp^n)}, \tag{A.10}$$

and (A.5) gives the next estimate of the solution. At this point one applies the global mass conserving limiting algorithm (2.6) to $\{\phi^{n+1}\}_{i \in \mathcal{V}}$ with the global mass $\mathcal{M} := \sum_{i \in \mathcal{V}} m_i \phi_i^{n+1}$ and $\phi_i^{\min} = 0$ (and one can also enforce ϕ_i^{\max} if it happens that the maximum is a priori known). The algorithms stop when $\|p^{n+1}\|/\|\phi^{n+1}\|$ is smaller than some tolerance.

The final and key part of the algorithm is the conservative limiting (local and global). For each angle $k \in \mathcal{L}$, we proceed as follows: (1) For every $i \in \mathcal{V}$, we use the method of characteristics explained in §3 to compute the lower and upper bounds on the angular intensity, $\{\Psi_{i,k}^{\min}, \Psi_{i,k}^{\max}\}_{i \in \mathcal{V}}$. (2) We then apply the local mass conserving algorithm (2.2)–(2.5) to $\{\Psi_{i,k}\}_{i \in \mathcal{V}}$ using the bounds $\{\Psi_{i,k}^{\min}, \Psi_{i,k}^{\max}\}_{i \in \mathcal{V}}$ computed above. (3) Finally we apply the global mass conserving algorithm (2.6) to $\{\overline{\Psi}_{i,k}\}_{i \in \mathcal{V}}$ with the global mass $\mathcal{M} := \sum_{i \in \mathcal{V}} m_i \overline{\Psi}_{i,k}, \overline{\Psi}_{i,k}^{\min} = 0$ (one can also enforce $\overline{\Psi}_{i,k}^{\max}$ if the maximum happens to be a priori known).

Algorithm A.1 Preconditioned source iteration.

Require: Tolerance $\delta > 0$, upper bounds $\{\Phi_i^{\max}\}_{i \in \mathcal{V}}$ (optional).

<p>1: $n = 0, e = \infty$.</p> <p>2: $\phi^0 = \Phi(q), p^0 = \phi^0 - \overline{\Psi^0(\phi^0)} + \overline{\Psi^*} \leftarrow$ solution of (A.3).</p> <p>3: while $e > \delta$ do</p> <p style="padding-left: 20px;">$dp^n \leftarrow p^n - \overline{\Psi^0(p^n)}$</p> <p style="padding-left: 20px;">$r^n \leftarrow \Phi(\sigma^s p^n)$, solution of (A.1) with source $\sigma^s p^n$.</p> <p style="padding-left: 20px;">$dr^n \leftarrow r^n - \overline{\Psi^0(r^n)}$</p> <p>4: $\beta^n \leftarrow \frac{\ p^n\ ^2 (\ dp^n\ ^2 (p^n, dr^n) - (dp^n, dr^n)(p^n, dp^n))}{2(d p^n, d r^n)(p^n, d p^n)(p^n, d r^n) - \ d p^n\ ^2 (p^n, d r^n)^2 - \ d r^n\ ^2 (p^n, d p^n)^2}$.</p> <p style="padding-left: 20px;">$\alpha^n \leftarrow \frac{-\ p^n\ ^2 - \beta^n (p^n, dr^n)}{(p^n, dp^n)}$.</p> <p>5: $p^{n+1} \leftarrow p^n + \alpha^n dp^n + \beta^n dr^n$.</p> <p>6: $\phi^{n+1} \leftarrow \phi^n + \alpha^n p^n + \beta^n r^n$.</p> <p>7: $\mathcal{M} := \sum_{i \in \mathcal{V}} m_i \Phi_i^{n+1}$</p> <p style="padding-left: 20px;">$\Phi_i^{n+1, \min} = 0$ for all $i \in \mathcal{V}$.</p> <p style="padding-left: 20px;">Optional: $\Phi_i^{n+1, \max} \leftarrow \Phi_i^{\max}$ for all $i \in \mathcal{V}$.</p> <p style="padding-left: 20px;">$\{\Phi_i^{n+1}\}_{i \in \mathcal{V}} \leftarrow$ (2.6) with above parameters.</p> <p>8: $e \leftarrow \frac{\ \phi^{n+1} - \overline{\phi^{n+1}}\ }{\ \phi^{n+1}\ }$.</p> <p>9: $n \leftarrow n + 1$</p> <p>10: end while</p> <p>11: $\{\Psi_{i,k}\}_{i,k \in \mathcal{V} \times \mathcal{L}} \leftarrow \Psi^0(\phi^n) + \Psi^*$.</p> <p>12: return $\{\Psi_{i,k}\}_{i,k \in \mathcal{V} \times \mathcal{L}}$</p>	<p>\triangleright initialize iteration count and error</p> <p>\triangleright Initialize scalar flux and residual</p> <p style="padding-left: 20px;">\triangleright iteratively update $\Psi^0(\phi)$</p> <p style="padding-left: 40px;">\triangleright set dp^n</p> <p style="padding-left: 40px;">\triangleright set r^n</p> <p style="padding-left: 40px;">\triangleright set dr^n</p> <p>\triangleright optimize β^n</p> <p>\triangleright optimize α^n</p> <p style="padding-left: 20px;">\triangleright update p^n</p> <p style="padding-left: 20px;">\triangleright update ϕ^n</p> <p>\triangleright compute the global mass</p> <p style="padding-left: 20px;">\triangleright set lower bound</p> <p style="padding-left: 20px;">\triangleright set upper bound</p> <p>\triangleright post-processing by global limiter</p> <p>\triangleright estimate error at current step</p> <p style="padding-left: 20px;">\triangleright go to next iteration</p> <p>\triangleright Compute final solution</p>
--	--

Algorithm A.2 Conservative limiting.

Require: Numerical solution $\{\Psi_{i,k}\}_{i,k \in \mathcal{V} \times \mathcal{L}}$, Iteration number it^{\max} .

```

1: for  $it \in \{1 : it^{\max}\}$  do
2:   for  $(i, k) \in \mathcal{V} \times \mathcal{L}$  do
3:      $\Psi_{i,k}^{\max}, \Psi_{i,k}^{\min} \leftarrow (3.3)$  with scattering source defined with  $\overline{\Psi}_i$ .
        $\Psi_{i,k}^{\max}, \Psi_{i,k}^{\min} \leftarrow$  Relaxation Algorithm A.3.
        $\Psi_{i,k} \leftarrow (2.2)$ –(2.5) with above bounds.
4:   end for
5: end for
6:  $\mathcal{M}_k := \sum_{i \in \mathcal{V}} m_i \Psi_{i,k}$  for all  $k \in \mathcal{L}$ .
    $\Psi_{i,k}^{\min} = 0$  for all  $(i, k) \in \mathcal{V} \times \mathcal{L}$ .
   Optional:  $\Psi_{i,k}^{\max} \leftarrow$  input data, for all  $(i, k) \in \mathcal{V} \times \mathcal{L}$ .
    $\{\Psi_{i,k}\}_{i,k \in \mathcal{V} \times \mathcal{L}} \leftarrow (2.6)$ , with above bounds.
7: return  $\{\Psi_{i,k}\}_{i,k \in \mathcal{V} \times \mathcal{L}}$ 

```

▷ iteratively apply local limiter
 ▷ loop on each DoF
 ▷ estimate local bounds
 ▷ relax local bounds
 ▷ Local conservative limiting
 ▷ set global mass for each angular
 ▷ set lower bounds for global limiter
 ▷ set upper bounds for global limiter
 ▷ apply global conservative limiter

Algorithm A.3 Relaxation.

Require: Bounds $\{u_i^{\max}, u_i^{\min}\}_{i \in \mathcal{V}}$, $\{u_i\}_{i \in \mathcal{V}}$, FEM basis $\{\varphi_i\}_{i \in \mathcal{V}}$, u^{\min}, u^{\max} .

```

1:  $\beta_{ij} := \int_D \nabla \varphi_i \cdot \nabla \varphi_j \, dx$ .
2:  $\alpha_i \leftarrow \frac{\sum_{j \in \mathcal{I}(i)^*} \beta_{ij}(u_i - u_j)}{\sum_{j \in \mathcal{I}(i)^*} d_{ij}}$  for each  $i \in \mathcal{V}$ .
3: for  $i \in \mathcal{V}$  do
4:    $\beta_i \leftarrow \alpha_i$ .
5:   for  $j \in \mathcal{I}(i)^*$  do
6:     if  $\beta_i \alpha_j \leq 0$  then
7:        $\beta_i \leftarrow 0$ .
8:     else if  $|\beta_i| > |\alpha_j|$  then
9:        $\beta_i \leftarrow \alpha_j$ .
10:    end if
11:  end for
12: end for
13:  $u_i^{\max} \leftarrow \min(u_i^{\max} + |\beta_i|, u^{\max})$  for each  $i \in \mathcal{V}$ .
14:  $u_i^{\min} \leftarrow \max(u_i^{\min} - |\beta_i|, u^{\min})$  for each  $i \in \mathcal{V}$ .
15: return  $\{u_i^{\max}, u_i^{\min}\}_{i \in \mathcal{V}}$ 

```

▷ set stiffness matrix
 ▷ estimate second-order increment
 ▷ tighten bounds by minmod
 ▷ initialize relaxation parameter
 ▷ minmod process
 ▷ relax upper bound
 ▷ relax lower bound

Data availability

Data will be made available on request.

References

[1] M.L. Adams, Discontinuous finite element transport solutions in thick diffusive problems, *Nucl. Sci. Eng.* 137 (3) (2001) 298–333.

[2] M. Ancellin, B. Després, S. Jaouen, Extension of generic two-component VOF interface advection schemes to an arbitrary number of components, *J. Comput. Phys.* 473 (2023) 111721.

[3] P. Bochev, D. Ridzal, G. Scovazzi, M. Shashkov, Formulation, analysis and numerical study of an optimization-based conservative interpolation (remap) of scalar fields for arbitrary Lagrangian–Eulerian methods, *J. Comput. Phys.* 230 (13) (2011) 5199–5225.

[4] P. Bochev, D. Ridzal, K. Peterson, Optimization-based remap and transport: a divide and conquer strategy for feature-preserving discretizations, *J. Comput. Phys.* 257 (2014) 1113–1139.

[5] E. Burman, P. Hansbo, Edge stabilization for Galerkin approximations of convection-diffusion-reaction problems, *Comput. Methods Appl. Mech. Eng.* 193 (15–16) (2004) 1437–1453.

[6] S. Chandrasekhar, *Radiative Transfer*, Oxford University Press, 1950.

[7] J. Douglas, T. Dupont, Interior penalty procedures for elliptic and parabolic Galerkin methods, in: R. Glowinski, J.L. Lions (Eds.), *Computing Methods in Applied Sciences*, Springer Berlin Heidelberg, Berlin, Heidelberg, 1976, pp. 207–216.

[8] L. Gosse, G. Toscani, An asymptotic-preserving well-balanced scheme for the hyperbolic heat equations, *C. R. Math. Acad. Sci. Paris* 334 (4) (2002) 337–342.

[9] J.-L. Guermont, G. Kanschat, Asymptotic analysis of upwind discontinuous Galerkin approximation of the radiative transport equation in the diffusive limit, *SIAM J. Numer. Anal.* 48 (1) (2010) 53–78.

[10] J.-L. Guermont, 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.

[11] J.-L. Guermont, B. Popov, J. Ragusa, Positive and asymptotic preserving approximation of the radiation transport equation, *SIAM J. Numer. Anal.* 58 (1) (2020) 519–540.

[12] S. Hamilton, M. Benzi, J. Warsa, *Negative Flux Fixups in Discontinuous Finite Element s_n Transport*, Saratoga Springs, NY, 2009.

[13] A. Harten, S. Osher, Uniformly high-order accurate nonoscillatory schemes. I, *SIAM J. Numer. Anal.* 24 (2) (1987) 279–309.

[14] E.W. Larsen, On numerical solutions of transport problems in the diffusion limit, *Nucl. Sci. Eng.* 83 (1) (1983) 90–99.

[15] E.W. Larsen, J.E. Morel, W.F. Miller Jr., Asymptotic solutions of numerical transport problems in optically thick, diffusive regimes, *J. Comput. Phys.* 69 (2) (1987) 283–324.

[16] K. Lathrop, Spatial differencing of the transport equation: positivity vs. accuracy, *J. Comput. Phys.* 4 (4) (1969) 475–498.

[17] E.E. Lewis, W.F. Miller, *Computational Methods of Neutron Transport*, American Nuclear Society, 1993.

[18] C. Liu, B. Riviere, J. Shen, X. Zhang, A simple and efficient convex optimization based bound-preserving high order accurate limiter for Cahn–Hilliard–Navier–Stokes system, *SIAM J. Sci. Comput.* 46 (3) (2024) A1923–A1948.

[19] P.G. Maginot, J.E. Morel, J.C. Ragusa, A non-negative moment-preserving spatial discretization scheme for the linearized Boltzmann transport equation in 1-D and 2-D Cartesian geometries, *J. Comput. Phys.* 231 (20) (2012) 6801–6826.

- [20] P.G. Maginot, J.C. Ragusa, J.E. Morel, Nonnegative methods for bilinear discontinuous differencing of the sn equations on quadrilaterals, *Nucl. Sci. Eng.* 185 (1) (2017) 53–69.
- [21] F. Malvagi, G.C. Pomraning, Initial and boundary conditions for diffusive linear transport problems, *J. Math. Phys.* 32 (3) (1991) 805–820.
- [22] S. Olivier, W. Pazner, T.S. Haut, B.C. Yee, A family of independent variable Eddington factor methods with efficient preconditioned iterative solvers, *J. Comput. Phys.* 473 (2023) 111747.
- [23] Z. Peng, R.G. McClarren, A sweep-based low-rank method for the discrete ordinate transport equation, *J. Comput. Phys.* 473 (2023) 111748.
- [24] K. Peterson, P. Bochev, D. Ridzal, Optimization-based, property-preserving algorithm for passive tracer transport, *Comput. Math. Appl.* 159 (2024) 267–286.
- [25] R. Sanchez, N.J. McCormick, A review of neutron transport approximations, *Nucl. Sci. Eng.* 80 (4) (1982) 481–535.
- [26] B. Schmidtman, R. Abgrall, M. Torrilhon, On third-order limiter functions for finite volume methods, *Bull. Braz. Math. Soc. (N. S.)* 47 (2) (2016) 753–764.
- [27] B.S. Southworth, M. Holec, T.S. Haut, Diffusion synthetic acceleration for heterogeneous domains, compatible with voids, *Nucl. Sci. Eng.* 195 (2) (2021) 119–136.
- [28] B.C. Yee, S.S. Olivier, T.S. Haut, M. Holec, V.Z. Tomov, P.G. Maginot, A quadratic programming flux correction method for high-order dg discretizations of sn transport, *J. Comput. Phys.* 419 (2020) 109696.
- [29] S.T. Zalesak, Fully multidimensional flux-corrected transport algorithms for fluids, *J. Comput. Phys.* 31 (3) (1979) 335–362.