

1 **ACTIVE FLUX METHODS FOR HYPERBOLIC CONSERVATION**
2 **LAWS – FLUX VECTOR SPLITTING AND BOUND-PRESERVATION***

3 JUNMING DUAN[†], WASILIJ BARSUKOW[‡], AND CHRISTIAN KLINGENBERG[§]

4 **Abstract.** The active flux (AF) method is a compact high-order finite volume method that
5 simultaneously evolves cell averages and point values at cell interfaces. Within the method of lines
6 framework, the existing Jacobian splitting-based point value update incorporates the upwind idea
7 but suffers from a stagnation issue for nonlinear problems due to inaccurate estimation of the up-
8 wind direction and also from a mesh alignment issue partially resulting from decoupled point value
9 updates. This paper proposes to use flux vector splitting for the point value update, offering a
10 natural and uniform remedy to those two issues. To improve robustness, this paper also develops
11 bound-preserving (BP) AF methods for hyperbolic conservation laws. Two cases are considered:
12 preservation of the maximum principle for the scalar case, and preservation of positive density and
13 pressure for the compressible Euler equations. The update of the cell average is rewritten as a convex
14 combination of the original high-order fluxes and robust low-order (local Lax-Friedrichs or Rusanov)
15 fluxes, and the desired bounds are enforced by choosing the right amount of low-order fluxes. A
16 similar blending strategy is used for the point value update. In addition, a shock sensor-based lim-
17 iting is proposed to enhance the convex limiting for the cell average, which can suppress oscillations
18 well. Several challenging tests are conducted to verify the robustness and effectiveness of the BP AF
19 methods, including flow past a forward-facing step and high Mach number jets.

20 **Key words.** Hyperbolic conservation laws, active flux, flux vector splitting, bound-preserving,
21 convex limiting, shock sensor

22 **MSC codes.** 65M08, 65M12, 65M20, 35L65

23 **1. Introduction.** This paper focuses on the development of robust active flux
24 (AF) methods for hyperbolic conservation laws. The AF method is a new finite vol-
25 ume method [17, 16, 18, 36], that was inspired by [39]. Apart from cell averages, it
26 incorporates additional degrees of freedom (DoFs) as point values located at the cell
27 interfaces, evolved simultaneously with the cell average. The original AF method em-
28 ploys a globally continuous representation of the numerical solution using a piecewise
29 quadratic reconstruction, leading naturally to a third-order accurate method with a
30 compact stencil. The introduction of point values at the cell interfaces avoids the us-
31 age of Riemann solvers as in usual Godunov methods, because the numerical solution
32 is continuous across the cell interface and the numerical flux is available directly.

33 The independence of the point value update adds flexibility to the AF methods.
34 Based on the evolution of the point value, there are generally two kinds of AF methods.
35 The original one uses exact or approximate evolution operators and Simpson’s rule for
36 flux quadrature in time, i.e., it does not require time integration methods like Runge-
37 Kutta methods. Exact evolution operators have been studied for linear equations
38 in [8, 19, 18, 39]. Approximate evolution operators have been explored for Burgers’

*Submitted to the editors DATE.

Funding: JD was supported by an Alexander von Humboldt Foundation Research Fellowship CHN-1234352-HFST-P. CK and WB acknowledge funding by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) within *SPP 2410 Hyperbolic Balance Laws in Fluid Mechanics: Complexity, Scales, Randomness (CoScaRa)*, project number 525941602.

[†]Corresponding author. Institute of Mathematics, University of Würzburg, Emil-Fischer-Straße 40, 97074 Würzburg, Germany (junming.duan@uni-wuerzburg.de).

[‡]Institut de Mathématiques de Bordeaux (IMB), CNRS UMR 5251, University of Bordeaux, 33405 Talence, France (wasilij.barsukow@math.u-bordeaux.fr).

[§]Institute of Mathematics, University of Würzburg, Emil-Fischer-Straße 40, 97074 Würzburg, Germany (christian.klingenberg@uni-wuerzburg.de).

39 equation [17, 16, 36, 5], the compressible Euler equations in one spatial dimension
 40 [17, 27, 5], and hyperbolic balance laws [7, 6], etc. One of the advantages of the AF
 41 method over standard finite volume methods is its structure-preserving property. For
 42 instance, it preserves the vorticity and stationary states for multi-dimensional acoustic
 43 equations [8], and it is naturally well-balanced for acoustics with gravity [7].

44 Since it may not be convenient to derive exact or approximate evolution operators
 45 for nonlinear systems, especially in multi-dimensions, another kind of generalized AF
 46 method was presented in [1, 2, 3]. A method of lines was used, where the cell average
 47 and point value updates are written in semi-discrete form and advanced in time with
 48 time integration methods. In the point values update, the Jacobian matrix is split
 49 based on the sign of the eigenvalues (Jacobian splitting (JS)), and upwind-biased stencils
 50 are used to compute the approximation of derivatives. There are some deficiencies
 51 of the JS-based AF methods, e.g., the **stagnation** issue [27] for nonlinear problems,
 52 **and mesh alignment issue in 2D to be introduced in Subsection 3.2**. Some remedies
 53 are suggested for the stagnation issue, e.g., using discontinuous reconstruction [27] or
 54 evaluating the upwind direction using more neighboring information [5].

55 Solutions to hyperbolic conservation laws often stay in an *admissible state set* \mathcal{G} ,
 56 also called the invariant domain. For instance, the solutions to initial value problems of
 57 scalar conservation laws satisfy a strict maximum principle (MP) [14]. Physically, both
 58 the density and pressure in the solutions to the compressible Euler equations should
 59 stay positive. It is desired to conceive so-called bound-preserving (BP) methods, i.e.,
 60 those guaranteeing that the numerical solutions at a later time will stay in \mathcal{G} , if the
 61 initial numerical solutions belong to \mathcal{G} . The BP property of numerical methods is very
 62 important for both theoretical analysis and numerical stability. Many BP methods
 63 have been developed in the past few decades, e.g., a series of works by Shu and
 64 collaborators [46, 28, 43], a recent general framework on BP methods [42], and the
 65 convex limiting approach [21, 25, 30], which can be traced back to the flux-corrected
 66 transport (FCT) schemes for scalar conservation laws [13, 23, 33, 31]. The previous
 67 studies on the AF methods pay limited attention to high-speed flows, or problems
 68 involving strong discontinuities, with some efforts on the limiting for the point value
 69 update, see e.g. [5, 27, 10]. **Although those limitings can reduce oscillations, the**
 70 **new cell average may violate the bound even for linear advection [5, 27], and it is**
 71 **not straightforward to extend them to the multi-dimensional case. In [10, 9], the**
 72 **authors proposed to adopt a discontinuous reconstruction based on the scaling limiter**
 73 **[46]. The flux is computed based on the limited point values, resulting in BP AF**
 74 **methods for scalar conservation laws.** In a very recent paper, the MOOD [11] based
 75 stabilization was adopted to achieve the BP property [4] in an a posteriori fashion.

76 This paper presents a new way for the point value update to cure the **stagnation**
 77 **and mesh alignment issues**, develops suitable BP limitings for the AF methods, **and**
 78 **also proposes a shock sensor-based limiting to further suppress oscillations**. The main
 79 contributions and findings in this work can be summarized as follows.

80 **i).** We propose to employ the flux vector splitting (FVS) for the point value update,
 81 **which can cure both the stagnation and the mesh alignment issues effectively, because**
 82 **the FVS couples the neighboring information in a uniform and natural way.** The AF
 83 method based on the FVS is also shown to give better results than the JS, especially
 84 the local Lax-Friedrichs (LLF) FVS, in terms of the CFL number and shock-capturing
 85 ability.

86 **ii).** We develop BP limitings for both the cell average and point value by blending the
 87 high-order AF methods with the first-order LLF method in a convex combination. The
 88 main idea is to retain as much as possible of the high-order method while guaranteeing

89 the numerical solutions to be BP, and the blending coefficients are computed by
 90 enforcing the bounds. We show that using a suitable time step size and BP limitings,
 91 the BP AF methods satisfy the MP for scalar conservation laws, and preserve positive
 92 density and pressure for the compressible Euler equations.

93 **iii).** We design a shock sensor-based limiting, which helps to reduce oscillations by
 94 detecting shock strength. It is shown to strongly improve the shock-capturing ability
 95 in the numerical tests.

96 **iv).** Several challenging numerical tests are used to demonstrate the robustness and
 97 effectiveness of our BP AF methods. Moreover, for the forward-facing step problem,
 98 our BP AF method captures small-scale features better compared to the third-order
 99 DG method with the TVB limiter on the same mesh resolution, while using fewer
 100 DoFs, demonstrating its efficiency and potential for high Mach number flows.

101 The remainder of this paper is organized as follows. Section 2 introduces the 1D
 102 AF methods based on the FVS for the point value update. Section 3 extends our
 103 FVS-based AF methods to the 2D case. To design BP methods, Section 4 describes
 104 our convex limiting approach for the cell average, and the limiting for the point value.
 105 The shock sensor-based limiting is also proposed in Section 4 to suppress oscillations.
 106 The 1D limitings can be reduced from the 2D case, and more details are given in
 107 Section SM3 in the supplementary material. Some numerical tests are conducted
 108 in Section 5 to experimentally demonstrate the accuracy, BP properties, and shock-
 109 capturing ability of the methods. Section 6 concludes the paper with final remarks.

110 **2. 1D active flux methods.** This section presents the generalized AF methods
 111 using the method of lines for the 1D hyperbolic conservation laws

$$112 \quad (2.1) \quad \mathbf{U}_t + \mathbf{F}(\mathbf{U})_x = 0, \quad \mathbf{U}(x, 0) = \mathbf{U}_0(x),$$

113 where $\mathbf{U} \in \mathbb{R}^m$ is the vector of m conservative variables, $\mathbf{F} \in \mathbb{R}^m$ is the flux function,
 114 and $\mathbf{U}_0(x)$ is assumed to be initial data of bounded variation. Two cases are of
 115 particular interest. The first is a scalar conservation law ($m = 1$)

$$116 \quad (2.2) \quad u_t + f(u)_x = 0, \quad u(x, 0) = u_0(x).$$

117 The second case is that of compressible Euler equations of gas dynamics with $\mathbf{U} =$
 118 $(\rho, \rho v, E)^\top$ and $\mathbf{F} = (\rho v, \rho v^2 + p, (E + p)v)^\top$, where ρ denotes the density, v the velocity,
 119 p the pressure, and $E = \frac{1}{2}\rho v^2 + \rho e$ the total energy with e the specific internal energy.
 120 The perfect gas equation of state (EOS) $p = (\gamma - 1)\rho e$ is used to close the system with
 121 the adiabatic index $\gamma > 1$. Note that this paper uses bold symbols to refer to vectors
 122 and matrices, such that they are easier to distinguish from scalars.

123 Assume that a 1D computational domain is divided into N cells $I_i = [x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}]$
 124 with cell centers $x_i = (x_{i-\frac{1}{2}} + x_{i+\frac{1}{2}})/2$ and cell sizes $\Delta x_i = x_{i+\frac{1}{2}} - x_{i-\frac{1}{2}}$, $i = 1, \dots, N$. The
 125 DoFs of the AF methods are the approximations to cell averages of the conservative
 126 variable as well as point values at the cell interfaces, allowing some freedom in the
 127 choice of the point values, e.g. conservative variables, primitive variables, entropy
 128 variables, etc. This paper only considers using the conservative variables, and the
 129 DoFs are denoted by

$$130 \quad (2.3) \quad \bar{\mathbf{U}}_i(t) = \frac{1}{\Delta x_i} \int_{I_i} \mathbf{U}_h(x, t) \, dx, \quad \mathbf{U}_{i+\frac{1}{2}}(t) = \mathbf{U}_h(x_{i+\frac{1}{2}}, t),$$

131 where $\mathbf{U}_h(x, t)$ is the numerical solution. The cell average is updated by integrating

132 (2.1) over I_i in the following semi-discrete finite volume manner

$$133 \quad (2.4) \quad \frac{d\bar{U}_i}{dt} = -\frac{1}{\Delta x_i} \left[\mathbf{F}(U_{i+\frac{1}{2}}) - \mathbf{F}(U_{i-\frac{1}{2}}) \right].$$

134 Thus, the accuracy of (2.4) is determined by the approximation accuracy of the point
135 values. It was so far (e.g. in [2]) considered sufficient to update the point values with
136 any finite-difference-like formula

$$137 \quad (2.5) \quad \frac{dU_{i+\frac{1}{2}}}{dt} = -\mathcal{R} \left(U_{i+\frac{1}{2}-l_1}(t), \bar{U}_{i+1-l_1}(t), \dots, \bar{U}_{i+l_2}(t), U_{i+\frac{1}{2}+l_2}(t) \right), \quad l_1, l_2 \geq 0,$$

138 with \mathcal{R} a consistent approximation of $\partial \mathbf{F} / \partial x$ at $x_{i+\frac{1}{2}}$, as long as it gave rise to a
139 stable method. This paper explores further conditions on \mathcal{R} for nonlinear problems.

140 **2.1. Stagnation issue when using Jacobian splitting.** Let us first briefly
141 describe the point value update based on the JS [2], which reads

$$142 \quad (2.6) \quad \frac{dU_{i+\frac{1}{2}}}{dt} = - \left[\mathbf{J}^+(U_{i+\frac{1}{2}}) \mathbf{D}_{i+\frac{1}{2}}^+(U) + \mathbf{J}^-(U_{i+\frac{1}{2}}) \mathbf{D}_{i+\frac{1}{2}}^-(U) \right],$$

143 where the splitting of the Jacobian matrix $\mathbf{J} = \mathbf{J}^+ + \mathbf{J}^-$ is defined as

$$144 \quad \mathbf{J}^\pm = \mathbf{R} \mathbf{\Lambda}^\pm \mathbf{R}^{-1}, \quad \mathbf{\Lambda}^\pm = \text{diag}\{\lambda_1^\pm, \dots, \lambda_m^\pm\},$$

145 based on the eigendecomposition $\partial \mathbf{F} / \partial \mathbf{U} = \mathbf{R} \mathbf{\Lambda} \mathbf{R}^{-1}$, $\mathbf{\Lambda} = \text{diag}\{\lambda_1, \dots, \lambda_m\}$, where
146 $\lambda_1, \dots, \lambda_m$ are the eigenvalues, with the columns of \mathbf{R} the corresponding eigenvectors,
147 and $a^+ = \max\{a, 0\}$, $a^- = \min\{a, 0\}$. To derive the approximation of the derivatives
148 in (2.6), one can first obtain a high-order reconstruction for \mathbf{U} in the upwind cell,
149 and then differentiate the reconstructed polynomial. As an example, a parabolic
150 reconstruction in cell I_i is

$$151 \quad \mathbf{U}_{\text{para},1}(x) = -3(2\bar{U}_i - U_{i-\frac{1}{2}} - U_{i+\frac{1}{2}}) \frac{x^2}{\Delta x_i^2} + (U_{i+\frac{1}{2}} - U_{i-\frac{1}{2}}) \frac{x}{\Delta x_i} \\ 152 \quad (2.7) \quad + \frac{1}{4}(6\bar{U}_i - U_{i-\frac{1}{2}} - U_{i+\frac{1}{2}})$$

153 satisfying $\mathbf{U}_{\text{para},1}(\pm \Delta x_i / 2) = U_{i\pm\frac{1}{2}}$, $\frac{1}{\Delta x_i} \int_{-\Delta x_i/2}^{\Delta x_i/2} \mathbf{U}_{\text{para},1}(x) dx = \bar{U}_i$. Then the deriva-
154 tives are

$$155 \quad (2.8a) \quad \mathbf{D}_{i+\frac{1}{2}}^+(U) = \mathbf{U}'_{\text{para},1}(\Delta x_i / 2) = \frac{1}{\Delta x_i} \left(2U_{i-\frac{1}{2}} - 6\bar{U}_i + 4U_{i+\frac{1}{2}} \right),$$

$$156 \quad (2.8b) \quad \mathbf{D}_{i+\frac{1}{2}}^-(U) = \frac{1}{\Delta x_{i+1}} \left(-4U_{i+\frac{1}{2}} + 6\bar{U}_{i+1} - 2U_{i+\frac{3}{2}} \right).$$

157 One of the deficiencies of using the JS is the **stagnation** issue that appears in
158 certain setups for nonlinear problems, as observed in [27, 5]. **As shown in Example 5.1**
159 **for Burgers' equation**, the numerical solution based on the JS without limiting gives
160 a spike in the cell average at the initial discontinuity $x = 0.2$, which grows linearly
161 in time. The reason for this behavior is the inaccurate estimation of the upwind
162 direction at the cell interface, **required to split the Jacobian in (2.6)**. In this example,
163 there are two successive point values with different signs near the initial discontinuity:
164 $u_{i-\frac{1}{2}} = 2$, $u_{i+\frac{1}{2}} = -1$. At the cell interface $x_{i-\frac{1}{2}}$ or $x_{i+\frac{1}{2}}$, depending on the details of

165 initialization, the upwind discretization in (2.8) only uses the data from the left or
 166 right, leading to zero derivatives, thus the point values $u_{i-\frac{1}{2}}$ and $u_{i+\frac{1}{2}}$ stay unchanged.
 167 However, according to the update of the cell average (2.4), \bar{u}_i increases gradually
 168 (which is the observed spike). Proposed solutions to handle the stagnation issue
 169 involve estimating the Jacobian not only at the relevant cell interface, but also at the
 170 neighboring interfaces, and to select a better upwind direction (e.g. [5]), or achieve
 171 the same by blending (e.g. [9]). As will be shown below, using FVS instead of the JS
 172 naturally has a similar effect.

173 **2.2. Point value update using flux vector splitting.** In this paper, we pro-
 174 pose to use the FVS for the point value update, which was originally used to identify
 175 the upwind directions, and is simpler and somewhat more efficient than Godunov-type
 176 methods for solving hyperbolic systems [38]. The FVS for the point value update reads
 177

$$178 \quad (2.9) \quad \frac{dU_{i+\frac{1}{2}}}{dt} = -[\tilde{D}^+ F^+(U) + \tilde{D}^- F^-(U)]_{i+\frac{1}{2}},$$

179 where the flux F is split into the positive and negative parts $F = F^+ + F^-$ satisfying

$$180 \quad (2.10) \quad \lambda\left(\frac{\partial F^+}{\partial U}\right) \geq 0, \quad \lambda\left(\frac{\partial F^-}{\partial U}\right) \leq 0,$$

181 i.e., all the eigenvalues of $\frac{\partial F^+}{\partial U}$ and $\frac{\partial F^-}{\partial U}$ are non-negative and non-positive, respectively.
 182 Different FVS can be adopted as long as they satisfy the constraint (2.10), to be
 183 discussed later. Finite difference formulae to approximate the flux derivatives are
 184 obtained as follows. From the reconstruction of U (2.7), one can evaluate the flux
 185 F , and also the split fluxes F^+ pointwise. We compute them at the endpoints of the
 186 cell and in the middle. Then a parabolic reconstruction for, say, F^+ in the cell I_i is
 187 obtained as follows

$$188 \quad F_{\text{para},2}^+(x) = 2(F_{i-\frac{1}{2}}^+ - 2F_i^+ + F_{i+\frac{1}{2}}^+) \frac{x^2}{\Delta x_i^2} + (F_{i+\frac{1}{2}}^+ - F_{i-\frac{1}{2}}^+) \frac{x}{\Delta x_i} + F_i^+,$$

189 satisfying $F_{\text{para},2}^+(\pm\Delta x_i/2) = F_{i\pm\frac{1}{2}}^+ = F^+(U_{i\pm\frac{1}{2}})$, and $F_{\text{para},2}^+(0) = F_i^+ = F^+(U_i)$. The
 190 cell-centered point value is $U_i = (-U_{i-\frac{1}{2}} + 6\bar{U}_i - U_{i+\frac{1}{2}})/4$. Then the discrete derivatives
 191 are

$$192 \quad (2.11a) \quad (\tilde{D}^+ F^+)_{i+\frac{1}{2}} = (F_{\text{para},2}^+)'(\Delta x_i/2) = \frac{1}{\Delta x_i} (F_{i-\frac{1}{2}}^+ - 4F_i^+ + 3F_{i+\frac{1}{2}}^+),$$

$$193 \quad (2.11b) \quad (\tilde{D}^- F^-)_{i+\frac{1}{2}} = \frac{1}{\Delta x_{i+1}} (-3F_{i+\frac{1}{2}}^- + 4F_{i+1}^- - F_{i+\frac{3}{2}}^-).$$

194 These finite differences are third-order accurate. While the reconstructions of both
 195 U and F are parabolic, the coefficients in the formula (2.11) differ from that in [2]
 196 because (2.11) uses the cell-centered value rather than the cell average.

197 The FVS-based point value update borrows the information from the neighbors
 198 naturally, and can eliminate the generation of the spike effectively, as shown in Fig-
 199 ure 3, similar to the idea of the remedy in [5]. Note that we still use the original
 200 continuous reconstruction in the AF methods. We remark that, in AF methods, it is
 201 not clear how to define the point values at discontinuities, thus there may be other
 202 methods to fix the stagnation issue.

203 **2.2.1. Local Lax-Friedrichs flux vector splitting.** The first FVS we consider
 204 is the LLF FVS, defined as

$$205 \quad \mathbf{F}^\pm = \frac{1}{2}(\mathbf{F}(\mathbf{U}) \pm \alpha \mathbf{U}),$$

206 where the choice of α should fulfill (2.10) across the spatial stencil. In our implemen-
 207 tation, it is determined by

$$208 \quad (2.12) \quad \alpha_{i+\frac{1}{2}} = \max_r \{\varrho(\mathbf{U}_r)\}, \quad r \in \left\{i - \frac{1}{2}, i, i + \frac{1}{2}, i + 1, i + \frac{3}{2}\right\},$$

209 where ϱ is the spectral radius of $\partial \mathbf{F} / \partial \mathbf{U}$. One can also choose α to be the maxi-
 210 mal spectral radius in the whole domain, corresponding to the (global) LF splitting.
 211 Note, however, that a larger α generally leads to a smaller time step size and more
 212 dissipation.

213 **2.2.2. Upwind flux vector splitting.** One can also split the Jacobian matrix
 214 based on each characteristic field,

$$215 \quad (2.13) \quad \mathbf{F}^\pm = \frac{1}{2}(\mathbf{F}(\mathbf{U}) \pm |\mathbf{J}|\mathbf{U}), \quad |\mathbf{J}| = \mathbf{R}(\mathbf{\Lambda}^+ - \mathbf{\Lambda}^-)\mathbf{R}^{-1}.$$

216 Note that we evaluate the Jacobian at three locations in the cell I_i to get corresponding
 217 \mathbf{F}^\pm . For linear systems, one has $\mathbf{F} = \mathbf{J}\mathbf{U}$, so (2.13) reduces to the JS, because in this
 218 case

$$219 \quad \mathbf{F}^\pm = \frac{1}{2}(\mathbf{J} \pm |\mathbf{J}|)\mathbf{U} = \mathbf{R}\mathbf{\Lambda}^\pm \mathbf{R}^{-1}\mathbf{U} = \mathbf{J}^\pm \mathbf{U},$$

220 with \mathbf{J}^\pm a constant matrix so that $\tilde{\mathbf{D}}^\pm \mathbf{F}^\pm(\mathbf{U}) = \mathbf{J}^\pm \tilde{\mathbf{D}}^\pm \mathbf{U}$, which is the same as $\mathbf{J}^\pm \mathbf{D}^\pm \mathbf{U}$
 221 if \mathbf{D}^+ and $\tilde{\mathbf{D}}^+$ are derived from the same reconstructed polynomial. In other words, for
 222 linear systems, the AF methods using this FVS are the same as the original JS-based
 223 AF methods.

224 Such an FVS is also known as the Steger-Warming (SW) FVS [37] for the Euler
 225 equations, since the ‘‘homogeneity property’’ $\mathbf{F} = \mathbf{J}\mathbf{U}$ holds [38]. One can write down
 226 the SW FVS explicitly

$$227 \quad \mathbf{F}^\pm = \begin{bmatrix} \frac{\rho}{2\gamma} \alpha^\pm \\ \frac{\rho}{2\gamma} (\alpha^\pm v + a(\lambda_2^\pm - \lambda_3^\pm)) \\ \frac{\rho}{2\gamma} \left(\frac{1}{2} \alpha^\pm v^2 + av(\lambda_2^\pm - \lambda_3^\pm) + \frac{a^2}{\gamma-1} (\lambda_2^\pm + \lambda_3^\pm) \right) \end{bmatrix},$$

228 where $\lambda_1 = v$, $\lambda_2 = v + a$, $\lambda_3 = v - a$, $\alpha^\pm = 2(\gamma - 1)\lambda_1^\pm + \lambda_2^\pm + \lambda_3^\pm$, and $a = \sqrt{\gamma p / \rho}$ is the
 229 sound speed.

230 *Remark 2.1.* It should be noted that \mathbf{F}^\pm in this FVS may not be differentiable
 231 with respect to \mathbf{U} for nonlinear systems (e.g. Euler), as the splitting is based on the
 232 absolute value. In [40], the mass flux of \mathbf{F}^\pm is shown to be not differentiable, which
 233 might explain the accuracy degradation in Example 5.7.

234 **2.2.3. Van Leer-Hanel flux vector splitting for the Euler equations.**

235 Another popular FVS for the Euler equations was proposed by van Leer [40], and
 236 improved by [26]. The flux can be split based on the Mach number $M = v/a$ as

$$237 \quad \mathbf{F} = \begin{bmatrix} \rho a M \\ \rho a^2 (M^2 + \frac{1}{\gamma}) \\ \rho a^3 M (\frac{1}{2} M^2 + \frac{1}{\gamma-1}) \end{bmatrix} = \mathbf{F}^+ + \mathbf{F}^-, \quad \mathbf{F}^\pm = \begin{bmatrix} \pm \frac{1}{4} \rho a (M \pm 1)^2 \\ \pm \frac{1}{4} \rho a (M \pm 1)^2 v + p^\pm \\ \pm \frac{1}{4} \rho a (M \pm 1)^2 H \end{bmatrix},$$

238 with the enthalpy $H = (E + p)/\rho$, and the pressure splitting $p^\pm = \frac{1}{2}(1 \pm \gamma M)p$. This
 239 FVS gives a quadratic differentiable splitting with respect to the Mach number.

240 *Remark 2.2.* Different FVS may lead to different stability conditions but it is
 241 difficult to perform the analysis theoretically. We provide experimental CFL numbers
 242 for different FVS in some 1D tests. Our numerical tests in Section 5 show that the AF
 243 methods based on the FVS generally give better results than the JS, and the LLF FVS
 244 is the best among all the three FVS in terms of the CFL number and non-oscillatory
 245 property for high-speed flows involving strong discontinuities.

246 **2.3. Time discretization.** The fully-discrete scheme is obtained by using the
 247 SSP-RK3 method [20]

$$\begin{aligned}
 \mathbf{U}^* &= \mathbf{U}^n + \Delta t^n \mathbf{L}(\mathbf{U}^n), \\
 \mathbf{U}^{**} &= \frac{3}{4}\mathbf{U}^n + \frac{1}{4}(\mathbf{U}^* + \Delta t^n \mathbf{L}(\mathbf{U}^*)), \\
 \mathbf{U}^{n+1} &= \frac{1}{3}\mathbf{U}^n + \frac{2}{3}(\mathbf{U}^{**} + \Delta t^n \mathbf{L}(\mathbf{U}^{**})),
 \end{aligned}
 \tag{2.14}$$

249 where \mathbf{L} is the right-hand side of the semi-discrete schemes (2.4) or (2.5). The time
 250 step size is determined by the usual CFL condition

$$\Delta t^n = \frac{C_{\text{CFL}}}{\max_i \{\varrho(\bar{\mathbf{U}}_i)/\Delta x_i\}}.
 \tag{2.15}$$

252 **3. 2D active flux methods on Cartesian meshes.** This section presents the
 253 generalized AF methods using the method of lines for the 2D hyperbolic conservation
 254 laws

$$\mathbf{U}_t + \mathbf{F}_1(\mathbf{U})_x + \mathbf{F}_2(\mathbf{U})_y = 0, \quad \mathbf{U}(x, y, 0) = \mathbf{U}_0(x, y).
 \tag{3.1}$$

256 We will consider the scalar conservation law

$$u_t + f_1(u)_x + f_2(u)_y = 0, \quad u(x, y, 0) = u_0(x, y),
 \tag{3.2}$$

258 and the Euler equations with $\mathbf{U} = (\rho, \rho\mathbf{v}, E)^\top$, $\mathbf{F}_1 = (\rho v_1, \rho v_1^2 + p, \rho v_1 v_2, (E + p)v_1)^\top$,
 259 $\mathbf{F}_2 = (\rho v_2, \rho v_1 v_2, \rho v_2^2 + p, (E + p)v_2)^\top$, where $\mathbf{v} = (v_1, v_2)$ is the velocity vector, and the
 260 other notations are the same as for 1D in Section 2. The SSP-RK3 method is used to
 261 obtain the fully-discrete method.

262 Assume that a 2D computational domain is divided into $N_1 \times N_2$ cells, $I_{i,j} =$
 263 $[x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}] \times [y_{j-\frac{1}{2}}, y_{j+\frac{1}{2}}]$ with the cell sizes $\Delta x_i = x_{i+\frac{1}{2}} - x_{i-\frac{1}{2}}$, $\Delta y_j = y_{j+\frac{1}{2}} - y_{j-\frac{1}{2}}$,
 264 and cell centers $(x_i, y_j) = \left(\frac{1}{2}(x_{i-\frac{1}{2}} + x_{i+\frac{1}{2}}), \frac{1}{2}(y_{j-\frac{1}{2}} + y_{j+\frac{1}{2}})\right)$, $i = 1, \dots, N_1$, $j = 1, \dots, N_2$.
 265 The DoFs consist of the cell average $\bar{\mathbf{U}}_{i,j}(t) = \frac{1}{\Delta x_i \Delta y_j} \int_{I_{i,j}} \mathbf{U}_h(x, y, t) \, dx dy$, the face-
 266 centered values $\mathbf{U}_{i+\frac{1}{2},j}(t) = \mathbf{U}_h(x_{i+\frac{1}{2}}, y_j, t)$, $\mathbf{U}_{i,j+\frac{1}{2}}(t) = \mathbf{U}_h(x_i, y_{j+\frac{1}{2}}, t)$, and the value
 267 at the corner $\mathbf{U}_{i+\frac{1}{2},j+\frac{1}{2}}(t) = \mathbf{U}_h(x_{i+\frac{1}{2}}, y_{j+\frac{1}{2}}, t)$, where $\mathbf{U}_h(x, y, t)$ is the numerical so-
 268 lution. A sketch of the DoFs for the third-order AF method (for the scalar case) is
 269 given in Figure 1.

270 The cell average is evolved as follows

$$\frac{d\bar{\mathbf{U}}_{i,j}}{dt} = -\frac{1}{\Delta x_i} \left(\widehat{\mathbf{F}}_{i+\frac{1}{2},j} - \widehat{\mathbf{F}}_{i-\frac{1}{2},j} \right) - \frac{1}{\Delta y_j} \left(\widehat{\mathbf{F}}_{i,j+\frac{1}{2}} - \widehat{\mathbf{F}}_{i,j-\frac{1}{2}} \right),
 \tag{3.3}$$

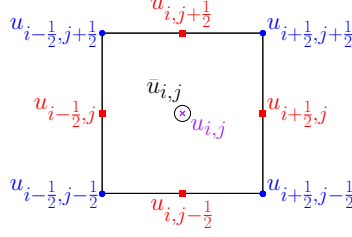


Fig. 1: The DoFs for the third-order AF method: cell average (circle), face-centered values (squares), values at corners (dots). Note that the cell-centered point value $u_{i,j}$ (cross) is used in constructing the scheme, but does not belong to the DoFs.

where $\widehat{\mathbf{F}}_{i+\frac{1}{2},j}$ and $\widehat{\mathbf{F}}_{i,j+\frac{1}{2}}$ are the numerical fluxes

(3.4)

$$\widehat{\mathbf{F}}_{i+\frac{1}{2},j} = \frac{1}{\Delta y_j} \int_{y_{j-\frac{1}{2}}}^{y_{j+\frac{1}{2}}} \mathbf{F}_1(\mathbf{U}_h(x_{i+\frac{1}{2}}, y)) \, dy, \quad \widehat{\mathbf{F}}_{i,j+\frac{1}{2}} = \frac{1}{\Delta x_i} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \mathbf{F}_2(\mathbf{U}_h(x, y_{j+\frac{1}{2}})) \, dx.$$

To achieve third-order accuracy, one can use Simpson's rule

$$\widehat{\mathbf{F}}_{i+\frac{1}{2},j} = \frac{1}{6} \left(\mathbf{F}_1(\mathbf{U}_{i+\frac{1}{2},j-\frac{1}{2}}) + 4\mathbf{F}_1(\mathbf{U}_{i+\frac{1}{2},j}) + \mathbf{F}_1(\mathbf{U}_{i+\frac{1}{2},j+\frac{1}{2}}) \right).$$

3.1. Point value update using flux vector splitting. For the evolution of the point values, consider the following general form

$$\frac{d\mathbf{U}_\sigma}{dt} = -\mathcal{R}(\overline{\mathbf{U}}_c(t), \mathbf{U}_{\sigma'}(t)), \quad c \in \mathcal{C}(\sigma), \sigma' \in \Sigma(\sigma),$$

where \mathcal{R} is a consistent approximation of $\partial \mathbf{F}_1 / \partial x + \partial \mathbf{F}_2 / \partial y$ at the point σ , $\mathcal{C}(\sigma)$ and $\Sigma(\sigma)$ are the spatial stencils containing the cell averages and point values, respectively. One can use the JS in [3], or employ the FVS for the point value update. E.g. for the point value at the corner $(x_{i+\frac{1}{2}}, y_{j+\frac{1}{2}})$ the FVS-based update reads

$$\frac{d\mathbf{U}_{i+\frac{1}{2},j+\frac{1}{2}}}{dt} = - \sum_{\ell=1}^2 \left[\widetilde{\mathbf{D}}_\ell^+ \mathbf{F}_\ell^+(\mathbf{U}) + \widetilde{\mathbf{D}}_\ell^- \mathbf{F}_\ell^-(\mathbf{U}) \right]_{i+\frac{1}{2},j+\frac{1}{2}},$$

where the fluxes are split as $\mathbf{F}_\ell = \mathbf{F}_\ell^+ + \mathbf{F}_\ell^-$, $\lambda \left(\frac{\partial \mathbf{F}_\ell^+}{\partial U} \right) \geq 0$, $\lambda \left(\frac{\partial \mathbf{F}_\ell^-}{\partial U} \right) \leq 0$. The explicit expressions of the 2D FVS used in this paper can be found in the supplementary material Section SM1. The finite difference operators $\widetilde{\mathbf{D}}_1^\pm$ and $\widetilde{\mathbf{D}}_2^\pm$ can be obtained similarly to Subsection 2.2. For third-order accuracy, starting with a bi-parabolic reconstruction of \mathbf{U} and computing a bi-parabolic interpolation of \mathbf{F}_ℓ^\pm , one thus obtains $\widetilde{\mathbf{D}}_1^\pm$ in the x -direction as

$$\left(\widetilde{\mathbf{D}}_1^+ \mathbf{F}_1^+ \right)_{i+\frac{1}{2},j+\frac{1}{2}} = \frac{1}{\Delta x_i} \left((\mathbf{F}_1^+)_{i-\frac{1}{2},j+\frac{1}{2}} - 4(\mathbf{F}_1^+)_{i,j+\frac{1}{2}} + 3(\mathbf{F}_1^+)_{i+\frac{1}{2},j+\frac{1}{2}} \right),$$

$$\left(\widetilde{\mathbf{D}}_1^- \mathbf{F}_1^- \right)_{i+\frac{1}{2},j+\frac{1}{2}} = \frac{1}{\Delta x_{i+1}} \left(-3(\mathbf{F}_1^-)_{i+\frac{1}{2},j+\frac{1}{2}} + 4(\mathbf{F}_1^-)_{i+1,j+\frac{1}{2}} - (\mathbf{F}_1^-)_{i+\frac{3}{2},j+\frac{1}{2}} \right).$$

For the face-centered point value at $(x_{i+\frac{1}{2}}, y_j)$, the FVS-based update reads

$$\frac{d\mathbf{U}_{i+\frac{1}{2},j}}{dt} = - \left[\widetilde{\mathbf{D}}_1^+ \mathbf{F}_1^+(\mathbf{U}) + \widetilde{\mathbf{D}}_1^- \mathbf{F}_1^-(\mathbf{U}) \right]_{i+\frac{1}{2},j} - \left(\widetilde{\mathbf{D}}_2 \mathbf{F}_2(\mathbf{U}) \right)_{i+\frac{1}{2},j},$$

294 where

$$\begin{aligned}
295 \quad & (\tilde{\mathbf{D}}_1^+ \mathbf{F}_1^+)_{i+\frac{1}{2},j} = \frac{1}{\Delta x_i} \left((\mathbf{F}_1^+)_{i-\frac{1}{2},j} - 4(\mathbf{F}_1^+)_{i,j} + 3(\mathbf{F}_1^+)_{i+\frac{1}{2},j} \right), \\
296 \quad & (\tilde{\mathbf{D}}_1^- \mathbf{F}_1^-)_{i+\frac{1}{2},j} = \frac{1}{\Delta x_{i+1}} \left(-3(\mathbf{F}_1^-)_{i+\frac{1}{2},j} + 4(\mathbf{F}_1^-)_{i+1,j} - (\mathbf{F}_1^-)_{i+\frac{3}{2},j} \right), \\
297 \quad & (\tilde{\mathbf{D}}_2 \mathbf{F}_2)_{i+\frac{1}{2},j} = \frac{1}{\Delta y_j} \left((\mathbf{F}_2)_{i+\frac{1}{2},j+\frac{1}{2}} - (\mathbf{F}_2)_{i+\frac{1}{2},j-\frac{1}{2}} \right),
\end{aligned}$$

298 and the cell-centered point value is computed from the bi-parabolic reconstruction [3]
299 as

$$\begin{aligned}
300 \quad (3.11) \quad & \mathbf{U}_{i,j} = \frac{1}{16} \left[36\bar{\mathbf{U}}_{i,j} - 4 \left(\mathbf{U}_{i-\frac{1}{2},j} + \mathbf{U}_{i+\frac{1}{2},j} + \mathbf{U}_{i,j-\frac{1}{2}} + \mathbf{U}_{i,j+\frac{1}{2}} \right) \right. \\
301 \quad (3.12) \quad & \left. - \left(\mathbf{U}_{i-\frac{1}{2},j-\frac{1}{2}} + \mathbf{U}_{i+\frac{1}{2},j-\frac{1}{2}} + \mathbf{U}_{i-\frac{1}{2},j+\frac{1}{2}} + \mathbf{U}_{i+\frac{1}{2},j+\frac{1}{2}} \right) \right].
\end{aligned}$$

302 The update for the point value at $(x_i, y_{j+\frac{1}{2}})$ is omitted here, which is similar to (3.9).

303 **3.2. Mesh alignment issue when using Jacobian splitting.** The mesh
304 alignment issue was observed for the fully-discrete AF methods in [34], where the
305 convergence rate reduces to 2 for the linear advection problem, when the advection
306 velocity is aligned with the grid. For the generalized AF methods based on the JS,
307 such an issue is also observed. Consider [Example 5.8](#), where we solve a quasi-2D Sod
308 shock tube along the x -direction on a 100×2 uniform mesh. As shown in [Figure 10](#),
309 the density based on the JS shows large deviations between the contact discontinuity
310 and shock wave. From [Figure 11](#), it can be seen that the solutions of the DoFs
311 at the corner $(x_{i+\frac{1}{2}}, y_{j+\frac{1}{2}})$ and horizontal face $(x_i, y_{j+\frac{1}{2}})$ are decoupled from that at
312 the vertical face $(x_{i+\frac{1}{2}}, y_j)$ and cell averages. The reason is complicated because the
313 mesh alignment issue seems to be caused by the decoupled point value update and its
314 interaction with the JS.

315 **3.3. Boundary treatment.** The numerical boundary conditions can be imple-
316 mented using ghost cells as usual finite volume methods. Take the reflective boundary
317 for the Euler equations as an example. Let $x = x_{N_1-\frac{1}{2}}$ be the boundary, then the cell
318 averages and point values in the ghost cell $I_{N_1,j}$ are given by

$$\begin{aligned}
319 \quad & \bar{\mathbf{U}}_{N_1,j} = \mathcal{M}(\bar{\mathbf{U}}_{N_1-1,j}), \quad \mathbf{U}_{N_1+\frac{1}{2},j} = \mathcal{M}(\mathbf{U}_{N_1-\frac{3}{2},j}), \\
320 \quad & \mathbf{U}_{N_1,j-\frac{1}{2}} = \mathcal{M}(\mathbf{U}_{N_1-1,j-\frac{1}{2}}), \quad \mathbf{U}_{N_1+\frac{1}{2},j-\frac{1}{2}} = \mathcal{M}(\mathbf{U}_{N_1-\frac{3}{2},j-\frac{1}{2}}),
\end{aligned}$$

321 where \mathcal{M} reverses the sign of the ρv_1 component while keeping others unchanged.
322 Then the point value update at the boundary can be computed in the same way as the
323 interior points, but the numerical flux on the boundary for the cell average is computed
324 through the LLF flux as suggested in [4]. For instance, the flux $\mathbf{F}_1(\mathbf{U}_{N_1-\frac{1}{2},j-\frac{1}{2}})$ in the
325 right-hand side of (3.5) is replaced by $\widehat{\mathbf{F}}_1^{\text{LLF}}(\mathbf{U}_{N_1-1,j-\frac{1}{2}}, \mathcal{M}(\mathbf{U}_{N_1-1,j-\frac{1}{2}}))$.

326 **4. 2D bound-preserving active flux methods.** In this paper, the admissible
327 state set \mathcal{G} is assumed to be convex. Two cases are considered. For the scalar con-
328 servation law (3.2), its solutions satisfy a strict maximum principle (MP) [14], i.e.,
329

$$330 \quad (4.1) \quad \mathcal{G} = \{u \mid m_0 \leq u \leq M_0\}, \quad m_0 = \min_{x,y} u_0(x,y), \quad M_0 = \max_{x,y} u_0(x,y).$$

331 For the compressible Euler equations, the admissible state set is

$$332 \quad (4.2) \quad \mathcal{G} = \left\{ \mathbf{U} = (\rho, \rho \mathbf{v}, E) \mid \rho > 0, p = (\gamma - 1) \left(E - \|\rho \mathbf{v}\|^2 / (2\rho) \right) > 0 \right\},$$

333 which is convex, see e.g. [47].

334 **DEFINITION 4.1.** *An AF method is called bound-preserving (BP) if starting from*
 335 *cell averages and point values in the admissible state set \mathcal{G} , the cell averages and point*
 336 *values remain in \mathcal{G} at the next time step.*

337 Note that to avoid the effect of the round-off error, we need to choose the desired lower
 338 bounds for the density and pressure. In the numerical tests, we will enforce $\rho \geq \varepsilon^\rho$,
 339 $p \geq \varepsilon^p$ with $\varepsilon^\rho, \varepsilon^p$ to be defined later. Since the DoFs in the AF methods include both
 340 cell averages and point values, it is necessary to design suitable BP limitings for both
 341 of them to achieve the BP property. The limiting for the cell average has not been
 342 addressed much in the literature, except for a very recent work [4]. The 1D limitings
 343 can be reduced from this Section, given in Section SM3 in the supplementary material.

344 **4.1. Convex limiting for the cell average.** This section presents a convex
 345 limiting approach to achieve the BP property of the cell average update. The basic
 346 idea of the convex limiting approaches [21, 25, 30] is to enforce the preservation of local
 347 or global bounds by constraining individual numerical fluxes. The BP or invariant
 348 domain-preserving (IDP) properties of flux-limited approximations are shown using
 349 representations in terms of intermediate states that stay in convex admissible state
 350 sets [21, 24]. The low-order scheme is chosen as the first-order LLF scheme

$$351 \quad \bar{\mathbf{U}}_{i,j}^L = \bar{\mathbf{U}}_{i,j}^n - \mu_{1,i} \left(\widehat{\mathbf{F}}_{i+\frac{1}{2},j}^L - \widehat{\mathbf{F}}_{i-\frac{1}{2},j}^L \right) - \mu_{2,j} \left(\widehat{\mathbf{F}}_{i,j+\frac{1}{2}}^L - \widehat{\mathbf{F}}_{i,j-\frac{1}{2}}^L \right),$$

352 where $\widehat{\mathbf{F}}_{i+\frac{1}{2},j}^L$ and $\widehat{\mathbf{F}}_{i,j+\frac{1}{2}}^L$ are the LLF fluxes. Take the x -direction as an example,

$$353 \quad (4.3) \quad \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^L := \widehat{\mathbf{F}}_1^{\text{LLF}}(\bar{\mathbf{U}}_{i,j}^n, \bar{\mathbf{U}}_{i+1,j}^n)$$

$$354 \quad = \frac{1}{2} \left(\mathbf{F}_1(\bar{\mathbf{U}}_{i,j}^n) + \mathbf{F}_1(\bar{\mathbf{U}}_{i+1,j}^n) \right) - \frac{(\alpha_1)_{i+\frac{1}{2},j}}{2} \left(\bar{\mathbf{U}}_{i+1,j}^n - \bar{\mathbf{U}}_{i,j}^n \right),$$

$$355 \quad (\alpha_1)_{i+\frac{1}{2},j} = \max\{\varrho_1(\bar{\mathbf{U}}_{i,j}^n), \varrho_1(\bar{\mathbf{U}}_{i+1,j}^n)\},$$

$$356 \quad \mu_{1,i} = \Delta t^n / \Delta x_i,$$

357 where ϱ_1 is the spectral radius of $\partial \mathbf{F}_1 / \partial \mathbf{U}$. Note that here $\alpha_{i+\frac{1}{2},j}$ is not the same as
 358 the one in the LLF FVS. Following [22], the first-order LLF scheme can be rewritten
 359 as

$$360 \quad \bar{\mathbf{U}}_{i,j}^L = \left[1 - \mu_{1,i} \left((\alpha_1)_{i-\frac{1}{2},j} + (\alpha_1)_{i+\frac{1}{2},j} \right) - \mu_{2,j} \left((\alpha_2)_{i,j-\frac{1}{2}} + (\alpha_2)_{i,j+\frac{1}{2}} \right) \right] \bar{\mathbf{U}}_{i,j}^n$$

$$361 \quad + \mu_{1,i} (\alpha_1)_{i-\frac{1}{2},j} \tilde{\mathbf{U}}_{i-\frac{1}{2},j} + \mu_{1,i} (\alpha_1)_{i+\frac{1}{2},j} \tilde{\mathbf{U}}_{i+\frac{1}{2},j}$$

$$362 \quad (4.4) \quad + \mu_{2,j} (\alpha_2)_{i,j-\frac{1}{2}} \tilde{\mathbf{U}}_{i,j-\frac{1}{2}} + \mu_{2,j} (\alpha_2)_{i,j+\frac{1}{2}} \tilde{\mathbf{U}}_{i,j+\frac{1}{2}},$$

363 with four intermediate states, and the explicit expressions in the x -direction are

$$364 \quad (4.5) \quad \tilde{\mathbf{U}}_{i\pm\frac{1}{2},j} = \frac{1}{2} \left(\bar{\mathbf{U}}_{i,j}^n + \bar{\mathbf{U}}_{i\pm 1,j}^n \right) \mp \frac{1}{2(\alpha_1)_{i\pm\frac{1}{2},j}} \left[\mathbf{F}_1(\bar{\mathbf{U}}_{i,j}^n) - \mathbf{F}_1(\bar{\mathbf{U}}_{i\pm 1,j}^n) \right].$$

365 The proofs of $\tilde{\mathbf{U}}_{i\pm\frac{1}{2},j}, \tilde{\mathbf{U}}_{i,j\pm\frac{1}{2}} \in \mathcal{G}$ are given in the supplementary material Section SM2,
 366 for the scalar case and Euler equations.

367 LEMMA 4.2. *If the time step size Δt^n satisfies*

$$368 \quad (4.6) \quad \Delta t^n \leq \frac{1}{2} \min \left\{ \frac{\Delta x_i}{(\alpha_1)_{i-\frac{1}{2},j} + (\alpha_1)_{i+\frac{1}{2},j}}, \frac{\Delta y_j}{(\alpha_2)_{i,j-\frac{1}{2}} + (\alpha_2)_{i,j+\frac{1}{2}}} \right\},$$

369 *then (4.4) is a convex combination, and the first-order LLF scheme is BP.*

370 *The proof (see e.g. [22, 35]) relies on $\bar{U}_{i,j}^n, \tilde{U}_{i\pm\frac{1}{2},j}, \tilde{U}_{i,j\pm\frac{1}{2}} \in \mathcal{G}$ and the convexity of \mathcal{G} .*

371 Upon defining the anti-diffusive flux $\Delta \widehat{F}_{i\pm\frac{1}{2},j} = \widehat{F}_{i\pm\frac{1}{2},j}^H - \widehat{F}_{i\pm\frac{1}{2},j}^L$, and $\widehat{F}_{i\pm\frac{1}{2},j}^H$ is given
372 in (3.4), a forward-Euler step applied to the semi-discrete high-order scheme for the
373 cell average (3.3) can be written as

$$\begin{aligned} 374 \quad \bar{U}_{i,j}^H &= \bar{U}_{i,j}^n - \mu_{1,i} \left(\widehat{F}_{i+\frac{1}{2},j}^L - \widehat{F}_{i-\frac{1}{2},j}^L \right) - \mu_{2,j} \left(\widehat{F}_{i,j+\frac{1}{2}}^L - \widehat{F}_{i,j-\frac{1}{2}}^L \right) \\ 375 &\quad - \mu_{1,i} \left(\Delta \widehat{F}_{i+\frac{1}{2},j} - \Delta \widehat{F}_{i-\frac{1}{2},j} \right) - \mu_{2,j} \left(\Delta \widehat{F}_{i,j+\frac{1}{2}} - \Delta \widehat{F}_{i,j-\frac{1}{2}} \right) \\ 376 &= \left[1 - \mu_{1,i} \left((\alpha_1)_{i-\frac{1}{2},j} + (\alpha_1)_{i+\frac{1}{2},j} \right) - \mu_{2,j} \left((\alpha_2)_{i,j-\frac{1}{2}} + (\alpha_2)_{i,j+\frac{1}{2}} \right) \right] \bar{U}_{i,j}^n \\ 377 &\quad + \mu_{1,i} (\alpha_1)_{i-\frac{1}{2},j} \tilde{U}_{i-\frac{1}{2},j}^{H,+} + \mu_{1,i} (\alpha_1)_{i+\frac{1}{2},j} \tilde{U}_{i+\frac{1}{2},j}^{H,-} \\ 378 \quad (4.7) &\quad + \mu_{2,j} (\alpha_2)_{i,j-\frac{1}{2}} \tilde{U}_{i,j-\frac{1}{2}}^{H,+} + \mu_{2,j} (\alpha_2)_{i,j+\frac{1}{2}} \tilde{U}_{i,j+\frac{1}{2}}^{H,-}, \end{aligned}$$

379 with the high-order intermediate states

$$380 \quad \tilde{U}_{i\pm\frac{1}{2},j}^{H,\mp} := \tilde{U}_{i\pm\frac{1}{2},j} \mp \frac{\Delta \widehat{F}_{i\pm\frac{1}{2},j}^H}{(\alpha_1)_{i\pm\frac{1}{2},j}}, \quad \tilde{U}_{i,j\pm\frac{1}{2}}^{H,\mp} := \tilde{U}_{i,j\pm\frac{1}{2}} \mp \frac{\Delta \widehat{F}_{i,j\pm\frac{1}{2}}^H}{(\alpha_2)_{i,j\pm\frac{1}{2}}}.$$

381 With the low-order scheme (4.4) and high-order scheme (4.7) having the same abstract
382 form, one can blend them to define the limited scheme for the cell average as

$$\begin{aligned} 383 \quad \bar{U}_{i,j}^{\text{Lim}} &= \left[1 - \mu_{1,i} \left((\alpha_1)_{i-\frac{1}{2},j} + (\alpha_1)_{i+\frac{1}{2},j} \right) - \mu_{2,j} \left((\alpha_2)_{i,j-\frac{1}{2}} + (\alpha_2)_{i,j+\frac{1}{2}} \right) \right] \bar{U}_{i,j}^n \\ 384 &\quad + \mu_{1,i} (\alpha_1)_{i-\frac{1}{2},j} \tilde{U}_{i-\frac{1}{2},j}^{\text{Lim},+} + \mu_{1,i} (\alpha_1)_{i+\frac{1}{2},j} \tilde{U}_{i+\frac{1}{2},j}^{\text{Lim},-} \\ 385 \quad (4.8) &\quad + \mu_{2,j} (\alpha_2)_{i,j-\frac{1}{2}} \tilde{U}_{i,j-\frac{1}{2}}^{\text{Lim},+} + \mu_{2,j} (\alpha_2)_{i,j+\frac{1}{2}} \tilde{U}_{i,j+\frac{1}{2}}^{\text{Lim},-}, \end{aligned}$$

386 where the limited intermediate states are

$$\begin{aligned} 387 \quad (4.9) \quad \tilde{U}_{i\pm\frac{1}{2},j}^{\text{Lim},\mp} &= \tilde{U}_{i\pm\frac{1}{2},j} \mp \frac{\Delta \widehat{F}_{i\pm\frac{1}{2},j}^{\text{Lim}}}{(\alpha_1)_{i\pm\frac{1}{2},j}} := \tilde{U}_{i\pm\frac{1}{2},j} \mp \frac{\theta_{i\pm\frac{1}{2},j} \Delta \widehat{F}_{i\pm\frac{1}{2},j}^H}{(\alpha_1)_{i\pm\frac{1}{2},j}}, \\ \tilde{U}_{i,j\pm\frac{1}{2}}^{\text{Lim},\mp} &= \tilde{U}_{i,j\pm\frac{1}{2}} \mp \frac{\Delta \widehat{F}_{i,j\pm\frac{1}{2}}^{\text{Lim}}}{(\alpha_2)_{i,j\pm\frac{1}{2}}} := \tilde{U}_{i,j\pm\frac{1}{2}} \mp \frac{\theta_{i,j\pm\frac{1}{2}} \Delta \widehat{F}_{i,j\pm\frac{1}{2}}^H}{(\alpha_2)_{i,j\pm\frac{1}{2}}}, \end{aligned}$$

388 and $\theta_{i\pm\frac{1}{2},j}, \theta_{i,j\pm\frac{1}{2}} \in [0, 1]$ are the blending coefficients. The limited scheme (4.8) re-
389 duces to the first-order LLF scheme if $\theta_{i\pm\frac{1}{2},j} = \theta_{i,j\pm\frac{1}{2}} = 0$, and recovers the high-order
390 AF scheme (3.3) when $\theta_{i\pm\frac{1}{2},j} = \theta_{i,j\pm\frac{1}{2}} = 1$.

391 PROPOSITION 4.3. *If the cell average at the last time step $\bar{U}_{i,j}^n$ and the limited*
392 *intermediate states $\tilde{U}_{i\pm\frac{1}{2},j}^{\text{Lim},\mp}, \tilde{U}_{i,j\pm\frac{1}{2}}^{\text{Lim},\mp}$ belong to the admissible state set \mathcal{G} , then the limited*
393 *average update (4.8) is BP, i.e., $\bar{U}_{i,j}^{\text{Lim}} \in \mathcal{G}$, under the CFL condition (4.6). If the SSP-*
394 *RK3 (2.14) is used for the time integration, the high-order scheme is also BP.*

395 *Proof.* Under the constraint (4.6), the limited cell average update $\bar{U}_{i,j}^{\text{Lim}}$ is a convex
 396 combination of $\bar{U}_{i,j}^n$, $\tilde{U}_{i\pm\frac{1}{2},j}^{\text{Lim},\mp}$, and $\tilde{U}_{i,j\pm\frac{1}{2}}^{\text{Lim},\mp}$, thus it belongs to \mathcal{G} due to the convexity of \mathcal{G} .
 397 Because the SSP-RK3 is a convex combination of forward-Euler stages, the high-order
 398 scheme equipped with the SSP-RK3 is also BP according to the convexity. \square

399 *Remark 4.4.* The scheme (4.8) is conservative as it amounts to using the x -
 400 directional numerical flux $\hat{F}_{i+\frac{1}{2},j}^{\text{L}} + \theta_{i+\frac{1}{2},j} \Delta \hat{F}_{i+\frac{1}{2},j} = \theta_{i+\frac{1}{2},j} \hat{F}_{i+\frac{1}{2},j}^{\text{H}} + (1 - \theta_{i+\frac{1}{2},j}) \hat{F}_{i+\frac{1}{2},j}^{\text{L}}$,
 401 which is a convex combination of the high-order and low-order fluxes.

402 *Remark 4.5.* It should be noted that the time step size (4.6) is determined based
 403 on the solutions at t^n . If the constraint is not satisfied at the later stage of the
 404 SSP-RK3, the BP property may not be achieved because (4.8) is no longer a convex
 405 combination. In our implementation, we start from the usual CFL condition (2.15).
 406 Then, if the high-order AF solutions need BP limitings and (4.5) is not BP or (4.6) is
 407 not satisfied, the numerical solutions are set back to the last time step, and we rerun
 408 with a halved time step size until (4.5) is BP and the constraint (4.6) is satisfied. This
 409 is a typical implementation in other BP methods, e.g. [45].

410 The remaining task is to determine the coefficients at each interface $\theta_{i\pm\frac{1}{2},j}, \theta_{i,j\pm\frac{1}{2}}$
 411 such that $\tilde{U}_{i\pm\frac{1}{2},j}^{\text{Lim},\mp}, \tilde{U}_{i,j\pm\frac{1}{2}}^{\text{Lim},\mp} \in \mathcal{G}$ and stay as close as possible to the high-order solu-
 412 tions $\tilde{U}_{i\pm\frac{1}{2},j}^{\text{H}}, \tilde{U}_{i,j\pm\frac{1}{2}}^{\text{H}}$, i.e., the goal is to find the largest $\theta_{i\pm\frac{1}{2},j}, \theta_{i,j\pm\frac{1}{2}} \in [0, 1]$ such that
 413 $\tilde{U}_{i\pm\frac{1}{2},j}^{\text{Lim},\mp}, \tilde{U}_{i,j\pm\frac{1}{2}}^{\text{Lim},\mp} \in \mathcal{G}$. The explanations will be given for the x -direction.

414 **4.1.1. Application to scalar conservation laws.** This section is devoted to
 415 applying the convex limiting approach to scalar conservation laws (3.2), such that the
 416 limited cell averages (4.8) satisfy the MP $u_{i,j}^{\text{min}} \leq \bar{u}_{i,j}^{\text{Lim}} \leq u_{i,j}^{\text{max}}$, where $u_{i,j}^{\text{min}} = \min \mathcal{N}$,
 417 $u_{i,j}^{\text{max}} = \max \mathcal{N}$, and \mathcal{N} will be defined later. According to the convex decomposition,
 418 the blending coefficient $\theta_{i+\frac{1}{2},j} \in [0, 1]$ or $\Delta \hat{f}_{i+\frac{1}{2},j}^{\text{Lim}}$ should be chosen such that $u_{i,j}^{\text{min}} \leq$
 419 $\tilde{u}_{i+\frac{1}{2},j}^{\text{Lim},-} \leq u_{i,j}^{\text{max}}$, $u_{i+1,j}^{\text{min}} \leq \tilde{u}_{i+\frac{1}{2},j}^{\text{Lim},+} \leq u_{i+1,j}^{\text{max}}$. Solving the first condition, i.e. $u_{i,j}^{\text{min}} \leq \tilde{u}_{i+\frac{1}{2},j} -$
 420 $\Delta \hat{f}_{i+\frac{1}{2},j}^{\text{Lim}} / \alpha_{i+\frac{1}{2},j} \leq u_{i,j}^{\text{max}}$, one has $\Delta \hat{f}_{i+\frac{1}{2},j}^{\text{Lim}} \leq \alpha_{i+\frac{1}{2},j} (\tilde{u}_{i+\frac{1}{2},j} - u_{i,j}^{\text{min}})$ if $\Delta \hat{f}_{i+\frac{1}{2},j} \geq 0$, or
 421 $\Delta \hat{f}_{i+\frac{1}{2},j}^{\text{Lim}} \geq \alpha_{i+\frac{1}{2},j} (\tilde{u}_{i+\frac{1}{2},j} - u_{i,j}^{\text{max}})$ if $\Delta \hat{f}_{i+\frac{1}{2},j} < 0$. Solving the second condition $u_{i+1,j}^{\text{min}} \leq$
 422 $\tilde{u}_{i+\frac{1}{2},j}^{\text{Lim},+} \leq u_{i+1,j}^{\text{max}}$ in the same way and combining the two sets of results yields

$$423 \quad \Delta \hat{f}_{i+\frac{1}{2},j}^{\text{Lim}} = \begin{cases} \min \{ \Delta \hat{f}_{i+\frac{1}{2},j}, \Delta \hat{f}_{i+\frac{1}{2},j}^+ \}, & \text{if } \Delta \hat{f}_{i+\frac{1}{2},j} \geq 0, \\ \max \{ \Delta \hat{f}_{i+\frac{1}{2},j}, \Delta \hat{f}_{i+\frac{1}{2},j}^- \}, & \text{otherwise,} \end{cases}$$

$$424 \quad \Delta \hat{f}_{i+\frac{1}{2},j}^+ = (\alpha_1)_{i+\frac{1}{2},j} \min \{ \tilde{u}_{i+\frac{1}{2},j} - u_{i,j}^{\text{min}}, u_{i+1,j}^{\text{max}} - \tilde{u}_{i+\frac{1}{2},j} \},$$

$$425 \quad \Delta \hat{f}_{i+\frac{1}{2},j}^- = (\alpha_1)_{i+\frac{1}{2},j} \max \{ u_{i+1,j}^{\text{min}} - \tilde{u}_{i+\frac{1}{2},j}, \tilde{u}_{i+\frac{1}{2},j} - u_{i,j}^{\text{max}} \}.$$

426 Finally, the limited numerical flux is

$$427 \quad (4.10) \quad \hat{f}_{i+\frac{1}{2},j}^{\text{Lim}} = \hat{f}_{i+\frac{1}{2},j}^{\text{L}} + \Delta \hat{f}_{i+\frac{1}{2},j}^{\text{Lim}}.$$

428 If considering the global MP, $\mathcal{N} = \cup_{i,j,\sigma} \{ \bar{u}_{i,j}^n, u_{\sigma}^n \}$. One can also enforce the local MP,
 429 which helps to suppress spurious oscillations [21, 31, 22], by choosing

$$430 \quad \mathcal{N} = \left\{ \bar{u}_{i,j}^n, \tilde{u}_{i-\frac{1}{2},j}, \tilde{u}_{i+\frac{1}{2},j}, \tilde{u}_{i,j-\frac{1}{2}}, \tilde{u}_{i,j+\frac{1}{2}}, \bar{u}_{i-1,j}^n, \bar{u}_{i+1,j}^n, \bar{u}_{i,j-1}^n, \bar{u}_{i,j+1}^n \right\},$$

431 which includes the intermediate states and neighboring cell averages.

432 **4.1.2. Application to the compressible Euler equations.** This section aims
 433 at enforcing the positivity of density and pressure. To avoid the effect of the round-
 434 off error, we need to choose the desired lower bounds. Denote the lowest density and
 435 pressure in the domain by

$$436 \quad (4.11) \quad \varepsilon^\rho := \min\{\bar{\mathbf{U}}_{i,j}^{n,\rho}, \mathbf{U}_\sigma^{n,\rho}\}, \quad \varepsilon^p := \min\{p(\bar{\mathbf{U}}_{i,j}^n), p(\mathbf{U}_\sigma^n)\},$$

437 where $\mathbf{U}^{*\cdot\rho}$ and $p(\mathbf{U}^*)$ denote the density component and pressure recovered from \mathbf{U}^* ,
 438 respectively, and σ denotes the locations of point values in the DoFs. The limiting
 439 (4.9) is feasible if the constraints are satisfied by the first-order LLF intermediate
 440 states (4.5), thus the lower bounds can be defined as

$$441 \quad \varepsilon_{i,j}^\rho := \min\{10^{-13}, \varepsilon^\rho, \tilde{\mathbf{U}}_{i-\frac{1}{2},j}^\rho, \tilde{\mathbf{U}}_{i+\frac{1}{2},j}^\rho, \tilde{\mathbf{U}}_{i,j-\frac{1}{2}}^\rho, \tilde{\mathbf{U}}_{i,j+\frac{1}{2}}^\rho\},$$

$$442 \quad \varepsilon_{i,j}^p := \min\{10^{-13}, \varepsilon^p, p(\tilde{\mathbf{U}}_{i-\frac{1}{2},j}), p(\tilde{\mathbf{U}}_{i+\frac{1}{2},j}), p(\tilde{\mathbf{U}}_{i,j-\frac{1}{2}}), p(\tilde{\mathbf{U}}_{i,j+\frac{1}{2}})\}.$$

443
 444 i) **Positivity of density.** The first step is to impose the density positivity
 445 $\tilde{\mathbf{U}}_{i+\frac{1}{2},j}^{\text{Lim},\pm,\rho} \geq \bar{\varepsilon}_{i+\frac{1}{2},j}^\rho := \min\{\varepsilon_{i,j}^\rho, \varepsilon_{i+1,j}^\rho\}$. Similarly to the derivation of the scalar case, the
 446 corresponding density component of the limited anti-diffusive flux is

$$447 \quad \Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},\rho} = \begin{cases} \min\left\{\Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^\rho, (\alpha_1)_{i+\frac{1}{2},j} \left(\tilde{\mathbf{U}}_{i+\frac{1}{2},j}^\rho - \bar{\varepsilon}_{i+\frac{1}{2},j}^\rho\right)\right\}, & \text{if } \Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^\rho \geq 0, \\ \max\left\{\Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^\rho, (\alpha_1)_{i+\frac{1}{2},j} \left(\bar{\varepsilon}_{i+\frac{1}{2},j}^\rho - \tilde{\mathbf{U}}_{i+\frac{1}{2},j}^\rho\right)\right\}, & \text{otherwise.} \end{cases}$$

448 Then the density component of the limited numerical flux is $\widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},*\cdot\rho} = \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{L},\rho} + \Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},\rho}$,
 449 with the other components remaining the same as $\widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{H}}$.

450 ii) **Positivity of pressure.** The second step is to enforce pressure positivity
 451 $p(\tilde{\mathbf{U}}_{i+\frac{1}{2},j}^{\text{Lim},\pm}) \geq \bar{\varepsilon}_{i+\frac{1}{2},j}^p := \min\{\varepsilon_{i,j}^p, \varepsilon_{i+1,j}^p\}$. Since

$$452 \quad \tilde{\mathbf{U}}_{i+\frac{1}{2},j}^{\text{Lim},\pm} = \tilde{\mathbf{U}}_{i+\frac{1}{2},j} \pm \frac{\theta_{i+\frac{1}{2},j} \Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},*}}{\alpha_{i+\frac{1}{2},j}}, \quad \Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},*} = \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},*} - \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{L}},$$

453 the constraints lead to two inequalities after some algebraic operations

$$454 \quad (4.12) \quad A_{i+\frac{1}{2},j} \theta_{i+\frac{1}{2},j}^2 \pm B_{i+\frac{1}{2},j} \theta_{i+\frac{1}{2},j} \leq C_{i+\frac{1}{2},j},$$

455 with the coefficients (the subscript $(\cdot)_{i+\frac{1}{2},j}$ is omitted in the right-hand side)

$$456 \quad A_{i+\frac{1}{2},j} = \frac{1}{2} \left\| \Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},*\cdot\rho v} \right\|_2^2 - \Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},*\cdot\rho} \Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},*E},$$

$$457 \quad B_{i+\frac{1}{2},j} = \alpha_1 \left(\Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},*\cdot\rho} \tilde{\mathbf{U}}^E + \tilde{\mathbf{U}}^\rho \Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},*E} - \Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},*\cdot\rho v} \cdot \tilde{\mathbf{U}}^{\rho v} - \bar{\varepsilon} \Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},*\cdot\rho} \right),$$

$$458 \quad C_{i+\frac{1}{2},j} = \alpha_1^2 \left(\tilde{\mathbf{U}}^\rho \tilde{\mathbf{U}}^E - \frac{1}{2} \left\| \tilde{\mathbf{U}}^{\rho v} \right\|_2^2 - \bar{\varepsilon} \tilde{\mathbf{U}}^\rho \right), \quad \bar{\varepsilon} = \bar{\varepsilon}_{i+\frac{1}{2},j}^p / (\gamma - 1).$$

459 Following [30], the inequalities (4.12) hold under the linear sufficient condition

$$460 \quad \left(\max\{0, A_{i+\frac{1}{2},j}\} + |B_{i+\frac{1}{2},j}| \right) \theta_{i+\frac{1}{2},j} \leq C_{i+\frac{1}{2},j},$$

461 if making use of $\theta_{i+\frac{1}{2},j}^2 \leq \theta_{i+\frac{1}{2},j}$, $\theta_{i+\frac{1}{2},j} \in [0,1]$. Thus the coefficient can be chosen as

$$462 \quad \theta_{i+\frac{1}{2},j} = \min \left\{ 1, \frac{C_{i+\frac{1}{2},j}}{\max\{0, A_{i+\frac{1}{2},j}\} + |B_{i+\frac{1}{2},j}|} \right\},$$

463 and the final limited numerical flux is

$$464 \quad (4.13) \quad \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},**} = \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{L}} + \theta_{i+\frac{1}{2},j} \Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},*}.$$

465 **4.1.3. Shock sensor-based limiting.** Spurious oscillations are observed, espe-
466 cially near strong shock waves, if only the BP limitings are employed, see [Exam-](#)
467 [ple 5.10](#). We propose to further limit the numerical fluxes using another parameter
468 $\theta_{i+\frac{1}{2},j}^s$ based on shock sensors. Consider the Jameson's shock sensor in [29],

$$469 \quad (\varphi_1)_{i,j} = \frac{|\bar{p}_{i+1,j} - 2\bar{p}_{i,j} + \bar{p}_{i-1,j}|}{|\bar{p}_{i+1,j} + 2\bar{p}_{i,j} + \bar{p}_{i-1,j}|},$$

470 and a modified Ducros' shock sensor [15]

$$471 \quad (\varphi_2)_{i,j} = \max \left\{ \frac{-(\nabla \cdot \bar{\mathbf{v}})_{i,j}}{\sqrt{(\nabla \cdot \bar{\mathbf{v}})_{i,j}^2 + (\nabla \times \bar{\mathbf{v}})_{i,j}^2 + 10^{-40}}}, 0 \right\},$$

472 where

$$473 \quad (\nabla \cdot \bar{\mathbf{v}})_{i,j} \approx \frac{2((\bar{v}_1)_{i+1,j} - (\bar{v}_1)_{i-1,j})}{\Delta x_i + \Delta x_{i+1}} + \frac{2((\bar{v}_2)_{i,j+1} - (\bar{v}_2)_{i,j-1})}{\Delta y_j + \Delta y_{j+1}},$$

$$474 \quad (\nabla \times \bar{\mathbf{v}})_{i,j} \approx \frac{2((\bar{v}_2)_{i+1,j} - (\bar{v}_2)_{i-1,j})}{\Delta x_i + \Delta x_{i+1}} - \frac{2((\bar{v}_1)_{i,j+1} - (\bar{v}_1)_{i,j-1})}{\Delta y_j + \Delta y_{j+1}},$$

475 with $\bar{v}_{i,j}$ and $\bar{p}_{i,j}$ the velocity and pressure recovered from the cell average $\bar{U}_{i,j}$. We
476 consider the sign of the velocity divergence, such that the shock waves can be located
477 better. The blending coefficient is designed as

$$478 \quad \theta_{i+\frac{1}{2},j}^s = \exp(-\kappa(\varphi_1)_{i+\frac{1}{2},j}(\varphi_2)_{i+\frac{1}{2},j}) \in (0,1],$$

$$479 \quad (\varphi_s)_{i+\frac{1}{2},j} = \max\{(\varphi_s)_{i,j}, (\varphi_s)_{i+1,j}\}, \quad s = 1, 2,$$

480 where the problem-dependent parameter κ adjusts the strength of the limiting, and
481 its optimal choice needs further investigation. The final limited numerical flux is

$$482 \quad (4.14) \quad \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim}} = \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{L}} + \theta_{i+\frac{1}{2},j}^s \Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},**},$$

483 with $\Delta \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},**} = \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},**} - \widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{L}}$, and $\widehat{\mathbf{F}}_{i+\frac{1}{2},j}^{\text{Lim},**}$ given in (4.13).

484 **4.2. Scaling limiter for point value.** To achieve the BP property, it is also
485 necessary to introduce BP limiting for the point value, because using the BP limiting
486 for cell average alone cannot guarantee the bounds, see [Example 5.5](#). As there is
487 no conservation requirement on the point value update, a simple scaling limiter [32]
488 is directly performed on the high-order solution rather than on the flux for the cell
489 average.

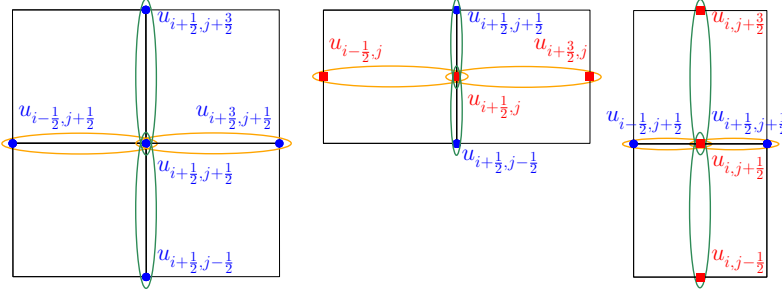


Fig. 2: The stencils for the first-order LLF schemes.

490 The first step is to define suitable first-order LLF schemes. The stencils are shown
 491 in Figure 2.

492 For the point value at the corner, one can choose

$$493 \quad \mathbf{U}_{i+\frac{1}{2}, j+\frac{1}{2}}^{\text{L}} = \mathbf{U}_{i+\frac{1}{2}, j+\frac{1}{2}}^{\text{n}} - \frac{2\Delta t^n}{\Delta x_i + \Delta x_{i+1}} \left(\widehat{\mathbf{F}}_{i+1, j+\frac{1}{2}}^{\text{L}} - \widehat{\mathbf{F}}_{i, j+\frac{1}{2}}^{\text{L}} \right)$$

$$494 \quad (4.15) \quad - \frac{2\Delta t^n}{\Delta y_j + \Delta y_{j+1}} \left(\widehat{\mathbf{F}}_{i+\frac{1}{2}, j+1}^{\text{L}} - \widehat{\mathbf{F}}_{i+\frac{1}{2}, j}^{\text{L}} \right),$$

495 with the LLF numerical fluxes

$$496 \quad \widehat{\mathbf{F}}_{i+1, j+\frac{1}{2}}^{\text{L}} := \widehat{\mathbf{F}}_1^{\text{LLF}}(\mathbf{U}_{i+\frac{1}{2}, j+\frac{1}{2}}^{\text{n}}, \mathbf{U}_{i+\frac{3}{2}, j+\frac{1}{2}}^{\text{n}}), \quad \widehat{\mathbf{F}}_{i+\frac{1}{2}, j+1}^{\text{L}} := \widehat{\mathbf{F}}_2^{\text{LLF}}(\mathbf{U}_{i+\frac{1}{2}, j+\frac{1}{2}}^{\text{n}}, \mathbf{U}_{i+\frac{1}{2}, j+\frac{3}{2}}^{\text{n}}).$$

497 Note that the x -directional LLF flux has been used in (4.3). For the vertical face-
 498 centered point value, we choose the first-order LLF scheme as

$$499 \quad (4.16) \quad \mathbf{U}_{i+\frac{1}{2}, j}^{\text{L}} = \mathbf{U}_{i+\frac{1}{2}, j}^{\text{n}} - \frac{2\Delta t^n}{\Delta x_i + \Delta x_{i+1}} \left(\widehat{\mathbf{F}}_{i+1, j}^{\text{L}} - \widehat{\mathbf{F}}_{i, j}^{\text{L}} \right) - \frac{\Delta t^n}{\Delta y_j} \left(\widehat{\mathbf{F}}_{i+\frac{1}{2}, j+\frac{1}{2}}^{\text{L}} - \widehat{\mathbf{F}}_{i+\frac{1}{2}, j-\frac{1}{2}}^{\text{L}} \right),$$

500 with the LLF numerical fluxes

$$501 \quad \widehat{\mathbf{F}}_{i+1, j}^{\text{L}} := \widehat{\mathbf{F}}_1^{\text{LLF}}(\mathbf{U}_{i+\frac{1}{2}, j}^{\text{n}}, \mathbf{U}_{i+\frac{3}{2}, j}^{\text{n}}), \quad \widehat{\mathbf{F}}_{i+\frac{1}{2}, j+\frac{1}{2}}^{\text{L}} := \widehat{\mathbf{F}}_2^{\text{LLF}}(\mathbf{U}_{i+\frac{1}{2}, j}^{\text{n}}, \mathbf{U}_{i+\frac{1}{2}, j+\frac{1}{2}}^{\text{n}}).$$

502 The LLF scheme for the face-centered value on the horizontal face can be chosen as

$$503 \quad (4.17) \quad \mathbf{U}_{i, j+\frac{1}{2}}^{\text{L}} = \mathbf{U}_{i, j+\frac{1}{2}}^{\text{n}} - \frac{\Delta t^n}{\Delta x_i} \left(\widehat{\mathbf{F}}_{i+\frac{1}{2}, j+\frac{1}{2}}^{\text{L}} - \widehat{\mathbf{F}}_{i-\frac{1}{2}, j+\frac{1}{2}}^{\text{L}} \right) - \frac{2\Delta t^n}{\Delta y_j + \Delta y_{j+1}} \left(\widehat{\mathbf{F}}_{i, j+1}^{\text{L}} - \widehat{\mathbf{F}}_{i, j}^{\text{L}} \right),$$

504 with similarly defined LLF numerical fluxes as for the vertical face.

505 Similarly to Lemma 4.2, it is straightforward to obtain the following Lemma.

506 LEMMA 4.6. *The LLF schemes (4.15)-(4.17) for the point value update are BP*
 507 *under the following time step size constraint*

$$508 \quad \Delta t^n \leq \frac{1}{2} \min \left\{ \frac{\Delta x_i + \Delta x_{i+1}}{2 \left((\alpha_1)_{i, j+\frac{1}{2}} + (\alpha_1)_{i+1, j+\frac{1}{2}} \right)}, \frac{\Delta y_j + \Delta y_{j+1}}{2 \left((\alpha_2)_{i+\frac{1}{2}, j} + (\alpha_2)_{i+\frac{1}{2}, j+1} \right)}, \right.$$

$$509 \quad \frac{\Delta x_i + \Delta x_{i+1}}{2 \left((\alpha_1)_{i, j} + (\alpha_1)_{i+1, j} \right)}, \frac{\Delta y_j}{\left(\alpha_2)_{i+\frac{1}{2}, j+\frac{1}{2}} + (\alpha_2)_{i+\frac{1}{2}, j-\frac{1}{2}} \right)},$$

$$510 \quad (4.18) \quad \left. \frac{\Delta x_i}{\left(\alpha_1)_{i+\frac{1}{2}, j+\frac{1}{2}} + (\alpha_1)_{i-\frac{1}{2}, j+\frac{1}{2}} \right)}, \frac{\Delta y_j + \Delta y_{j+1}}{2 \left((\alpha_2)_{i, j} + (\alpha_2)_{i, j+1} \right)} \right\},$$

511 where $(\alpha_1)_*$ and $(\alpha_2)_*$ are the viscosity coefficients in the LLF schemes.

512 The limited solution is obtained by blending the high-order AF scheme (3.6) with
513 the forward-Euler scheme and the LLF schemes (4.15)-(4.17) as $\mathbf{U}_\sigma^{\text{Lim}} = \theta_\sigma \mathbf{U}_\sigma^{\text{H}} + (1 -$
514 $\theta_\sigma) \mathbf{U}_\sigma^{\text{L}}$, such that $\mathbf{U}_\sigma^{\text{Lim}} \in \mathcal{G}$.

515 *Remark 4.7.* In the FVS, the cell-centered value obtained based on Simpson's rule
516 $\mathbf{U}_i = (-\mathbf{U}_{i-\frac{1}{2}} + 6\bar{\mathbf{U}}_i - \mathbf{U}_{i+\frac{1}{2}})/4$ in 1D or (3.11) in 2D is not a convex combination, thus
517 it is possible that $\mathbf{U}_i, \mathbf{U}_{i,j} \notin \mathcal{G}$. For the scalar case, it does not affect the BP property.
518 However, for the Euler equations, the computation of \mathbf{F}_i (resp. $(\mathbf{F}_\ell)_{i,j}$) requires that
519 $\mathbf{U}_i \in \mathcal{G}$ (resp. $\mathbf{U}_{i,j} \in \mathcal{G}$), thus the scaling limiter [45] is applied in the cell I_i (resp.
520 $I_{i,j}$), a procedure also mentioned in [10]. See more details in Remark 4.8.

521 **4.2.1. Application to scalar conservation laws.** This section enforces the
522 MP $u_\sigma^{\min} \leq u_\sigma^{\text{Lim}} \leq u_\sigma^{\max}$ using the scaling limiter [44]. The limited solution is

$$523 \quad (4.19) \quad u_\sigma^{\text{Lim}} = \theta_\sigma u_\sigma^{\text{H}} + (1 - \theta_\sigma) u_\sigma^{\text{L}},$$

524 with the coefficient

$$525 \quad \theta_\sigma = \min \left\{ 1, \left| \frac{u_\sigma^{\text{L}} - m_0}{u_\sigma^{\text{L}} - u_\sigma^{\text{H}}} \right|, \left| \frac{M_0 - u_\sigma^{\text{L}}}{u_\sigma^{\text{H}} - u_\sigma^{\text{L}}} \right| \right\}.$$

526

527 The bounds are determined by $u_\sigma^{\min} = \min \mathcal{N}$, $u_\sigma^{\max} = \max \mathcal{N}$, where the set \mathcal{N}
528 consists of all the DoFs in the domain, i.e., $\mathcal{N} = \cup_{i,j,\sigma} \{\bar{u}_{i,j}^n, u_\sigma^n\}$ for the global MP.
529 One can also consider the neighboring DoFs for the local MP. For the point value at
530 the corner $(x_{i+\frac{1}{2}}, y_{j+\frac{1}{2}})$, we choose

$$531 \quad \mathcal{N} = \left\{ u_{i+\frac{1}{2},j+\frac{1}{2}}^n, u_{i-\frac{1}{2},j+\frac{1}{2}}^n, u_{i+\frac{3}{2},j+\frac{1}{2}}^n, u_{i+\frac{1}{2},j-\frac{1}{2}}^n, u_{i+\frac{1}{2},j+\frac{3}{2}}^n \right\},$$

532 which should at least include all the DoFs that appeared in the first-order LLF scheme
533 (4.15). For the point value at the vertical face center $(x_{i+\frac{1}{2}}, y_j)$, similarly we choose

$$534 \quad \mathcal{N} = \left\{ u_{i+\frac{1}{2},j}^n, u_{i-\frac{1}{2},j}^n, u_{i+\frac{3}{2},j}^n, u_{i+\frac{1}{2},j-\frac{1}{2}}^n, u_{i+\frac{1}{2},j+\frac{1}{2}}^n \right\}.$$

535 For the point value at the horizontal face center $(x_i, y_{j+\frac{1}{2}})$, we choose

$$536 \quad \mathcal{N} = \left\{ u_{i,j+\frac{1}{2}}^n, u_{i,j-\frac{1}{2}}^n, u_{i,j+\frac{3}{2}}^n, u_{i-\frac{1}{2},j+\frac{1}{2}}^n, u_{i+\frac{1}{2},j+\frac{1}{2}}^n \right\}.$$

537

538 **4.2.2. Application to the compressible Euler equations.** The limiting con-
539 sists of two steps.

540 **i) Positivity of density.** First, the high-order solution $\mathbf{U}_\sigma^{\text{H}}$ is modified as $\mathbf{U}_\sigma^{\text{Lim},*}$,
541 such that its density component satisfies $\mathbf{U}_\sigma^{\text{Lim},*,\rho} \geq \varepsilon_\sigma^\rho := \min\{10^{-13}, \varepsilon^\rho, \mathbf{U}_\sigma^{\text{L},\rho}\}$ with ε^ρ
542 given in (4.11). Solving the inequality yields

$$543 \quad \theta_\sigma^* = \begin{cases} \frac{\mathbf{U}_\sigma^{\text{L},\rho} - \varepsilon_\sigma^\rho}{\mathbf{U}_\sigma^{\text{L},\rho} - \mathbf{U}_\sigma^{\text{H},\rho}}, & \text{if } \mathbf{U}_\sigma^{\text{H},\rho} < \varepsilon_\sigma^\rho, \\ 1, & \text{otherwise.} \end{cases}$$

544 Then the density component of the limited solution is $\mathbf{U}_\sigma^{\text{Lim},*,\rho} = \theta_\sigma^* \mathbf{U}_\sigma^{\text{H},\rho} + (1 - \theta_\sigma^*) \mathbf{U}_{i+\frac{1}{2}}^{\text{L},\rho}$,
545 with the other components remaining the same as $\mathbf{U}_\sigma^{\text{H}}$.

546 **ii) Positivity of pressure.** Then the limited solution $\mathbf{U}_\sigma^{\text{Lim},*}$ is modified as $\mathbf{U}_\sigma^{\text{Lim}}$,
 547 such that it gives positive pressure, i.e., $p(\mathbf{U}_\sigma^{\text{Lim}}) \geq \varepsilon_\sigma^p := \min\{10^{-13}, \varepsilon^p, p(\mathbf{U}_\sigma^{\text{L}})\}$, with
 548 ε^p given in (4.11). Let the final limited solution be

$$549 \quad (4.20) \quad \mathbf{U}_\sigma^{\text{Lim}} = \theta_\sigma^{**} \mathbf{U}_\sigma^{\text{Lim},*} + (1 - \theta_\sigma^{**}) \mathbf{U}_\sigma^{\text{L}}.$$

550 The pressure is a concave function of the conservative variables (see e.g. [46]), so that
 551 $p(\mathbf{U}_\sigma^{\text{Lim}}) \geq \theta_\sigma^{**} p(\mathbf{U}_\sigma^{\text{Lim},*}) + (1 - \theta_\sigma^{**}) p(\mathbf{U}_\sigma^{\text{L}})$ based on Jensen's inequality and $\mathbf{U}_\sigma^{\text{Lim},*,\rho} > 0$,
 552 $\mathbf{U}_\sigma^{\text{L},\rho} > 0$, $\theta_\sigma^{**} \in [0, 1]$. Thus the coefficient can be chosen as

$$553 \quad \theta_\sigma^{**} = \begin{cases} \frac{p(\mathbf{U}_\sigma^{\text{L}}) - \varepsilon_\sigma^p}{p(\mathbf{U}_\sigma^{\text{L}}) - p(\mathbf{U}_\sigma^{\text{Lim},*})}, & \text{if } p(\mathbf{U}_\sigma^{\text{Lim},*}) < \varepsilon_\sigma^p, \\ 1, & \text{otherwise.} \end{cases}$$

554 *Remark 4.8.* To compute the high-order FVS-based point value update, we should
 555 limit the cell-centered value \mathbf{U}_i in 1D (resp. $\mathbf{U}_{i,j}$ in 2D) at the beginning of each
 556 Runge-Kutta stage. For example, in 2D, we modify $\mathbf{U}_{i,j}$ as $\mathbf{U}_{i,j}^{\text{Lim}} = \theta_{i,j} \mathbf{U}_{i,j} + (1 -$
 557 $\theta_{i,j}) \overline{\mathbf{U}}_{i,j}$ such that

$$558 \quad \mathbf{U}_{i,j}^{\text{Lim},\rho} \geq \min\{10^{-13}, \overline{\mathbf{U}}_{i,j}^\rho\}, \quad p(\mathbf{U}_{i,j}^{\text{Lim}}) \geq \min\{10^{-13}, p(\overline{\mathbf{U}}_{i,j})\}.$$

559 The computation of $\theta_{i,j}$ is similar to the procedure in this section.

560 Let us summarize the main results of the BP AF methods in this paper.

561 **PROPOSITION 4.9.** *If the initial numerical solution $\overline{\mathbf{U}}_{i,j}^0, \mathbf{U}_\sigma^0 \in \mathcal{G}$ for all i, j, σ , and*
 562 *the time step size satisfies (4.6) and (4.18), then the AF methods (3.3)-(3.6) equipped*
 563 *with the SSP-RK3 (2.14) and the BP limitings*

- 564 • (4.10) and (4.19) preserve the maximum principle for scalar case;
- 565 • (4.13) and (4.20) preserve positive density and pressure for the Euler equations.

566 *Remark 4.10.* For uniform meshes, and if taking the maximal spectral radius of
 567 $\partial \mathbf{F}_1 / \partial \mathbf{U}$ and $\partial \mathbf{F}_2 / \partial \mathbf{U}$ in the domain as $\|\varrho_1\|_\infty$ and $\|\varrho_2\|_\infty$, the following CFL condition

$$568 \quad \Delta t^n \leq \frac{1}{4} \min \left\{ \frac{\Delta x}{\|\varrho_1\|_\infty}, \frac{\Delta y}{\|\varrho_2\|_\infty} \right\}$$

569 fulfills the time step size constraints (4.6) and (4.18).

570 **5. Numerical results.** This section presents some numerical tests to verify the
 571 accuracy, BP property, and shock-capturing ability of the proposed BP AF methods.
 572 The adiabatic index is $\gamma = 1.4$ for the Euler equations except for Example 5.11, where
 573 it is $5/3$. In the 2D plots, the numerical solutions are visualized on a refined mesh
 574 with half the mesh size, where the values at the grid points are the cell averages or
 575 point values on the original mesh. Note that the BP limitings naturally reduce some
 576 oscillations. Some additional tests are provided in Section SM4 in the supplementary
 577 material, including a 1D accuracy test for the Euler equations, double rarefaction
 578 problem, blast wave interaction problem using the power law reconstruction [5], and
 579 double Mach reflection problem.

580 *Example 5.1 (Self-steepening shock).* Consider the 1D Burgers' equation $u_t +$
 581 $(\frac{1}{2}u^2)_x = 0$ on the domain $[-1, 1]$ with periodic boundary conditions. The initial
 582 condition is a square wave, $u_0(x) = 2$ if $|x| < 0.2$, otherwise $u_0(x) = -1$.

583 **Figure 3** shows the cell averages and point values at $T = 0.5$ based on different
 584 point value updates with 200 cells. The CFL number is 0.2. The spike generation
 585 when using the JS has been observed in [27], and the reason is also discussed in
 586 **Subsection 2.2**. Such an issue cannot be eliminated by our BP limitings alone, but
 587 can be cured by additionally using the FVS for the point value update. The numerical
 588 solutions based on the FVS agree well with the reference solution when the limitings
 589 are activated.

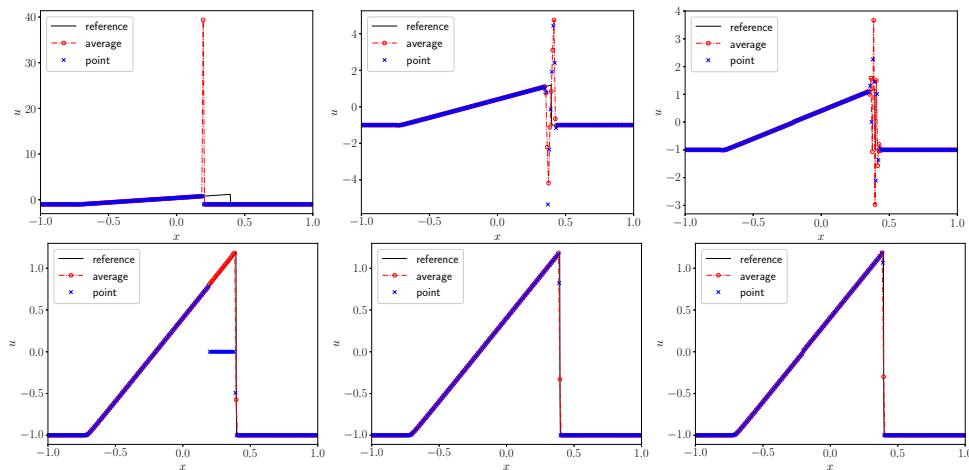


Fig. 3: Example 5.1, self-steepening shock for the Burgers' equation. The numerical solutions computed without limiting (top row) and with the BP limitings imposing the local MP (bottom row). From left to right: JS, LLF, and upwind FVS.

590 **Example 5.2** (LeBlanc shock tube). This is a Riemann problem with an extremely
 591 large initial pressure ratio. This test is solved until $T = 5 \times 10^{-6}$ on the domain $[0, 1]$
 592 with the initial data $(\rho, v, p) = (2, 0, 10^9)$ if $x < 0.5$, otherwise $(\rho, v, p) = (10^{-3}, 0, 1)$.

593 Without the BP limitings, the simulation will stop due to negative density or
 594 pressure. **Figure 4** shows the density computed on a uniform mesh of 6000 cells with
 595 the BP limitings and shock sensor-based limiting. **Note that, the numerical methods**
 596 **typically need a small mesh size to accurately obtain the right location of the shock**
 597 **wave.** The CFL number is 0.4 for the JS, LLF, and SW FVS, and 0.1 for the VH FVS
 598 for stability. The numerical solutions agree well with the exact solution with only a
 599 few undershoots at the discontinuities.

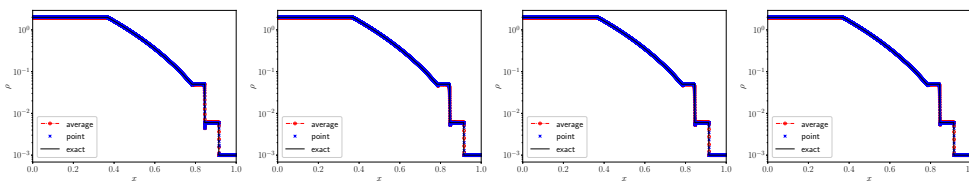


Fig. 4: Example 5.2, LeBlanc Riemann problem. The density computed with the BP limitings and the shock sensor-based limiting ($\kappa = 10$) on a uniform mesh of 6000 cells. From left to right: JS, LLF, SW, and VH FVS.

600 *Example 5.3* (Blast wave interaction [41]). This test describes the interaction of
 601 two strong shocks in the domain $[0, 1]$ with reflective boundary conditions. The test
 602 is solved until $T = 0.038$.

603 Due to the low-pressure region, the schemes blow up without the BP limitings.
 604 **Figure 5** shows the density plots obtained by using the BP limitings and shock sensor-
 605 based limiting on a uniform mesh of 800 cells. The CFL number is 0.4 for the JS,
 606 LLF, and SW FVS, and 0.36 for the VH FVS. The numerical solutions agree well the
 607 reference solution with a few overshoots/undershoots.

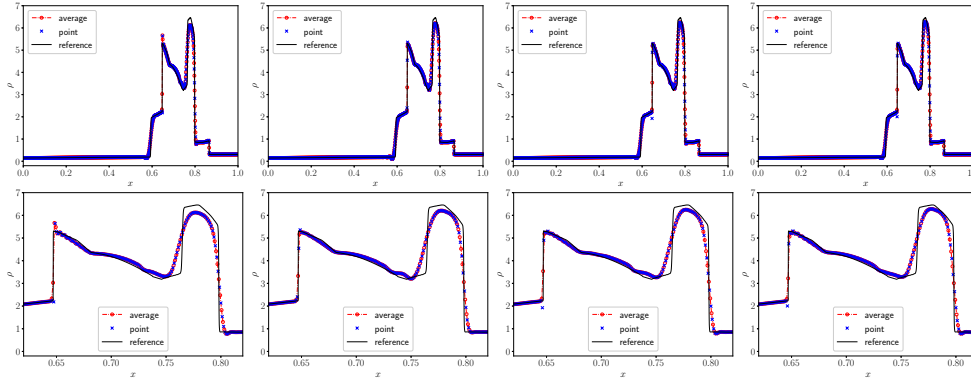


Fig. 5: **Example 5.3**, blast wave interaction. The density computed with the BP limitings and the shock sensor-based limiting ($\kappa = 1$). The corresponding enlarged views in $x \in [0.62, 0.82]$ are shown in the bottom row.

608 *Remark 5.4*. In the numerical tests, the maximal CFL numbers for stability are
 609 obtained **experimentally**. Note that the reduction of the CFL numbers is due to
 610 different stability bounds for different point value updates, and is not related to the
 611 BP property. The study of such an issue is beyond the scope of this paper, which will
 612 be explored in the future.

613 *Example 5.5* (2D advection equation). This test solves $u_t + u_x + u_y = 0$, on the
 614 periodic domain $[0, 1] \times [0, 1]$ with the following initial data

$$615 \quad u_0(x, y) = \begin{cases} 1 - |5r|, & \text{if } r = \sqrt{(x - 0.3)^2 + (y - 0.3)^2} < 0.2, \\ 1, & \text{if } \max\{|x - 0.7|, |y - 0.7|\} < 0.2, \\ 0, & \text{otherwise.} \end{cases}$$

616 For the advection equation, the JS and LLF FVS are equivalent. The results on
 617 the uniform 100×100 mesh obtained without and with BP limitings at $T = 2$ are
 618 presented in **Figure 6**. The BP limitings suppress the overshoots and undershoots
 619 well near the discontinuities. **Table 1** lists whether the numerical solutions stay in the
 620 bound $[0, 1]$. The bound is only preserved when both the BP limitings for the cell
 621 average and point value are activated, demonstrating that it is necessary to use the
 622 two kinds of BP limitings simultaneously.

623 *Example 5.6* (2D Burgers' equation). We solve $u_t + \left(\frac{1}{2}u^2\right)_x + \left(\frac{1}{2}u^2\right)_y = 0$ on the
 624 periodic domain $[0, 1] \times [0, 1]$, with the initial condition $u_0(x, y) = 0.5 + \sin(2\pi(x + y))$.
 625 This test is solved until $T = 0.3$, when the shock waves have emerged.

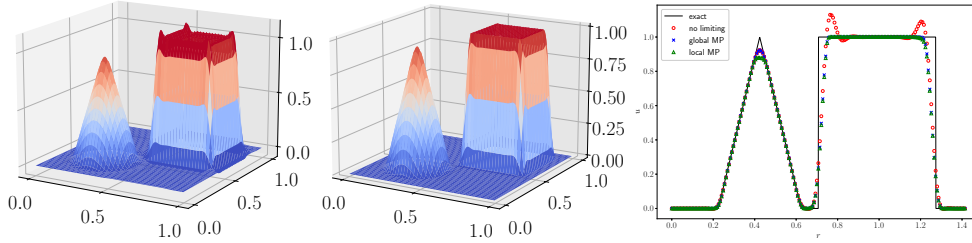


Fig. 6: [Example 5.5](#), 2D advection equation. From left to right: without any limiting, with BP limitings imposing the global MP, cut-line along $y = x$.

cell average \ point value	no limiting	global MP	local MP
	no limiting	✗	✗
global MP	✗	✓	✓
local MP	✗	✓	✓

Table 1: [Example 5.5](#), 2D advection equation. We list whether the numerical solutions stay in the bound $[0, 1]$ with different limitings.

626 [Figure 7](#) plots the solutions using the LLF FVS on the uniform 100×100 mesh
 627 without and with limitings. The oscillations near the shock waves are suppressed
 628 well when the limitings are activated, and the numerical solutions agree well with the
 629 reference solution. The blending coefficients $\theta_{i+\frac{1}{2},j}$, $\theta_{i,j+\frac{1}{2}}$ for the cell average and θ_σ
 630 for the point value when using the global MP are also presented in [Figure 8](#), verifying
 631 that the limitings are only locally activated near the shock waves.

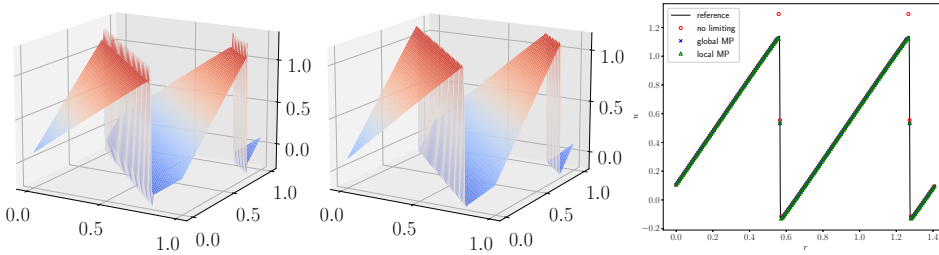


Fig. 7: [Example 5.6](#), 2D Burgers' equation. From left to right: without limiting, with BP limitings imposing the global MP, cut-line along $y = x$.

632 *Example 5.7* (2D isentropic vortex). The domain is $[-5, 5] \times [-5, 5]$ with periodic
 633 boundary conditions and the initial condition is

$$634 \quad \rho = T_0^{\frac{1}{\gamma-1}}, \quad (v_1, v_2) = (1, 1) + k_0(y, -x), \quad p = T_0 \rho, \quad k_0 = \frac{\epsilon}{2\pi} e^{0.5(1-r^2)}, \quad T_0 = 1 - \frac{\gamma-1}{2\gamma} k_0^2,$$

635 where $r^2 = x^2 + y^2$, and $\epsilon = 10.0828$ is the vortex strength. The lowest initial density
 636 and pressure are around 7.83×10^{-15} and 1.78×10^{-20} , respectively, so that the BP

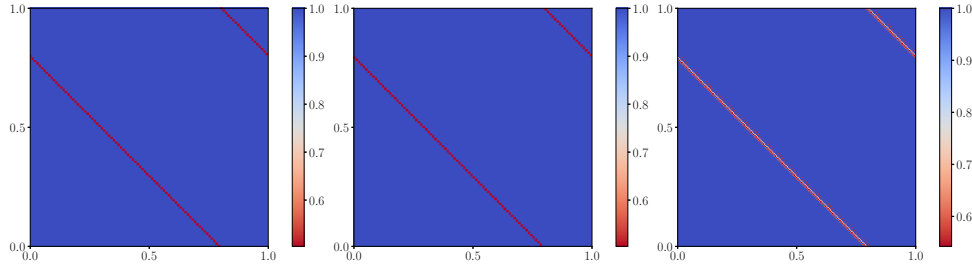


Fig. 8: [Example 5.6](#), 2D Burgers' equation. The blending coefficients in the limitings. From left to right: $\theta_{i+\frac{1}{2},j}$ and $\theta_{i,j+\frac{1}{2}}$ for the cell average, θ_σ for the point value.

637 limitings are necessary to run this test case. The problem is solved until $T = 1$.

638 Figure 9 shows the errors and corresponding convergence rates of the conservative
 639 variables in the ℓ^1 norm with the CFL number 0.2. The BP AF methods based on
 640 the JS, LLF, and VH FVS achieve the third-order accuracy, which is not affected by
 641 the BP limitings. The convergence rate based on the SW FVS reduces to around 2,
 642 due to the non-smoothness of the SW FVS as mentioned in [Remark 2.1](#).

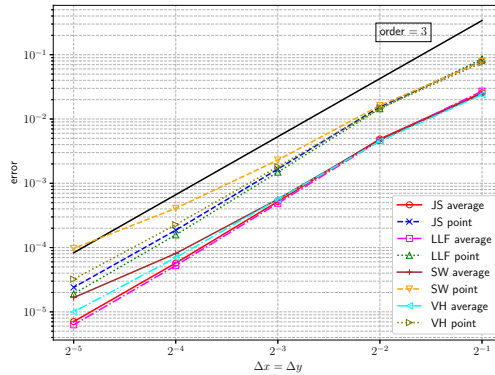


Fig. 9: [Example 5.7](#), 2D isentropic vortex problem. The errors and convergence rates.

643 *Example 5.8* (Quasi-2D Sod shock tube). This test solves the Sod shock tube
 644 problem along the x -direction on the domain $[0, 1] \times [0, 1]$ with a 100×2 uniform mesh
 645 until $T = 0.2$. The initial condition is $(\rho, v_1, v_2, p) = (1, 0, 0, 1)$ if $x < 0.5$, otherwise
 646 $(\rho, v_1, v_2, p) = (0.125, 0, 0, 0.1)$.

647 The density plots obtained by using different ways for the point value update
 648 without and with the shock sensor ($\kappa = 1$) are shown in [Figure 10](#). The density based
 649 on the JS shows large deviations between the contact discontinuity and shock wave,
 650 which cannot be reduced by the limiting. Seen from [Figure 11](#), the solutions belonging
 651 to the DoFs for different point values are decoupled, known as the mesh alignment
 652 issue, and has been explained in [Subsection 3.2](#). The results of all the FVS-based
 653 methods agree well with the exact solution when the limiting is activated. The FVS-
 654 based AF methods are more advantageous in simulations since they can cure both
 655 the stagnation and mesh alignment issues. To save space, in the following tests, we
 656 only show the results obtained using the LLF FVS.

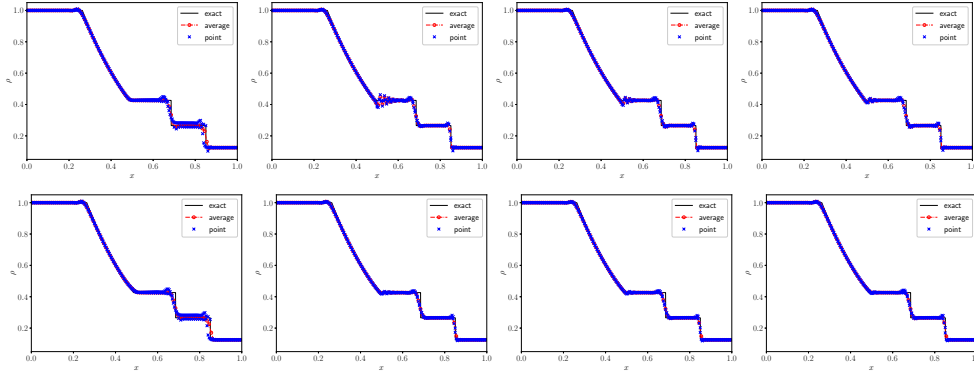


Fig. 10: [Example 5.8](#), quasi-2D Sod shock tube. The density are computed without (top row) and with the shock sensor-based limiting ($\kappa = 1$, bottom row). From left to tight: JS, LLF, SW, and VH FVS.

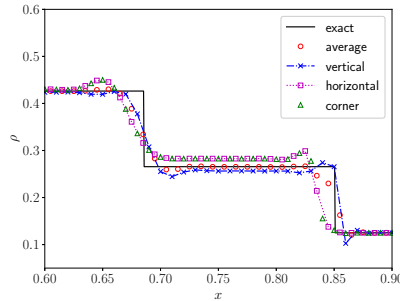


Fig. 11: [Example 5.8](#), quasi-2D Sod shock tube. Based on the JS, the solutions that belong to different kinds of DoFs are decoupled.

657 *Example 5.9* (Sedov blast wave). The domain is $[-1.1, 1.1] \times [-1.1, 1.1]$ with out-
 658 flow boundary conditions. The initial density is one, velocity is zero, and the total
 659 energy is 10^{-12} everywhere except that for the centered cell, the total energy of the cell
 660 average and the point values on its faces are $\frac{0.979264}{\Delta x \Delta y}$ with $\Delta x = 2.2/N_1$, $\Delta y = 2.2/N_2$,
 661 which is used to emulate a δ -function at the center.

662 This test is solved until $T = 1$ and the BP limitings are necessary, otherwise,
 663 the simulation fails due to negative pressure. The density plots obtained with the
 664 shock sensor ($\kappa = 0.5$) are shown in [Figure 12](#). The circular shock wave is well-
 665 captured and the numerical solutions converge to the exact solution without spurious
 666 oscillations. The blending coefficients based on the shock sensor are presented in
 667 [Figure 13](#), indicating that the limiting is locally activated.

668 *Example 5.10* (A Mach 3 wind tunnel with a forward-facing step). The initial
 669 condition is a Mach 3 flow $(\rho, v_1, v_2, p) = (1.4, 3, 0, 1)$. The computational domain is
 670 $[0, 3] \times [0, 1]$ and the step is of height 0.2 located from $x = 0.6$ to $x = 3$. The inflow
 671 and outflow boundary conditions are applied at the entrance ($x = 0$) and exit ($x = 3$),
 672 respectively, and the reflective boundary conditions are imposed at other boundaries.

673 The density computed by the BP AF method without and with the shock sensor-

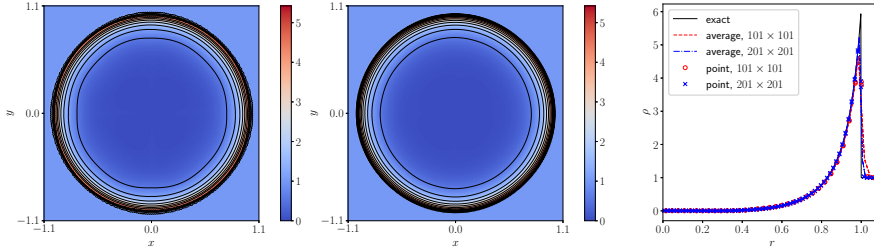


Fig. 12: [Example 5.9](#), 2D Sedov blast wave. The density plots computed by the BP AF method. From left to right: 10 equally spaced contour lines from 0 to 5.423 on the uniform 101×101 and 201×201 meshes, respectively, cut-line along $y = x$.

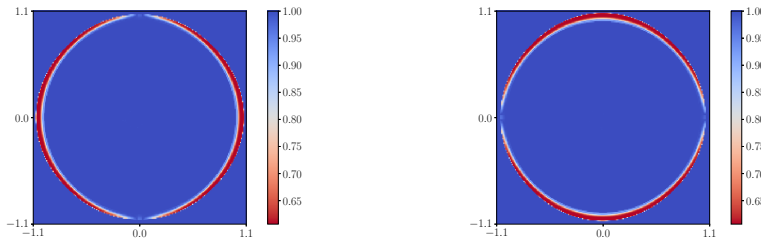


Fig. 13: [Example 5.9](#), 2D Sedov blast wave. The shock sensor-based blending coefficients $\theta^s_{i+\frac{1}{2},j}$ (left) and $\theta^s_{i,j+\frac{1}{2}}$ (right) on the 201×201 uniform mesh.

674 based limiting at $T = 4$ are shown in [Figure 14](#), and the blending coefficients $\theta^s_{i+\frac{1}{2},j}$,
675 $\theta^s_{i,j+\frac{1}{2}}$ are presented in [Figure 15](#). If only the BP limitings are used, there are oscillations in the numerical solutions, but the BP property is not violated. The numerical
676 solutions can be improved by our shock sensor-based limiting. Our BP AF method
677 can capture the main features and well-developed Kelvin–Helmholtz roll-ups that origi-
678 nate from the triple point. The noise after the shock waves is reduced by the shock
679 sensor-based limiting, while the roll-ups are preserved well. Compared to the results
680 obtained by the third-order P^2 DG method with the TVB limiter [12], the vortices are
681 better captured with the same mesh size $\Delta x = \Delta y = 1/160, 1/320$. Note that the AF
682 method uses fewer DoFs, showing its efficiency and potential for high Mach number
683 flows.
684

685 *Example 5.11* (High Mach number astrophysical jets). This test follows the setup
686 in [45]. The first case considers a Mach 80 jet on a computational domain $[0, 2] \times$
687 $[-0.5, 0.5]$, initially filled with ambient gas with $(\rho, v_1, v_2, p) = (0.5, 0, 0, 0.4127)$. A jet
688 is injected into the domain with $(\rho, v_1, v_2, p) = (5, 30, 0, 0.4127)$ at the left boundary
689 when $|y| < 0.05$. The free boundary conditions are applied on other boundaries. The
690 second case considers a Mach 2000 jet on a computational domain $[0, 1] \times [-0.25, 0.25]$.
691 The initial condition and boundary conditions are the same as the first case except
692 that the state of the jet is $(\rho, v_1, v_2, p) = (5, 800, 0, 0.4127)$. The adiabatic index is
693 $\gamma = 5/3$, and the output time is 0.07 and 0.001 for the two cases, respectively.

694 The numerical solutions obtained by the BP AF methods with the shock sensor
695 on the uniform 400×200 mesh are shown in [Figure 16](#). The main flow structures and
696 small-scale features are captured well, comparable to those in [45].

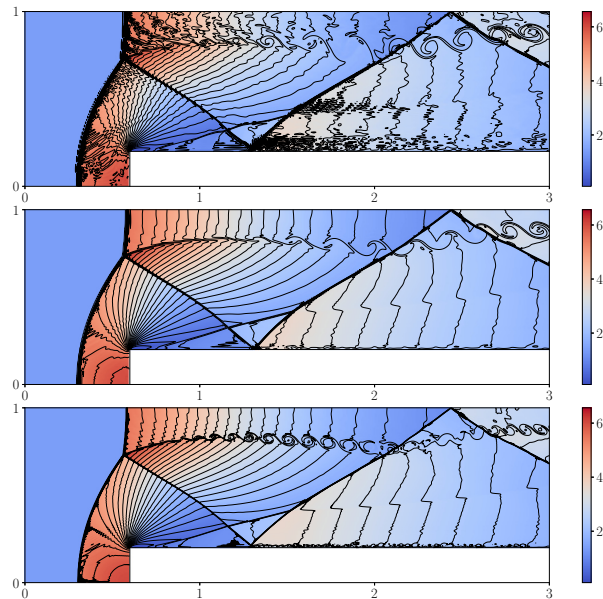


Fig. 14: [Example 5.10](#), forward-facing step problem. 30 equally spaced contour lines of the density from 0.098 to 6.566. From top to bottom: 480×160 mesh without shock sensor, 480×160 mesh with $\kappa = 1$, 960×320 mesh with $\kappa = 1$.

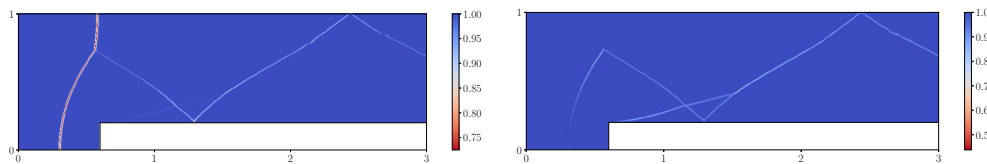


Fig. 15: [Example 5.10](#), forward-facing step problem. The blending coefficients $\theta_{i+\frac{1}{2},j}^s$ (left) and $\theta_{i,j+\frac{1}{2}}^s$ (right) based on the shock sensor with $\kappa = 1$ on the 960×320 mesh.

697 **6. Conclusion.** In the active flux (AF) methods, it is pivotal to design suitable
 698 point values update at cell interfaces, to achieve stability and high-order accuracy. The
 699 point value update based on the Jacobian splitting (JS) may lead to the stagnation and
 700 mesh alignment issues. This paper proposed to use the flux vector splitting (FVS) for
 701 the point value update instead of the JS, which keeps the continuous reconstruction as
 702 the original AF methods, and offers a natural and uniform remedy to those two issues.
 703 To further improve the robustness of the AF methods, this paper developed bound-
 704 preserving (BP) AF methods for hyperbolic conservation laws, achieved by blending
 705 the high-order AF methods with the first-order local Lax-Friedrichs (LLF) or Rusanov
 706 methods for both the cell average and point value updates, where the convex limiting
 707 and scaling limiter were employed, respectively. The shock sensor-based limiting was
 708 proposed to further improve the shock-capturing ability. The challenging numerical
 709 tests verified the robustness and effectiveness of our BP AF methods, and also showed
 710 that the LLF FVS is generally superior to others in terms of the CFL number and non-
 711 oscillatory property. Moreover, for the forward-facing step problem, the present FVS-

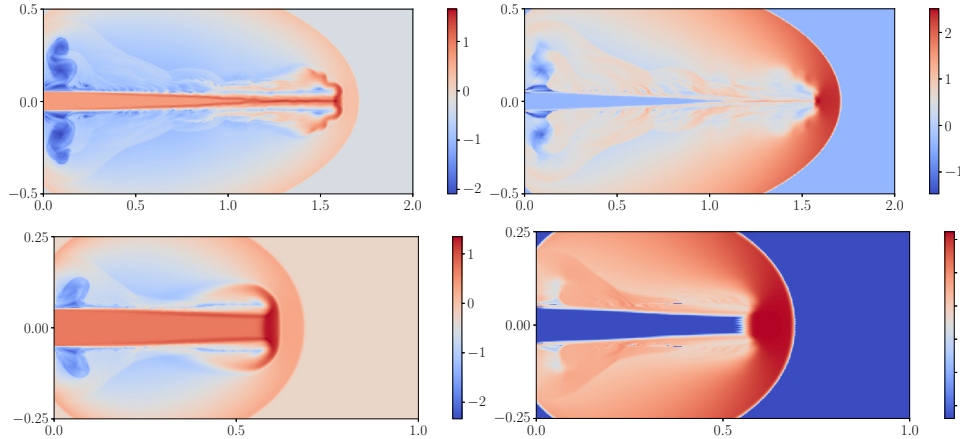


Fig. 16: **Example 5.11**, the Mach 80 jet (top row) and Mach 2000 jet (bottom row). $\log_{10} \rho$ (left) and $\log_{10} p$ (right) obtained with the BP limitings and shock sensor-based limiting ($\kappa = 1$ for Mach 80 and 10 for Mach 2000, respectively).

712 based BP AF method was able to capture small-scale features better compared to the
 713 third-order discontinuous Galerkin method with the TVB limiter on the same mesh
 714 resolution [12], while using fewer degrees of freedom, demonstrating the efficiency and
 715 potential of our BP AF method for high Mach number flows.

716 **Acknowledgments.** We acknowledge helpful discussions with Praveen Chan-
 717 drashekar at TIFR-CAM Bangalore on the Ducros' shock sensor.

718

REFERENCES

- 719 [1] R. ABGRALL, *A combination of residual distribution and the active flux formulations or a*
 720 *new class of schemes that can combine several writings of the same hyperbolic problem:*
 721 *Application to the 1D Euler equations*, Commun. Appl. Math. Comput., 5 (2023), pp. 370–
 722 402.
- 723 [2] R. ABGRALL AND W. BARSUKOW, *Extensions of active flux to arbitrary order of accuracy*,
 724 ESAIM: Math. Model. Numer. Anal., 57 (2023), pp. 991–1027.
- 725 [3] R. ABGRALL, W. BARSUKOW, AND C. KLINGENBERG, *The active flux method for the Euler*
 726 *equations on Cartesian grids*, Oct. 2023. arXiv:2310.00683.
- 727 [4] R. ABGRALL, J. LIN, AND Y. LIU, *Active flux for triangular meshes for compressible flows*
 728 *problems*, Dec. 2023.
- 729 [5] W. BARSUKOW, *The active flux scheme for nonlinear problems*, J. Sci. Comput., 86 (2021),
 730 p. 3.
- 731 [6] W. BARSUKOW AND J. P. BERBERICH, *A well-balanced active flux method for the shallow*
 732 *water equations with wetting and drying*, Commun. Appl. Math. Comput., (2023).
- 733 [7] W. BARSUKOW, J. P. BERBERICH, AND C. KLINGENBERG, *On the active flux scheme for*
 734 *hyperbolic PDEs with source terms*, SIAM J. Sci. Comput., 43 (2021), pp. A4015–A4042.
- 735 [8] W. BARSUKOW, J. HOHM, C. KLINGENBERG, AND P. L. ROE, *The active flux scheme on*
 736 *Cartesian grids and its low Mach number limit*, J. Sci. Comput., 81 (2019), pp. 594–622.
- 737 [9] E. CHUDZIK AND C. HELZEL, *A Review of Cartesian Grid Active Flux Methods for Hyperbolic*
 738 *Conservation Laws*, in Finite Volumes for Complex Applications X—Volume 1, Elliptic
 739 and Parabolic Problems, E. Franck, J. Fuhrmann, V. Michel-Dansac, and L. Navoret, eds.,
 740 Cham, 2023, Springer Nature Switzerland, pp. 93–109.
- 741 [10] E. CHUDZIK, C. HELZEL, AND D. KERKMANN, *The Cartesian grid active flux method: Linear*
 742 *stability and bound preserving limiting*, Appl. Math. Comput., 393 (2021), p. 125501.
- 743 [11] S. CLAIN, S. DIOT, AND R. LOUBÈRE, *A high-order finite volume method for systems of con-*

- 744 *servation laws—Multi-dimensional Optimal Order Detection (MOOD)*, J. Comput. Phys.,
 745 230 (2011), pp. 4028–4050.
- 746 [12] B. COCKBURN AND C. W. SHU, *Runge-Kutta discontinuous Galerkin methods for convection-*
 747 *dominated problems*, J. Sci. Comput., 16 (2001), pp. 173–261.
- 748 [13] C. J. COTTER AND D. KUZMIN, *Embedded discontinuous Galerkin transport schemes with*
 749 *localised limiters*, J. Comput. Phys., 311 (2016), pp. 363–373.
- 750 [14] C. M. DAFERMOS, *Hyperbolic Conservation Laws in Continuum Physics*, Springer Berlin Hei-
 751 delberg, 2000.
- 752 [15] F. DUCROS, V. FERRAND, F. NICOUD, C. WEBER, D. DARRACQ, C. GACHERIEU, AND
 753 T. POINSOT, *Large-eddy simulation of the shock/turbulence interaction*, Journal of Com-
 754 putational Physics, 152 (1999), pp. 517–549.
- 755 [16] T. EYMANN AND P. ROE, *Active flux schemes*, in 49th AIAA Aerospace Sciences Meeting
 756 including the New Horizons Forum and Aerospace Exposition, Orlando, Florida, Jan. 2011,
 757 American Institute of Aeronautics and Astronautics.
- 758 [17] T. EYMANN AND P. ROE, *Active flux schemes for systems*, in 20th AIAA Computational Fluid
 759 Dynamics Conference, Fluid Dynamics and Co-located Conferences, American Institute of
 760 Aeronautics and Astronautics, June 2011.
- 761 [18] T. A. EYMANN AND P. L. ROE, *Multidimensional active flux schemes*, in 21st AIAA Computa-
 762 tional Fluid Dynamics Conference, Fluid Dynamics and Co-located Conferences, American
 763 Institute of Aeronautics and Astronautics, June 2013.
- 764 [19] D. FAN AND P. L. ROE, *Investigations of a new scheme for wave propagation*, in 22nd AIAA
 765 Computational Fluid Dynamics Conference, American Institute of Aeronautics and Astro-
 766 nautics, 2015.
- 767 [20] S. GOTTLIEB, C.-W. SHU, AND E. TADMOR, *Strong Stability-Preserving High-Order Time*
 768 *Discretization Methods*, SIAM Rev., 43 (2001), pp. 89–112.
- 769 [21] J.-L. GUERMOND, M. NAZAROV, B. POPOV, AND I. TOMAS, *Second-order invariant domain*
 770 *preserving approximation of the Euler equations using convex limiting*, SIAM J. Sci. Com-
 771 put., 40 (2018), pp. A3211–A3239.
- 772 [22] J.-L. GUERMOND AND B. POPOV, *Invariant domains and first-order continuous finite element*
 773 *approximation for hyperbolic systems*, SIAM J. Numer. Anal., 54 (2016), pp. 2466–2489.
- 774 [23] J.-L. GUERMOND AND B. POPOV, *Invariant domains and second-order continuous finite ele-*
 775 *ment approximation for scalar conservation equations*, SIAM J. Numer. Anal., 55 (2017),
 776 pp. 3120–3146.
- 777 [24] J.-L. GUERMOND, B. POPOV, AND I. TOMAS, *Invariant domain preserving discretization-*
 778 *independent schemes and convex limiting for hyperbolic systems*, Comput. Method. Appl.
 779 M., 347 (2019), pp. 143–175.
- 780 [25] H. HAJDUK, *Monolithic convex limiting in discontinuous Galerkin discretizations of hyperbolic*
 781 *conservation laws*, Comput. Math. Appl., 87 (2021), pp. 120–138.
- 782 [26] D. HÄNEL, R. SCHWANE, AND G. SEIDER, *On the accuracy of upwind schemes for the solution*
 783 *of the Navier-Stokes equations*, Fluid Dynamics and Co-located Conferences, American
 784 Institute of Aeronautics and Astronautics, June 1987.
- 785 [27] C. HELZEL, D. KERKMANN, AND L. SCANDURRA, *A new ADER method inspired by the active*
 786 *flux method*, J. Sci. Comput., 80 (2019), pp. 1463–1497.
- 787 [28] X. Y. HU, N. A. ADAMS, AND C.-W. SHU, *Positivity-preserving method for high-order con-*
 788 *servative schemes solving compressible Euler equations*, J. Comput. Phys., 242 (2013),
 789 pp. 169–180.
- 790 [29] A. JAMESON, W. SCHMIDT, AND E. TURKEL, *Solutions of the Euler equations by finite volume*
 791 *methods using Runge-Kutta time-stepping schemes*, AIAA J., 1259 (1981).
- 792 [30] D. KUZMIN, *Monolithic convex limiting for continuous finite element discretizations of hyper-*
 793 *bolic conservation laws*, Computer Methods in Applied Mechanics and Engineering, 361
 794 (2020), p. 112804.
- 795 [31] D. KUZMIN, R. LÖHNER, AND S. TUREK, eds., *Flux-Corrected Transport: Principles, Algo-*
 796 *rithms, and Applications*, Scientific Computation, Springer Netherlands, Dordrecht, 2012.
- 797 [32] X.-D. LIU AND S. OSHER, *Nonoscillatory high order accurate self-similar maximum principle*
 798 *satisfying shock capturing schemes I*, SIAM J. Numer. Anal., 33 (1996), pp. 760–779.
- 799 [33] C. LOHMANN, D. KUZMIN, J. N. SHADID, AND S. MABUZA, *Flux-corrected transport algo-*
 800 *rithms for continuous Galerkin methods based on high order Bernstein finite elements*, J.
 801 Comput. Phys., 344 (2017), pp. 151–186.
- 802 [34] J. MAENG, *On the Advective Component of Active Flux Schemes for Nonlinear Hyperbolic*
 803 *Conservation Laws*, PhD thesis, 2017.
- 804 [35] B. PERTHAME AND C.-W. SHU, *On positivity preserving finite volume schemes for Euler*
 805 *equations*, Numer. Math., 73 (1996), pp. 119–130.

- 806 [36] P. ROE, *Is discontinuous reconstruction really a good idea?*, J. Sci. Comput., 73 (2017),
807 pp. 1094–1114.
- 808 [37] J. L. STEGER AND R. F. WARMING, *Flux vector splitting of the inviscid gasdynamic equations*
809 *with application to finite-difference methods*, J. Comput. Phys., 40 (1981), pp. 263–293.
- 810 [38] E. F. TORO, *Riemann Solvers and Numerical Methods for Fluid Dynamics*, Springer Berlin
811 Heidelberg, 2009.
- 812 [39] B. VAN LEER, *Towards the ultimate conservative difference scheme. IV. A new approach to*
813 *numerical convection*, J. Comput. Phys., 23 (1977), pp. 276–299.
- 814 [40] B. VAN LEER, *Flux-vector splitting for the Euler equations*, in Eighth International Conference
815 on Numerical Methods in Fluid Dynamics, E. Krause, ed., Lecture Notes in Physics, Berlin,
816 Heidelberg, 1982, Springer, pp. 507–512.
- 817 [41] P. WOODWARD AND P. COLELLA, *The numerical simulation of two-dimensional fluid flow*
818 *with strong shocks*, J. Comput. Phys., 54 (1984), pp. 115–173.
- 819 [42] K. WU AND C.-W. SHU, *Geometric quasilinearization framework for analysis and design of*
820 *bound-preserving schemes*, SIAM Rev., 65 (2023), pp. 1031–1073.
- 821 [43] Z. XU, *Parametrized maximum principle preserving flux limiters for high order schemes solving*
822 *hyperbolic conservation laws: one-dimensional scalar problem*, Math. Comput., 83 (2014),
823 pp. 2213–2238.
- 824 [44] X. ZHANG AND C.-W. SHU, *On maximum-principle-satisfying high order schemes for scalar*
825 *conservation laws*, J. Comput. Phys., 229 (2010), pp. 3091–3120.
- 826 [45] X. ZHANG AND C.-W. SHU, *On positivity-preserving high order discontinuous Galerkin*
827 *schemes for compressible Euler equations on rectangular meshes*, J. Comput. Phys., 229
828 (2010), pp. 8918–8934.
- 829 [46] X. ZHANG AND C.-W. SHU, *Maximum-principle-satisfying and positivity-preserving high-order*
830 *schemes for conservation laws: survey and new developments*, Proceedings of the Royal
831 Society A: Mathematical, Physical and Engineering Sciences, 467 (2011), pp. 2752–2776.
- 832 [47] X. ZHANG AND C.-W. SHU, *Positivity-preserving high order discontinuous Galerkin schemes*
833 *for compressible Euler equations with source terms*, J. Comput. Phys., 230 (2011),
834 pp. 1238–1248.