

1     **ACTIVE FLUX METHODS FOR HYPERBOLIC CONSERVATION**  
2     **LAWS – FLUX VECTOR SPLITTING AND**  
3     **BOUND-PRESERVATION: ONE-DIMENSIONAL CASE \***

4     JUNMING DUAN<sup>†</sup>, WASILIJ BARSUKOW<sup>‡</sup>, AND CHRISTIAN KLINGENBERG<sup>†</sup>

5     **Abstract.** The active flux (AF) method is a compact high-order finite volume method that  
6     evolves cell averages and point values at cell interfaces independently. Within the method of lines  
7     framework, the point value can be updated based on Jacobian splitting (JS), incorporating the up-  
8     wind idea. However, such JS-based AF methods encounter transonic issues for nonlinear problems  
9     due to inaccurate upwind direction estimation. This paper proposes to use flux vector splitting for  
10    the point value update, offering a natural and uniform remedy to the transonic issue. To improve  
11    robustness, this paper also develops bound-preserving (BP) AF methods for one-dimensional hyper-  
12    bolic conservation laws. Two cases are considered: preservation of the maximum principle for the  
13    scalar case, and preservation of positive density and pressure for the compressible Euler equations.  
14    The update of the cell average in high-order AF methods is rewritten as a convex combination of us-  
15    ing the original high-order fluxes and robust low-order (local Lax-Friedrichs or Rusanov) fluxes, and  
16    the desired bounds are enforced by choosing the right amount of low-order fluxes. A similar blending  
17    strategy is used for the point value update. Several challenging benchmark tests are conducted to  
18    verify the accuracy, BP properties, and shock-capturing ability of the methods.

19    **Key words.** Hyperbolic conservation laws, finite volume method, active flux, flux vector split-  
20    ting, bound-preserving, convex limiting, scaling limiter

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

22    **1. Introduction.** This paper is concerned with solving systems of hyperbolic  
23    conservation laws

24    (1.1)     
$$\frac{\partial \mathbf{U}(x, t)}{\partial t} + \frac{\partial \mathbf{F}(\mathbf{U})}{\partial x} = 0, \quad \mathbf{U}(x, 0) = \mathbf{U}_0(x), \quad (x, t) \in \mathbb{R} \times \mathbb{R}^+,$$

25    where  $\mathbf{U} \in \mathbb{R}^m$  is the vector of  $m$  conservative variables,  $\mathbf{F} \in \mathbb{R}^m$  is the physical flux,  
26    and  $\mathbf{U}_0(x)$  is assumed to be initial data of bounded variation. In this paper, we would  
27    like to consider two cases. The first is a scalar conservation law ( $m = 1$ )

28    (1.2)     
$$\frac{\partial u}{\partial t} + \frac{\partial f(u)}{\partial x} = 0, \quad u(x, 0) = u_0(x).$$

29    The second case is that of compressible Euler equations of gas dynamics with  $\mathbf{U} =$   
30     $(\rho, \rho v, E)^\top$  and  $\mathbf{F} = (\rho v, \rho v^2 + p, (E + p)v)^\top$ , i.e.,

31    (1.3)     
$$\frac{\partial}{\partial t} \begin{pmatrix} \rho \\ \rho v \\ E \end{pmatrix} + \frac{\partial}{\partial x} \begin{pmatrix} \rho v \\ \rho v^2 + p \\ (E + p)v \end{pmatrix} = \mathbf{0}, \quad (\rho, v, p)(x, 0) = (\rho_0, v_0, p_0).$$

32    Here  $\rho$  denotes the density,  $v$  the velocity,  $p$  the pressure, and  $E = \frac{1}{2}\rho v^2 + \rho e$  the  
33    total energy with  $e$  the specific internal energy. The system (1.3) should be closed by

---

\*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.

<sup>†</sup>Department of Mathematics, University of Würzburg, Germany ([junming.duan@uni-wuerzburg.de](mailto:junming.duan@uni-wuerzburg.de), [christian.klingenberg@uni-wuerzburg.de](mailto:christian.klingenberg@uni-wuerzburg.de)).

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

34 an equation of state (EOS). This paper considers the perfect gas EOS,  $p = (\gamma - 1)\rho e$ ,  
 35 with the adiabatic index  $\gamma > 1$ . Note that this paper uses bold symbols to refer to  
 36 vectors and matrices, such that they are easier to distinguish from scalars.

37 The active flux (AF) method is a new finite volume method [13, 12, 14, 34], that  
 38 Roe took inspiration by [39]. Apart from cell averages, it incorporates additional de-  
 39 grees of freedom as point values located at the cell interfaces, evolved independently  
 40 from the cell average. The original AF method gives a global continuous represen-  
 41 tation of the numerical solution using a piecewise quadratic reconstruction, leading  
 42 naturally to a third-order accurate method with a compact stencil. The introduc-  
 43 tion of point values at the cell interfaces avoids the usage of Riemann solvers as in  
 44 usual Godunov methods, because the numerical solution is continuous across the cell  
 45 interface and the numerical flux for the cell average update is available directly.

46 The independence of the point value update adds flexibility to the AF methods.  
 47 Based on the evolution of the point value, there are generally two kinds of AF methods.  
 48 The original one uses exact or approximate evolution operators and Simpson's rule for  
 49 flux quadrature in time, i.e. it does not require time integration methods like Runge-  
 50 Kutta methods. Exact evolution operators have been studied for linear equations  
 51 in [7, 15, 14, 39]. Approximate evolution operators have been explored for Burgers'  
 52 equation [13, 12, 34, 4], the compressible Euler equations in one spatial dimension  
 53 [13, 25, 4], and hyperbolic balance laws [6, 5], etc. One of the advantages of the AF  
 54 method over standard finite volume methods is its structure-preserving property. For  
 55 instance, it preserves the vorticity and stationary states for multi-dimensional acoustic  
 56 equations [7], and it is naturally well-balanced for acoustics with gravity [6].

57 Since it may not be convenient to derive exact or approximate evolution operators  
 58 for nonlinear systems, especially in multi-dimensions, another kind of generalized AF  
 59 method was presented in [1, 2]. A method of lines was used, where the cell average and  
 60 point value updates are written in semi-discrete form and advanced in time with time  
 61 integration methods. In the point values update, the Jacobian matrix is split based on  
 62 the sign of the eigenvalues (Jacobian splitting (JS)), and upwind-biased stencils are  
 63 used to compute the approximation of derivatives. There are some deficiencies of the  
 64 JS when used for the AF methods, e.g., the transonic issue [25] for nonlinear problems,  
 65 leading to spikes in the cell average. Some remedies are suggested in the literature,  
 66 e.g., using discontinuous reconstruction [25] or evaluating the upwind direction using  
 67 more information from the neighbors [4].

68 Solutions to hyperbolic systems (1.1) often stay in an *admissible state set*  $\mathcal{G}$ , also  
 69 called the invariant domain. For instance, the solutions to initial value problems of  
 70 scalar conservation laws (1.2) satisfy a strict maximum principle (MP) [11], i.e.,

$$71 \quad (1.4) \quad \mathcal{G} = \{u \mid m_0 \leq u \leq M_0\}, \quad m_0 = \min_x u_0(x), \quad M_0 = \max_x u_0(x).$$

72 Physically, both the density and pressure in the solutions to the compressible Euler  
 73 equations (1.3) should stay positive, i.e.,

$$74 \quad (1.5) \quad \mathcal{G} = \left\{ \mathbf{U} = (\rho, \rho v, E) \mid \rho > 0, \quad p = (\gamma - 1) \left( E - \frac{(\rho v)^2}{2\rho} \right) > 0 \right\}.$$

75 Throughout this paper, it is assumed that  $\mathcal{G}$  is a *convex* set, which is obvious for the  
 76 scalar case (1.4) and can be verified for the Euler equations (1.5), see e.g. [46]. It is  
 77 desirable to conceive so-called bound-preserving (BP) methods, i.e., those guaranteeing  
 78 that the numerical solutions at a later time will stay in  $\mathcal{G}$ , if the initial numerical solu-  
 79 tions belong to  $\mathcal{G}$ . The BP property of numerical methods is very important for both

80 theoretical analysis and numerical stability. Many BP methods have been developed  
 81 in the past few decades, e.g., a series of works by Shu and collaborators [45, 26, 43],  
 82 a recent general framework on BP methods [42], and the convex limiting approach  
 83 [17, 22, 29], which can be traced back to the flux-corrected transport (FCT) schemes  
 84 for scalar conservation laws [10, 20, 32, 30]. The previous studies on the AF methods  
 85 pay limited attention to high-speed flows, or problems containing strong discontinu-  
 86 ities, with some efforts on the limiting for the point value update, see e.g. [4, 8].  
 87 However, those limitings are not enough to guarantee the BP property, as shown in  
 88 our numerical tests. In a very recent paper, the MOOD [9] based stabilization was  
 89 adopted to achieve the BP property [3] in an a posteriori fashion.

90 This paper presents a new way for the point value update to cure the transonic  
 91 issue and develops suitable BP limiting strategies for the AF methods. The main  
 92 contributions and findings in this work can be summarized as follows.

93 i). We propose to employ the flux vector splitting (FVS) methods for the point value  
 94 update to cure the transonic issue, since it borrows information from the neighbors  
 95 naturally and uniformly. The FVS was originally used to identify the upwind direc-  
 96 tions, which is simpler and somewhat more efficient than Godunov-type methods for  
 97 solving hyperbolic systems [38]. In our numerical tests, the FVS is also shown to  
 98 give better results than the JS, especially the local Lax-Friedrichs (LLF) or Rusanov  
 99 FVS, in terms of the CFL number and shock-capturing ability. The FVS can also  
 100 cure some defects in two dimensions observed in the JS, which will be shown in our  
 101 future companion paper.

102 ii). We design BP limitings for both the update of the cell average and the point value  
 103 by blending the high-order AF methods with the first-order LLF method in a convex  
 104 combination. The convex limiting [17, 22, 29] and the scaling limiter [31] are applied  
 105 to the cell average and point value updates, respectively. The main idea is to retain as  
 106 much as possible of the high-order method while guaranteeing the numerical solutions  
 107 to be BP, and the blending coefficients are computed by enforcing the bounds. We  
 108 show that using a suitable time step size and BP limitings, the numerical solutions  
 109 of the BP AF methods satisfy the MP for scalar conservation laws, and give positive  
 110 density and pressure for the compressible Euler equations.

111 iii). Several challenging test cases such as the LeBlanc and double rarefaction Rie-  
 112 mann problems, the Sedov point blast wave, and blast wave interaction problems are  
 113 conducted to demonstrate the BP properties and the shock-capturing ability, which  
 114 are rare in the literature for the AF methods.

115 The remainder of this paper is structured as follows. Section 2 introduces the  
 116 AF methods based on the JS or FVS for the point value update, and the power  
 117 law reconstruction for limiting the derivatives in the point value update. To design  
 118 BP methods, Section 3 describes our convex limiting approach for the cell average,  
 119 while Section 4 deals with the limiting for the point value. Some numerical tests are  
 120 conducted in Section 5 to experimentally demonstrate the accuracy, BP properties,  
 121 and shock-capturing ability of the methods. Section 6 concludes the paper with final  
 122 remarks and future directions.

123 **2. 1D active flux methods for hyperbolic conservation laws.** This section  
 124 presents the 1D semi-discrete AF methods for the hyperbolic conservation laws (1.1),  
 125 based on the JS [2] or FVS for the point value update. The fully-discrete methods  
 126 are obtained using Runge-Kutta methods.

127 Assume that a 1D computational domain is divided into  $N$  cells  $I_i = [x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}]$   
 128 with cell centers  $x_i = (x_{i-\frac{1}{2}} + x_{i+\frac{1}{2}})/2$  and cell sizes  $\Delta x_i = x_{i+\frac{1}{2}} - x_{i-\frac{1}{2}}$ ,  $i = 1, \dots, N$ .

129 The degrees of freedom of the AF methods are the approximations to cell averages  
 130 of the conservative variable as well as point values at the cell interfaces, allowing  
 131 some freedom in the choice of the point values, e.g. conservative variables, primitive  
 132 variables, entropy variables, etc. This paper only considers using the conservative  
 133 variables, and the degrees of freedom are denoted by

$$134 \quad (2.1) \quad \bar{\mathbf{U}}_i(t) = \frac{1}{\Delta x_i} \int_{I_i} \mathbf{U}(x, t) \, dx, \quad \mathbf{U}_{i+\frac{1}{2}}(t) = \mathbf{U}(x_{i+\frac{1}{2}}, t).$$

135 The cell average is updated by integrating (1.1) over  $I_i$  in the following semi-discrete  
 136 finite volume manner

$$137 \quad (2.2) \quad \frac{d\bar{\mathbf{U}}_i}{dt} = -\frac{1}{\Delta x_i} \left[ \mathbf{F}(\mathbf{U}_{i+\frac{1}{2}}) - \mathbf{F}(\mathbf{U}_{i-\frac{1}{2}}) \right].$$

138 Thus, the accuracy of (2.2) is determined by the approximation accuracy of the point  
 139 values. It was so far (e.g. in [2]) considered sufficient to update the point values with  
 140 any finite-difference-like formula

$$141 \quad (2.3) \quad \frac{d\mathbf{U}_{i+\frac{1}{2}}}{dt} = -\mathcal{R} \left( \mathbf{U}_{i+\frac{1}{2}-l_1}(t), \bar{\mathbf{U}}_{i+1-l_1}(t), \dots, \bar{\mathbf{U}}_{i+l_2}(t), \mathbf{U}_{i+\frac{1}{2}+l_2}(t) \right), \quad l_1, l_2 \geq 0,$$

142 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  
 143 stable method. This paper explores further conditions on  $\mathcal{R}$  for nonlinear problems.

144 **2.1. Point value update using Jacobian splitting.** For smooth solutions,  
 145 we have an equivalent formulation in the form

$$146 \quad (2.4) \quad \frac{\partial \mathbf{U}}{\partial t} + \mathbf{J}(\mathbf{U}) \frac{\partial \mathbf{U}}{\partial x} = 0, \quad \mathbf{J}(\mathbf{U}) = \frac{\partial \mathbf{F}(\mathbf{U})}{\partial \mathbf{U}}.$$

147 Inspired by the upwind scheme, (2.4) can be discretized by the JS [1, 2] as follows

$$148 \quad (2.5) \quad \frac{d\mathbf{U}_{i+\frac{1}{2}}}{dt} = - \left[ \mathbf{J}^+(\mathbf{U}_{i+\frac{1}{2}}) \mathbf{D}_{i+\frac{1}{2}}^+(\mathbf{U}) + \mathbf{J}^-(\mathbf{U}_{i+\frac{1}{2}}) \mathbf{D}_{i+\frac{1}{2}}^-(\mathbf{U}) \right],$$

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

$$150 \quad \mathbf{J}^+ = \mathbf{R} \mathbf{\Lambda}^+ \mathbf{R}^{-1}, \quad \mathbf{J}^- = \mathbf{R} \mathbf{\Lambda}^- \mathbf{R}^{-1},$$

$$151 \quad \mathbf{\Lambda}^+ = \text{diag}\{\max(\lambda_1, 0), \dots, \max(\lambda_m, 0)\},$$

$$152 \quad \mathbf{\Lambda}^- = \text{diag}\{\min(\lambda_1, 0), \dots, \min(\lambda_m, 0)\},$$

153 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  
 154  $\lambda_1, \dots, \lambda_m$  are the eigenvalues, with the columns of  $\mathbf{R}$  the corresponding eigenvectors.

155 To derive the approximation of the derivatives in (2.5), one can first obtain a high-  
 156 order reconstruction for  $\mathbf{U}$  in the upwind cell, and then differentiate the reconstructed  
 157 polynomial. As an example, a parabolic reconstruction in cell  $i$  is

$$158 \quad \mathbf{U}_{\text{para},1}(x) = -3(2\bar{\mathbf{U}}_i - \mathbf{U}_{i-\frac{1}{2}} - \mathbf{U}_{i+\frac{1}{2}}) \frac{x^2}{\Delta x_i^2} + (\mathbf{U}_{i+\frac{1}{2}} - \mathbf{U}_{i-\frac{1}{2}}) \frac{x}{\Delta x_i}$$

$$159 \quad (2.6) \quad + \frac{1}{4}(6\bar{\mathbf{U}}_i - \mathbf{U}_{i-\frac{1}{2}} - \mathbf{U}_{i+\frac{1}{2}})$$

160 satisfying  $\mathbf{U}_{\text{para},1}(\pm\Delta x_i/2) = \mathbf{U}_{i\pm\frac{1}{2}}$ ,  $\frac{1}{\Delta x_i} \int_{-\Delta x_i/2}^{\Delta x_i/2} \mathbf{U}_{\text{para},1}(x) dx = \bar{\mathbf{U}}_i$ . Then the de-  
 161 rivatives are

$$162 \quad (2.7a) \quad \mathbf{D}_{i+\frac{1}{2}}^+(\mathbf{U}) = \mathbf{U}'_{\text{para},1}(\Delta x_i/2) = \frac{1}{\Delta x_i} \left( 2\mathbf{U}_{i-\frac{1}{2}} - 6\bar{\mathbf{U}}_i + 4\mathbf{U}_{i+\frac{1}{2}} \right),$$

$$163 \quad (2.7b) \quad \mathbf{D}_{i+\frac{1}{2}}^-(\mathbf{U}) = \frac{1}{\Delta x_{i+1}} \left( -4\mathbf{U}_{i+\frac{1}{2}} + 6\bar{\mathbf{U}}_{i+1} - 2\mathbf{U}_{i+\frac{3}{2}} \right).$$

164 They are third-order accurate. Higher-order extensions can be obtained by higher-  
 165 order finite difference formulae using a larger spatial stencil, see [2] for examples.

166 **2.2. Point value update using flux vector splitting.** One of the deficiencies  
 167 of using the JS is the transonic issue that appears for nonlinear problems, as observed  
 168 in [25, 4] and described in more detail next. Consider [Example 5.2](#), where we solve  
 169 Burgers' equation with a square wave as the initial data. [Figure 3](#) shows the cell  
 170 averages and point values based on the JS with 200 cells, as well as the reference  
 171 solution. The numerical solution based on the JS without limiting gives a spike at  
 172 the initial discontinuity  $x = 0.2$ , which grows linearly in time. The reason for this  
 173 behaviour is the inaccurate estimation of the upwind direction at the cell interface.  
 174 In this example, there are two successive point values with different initial data near  
 175 the initial discontinuity, denoted by  $u_{i-\frac{1}{2}} = 2$ ,  $u_{i+\frac{1}{2}} = -1$ , respectively. At the cell  
 176 interface  $x_{i-\frac{1}{2}}$  or  $x_{i+\frac{1}{2}}$ , the upwind discretization in (2.7) only uses the data from the  
 177 left or right, leading to zero derivatives, thus the point values  $u_{i-\frac{1}{2}}$  and  $u_{i+\frac{1}{2}}$  stay  
 178 unchanged. However, according to the update of the cell average (2.2),  $\bar{u}_i$  increases  
 179 gradually (which is the observed spike). This deficiency cannot be eliminated by  
 180 limitings, as one observes from [Figure 3](#). Some remedies have been proposed, such  
 181 as using discontinuous reconstruction [25] and an ‘‘entropy fix’’ that evaluates the  
 182 upwind direction not only at the corresponding cell interface but also with values  
 183 from its neighbors [4].

184 In this paper, we propose to use the FVS for the point value update, which  
 185 borrows the information from the neighbors naturally, still based on the continuous  
 186 reconstruction, and can eliminate the generation of the spike effectively, as shown in  
 187 [Figure 4](#). The FVS for the point value update reads

$$188 \quad (2.8) \quad \frac{d\mathbf{U}_{i+\frac{1}{2}}}{dt} = - \left[ \tilde{\mathbf{D}}^+ \mathbf{F}^+(\mathbf{U}) + \tilde{\mathbf{D}}^- \mathbf{F}^-(\mathbf{U}) \right]_{i+\frac{1}{2}},$$

189 where the flux  $\mathbf{F}$  is split into the positive and negative parts  $\mathbf{F} = \mathbf{F}^+ + \mathbf{F}^-$  satisfying  
 190

$$191 \quad (2.9) \quad \lambda \left( \frac{\partial \mathbf{F}^+}{\partial \mathbf{U}} \right) \geq 0, \quad \lambda \left( \frac{\partial \mathbf{F}^-}{\partial \mathbf{U}} \right) \leq 0,$$

192 i.e., all the eigenvalues of  $\frac{\partial \mathbf{F}^+}{\partial \mathbf{U}}$  and  $\frac{\partial \mathbf{F}^-}{\partial \mathbf{U}}$  are non-negative and non-positive, respec-  
 193 tively. Different FVS can be adopted as long as they satisfy the constraint (2.9),  
 194 to be discussed later. Finite difference formulae to approximate the flux derivatives  
 195 are obtained similarly to the computation of the derivatives in the JS. A parabolic  
 196 reconstruction of the flux can be obtained based on the three flux values as follows

$$197 \quad \mathbf{F}_{\text{para},2}(x) = 2(\mathbf{F}_{i-\frac{1}{2}} - 2\mathbf{F}_i + \mathbf{F}_{i+\frac{1}{2}}) \frac{x^2}{\Delta x_i^2} + (\mathbf{F}_{i+\frac{1}{2}} - \mathbf{F}_{i-\frac{1}{2}}) \frac{x}{\Delta x_i} + \mathbf{F}_i,$$

198 satisfying  $\mathbf{F}_{\text{para},2}(\pm\Delta x_i/2) = \mathbf{F}_{i\pm\frac{1}{2}}$ ,  $\mathbf{F}_{\text{para},2}(0) = \mathbf{F}_i$ , with  $\mathbf{F}_{i\pm\frac{1}{2}} = \mathbf{F}(\mathbf{U}_{i\pm\frac{1}{2}})$ , and the  
 199 cell-centered point value  $\mathbf{F}_i = \mathbf{F}(\mathbf{U}_i)$  is obtained by evaluating the reconstruction  
 200 of  $\mathbf{U}$ , i.e. according to Simpson's rule  $\mathbf{U}_i = (-\mathbf{U}_{i-\frac{1}{2}} + 6\bar{\mathbf{U}}_i - \mathbf{U}_{i+\frac{1}{2}})/4$ . Then the  
 201 derivatives are

$$202 \quad (2.10a) \quad \left(\tilde{\mathbf{D}}^+ \mathbf{F}^+\right)_{i+\frac{1}{2}} = \mathbf{F}'_{\text{para},2}(\Delta x_i/2) = \frac{1}{\Delta x_i} \left(\mathbf{F}_{i-\frac{1}{2}} - 4\mathbf{F}_i + 3\mathbf{F}_{i+\frac{1}{2}}\right),$$

$$203 \quad (2.10b) \quad \left(\tilde{\mathbf{D}}^- \mathbf{F}^-\right)_{i+\frac{1}{2}} = \frac{1}{\Delta x_{i+1}} \left(-3\mathbf{F}_{i+\frac{1}{2}} + 4\mathbf{F}_{i+1} - \mathbf{F}_{i+\frac{3}{2}}\right).$$

204 These finite differences are third-order accurate. While the reconstructions of both  $\mathbf{U}$   
 205 and  $\mathbf{F}$  are parabolic, the coefficients in the formula (2.10) differ from (2.7) because  
 206 (2.10) uses the cell-centered value rather than the cell average. Our numerical tests in  
 207 Section 5 show that the AF methods based on the FVS generally give better results  
 208 than the JS.

209 **2.2.1. Local Lax-Friedrichs flux vector splitting.** The first FVS we consider  
 210 is the LLF FVS, defined as

$$211 \quad \mathbf{F}^\pm = \frac{1}{2}(\mathbf{F}(\mathbf{U}) \pm \alpha \mathbf{U}),$$

212 where the choice of  $\alpha$  should fulfill (2.9) across the spatial stencil. In our implemen-  
 213 tation, it is determined by

$$214 \quad (2.11) \quad \alpha_{i+\frac{1}{2}} = \max_{r,\ell} \{|\lambda_\ell(\mathbf{U}_r)|\}, \quad r \in \left\{i - \frac{1}{2}, i, i + \frac{1}{2}, i + 1, u + \frac{3}{2}\right\}, \quad \ell = 1, \dots, m.$$

215 One can also choose  $\alpha$  to be the maximal absolute value of the eigenvalues in the whole  
 216 domain, corresponding to the (global) LF splitting. Note, however, that a larger  $\alpha$   
 217 generally leads to a smaller time step size and more dissipation.

218 **2.2.2. Upwind flux vector splitting.** One can also split the Jacobian matrix  
 219 based on each characteristic field,

$$220 \quad (2.12) \quad \mathbf{F}^\pm = \frac{1}{2}(\mathbf{F}(\mathbf{U}) \pm |\mathbf{J}|\mathbf{U}), \quad |\mathbf{J}| = \mathbf{R}(\mathbf{\Lambda}^+ - \mathbf{\Lambda}^-)\mathbf{R}^{-1}.$$

221 For linear systems, one has  $\mathbf{F} = \mathbf{J}\mathbf{U}$ , so (2.12) reduces to the JS. To be specific,

$$222 \quad \mathbf{F}^\pm = \frac{1}{2}(\mathbf{J} \pm |\mathbf{J}|)\mathbf{U} = \mathbf{R}\mathbf{\Lambda}^\pm\mathbf{R}^{-1}\mathbf{U} = \mathbf{J}^\pm\mathbf{U},$$

223 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  
 224  $\mathbf{J}^\pm \mathbf{D}^\pm \mathbf{U}$  if  $\mathbf{D}^+$  and  $\tilde{\mathbf{D}}^+$  are derived from the same reconstructed polynomial. In  
 225 other words, the AF methods using this FVS enjoy the same properties as the original  
 226 JS-based AF methods for linear systems.

227 Such an FVS is also known as the Steger-Warming (SW) FVS [36] for the Euler  
 228 equations (1.3), since the ‘‘homogeneity property’’ holds [38], i.e.,  $\mathbf{F} = \mathbf{J}\mathbf{U}$ . One can  
 229 write down the SW FVS explicitly

$$230 \quad \mathbf{F}^\pm = \begin{bmatrix} \frac{\rho}{2\gamma} \alpha^\pm \\ \frac{\rho}{2\gamma} (\alpha^\pm v + a(\lambda_2^\pm - \lambda_3^\pm)) \\ \frac{\rho}{2\gamma} \left(\frac{1}{2}\alpha^\pm v^2 + av(\lambda_2^\pm - \lambda_3^\pm) + \frac{a^2}{\gamma-1}(\lambda_2^\pm + \lambda_3^\pm)\right) \end{bmatrix},$$

231 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  $a = \sqrt{\gamma p/\rho}$   
 232 is the sound speed.

233 **2.2.3. Van Leer-Hänel flux vector splitting for the Euler equations.**

234 Another popular FVS for the Euler equations was proposed by Van Leer [40], and  
 235 improved by [23]. The flux can be split based on the Mach number  $M = v/a$  as

$$236 \quad \mathbf{F} = \begin{bmatrix} \rho a M \\ \rho a^2 (M^2 + \frac{1}{\gamma}) \\ \rho a^3 M (\frac{1}{2} M^2 + \frac{1}{\gamma-1}) \end{bmatrix} = \mathbf{F}^+ + \mathbf{F}^-, \quad \mathbf{F}^\pm = \begin{bmatrix} \pm \frac{1}{4} \rho a (M \pm 1)^2 \\ \pm \frac{1}{4} \rho a (M \pm 1)^2 v + p^\pm \\ \pm \frac{1}{4} \rho a (M \pm 1)^2 H \end{bmatrix},$$

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

239 **2.3. 1D power law reconstruction for point value update.** When the  
 240 numerical solutions contain discontinuities, the computation of the derivatives (2.7)  
 241 or (2.10) based on the parabolic reconstructions may cause oscillations. Thus, it is  
 242 reasonable to seek finite difference approximations based on differentiating a modified  
 243 reconstruction with improved monotonicity properties. This section only considers  
 244 the scalar case and can be extended to systems of equations in a component-wise  
 245 fashion.

246 The power law reconstruction proposed in [4] can be used to replace the original  
 247 parabolic reconstruction to achieve monotonicity on some occasions. It is shown in  
 248 Theorem 5 in [4] that the extremum is not avoidable in the cell  $I_i = [x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}]$   
 249 for continuous reconstructions if the cell average lies outside the range of the point  
 250 values  $(\bar{u}_i - u_{i-\frac{1}{2}})(u_{i+\frac{1}{2}} - \bar{u}_i) < 0$ . The parabola is monotone, and thus no action  
 251 is required when  $(2u_{i-\frac{1}{2}} + u_{i+\frac{1}{2}})/3 < \bar{u}_i < (u_{i-\frac{1}{2}} + 2u_{i+\frac{1}{2}})/3$  or  $(2u_{i-\frac{1}{2}} + u_{i+\frac{1}{2}})/3 >$   
 252  $\bar{u}_i > (u_{i-\frac{1}{2}} + 2u_{i+\frac{1}{2}})/3$ . Upon defining  $r = \frac{u_{i+1/2} - \bar{u}_i}{\bar{u}_i - u_{i-1/2}}$ , one can equivalently express  
 253 that the parabola is monotone when  $1/2 < r < 2$ . In both these cases, the parabolic  
 254 reconstruction is used, and the derivatives are obtained by (2.7) or (2.10). Otherwise,  
 255 the following power law reconstruction is used.

256 PROPOSITION 2.1 (Barsukow [4]). The power law reconstruction

$$257 \quad (2.13) \quad \begin{cases} u_{\text{pw1},1}(x) = u_{i-\frac{1}{2}} + (u_{i+\frac{1}{2}} - u_{i-\frac{1}{2}}) \left( \frac{x - x_i}{\Delta x_i} + \frac{1}{2} \right)^r, & \text{if } r > 2 \\ u_{\text{pw1},2}(x) = u_{i+\frac{1}{2}} - (u_{i+\frac{1}{2}} - u_{i-\frac{1}{2}}) \left( \frac{1}{2} - \frac{x - x_i}{\Delta x_i} \right)^{1/r}, & \text{if } 0 < r < 1/2 \end{cases}$$

258 is monotone and satisfies

$$259 \quad u_{\text{pw1},l}(x_{i-\frac{1}{2}}) = u_{i-\frac{1}{2}}, \quad u_{\text{pw1},l}(x_{i+\frac{1}{2}}) = u_{i+\frac{1}{2}}, \quad \frac{1}{\Delta x_i} \int_{I_i} u_{\text{pw1},l}(x) \, dx = \bar{u}_i, \quad l = 1, 2.$$

260 A comparison between the parabolic reconstruction (2.6) and power law recon-  
 261 struction (2.13) is given in Figure 1 with point values fixed as  $-1$  and  $1$  at the inter-  
 262 faces, and different cell averages  $\{-1.1, -0.8, -1/3, 0.1, 1/3, 0.8, 1.1\}$ . One can observe  
 263 monotone profiles for the power law reconstruction when the cell average lies between  
 264 the two point values. Based on (2.13), the derivatives can be computed directly

$$265 \quad \begin{cases} u'_{\text{pw1},1}(x) = \frac{u_{i+\frac{1}{2}} - u_{i-\frac{1}{2}}}{\Delta x_i} r \left( \frac{x - x_i}{\Delta x_i} + \frac{1}{2} \right)^{r-1}, & \text{if } r > 2, \\ u'_{\text{pw1},2}(x) = \frac{u_{i+\frac{1}{2}} - u_{i-\frac{1}{2}}}{\Delta x_i} \frac{1}{r} \left( \frac{1}{2} - \frac{x - x_i}{\Delta x_i} \right)^{1/r-1}, & \text{if } 0 < r < 1/2. \end{cases}$$

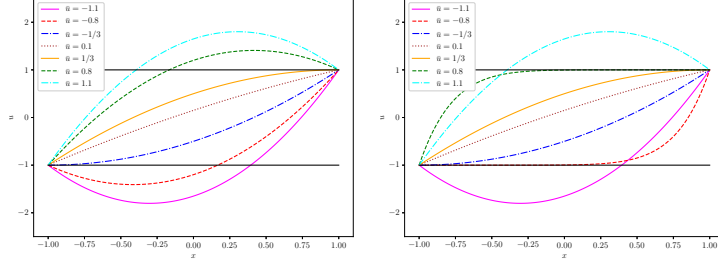


Fig. 1: The parabolic (2.6) and power law reconstruction (2.13) obtained with different cell averages  $\{-1.1, -0.8, -1/3, 0.1, 1/3, 0.8, 1.1\}$ , and fixed point values as  $-1$  and  $1$  at the left and right interfaces.

266 At the left interface, the derivative is

$$267 \quad (2.14) \quad \begin{cases} u'_{\text{pw}1,1}(x_{i-\frac{1}{2}}^+) = 0, & \text{if } r > 2, \\ u'_{\text{pw}1,2}(x_{i-\frac{1}{2}}^+) = \frac{u_{i+\frac{1}{2}} - u_{i-\frac{1}{2}}}{\Delta x_i} \frac{1}{r}, & \text{if } 0 < r < 1/2, \end{cases}$$

268 and at the right interface, the derivative is

$$269 \quad (2.15) \quad \begin{cases} u'_{\text{pw}1,1}(x_{i+\frac{1}{2}}^-) = \frac{u_{i+\frac{1}{2}} - u_{i-\frac{1}{2}}}{\Delta x_i} r, & \text{if } r > 2, \\ u'_{\text{pw}1,2}(x_{i+\frac{1}{2}}^-) = 0, & \text{if } 0 < r < 1/2. \end{cases}$$

270 To avoid computational issues, when  $r \notin [1/50, 50]$ , the parabolic reconstruction is  
271 adopted directly.

272 For the FVS, as the cell average of the flux can be obtained through Simpson's  
273 rule,  $\bar{f}_i = (f_{i-\frac{1}{2}} + 4f_i + f_{i+\frac{1}{2}})/6$ , the flux derivatives can be computed by (2.14)-(2.15).

274 *Remark 2.2.* In [2], it is mentioned that if the signs of the derivatives of the  
275 parabolic reconstruction and the first-order reconstruction are the same, then the  
276 parabolic reconstruction is adopted. This strategy is not employed in this paper as  
277 the numerical results may be worse.

278 **2.4. Time discretization.** The fully-discrete scheme is obtained by using the  
279 SSP-RK3 method [16]

$$280 \quad (2.16) \quad \begin{aligned} \mathbf{U}^* &= \mathbf{U}^n + \Delta t^n \mathbf{L}(\mathbf{U}^n), \\ \mathbf{U}^{**} &= \frac{3}{4} \mathbf{U}^n + \frac{1}{4} (\mathbf{U}^* + \Delta t^n \mathbf{L}(\mathbf{U}^*)), \\ \mathbf{U}^{n+1} &= \frac{1}{3} \mathbf{U}^n + \frac{2}{3} (\mathbf{U}^{**} + \Delta t^n \mathbf{L}(\mathbf{U}^{**})), \end{aligned}$$

281 where  $\mathbf{L}$  is the right-hand side of the semi-discrete schemes (2.2) or (2.3). The time  
282 step size is determined by the usual CFL condition

$$283 \quad (2.17) \quad \Delta t^n = \frac{C_{\text{CFL}}}{\max_{i,\ell} \{\lambda_\ell(\bar{\mathbf{U}}_i) / \Delta x_i\}}.$$



284 **3. Convex limiting for the cell average.** Although the power law recon-  
 285 struction [4] has been shown to effectively reduce spurious oscillations, the numerical  
 286 solutions may still violate certain bounds, e.g., the appearance of negative density or  
 287 pressure, leading to unphysical solutions or even causing the simulations to blow up.  
 288 Since the degrees of freedom in the AF methods include both cell averages and point  
 289 values, it is necessary to design suitable BP limitings for both of them to achieve the  
 290 BP property. The limiting for the cell average has not been addressed much in the  
 291 literature, except for a very recent work [3].

292 **DEFINITION 3.1.** *An AF method is called bound-preserving (BP) if starting from*  
 293 *cell averages and point values in the admissible state set  $\mathcal{G}$ , the cell averages and point*  
 294 *values remain in  $\mathcal{G}$  at the next time step.*

295 This section presents a convex limiting approach to achieve the BP property of  
 296 the cell average update. The basic idea of the convex limiting approaches [17, 22, 29]  
 297 is to enforce the preservation of local and global bounds by constraining individual  
 298 numerical fluxes. The BP or invariant domain-preserving (IDP) properties of flux-  
 299 limited approximations are shown using representations in terms of intermediate states  
 300 that stay in convex admissible state sets [17, 21]. The low-order scheme is chosen as  
 301 the first-order LLF scheme

$$302 \quad \bar{U}_i^L = \bar{U}_i^n - \mu_i \left( \hat{F}_{i+\frac{1}{2}}^L - \hat{F}_{i-\frac{1}{2}}^L \right), \quad \mu_i = \Delta t^n / \Delta x_i,$$

$$303 \quad \hat{F}_{i+\frac{1}{2}}^L = \frac{1}{2} \left( F(\bar{U}_i^n) + F(\bar{U}_{i+1}^n) \right) - \frac{\alpha_{i+\frac{1}{2}}}{2} \left( \bar{U}_{i+1}^n - \bar{U}_i^n \right),$$

304 where  $\alpha_{i+\frac{1}{2}}$  is an *upper bound* for the maximum wave speed of the Riemann problem  
 305 with the initial data  $U_i, U_{i+1}$ , whose estimation for scalar conservation laws and the  
 306 Euler equations can be found in [19] and [18], respectively. Note that here  $\alpha_{i+\frac{1}{2}}$   
 307 is not the same as the one in the LLF FVS (2.11). Following [19], the first-order LLF  
 308 scheme can be rewritten as

$$309 \quad (3.1) \quad \bar{U}_i^L = \left[ 1 - \mu_i \left( \alpha_{i-\frac{1}{2}} + \alpha_{i+\frac{1}{2}} \right) \right] \bar{U}_i^n + \mu_i \alpha_{i-\frac{1}{2}} \tilde{U}_{i-\frac{1}{2}} + \mu_i \alpha_{i+\frac{1}{2}} \tilde{U}_{i+\frac{1}{2}},$$

310 with the intermediate states defined as

$$311 \quad (3.2) \quad \tilde{U}_{i-\frac{1}{2}} := \frac{1}{2} \left( \bar{U}_{i-1}^n + \bar{U}_i^n \right) + \frac{1}{2\alpha_{i-\frac{1}{2}}} \left[ F(\bar{U}_{i-1}^n) - F(\bar{U}_i^n) \right],$$

$$\tilde{U}_{i+\frac{1}{2}} := \frac{1}{2} \left( \bar{U}_i^n + \bar{U}_{i+1}^n \right) + \frac{1}{2\alpha_{i+\frac{1}{2}}} \left[ F(\bar{U}_i^n) - F(\bar{U}_{i+1}^n) \right].$$

312 **Remark 3.2.** As  $\alpha_{i+\frac{1}{2}}$  is chosen to be larger than the leftmost and rightmost wave  
 313 speed, the intermediate state defined in (3.2) is indeed an average of the exact Riemann  
 314 solution [19], thus it belongs to  $\mathcal{G}$ . For systems, it is also the intermediate state of the  
 315 HLL solver [24]. Moreover, the intermediate state (3.2) preserves all *convex invariants*  
 316 (e.g., density and pressure positivity, and minimum entropy principle for the Euler  
 317 equations) of initial value problems for hyperbolic systems [19].

318 **LEMMA 3.3** (Guermont and Popov [19]). If the time step size  $\Delta t^n$  satisfies

$$319 \quad (3.3) \quad \Delta t^n \leq \frac{\Delta x_i}{\alpha_{i-\frac{1}{2}} + \alpha_{i+\frac{1}{2}}},$$

320 then (3.1) is a convex combination, and the first-order LLF scheme (3.1) is BP.

321 The proof relies on the fact that the intermediate state (3.2) stays in the admissible  
 322 state set  $\mathcal{G}$  and the convexity of  $\mathcal{G}$ .

323 Upon defining the anti-diffusive flux  $\Delta\widehat{\mathbf{F}}_{i\pm\frac{1}{2}}^{\text{H}} := \widehat{\mathbf{F}}_{i\pm\frac{1}{2}}^{\text{H}} - \widehat{\mathbf{F}}_{i\pm\frac{1}{2}}^{\text{L}}$  with  $\widehat{\mathbf{F}}_{i\pm\frac{1}{2}}^{\text{H}} :=$   
 324  $\mathbf{F}(\mathbf{U}_{i\pm\frac{1}{2}})$ , a forward-Euler step applied to the semi-discrete high-order scheme for  
 325 the cell average (2.2) can be written as

$$\begin{aligned} 326 \quad \overline{\mathbf{U}}_i^{\text{H}} &= \overline{\mathbf{U}}_i^n - \mu_i(\widehat{\mathbf{F}}_{i+\frac{1}{2}}^{\text{H}} - \widehat{\mathbf{F}}_{i-\frac{1}{2}}^{\text{H}}) = \overline{\mathbf{U}}_i^n - \mu_i(\widehat{\mathbf{F}}_{i+\frac{1}{2}}^{\text{L}} - \widehat{\mathbf{F}}_{i-\frac{1}{2}}^{\text{L}}) - \mu_i(\Delta\widehat{\mathbf{F}}_{i+\frac{1}{2}} - \Delta\widehat{\mathbf{F}}_{i-\frac{1}{2}}) \\ 327 \quad (3.4) \quad &=: \left[1 - \mu_i(\alpha_{i-\frac{1}{2}} + \alpha_{i+\frac{1}{2}})\right] \overline{\mathbf{U}}_i^n + \mu_i\alpha_{i-\frac{1}{2}}\widetilde{\mathbf{U}}_{i-\frac{1}{2}}^{\text{H}} + \mu_i\alpha_{i+\frac{1}{2}}\widetilde{\mathbf{U}}_{i+\frac{1}{2}}^{\text{H}}, \\ 328 \quad \widetilde{\mathbf{U}}_{i-\frac{1}{2}}^{\text{H}} &:= \left(\widetilde{\mathbf{U}}_{i-\frac{1}{2}} + \frac{\Delta\widehat{\mathbf{F}}_{i-\frac{1}{2}}}{\alpha_{i-\frac{1}{2}}}\right), \quad \widetilde{\mathbf{U}}_{i+\frac{1}{2}}^{\text{H}} := \left(\widetilde{\mathbf{U}}_{i+\frac{1}{2}} - \frac{\Delta\widehat{\mathbf{F}}_{i+\frac{1}{2}}}{\alpha_{i+\frac{1}{2}}}\right). \end{aligned}$$

329 With the low-order scheme (3.1) and high-order scheme (3.4) having the same  
 330 form one can now define the limited scheme for the cell average as

$$331 \quad (3.5) \quad \overline{\mathbf{U}}_i^{\text{Lim}} = \left[1 - \mu_i(\alpha_{i-\frac{1}{2}} + \alpha_{i+\frac{1}{2}})\right] \overline{\mathbf{U}}_i^n + \mu_i\alpha_{i-\frac{1}{2}}\widetilde{\mathbf{U}}_{i-\frac{1}{2}}^{\text{Lim},+} + \mu_i\alpha_{i+\frac{1}{2}}\widetilde{\mathbf{U}}_{i+\frac{1}{2}}^{\text{Lim},-},$$

332 with the limited intermediate states

$$\begin{aligned} 333 \quad \widetilde{\mathbf{U}}_{i-\frac{1}{2}}^{\text{Lim},+} &= \widetilde{\mathbf{U}}_{i-\frac{1}{2}} + \frac{\Delta\widehat{\mathbf{F}}_{i-\frac{1}{2}}^{\text{Lim}}}{\alpha_{i-\frac{1}{2}}} := \widetilde{\mathbf{U}}_{i-\frac{1}{2}} + \frac{\theta_{i-\frac{1}{2}}\Delta\widehat{\mathbf{F}}_{i-\frac{1}{2}}}{\alpha_{i-\frac{1}{2}}}, \\ 334 \quad \widetilde{\mathbf{U}}_{i+\frac{1}{2}}^{\text{Lim},-} &= \widetilde{\mathbf{U}}_{i+\frac{1}{2}} - \frac{\Delta\widehat{\mathbf{F}}_{i+\frac{1}{2}}^{\text{Lim}}}{\alpha_{i+\frac{1}{2}}} := \widetilde{\mathbf{U}}_{i+\frac{1}{2}} - \frac{\theta_{i+\frac{1}{2}}\Delta\widehat{\mathbf{F}}_{i+\frac{1}{2}}}{\alpha_{i+\frac{1}{2}}}, \end{aligned}$$

335 where the coefficients  $\theta_{i\pm\frac{1}{2}} \in [0, 1]$ .

336 PROPOSITION 3.4. If the cell average at the last time step  $\overline{\mathbf{U}}_i^n$  and the limited  
 337 intermediate states  $\widetilde{\mathbf{U}}_{i\pm\frac{1}{2}}^{\text{Lim},\mp}$  belong to the admissible state set  $\mathcal{G}$ , then the limited  
 338 average update (3.5) is BP, i.e.,  $\overline{\mathbf{U}}_i^{\text{Lim}} \in \mathcal{G}$ , under the CFL condition (3.3). If the  
 339 SSP-RK3 (2.16) is used for the time integration, the high-order scheme is also BP.

340 *Proof.* Under the constraint (3.3), the limited cell average update  $\overline{\mathbf{U}}_i^{\text{Lim}}$  is a convex  
 341 combination of  $\overline{\mathbf{U}}_i^n$  and  $\widetilde{\mathbf{U}}_{i\pm\frac{1}{2}}^{\text{Lim},\mp}$ , thus it belongs to  $\mathcal{G}$  due to the convexity of  $\mathcal{G}$ . Because  
 342 the SSP-RK3 is a convex combination of forward-Euler stages, the high-order scheme  
 343 equipped with the SSP-RK3 is also BP according to the convexity.  $\square$

344 *Remark 3.5.* The scheme (3.5) is conservative as it amounts to using the nu-  
 345 merical flux  $\widehat{\mathbf{F}}_{i+\frac{1}{2}}^{\text{L}} + \theta_{i+\frac{1}{2}}\Delta\widehat{\mathbf{F}}_{i+\frac{1}{2}} = \theta_{i+\frac{1}{2}}\widehat{\mathbf{F}}_{i+\frac{1}{2}}^{\text{H}} + (1 - \theta_{i+\frac{1}{2}})\widehat{\mathbf{F}}_{i+\frac{1}{2}}^{\text{L}}$ , which is a convex  
 346 combination of the high-order and low-order fluxes.

347 *Remark 3.6.* It should be noted that the time step size (3.3) is determined based  
 348 on the solutions at  $t^n$ . If the constraint is not satisfied at the later stage of the  
 349 SSP-RK3, the BP property may not be achieved because (3.5) is no longer a convex  
 350 combination. In our implementation, we start from the usual CFL condition (2.17).  
 351 Then, if the high-order AF states need BP limitings and (3.2) is not BP or (3.3) is not  
 352 satisfied, the numerical solutions are set back to the last time step, and we rerun with  
 353 a halved time step size until (3.2) is BP and the constraint (3.3) is satisfied. This is  
 354 also a typical implementation to save computational costs in other BP methods.

355 The remaining task is to determine the coefficients at each interface  $\theta_{i\pm\frac{1}{2}}$  such  
 356 that  $\tilde{U}_{i\pm\frac{1}{2}}^{\text{Lim},\mp} \in \mathcal{G}$  and stay as close as possible to the high-order states  $\tilde{U}_{i\pm\frac{1}{2}}^{\text{H}}$ , i.e., the  
 357 goal is to find the largest  $\theta_{i\pm\frac{1}{2}} \in [0, 1]$  such that  $\tilde{U}_{i\pm\frac{1}{2}}^{\text{Lim},\mp} \in \mathcal{G}$ .

358 **3.1. Application to scalar conservation laws.** This section is devoted to  
 359 applying the convex limiting approach to scalar conservation laws (1.2), such that the  
 360 numerical solutions satisfy the global or local MP. For the global MP, the blending  
 361 coefficient  $\theta_{i+\frac{1}{2}} \in [0, 1]$  should be chosen such that  $m_0 \leq \tilde{u}_{i+\frac{1}{2}}^{\text{Lim},\pm} \leq M_0$ , with  $m_0, M_0$   
 362 defined in (1.4), which gives

$$363 \quad \theta_{i+\frac{1}{2}} = \begin{cases} \min \left\{ 1, \frac{\alpha_{i+\frac{1}{2}}(\tilde{u}_{i+\frac{1}{2}} - m_0)}{\Delta \hat{f}_{i+\frac{1}{2}}}, \frac{\alpha_{i+\frac{1}{2}}(M_0 - \tilde{u}_{i+\frac{1}{2}})}{\Delta \hat{f}_{i+\frac{1}{2}}} \right\}, & \text{if } \Delta \hat{f}_{i+\frac{1}{2}} > 0, \\ \min \left\{ 1, \frac{\alpha_{i+\frac{1}{2}}(m_0 - \tilde{u}_{i+\frac{1}{2}})}{\Delta \hat{f}_{i+\frac{1}{2}}}, \frac{\alpha_{i+\frac{1}{2}}(\tilde{u}_{i+\frac{1}{2}} - M_0)}{\Delta \hat{f}_{i+\frac{1}{2}}} \right\}, & \text{if } \Delta \hat{f}_{i+\frac{1}{2}} < 0. \end{cases}$$

364 To avoid a small denominator, the limited anti-diffusive flux can be obtained directly,

$$365 \quad \Delta \hat{f}_{i+\frac{1}{2}}^{\text{Lim}} = \begin{cases} \min \left\{ \Delta \hat{f}_{i+\frac{1}{2}}, \alpha_{i+\frac{1}{2}}(\tilde{u}_{i+\frac{1}{2}} - m_0), \alpha_{i+\frac{1}{2}}(M_0 - \tilde{u}_{i+\frac{1}{2}}) \right\}, & \text{if } \Delta \hat{f}_{i+\frac{1}{2}} \geq 0, \\ \max \left\{ \Delta \hat{f}_{i+\frac{1}{2}}, \alpha_{i+\frac{1}{2}}(m_0 - \tilde{u}_{i+\frac{1}{2}}), \alpha_{i+\frac{1}{2}}(\tilde{u}_{i+\frac{1}{2}} - M_0) \right\}, & \text{otherwise.} \end{cases}$$

366 On the other hand, one can also enforce the local MP  $u_i^{\min} \leq \tilde{u}_{i+\frac{1}{2}}^{\text{Lim},-} \leq u_i^{\max}$ ,  
 367  $u_{i+1}^{\min} \leq \tilde{u}_{i+\frac{1}{2}}^{\text{Lim},+} \leq u_{i+1}^{\max}$ , which helps to suppress spurious oscillations and improve  
 368 shock-capturing ability. The corresponding limited anti-diffusive flux is

$$369 \quad \Delta \hat{f}_{i+\frac{1}{2}}^{\text{Lim}} = \begin{cases} \min \left\{ \Delta \hat{f}_{i+\frac{1}{2}}, \alpha_{i+\frac{1}{2}}(\tilde{u}_{i+\frac{1}{2}} - u_i^{\min}), \alpha_{i+\frac{1}{2}}(u_{i+1}^{\max} - \tilde{u}_{i+\frac{1}{2}}) \right\}, & \text{if } \Delta \hat{f}_{i+\frac{1}{2}} \geq 0, \\ \max \left\{ \Delta \hat{f}_{i+\frac{1}{2}}, \alpha_{i+\frac{1}{2}}(u_{i+1}^{\min} - \tilde{u}_{i+\frac{1}{2}}), \alpha_{i+\frac{1}{2}}(\tilde{u}_{i+\frac{1}{2}} - u_i^{\max}) \right\}, & \text{otherwise.} \end{cases}$$

370 The choice of the local bounds can be based on the intermediate states

$$371 \quad u_i^{\min} = \min \left\{ \tilde{u}_i^n, \tilde{u}_{i-\frac{1}{2}}, \tilde{u}_{i+\frac{1}{2}} \right\}, \quad u_i^{\max} = \max \left\{ \tilde{u}_i^n, \tilde{u}_{i-\frac{1}{2}}, \tilde{u}_{i+\frac{1}{2}} \right\}.$$

372 Finally, the numerical flux is

$$373 \quad (3.6) \quad \hat{f}_{i+\frac{1}{2}}^{\text{Lim}} = \hat{f}_{i+\frac{1}{2}}^{\text{L}} + \Delta \hat{f}_{i+\frac{1}{2}}^{\text{Lim}}.$$

374 **3.2. Application to the compressible Euler equations.** This section aims  
 375 at enforcing the strict positivity of density and pressure, i.e.,  $\rho > \varepsilon$ ,  $p > \varepsilon$ , with  $\varepsilon$  a  
 376 small positive number close to zero, chosen as  $10^{-13}$  in our numerical tests.

377 **3.2.1. Positivity of density.** The first step is to impose the density positivity  
 378  $\tilde{U}_{i+\frac{1}{2}}^{\text{Lim},\pm,\rho} > \varepsilon$ , where  $U^{*,\rho}$  denotes the density component of  $U^*$ . The corresponding  
 379 density component of the limited anti-diffusive flux is

$$380 \quad \Delta \hat{F}_{i+\frac{1}{2}}^{\text{Lim},*,\rho} = \begin{cases} \min \left\{ \Delta \hat{F}_{i+\frac{1}{2}}^{\rho}, \alpha_{i+\frac{1}{2}} \left( \tilde{U}_{i+\frac{1}{2}}^{\rho} - \varepsilon \right) \right\}, & \text{if } \Delta \hat{F}_{i+\frac{1}{2}}^{\rho} \geq 0, \\ \max \left\{ \Delta \hat{F}_{i+\frac{1}{2}}^{\rho}, \alpha_{i+\frac{1}{2}} \left( \varepsilon - \tilde{U}_{i+\frac{1}{2}}^{\rho} \right) \right\}, & \text{otherwise.} \end{cases}$$

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

383 **3.2.2. Positivity of pressure.** The second step is to enforce pressure positivity  
 384  $p(\tilde{U}_{i+\frac{1}{2}}^{\text{Lim},\pm}) > \varepsilon$ , where  $p(U^*)$  denotes the pressure recovered from  $U^*$ , with

$$385 \quad \tilde{U}_{i+\frac{1}{2}}^{\text{Lim},\pm} = \tilde{U}_{i+\frac{1}{2}} \pm \frac{\theta_{i+\frac{1}{2}} \Delta \hat{F}_{i+\frac{1}{2}}^{\text{Lim},*}}{\alpha_{i+\frac{1}{2}}}, \quad \Delta \hat{F}_{i+\frac{1}{2}}^{\text{Lim},*} = \hat{F}_{i+\frac{1}{2}}^{\text{Lim},*} - \hat{F}_{i+\frac{1}{2}}^{\text{L}}.$$

386 Such constraints lead to two inequalities

$$387 \quad (3.7) \quad \frac{A_{i+\frac{1}{2}}}{\alpha_{i+\frac{1}{2}}^2} \theta_{i+\frac{1}{2}}^2 \pm \frac{B_{i+\frac{1}{2}}}{\alpha_{i+\frac{1}{2}}} \theta_{i+\frac{1}{2}} < C_{i+\frac{1}{2}},$$

388 with the coefficients

$$389 \quad A_{i+\frac{1}{2}} = \frac{1}{2} \left( \Delta \hat{F}_{i+\frac{1}{2}}^{\text{Lim},*\rho\nu} \right)^2 - \Delta \hat{F}_{i+\frac{1}{2}}^{\text{Lim},*\rho} \Delta \hat{F}_{i+\frac{1}{2}}^{\text{Lim},*E},$$

$$390 \quad B_{i+\frac{1}{2}} = \alpha_{i+\frac{1}{2}} \left( \Delta \hat{F}_{i+\frac{1}{2}}^{\text{Lim},*\rho} \tilde{U}_{i+\frac{1}{2}}^E + \tilde{U}_{i+\frac{1}{2}}^\rho \Delta \hat{F}_{i+\frac{1}{2}}^{\text{Lim},*E} - \Delta \hat{F}_{i+\frac{1}{2}}^{\text{Lim},*\rho\nu} \tilde{U}_{i+\frac{1}{2}}^{\rho\nu} - \varepsilon \Delta \hat{F}_{i+\frac{1}{2}}^{\text{Lim},*\rho} \right),$$

$$391 \quad C_{i+\frac{1}{2}} = \alpha_{i+\frac{1}{2}}^2 \left( \tilde{U}_{i+\frac{1}{2}}^\rho \tilde{U}_{i+\frac{1}{2}}^E - \frac{1}{2} \left( \tilde{U}_{i+\frac{1}{2}}^{\rho\nu} \right)^2 - \varepsilon \tilde{U}_{i+\frac{1}{2}}^\rho \right).$$

392 Following [29], the inequalities (3.7) hold under the linear sufficient condition

$$393 \quad \max\{0, A_{i+\frac{1}{2}}\} + |B_{i+\frac{1}{2}}| \leq C_{i+\frac{1}{2}},$$

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

$$395 \quad \theta_{i+\frac{1}{2}} = \min \left\{ 1, \frac{C_{i+\frac{1}{2}}}{\max\{0, A_{i+\frac{1}{2}}\} + |B_{i+\frac{1}{2}}|} \right\},$$

396 and the final limited numerical flux is

$$397 \quad (3.8) \quad \hat{F}_{i+\frac{1}{2}}^{\text{Lim}} = \hat{F}_{i+\frac{1}{2}}^{\text{L}} + \theta_{i+\frac{1}{2}} \Delta \hat{F}_{i+\frac{1}{2}}^{\text{Lim},*}.$$

398 **4. Scaling limiter for point value.** To achieve the BP property, it is also necessary  
 399 to introduce BP limiting for the point value. As one will see in the numerical  
 400 tests in Section 5, using power law reconstruction or BP limiting for cell average,  
 401 individually or in combination, cannot guarantee the bounds. As there is no conser-  
 402 vation requirement on the point value update, a simple scaling limiter [31] is directly  
 403 performed on the high-order point values rather than on the flux for the cell average.

404 A first-order LLF scheme for the point value update can be

$$405 \quad (4.1) \quad U_{i+\frac{1}{2}}^{\text{L}} = U_{i+\frac{1}{2}}^n - \frac{2\Delta t^n}{\Delta x_i + \Delta x_{i+1}} \left( \hat{F}_{i+1}^{\text{L}}(U_{i+\frac{1}{2}}^n, U_{i+\frac{3}{2}}^n) - \hat{F}_i^{\text{L}}(U_{i-\frac{1}{2}}^n, U_{i+\frac{1}{2}}^n) \right),$$

406 with the numerical flux

$$407 \quad \hat{F}_i^{\text{L}}(U_{i-\frac{1}{2}}^n, U_{i+\frac{1}{2}}^n) = \frac{1}{2} \left( F(U_{i-\frac{1}{2}}^n) + F(U_{i+\frac{1}{2}}^n) \right) - \frac{\alpha_i}{2} \left( U_{i+\frac{1}{2}}^n - U_{i-\frac{1}{2}}^n \right),$$

$$408 \quad \alpha_i = \max\{\lambda(U_{i-\frac{1}{2}}^n), \lambda(U_{i+\frac{1}{2}}^n)\}.$$

409 Such an LLF scheme can be interpreted as a scheme on a staggered mesh if the point  
 410 value is viewed as the cell average on the staggered mesh. Based on the proof in [33],  
 411 it is straightforward to obtain the following Lemma.

412 LEMMA 4.1. The LLF scheme for the point value (4.1) is BP under the CFL  
413 condition

$$414 \quad (4.2) \quad \Delta t^n \leq \frac{\Delta x_i + \Delta x_{i+1}}{4\alpha_i}.$$

415 The limited state is obtained by blending the high-order AF scheme (2.3) with  
416 the forward Euler scheme and the LLF scheme (4.1) as  $\mathbf{U}_{i+\frac{1}{2}}^{\text{Lim}} = \theta_{i+\frac{1}{2}} \mathbf{U}_{i+\frac{1}{2}}^{\text{H}} + (1 -$   
417  $\theta_{i+\frac{1}{2}}) \mathbf{U}_{i+\frac{1}{2}}^{\text{L}}$ , such that  $\mathbf{U}_{i+\frac{1}{2}}^{\text{Lim}} \in \mathcal{G}$ .

418 *Remark 4.2.* In the FVS for the point value update, the cell-centered value  $\mathbf{U}_i$   
419 is used. It is possible that  $\mathbf{U}_i \notin \mathcal{G}$ , then it is set as  $\overline{\mathbf{U}}_i$  in such cases, which is a  
420 reasonable second-order approximation.

421 **4.1. Application to scalar conservation laws.** This section enforces the  
422 global MP  $m_0 \leq u_{i+\frac{1}{2}}^{\text{Lim}} \leq M_0$  by choosing the coefficient as

$$423 \quad \theta_{i+\frac{1}{2}} = \begin{cases} \frac{u_{i+\frac{1}{2}}^{\text{L}} - m_0}{u_{i+\frac{1}{2}}^{\text{L}} - u_{i+\frac{1}{2}}^{\text{H}}}, & \text{if } u_{i+\frac{1}{2}}^{\text{H}} < m_0, \\ M_0 - u_{i+\frac{1}{2}}^{\text{L}}, & \text{if } u_{i+\frac{1}{2}}^{\text{H}} > M_0. \\ \frac{u_{i+\frac{1}{2}}^{\text{H}} - u_{i+\frac{1}{2}}^{\text{L}}}{u_{i+\frac{1}{2}}^{\text{H}} - u_{i+\frac{1}{2}}^{\text{L}}}, & \end{cases}$$

424 The final limited state is

$$425 \quad (4.3) \quad u_{i+\frac{1}{2}}^{\text{Lim}} = \theta_{i+\frac{1}{2}} u_{i+\frac{1}{2}}^{\text{H}} + (1 - \theta_{i+\frac{1}{2}}) u_{i+\frac{1}{2}}^{\text{L}}.$$

426 **4.2. Application to the compressible Euler equations.** The limiting con-  
427 sists of two steps. First, the high-order state  $\mathbf{U}_{i+\frac{1}{2}}^{\text{H}}$  is modified as  $\mathbf{U}_{i+\frac{1}{2}}^{\text{Lim},*}$ , such that its  
428 density component satisfies  $\mathbf{U}_{i+\frac{1}{2}}^{\text{Lim},*\rho} > \varepsilon$ . Solving this inequality gives the coefficient

$$429 \quad \theta_{i+\frac{1}{2}}^* = \begin{cases} \frac{\mathbf{U}_{i+\frac{1}{2}}^{\text{L},\rho} - \varepsilon}{\mathbf{U}_{i+\frac{1}{2}}^{\text{L},\rho} - \mathbf{U}_{i+\frac{1}{2}}^{\text{H},\rho}}, & \text{if } \mathbf{U}_{i+\frac{1}{2}}^{\text{H},\rho} < \varepsilon, \\ 1, & \text{otherwise.} \end{cases}$$

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

432 Then the limited state  $\mathbf{U}_{i+\frac{1}{2}}^{\text{Lim},*}$  is modified as  $\mathbf{U}_{i+\frac{1}{2}}^{\text{Lim}}$ , such that it gives positive  
433 pressure, i.e.,  $p(\mathbf{U}_{i+\frac{1}{2}}^{\text{Lim}}) > \varepsilon$ . Let  $\mathbf{U}_{i+\frac{1}{2}}^{\text{Lim}} = \theta_{i+\frac{1}{2}}^{**} \mathbf{U}_{i+\frac{1}{2}}^{\text{Lim},*} + (1 - \theta_{i+\frac{1}{2}}^{**}) \mathbf{U}_{i+\frac{1}{2}}^{\text{L}}$ . Note that the  
434 pressure is a concave function (see e.g. [45]) of the conservative variables, such that

$$435 \quad p(\mathbf{U}_{i+\frac{1}{2}}^{\text{Lim}}) \geq \theta_{i+\frac{1}{2}}^{**} p(\mathbf{U}_{i+\frac{1}{2}}^{\text{Lim},*}) + (1 - \theta_{i+\frac{1}{2}}^{**}) p(\mathbf{U}_{i+\frac{1}{2}}^{\text{L}})$$

436 based on Jensen's inequality and  $\mathbf{U}_{i+\frac{1}{2}}^{\text{Lim},*\rho} > 0$ ,  $\mathbf{U}_{i+\frac{1}{2}}^{\text{L},\rho} > 0$ ,  $\theta_{i+\frac{1}{2}}^{**} \in [0, 1]$ . Thus a  
437 sufficient condition is

$$438 \quad \theta_{i+\frac{1}{2}}^{**} = \begin{cases} \frac{p(\mathbf{U}_{i+\frac{1}{2}}^{\text{L}}) - \varepsilon}{p(\mathbf{U}_{i+\frac{1}{2}}^{\text{L}}) - p(\mathbf{U}_{i+\frac{1}{2}}^{\text{Lim},*})}, & \text{if } p(\mathbf{U}_{i+\frac{1}{2}}^{\text{Lim},*}) < \varepsilon, \\ 1, & \text{otherwise.} \end{cases}$$

439 The final limited state is

$$440 \quad (4.4) \quad \mathbf{U}_{i+\frac{1}{2}}^{\text{Lim}} = \theta_{i+\frac{1}{2}}^{**} \mathbf{U}_{i+\frac{1}{2}}^{\text{Lim},*} + \left(1 - \theta_{i+\frac{1}{2}}^{**}\right) \mathbf{U}_{i+\frac{1}{2}}^{\text{L}}.$$

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

442 **THEOREM 4.3.** If the initial numerical solution  $\overline{\mathbf{U}}_i^0, \mathbf{U}_{i+\frac{1}{2}}^0 \in \mathcal{G}$  for all  $i$ , and the  
443 time step size satisfies (3.3) and (4.2), then the AF methods (2.2)-(2.3) equipped with  
444 the SSP-RK3 (2.16) and the BP limitings

- 445 • (3.6) and (4.3) preserve the maximum principle for scalar case;
- 446 • (3.8) and (4.4) preserve the density and pressure positivity for the Euler equations.

447 **5. Numerical results.** This section conducts some numerical tests to verify  
448 the accuracy of using the FVS for point value updates, the BP property, and the  
449 shock-capturing ability of our AF methods.

450 **5.1. Scalar conservation laws.** This section shows the results for the linear  
451 advection equation and the Burgers' equation, which demonstrate that the proposed  
452 limiting can preserve the MP and suppress oscillations well.

453 *Example 5.1* (Advection equation). Consider the 1D advection equation  $u_t + u_x =$   
454  $0$ , on the periodic domain  $[-1, 1]$  with the initial data [27]

$$455 \quad \begin{cases} \frac{1}{6} (G_1(x, \beta, z - \delta) + G_1(x, \beta, z + \delta) + 4G_1(x, \beta, z)), & \text{if } -0.8 \leq x \leq -0.6, \\ 1, & \text{if } -0.4 \leq x \leq -0.2, \\ 1 - |10(x - 0.1)|, & \text{if } 0 \leq x \leq 0.2, \\ \frac{1}{6} (G_2(x, \alpha, a - \delta) + G_2(x, \alpha, a + \delta) + 4G_2(x, \alpha, a)), & \text{if } 0.4 \leq x \leq 0.6, \\ 0, & \text{otherwise,} \end{cases}$$

456 where  $G_1(x, \beta, z) = \exp(-\beta(x - z)^2)$ ,  $G_2(x, \alpha, a) = \sqrt{\max(1 - \alpha^2(x - a)^2, 0)}$ , and the  
457 constants are  $a = -0.5, z = -0.7, \delta = 0.005, \alpha = 10, \beta = \ln 2 / (36\delta^2)$ . The problem is  
458 solved for one period, i.e., until  $T = 2$ .

459 For the advection equation, the JS and LLF FVS are equivalent. The maximal  
460 CFL number leading to a stable simulation is 0.41 without any limiting, and it reduces  
461 to 0.13 when only the power law reconstruction is activated, and it increases a little  
462 bit to 0.42 when only the BP limitings are used. When the power law reconstruction  
463 and the BP limitings are employed together, the maximal CFL number can be 0.4.  
464 The reduction of the CFL number with the power law reconstruction for semi-discrete  
465 AF has, in fact, not been noticed previously. Thus, in the following simulations we  
466 try not to use the power law reconstruction unless otherwise stated.

467 The results obtained with different limitings are shown in Figure 2, which are  
468 computed with 400 cells and the CFL number is 0.1. The ranges of the numerical  
469 solutions are listed in Table 1, considering both the cell averages and point values.  
470 One can observe that there are some oscillations near the discontinuities without  
471 any limiting, and that the power law reconstruction can eliminate the oscillations  
472 effectively but is still not BP. The activation of either the BP limiting for the cell  
473 average alone or the BP limiting for the point value alone also fails to preserve the  
474 bounds  $[0, 1]$ , as one can see from Table 1, as is the case when using both the BP  
475 limiting for the cell average and the power law reconstruction in the point value  
476 update. Only when a BP limiting is performed on both the cell average and the  
477 point value, the BP property is achieved, showing that using the two BP limitings

478 simultaneously is necessary for the preservation of the MP. Figure 2 also shows the  
 479 results obtained by imposing the global or local MP for the cell average, and global  
 480 MP for the point value (without power law reconstruction), indicating that the use of  
 481 local MP tends to dissipate the numerical solutions near the discontinuities and clip  
 482 maxima more than the global MP.

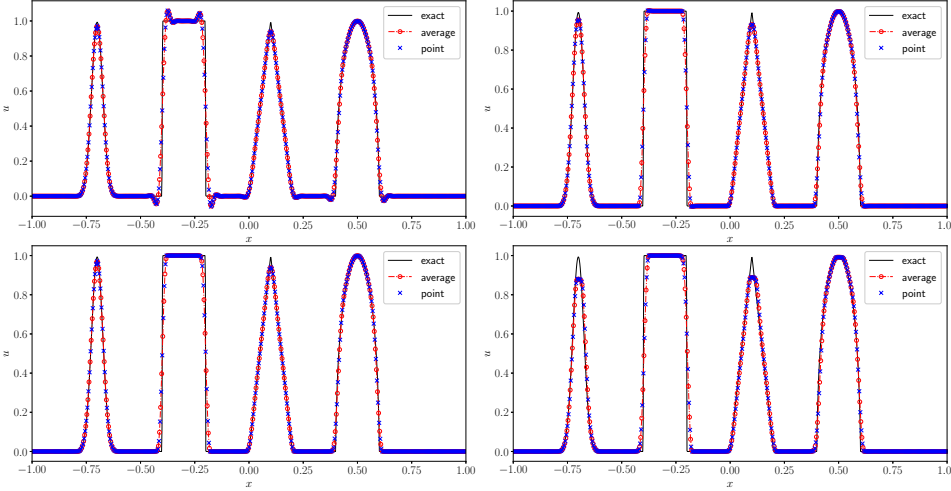


Fig. 2: Example 5.1, advection. The results are obtained without any limiting (upper left), with power law reconstruction (upper right), with BP limitings imposing global MP for the cell average and point value (lower left), with BP limitings imposing local and global MP for the cell average and point value (lower right).

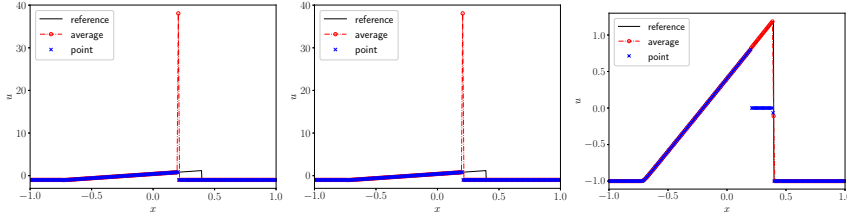
none	$[-5.9 \times 10^{-2}, 1 + 5.9 \times 10^{-2}]$	✗
PLR	$[-2.7 \times 10^{-3}, 1 + 2.6 \times 10^{-3}]$	✗
global MP for average	$[-1.7 \times 10^{-3}, 1 + 1.7 \times 10^{-3}]$	✗
local MP for average	$[-1.3 \times 10^{-3}, 1 + 1.3 \times 10^{-3}]$	✗
global MP for point	$[-3.0 \times 10^{-4}, 1 + 2.6 \times 10^{-4}]$	✗
PLR + global MP for average	$[-9.8 \times 10^{-6}, 1 + 2.7 \times 10^{-6}]$	✗
PLR + local MP for average	$[-1.4 \times 10^{-5}, 1 + 1.9 \times 10^{-5}]$	✗
global MP for average + global MP for point	$[0.0, 1.0]$	✓
local MP for average + global MP for point	$[0.0, 1 - 9.4 \times 10^{-13}]$	✓
PLR + global MP for average + global MP for point	$[0.0, 1 - 1.1 \times 10^{-16}]$	✓
PLR + local MP for average + global MP for point	$[0.0, 1 - 7.3 \times 10^{-14}]$	✓

Table 1: Example 5.1, advection. The ranges of the numerical solutions (including both the cell averages and the point values) obtained with different limitings after one period. “PLR” denotes the power law reconstruction.

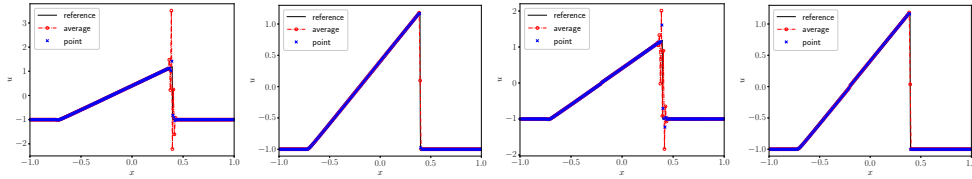
483 Example 5.2 (Self-steepening shock). Consider the 1D Burgers’ equation  $u_t +$   
 484  $(\frac{1}{2}u^2)_x = 0$  on the domain  $[-1, 1]$  with periodic boundary conditions. This test is  
 485 solved until  $T = 0.5$  with the initial condition as a square wave

$$u_0(x) = \begin{cases} 2, & \text{if } |x| < 0.2, \\ -1, & \text{otherwise.} \end{cases}$$

487 **Figures 3 and 4** plot the cell averages and point values based on different point  
 488 value updates with 200 cells, as well as the reference solution. The spike generation  
 489 has been observed in [25], and the reason is also discussed in **Subsection 2.2**. Such  
 490 spike generation cannot be eliminated by using the power law reconstruction, nor do  
 491 both BP limitings help to eliminate artefacts, as can be seen from **Figure 3**. The  
 492 numerical solutions based on the LLF or SW FVS are shown in **Figure 4**, in which no  
 493 spike appears. There are some oscillations near the discontinuity without limitings,  
 494 and the numerical solutions agree well with the reference solution when the limitings  
 495 are activated.



**Fig. 3: Example 5.2**, self-steepening shock for the Burgers' equation. The numerical solutions are based on the JS. From left to right: without limiting, with the power law reconstruction, with the BP limitings imposing local and global MP for the cell average and point value update, respectively.



**Fig. 4: Example 5.2**, self-steepening shock for the Burgers' equation. From left to right: the LLF FVS without limiting, the LLF FVS with limitings, the SW FVS without limiting, the SW FVS with limitings. The limitings consider the local and global MP for the cell average and point value updates, respectively.

496 **5.2. The compressible Euler equations.** This section shows some challeng-  
 497 ing tests, which require the BP property of the numerical methods in order to prevent  
 498 simulations from crashing at some time. The adiabatic index is  $\gamma = 1.4$  unless other-  
 499 wise stated. Note that the BP limiting naturally reduces some oscillations.

500 *Example 5.3* (1D accuracy test for the Euler equations). This test is used to  
 501 examine the accuracy of using different point value updates. The domain is  $[0, 1]$  with  
 502 periodic boundary conditions. Two manufactured solutions are constructed by adding  
 503 additional source terms  $\mathbf{S}$  to the Euler equations,

$$\begin{aligned}
 (5.1) \quad \rho &= 4 + 0.1s_1, \quad v_1 = s_1, \quad p = (6002 + 398c_2 + 305s_1 + 5s_3)/1000, \\
 \mathbf{S} &= (\pi(39c_1 + s_2)/5, \quad -\pi(905c_1 + 15c_3 - 776s_2)/125, \\
 &\quad \pi c_1(20421 + 1179c_2 + 2160s_1 + 20s_3)/500),
 \end{aligned}$$



507 and

$$\begin{aligned}
 508 \quad (5.2) \quad & \rho = 4 + 0.1s_1, \quad v_1 = 2 + 0.5s_1, \quad p = (12328 + 472c_2 - 5455s_1 + 5s_3)/4000, \\
 509 \quad & \mathbf{S} = (\pi(42c_1 + s_2)/10, \pi(4855c_1 - 15c_3 + 914s_2)/500, \\
 510 \quad & \pi c_1(14991 + 369c_2 - 2983s_1 + 5s_3)/1000),
 \end{aligned}$$

511 with  $s_k = \sin(2k\pi(x - t))$ ,  $c_k = \cos(2k\pi(x - t))$ ,  $k = 1, 2, 3$ . The source terms are  
 512 discretized by using Simpson’s rule for the cell average update. The problem is solved  
 513 until  $T = 0.4$ .

514 In this test, the maximal CFL number is around 0.18 for the VH FVS, while  
 515 around 0.43 for the JS, LLF, and SW FVS, thus we run the test with the same CFL  
 516 number as 0.18. Figure 5 shows the following errors and corresponding convergence  
 517 rates for the conservative variables in the  $\ell^1$  norm. It is seen that for the first exact  
 518 solution (5.1), the JS and all the FVS except for the SW FVS achieve the designed  
 519 third-order accuracy, while the SW FVS only gives second-order accuracy. Figure 6  
 520 plots the density and velocity profiles obtained by the SW FVS with 80 cells. One  
 521 can observe some defects in the density when the velocity is zero, similar to the “sonic  
 522 point glitch” in the literature [37]. For the second exact solution (5.2), the velocity  
 523 stays away from zero and no such issue appears. One possible reason is that the SW  
 524 FVS is based on the absolute value of the eigenvalues, which is not smooth when the  
 525 velocity is zero. Such an issue remains to be further explored in the future.

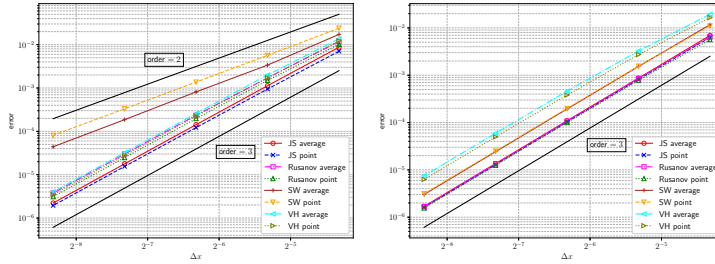


Fig. 5: Example 5.3, the accuracy tests for the 1D Euler equations based on the manufactured solutions (5.1) and (5.2) for the left and right plots, respectively.

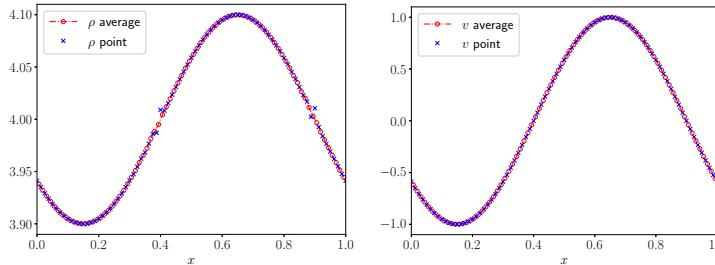


Fig. 6: Example 5.3, the density (left) and velocity (right) are obtained with the SW FVS and 80 cells for the 1D Euler equations based on the initial data (5.1).

526 *Example 5.4* (Double rarefaction problem). The exact solution to this problem  
 527 contains a vacuum, so that it is often used to verify the BP property of numerical  
 528 methods. The test is solved on a domain  $[0, 1]$  until  $T = 0.3$  with the initial data

$$529 \quad (\rho, v, p) = \begin{cases} (7, -1, 0.2), & \text{if } x < 0.5, \\ (7, 1, 0.2), & \text{otherwise.} \end{cases}$$

530 In this test, the AF method based on any kind of point value update mentioned  
 531 in this paper gives negative density or pressure without the BP limitations. **Figure 7**  
 532 shows the density computed with 400 cells and the BP limitations for the cell average  
 533 and point value updates. The power law reconstruction is not used in this test, and  
 534 the CFL number is 0.4 for all kinds of point value updates, except for 0.1 for the VH  
 535 FVS. One observes that the BP AF method gets good performance for this example.

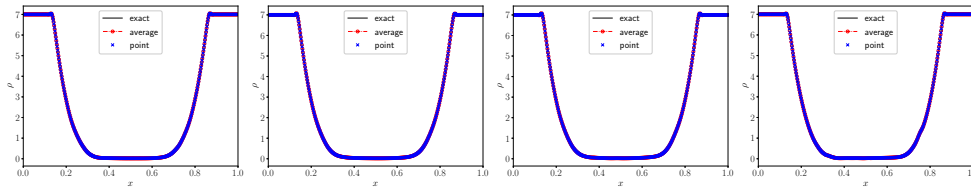


Fig. 7: **Example 5.4**, double rarefaction Riemann problem. The numerical solutions are computed with BP limitations for the cell average and point value updates on a uniform mesh of 400 cells. The power law reconstruction is not used. From left to right: JS, LLF, SW, and VH FVS.

536 *Example 5.5* (LeBlanc shock tube). This is a Riemann problem with an extremely  
 537 large initial pressure ratio. This test is solved until  $T = 5 \times 10^{-6}$  on a domain  $[0, 1]$   
 538 with the initial data

$$539 \quad (\rho, v, p) = \begin{cases} (2, 0, 10^9), & \text{if } x < 0.5, \\ (10^{-3}, 0, 1), & \text{otherwise.} \end{cases}$$

540 Without the BP limitations, the simulation will stop due to negative density or  
 541 pressure. **Figure 8** shows the density computed on a uniform mesh of 400 and 6000 cells  
 542 with the BP limitations for the cell average and point value updates. The CFL number  
 543 is 0.4 for the LLF and SW FVS, and 0.15 for the JS and VH FVS for stability when  
 544 the power law reconstruction is not used. It is seen that the numerical solutions on the  
 545 coarse mesh deviate from the exact solutions, which has also been observed in other  
 546 high-order BP methods, e.g., [44]. As the number of the mesh cells increases from 400  
 547 to 6000, one can observe from **Figure 8** that the numerical solutions converge to the  
 548 exact solutions with only a few overshoots/undershoots at the contact discontinuity.  
 549 The LLF and SW FVS give better results.

550 To verify whether the power law reconstruction can suppress spurious oscillations  
 551 and overshoots/undershoots, we rerun the test with the CFL number 0.1, and the  
 552 density profiles are shown in **Figure 9**. It is obvious that only reducing the CFL  
 553 number does not change the numerical solutions much except that the oscillations  
 554 near the contact discontinuity based on the VH FVS are damped. When the power  
 555 law reconstruction is activated, the overshoots/undershoots are reduced for the JS,

556 LLF, and SW FVS, while the VH FVS gives worse results even with a smaller CFL  
 557 number (e.g. 0.02, not shown here), which needs further investigation.

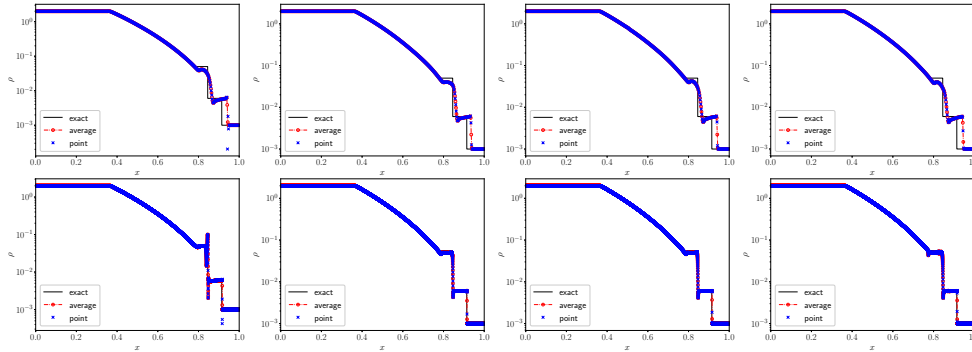


Fig. 8: **Example 5.5**, LeBlanc Riemann problem. The numerical solutions are computed with the BP limitations for the cell average and point value updates on a uniform mesh of 400 cells (top) and 6000 cells (bottom). The CFL number is 0.4 and the power law reconstruction is not used. From left to right: JS, LLF, SW, and VH FVS.

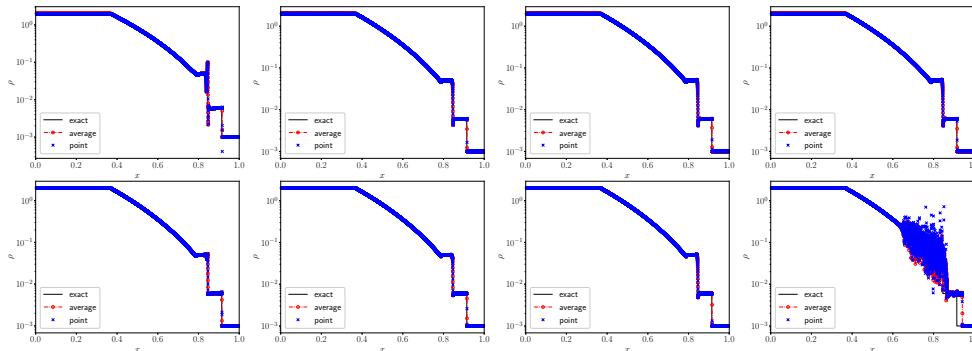


Fig. 9: **Example 5.5**, LeBlanc Riemann problem. The numerical solutions are computed with the BP limitations for the cell average and point value updates on a uniform mesh of 6000 cells. From left to right: JS, LLF, SW, and VH FVS. The CFL number is 0.1 and the power law reconstruction is not activated (top) and activated (bottom).

558 *Example 5.6* (Sedov problem). In this problem, a volume of uniform density and  
 559 temperature is initialized, and a large quantity of thermal energy is injected at the  
 560 center, developing into a blast wave that evolves in time in a self-similar fashion [35].  
 561 An exact analytical solution based on self-similarity arguments is available [28], which  
 562 contains very low density with strong shocks. The initial density is one, velocity is  
 563 zero, and total energy is  $10^{-12}$  everywhere except that in the center cell, the total  
 564 energy of the cell average and point values at two cell interfaces are  $3.2 \times 10^6 / \Delta x$   
 565 with  $\Delta x = 4/N$  with  $N$  the number of cells, which is used to emulate a  $\delta$ -function at  
 566 the center. The test is solved until  $T = 5 \times 10^{-6}$ .

567 This test is run with  $N = 801$  cells, and the density plots in the right half domain  
 568 are shown in **Figure 10**. The BP limitations are adopted for the cell average and point

569 value updates, while the power law reconstruction is not used. The maximal CFL  
 570 numbers for different point value updates to be stable are also listed in the caption,  
 571 i.e., 0.1 for the JS, 0.4, 0.3, and 0.25 for the LLF, SW, and VH FVS, respectively.  
 572 The numerical solutions obtained by the three FVS are nearly the same, while there  
 573 are some defects in the solution based on the JS. Thus the LLF FVS is superior to  
 574 others regarding the time step size and the shock-capturing ability.

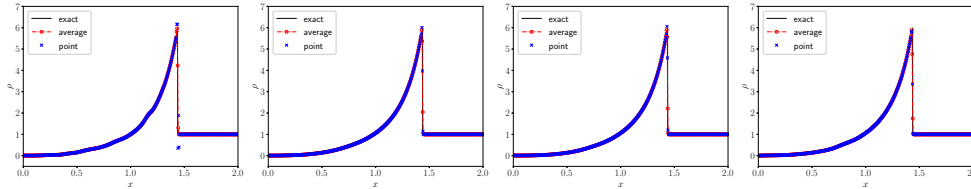


Fig. 10: [Example 5.6](#), Sedov problem. The numerical solutions are computed with the BP limitings for the cell average and point value updates on a uniform mesh of 801 cells, without the power law reconstruction. The CFL number is (from left to right): 0.1 for the JS, 0.4 for the LLF FVS, 0.3 for the SW FVS, 0.25 for the VH FVS.

575 *Example 5.7* (Blast wave interaction [\[41\]](#)). This test describes the interaction of  
 576 two strong shocks in the domain  $[0, 1]$  with reflective boundary conditions. The test  
 577 is solved until  $T = 0.038$ .

578 Due to the low-pressure region, the schemes blow up without the BP limitings.  
 579 [Figure 11](#) shows the density profiles and corresponding enlarged views in  
 580  $x \in [0.62, 0.82]$  obtained by using the BP limitings on a uniform mesh of 800 cells,  
 581 in which the power law reconstruction is not activated. It is seen that the numerical  
 582 solutions are close to the reference solution, although there are some oscillations in  
 583 the enlarged views. Then the power law reconstruction is additionally adopted to  
 584 see if it can suppress the oscillations. The results with the CFL number 0.1 and a  
 585 refined mesh of 1600 cells are shown in [Figure 12](#), from which one can observe that  
 586 the oscillations reduce, and the LLF FVS gives the best result.

587 *Remark 5.8.* In the numerical tests, the maximal CFL numbers for stability are  
 588 obtained by experiments. Note that the constraints [\(3.3\)](#) and [\(4.2\)](#) are used to guar-  
 589 antee the BP property, while the reduction of the CFL numbers is due to the stability  
 590 issue for different FVS and power law reconstruction.

591 **6. Conclusion.** In the active flux (AF) methods, the way how point values at  
 592 cell interfaces are updated is essential to achieve stability and high-order accuracy.  
 593 The point value update based on Jacobian splitting (JS) may lead to the so-called  
 594 transonic issue for nonlinear problems due to inaccurate estimation of the upwind di-  
 595 rection. This paper proposed to use the flux vector splitting (FVS) for the point value  
 596 update instead of the JS, which keeps the continuous reconstruction as the original AF  
 597 methods, and offers a natural and uniform remedy to the transonic issue. To further  
 598 improve the robustness of the AF methods, this paper developed bound-preserving  
 599 (BP) AF methods for general one-dimensional hyperbolic conservation laws, achieved  
 600 by blending the high-order AF methods with the first-order local Lax-Friedrichs (LLF)  
 601 or Rusanov methods for both the cell average and point value updates, where the con-  
 602 vex limiting and scaling limiter were employed, respectively. For scalar conservation  
 603 laws, the blending coefficient was determined based on the global or local maximum

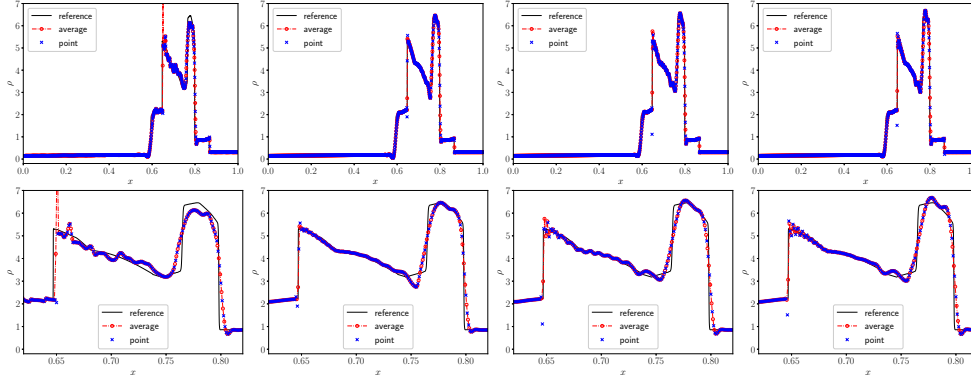


Fig. 11: [Example 5.7](#), blast wave interaction. The numerical solutions are computed with the BP limitings for the cell average and point value updates on a uniform mesh of 800 cells. The power law reconstruction is not used, and from left to right: the CFL number is 0.4, 0.4, 0.4, 0.35 for the JS, LLF, SW, and VH FVS, respectively. The corresponding enlarged views in  $[0.62, 0.82]$  are shown in the bottom row.

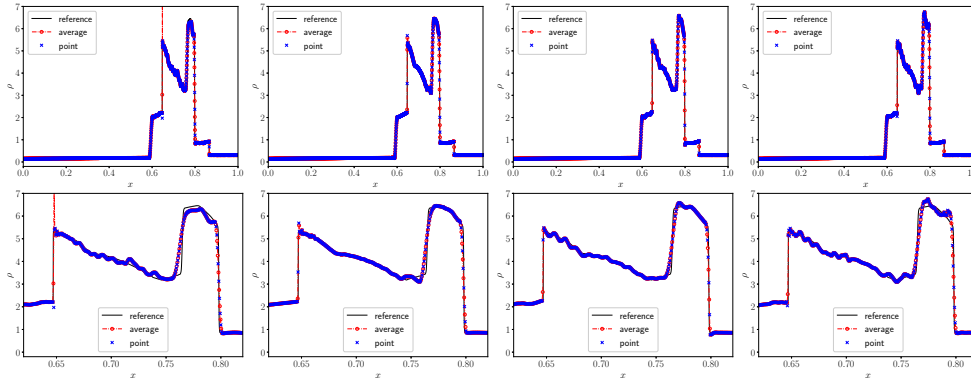


Fig. 12: [Example 5.7](#), blast wave interaction. The numerical solutions are computed with the power law reconstruction and the BP limitings for the cell average and point values update on a uniform mesh of 1600 cells. The CFL number is 0.1 for all the point value updates, and the corresponding enlarged views in  $[0.62, 0.82]$  are shown in the bottom row. From left to right: JS, LLF, SW, and VH FVS.

604 principle, while for the compressible Euler equations, it was obtained by enforcing  
 605 the positivity of density and pressure. Some challenging benchmark tests were con-  
 606 ducted based on different choices of the point value update, including the JS, LLF,  
 607 Steger-Warming, and Van Leer-Hänel FVS. The numerical results confirmed the ac-  
 608 curacy, BP property, and shock-capturing ability of our methods, and also showed  
 609 that the LLF FVS is generally superior to others in terms of the CFL number and  
 610 shock-capturing ability. Our future work will include, among others, extending the  
 611 current BP limitings to two-dimensional cases. We may also explore other ways to  
 612 further suppress oscillations for the Euler equations.

613

## REFERENCES

- 614 [1] R. ABGRALL, *A combination of residual distribution and the active flux formulations or a*  
615 *new class of schemes that can combine several writings of the same hyperbolic problem:*  
616 *Application to the 1D Euler equations*, Commun. Appl. Math. Comput., 5 (2023), pp. 370–  
617 402.
- 618 [2] R. ABGRALL AND W. BARSUKOW, *Extensions of active flux to arbitrary order of accuracy*,  
619 ESAIM: Math. Model. Numer. Anal., 57 (2023), pp. 991–1027.
- 620 [3] R. ABGRALL, J. LIN, AND Y. LIU, *Active flux for triangular meshes for compressible flows*  
621 *problems*, Dec. 2023, <https://arxiv.org/abs/2312.11271>.
- 622 [4] W. BARSUKOW, *The active flux scheme for nonlinear problems*, J. Sci. Comput., 86 (2021),  
623 p. 3.
- 624 [5] W. BARSUKOW AND J. P. BERBERICH, *A well-balanced active flux method for the shallow*  
625 *water equations with wetting and drying*, Commun. Appl. Math. Comput., (2023).
- 626 [6] W. BARSUKOW, J. P. BERBERICH, AND C. KLINGENBERG, *On the active flux scheme for*  
627 *hyperbolic PDEs with source terms*, SIAM J. Sci. Comput., 43 (2021), pp. A4015–A4042.
- 628 [7] W. BARSUKOW, J. HOHM, C. KLINGENBERG, AND P. L. ROE, *The active flux scheme on*  
629 *Cartesian grids and its low Mach number limit*, J. Sci. Comput., 81 (2019), pp. 594–622.
- 630 [8] E. CHUDZIK, C. HELZEL, AND D. KERKMANN, *The Cartesian grid active flux method: Linear*  
631 *stability and bound preserving limiting*, Appl. Math. Comput., 393 (2021), p. 125501.
- 632 [9] S. CLAIN, S. DIOT, AND R. LOUBÈRE, *A high-order finite volume method for systems of con-*  
633 *servations laws—Multi-dimensional Optimal Order Detection (MOOD)*, J. Comput. Phys.,  
634 230 (2011), pp. 4028–4050.
- 635 [10] C. J. COTTER AND D. KUZMIN, *Embedded discontinuous Galerkin transport schemes with*  
636 *localised limiters*, J. Comput. Phys., 311 (2016), pp. 363–373.
- 637 [11] C. M. DAFERMOS, *Hyperbolic Conservation Laws in Continuum Physics*, Springer Berlin Hei-  
638 delberg, 2000.
- 639 [12] T. EYMANN AND P. ROE, *Active flux schemes*, in 49th AIAA Aerospace Sciences Meeting  
640 including the New Horizons Forum and Aerospace Exposition, Orlando, Florida, Jan. 2011,  
641 American Institute of Aeronautics and Astronautics.
- 642 [13] T. EYMANN AND P. ROE, *Active flux schemes for systems*, in 20th AIAA Computational Fluid  
643 Dynamics Conference, Fluid Dynamics and Co-located Conferences, American Institute of  
644 Aeronautics and Astronautics, June 2011.
- 645 [14] T. A. EYMANN AND P. L. ROE, *Multidimensional active flux schemes*, in 21st AIAA Computa-  
646 tional Fluid Dynamics Conference, Fluid Dynamics and Co-located Conferences, American  
647 Institute of Aeronautics and Astronautics, June 2013.
- 648 [15] D. FAN AND P. L. ROE, *Investigations of a new scheme for wave propagation*, in 22nd AIAA  
649 Computational Fluid Dynamics Conference, American Institute of Aeronautics and Astro-  
650 nautics, 2015.
- 651 [16] S. GOTTLIEB, C.-W. SHU, AND E. TADMOR, *Strong Stability-Preserving High-Order Time*  
652 *Discretization Methods*, SIAM Rev., 43 (2001), pp. 89–112.
- 653 [17] J.-L. GUERMOND, M. NAZAROV, B. POPOV, AND I. TOMAS, *Second-order invariant domain*  
654 *preserving approximation of the Euler equations using convex limiting*, SIAM J. Sci. Com-  
655 put., 40 (2018), pp. A3211–A3239.
- 656 [18] J.-L. GUERMOND AND B. POPOV, *Fast estimation from above of the maximum wave speed in*  
657 *the Riemann problem for the Euler equations*, J. Comput. Phys., 321 (2016), pp. 908–926.
- 658 [19] J.-L. GUERMOND AND B. POPOV, *Invariant domains and first-order continuous finite element*  
659 *approximation for hyperbolic systems*, SIAM J. Numer. Anal., 54 (2016), pp. 2466–2489.
- 660 [20] J.-L. GUERMOND AND B. POPOV, *Invariant domains and second-order continuous finite ele-*  
661 *ment approximation for scalar conservation equations*, SIAM J. Numer. Anal., 55 (2017),  
662 pp. 3120–3146.
- 663 [21] J.-L. GUERMOND, B. POPOV, AND I. TOMAS, *Invariant domain preserving discretization-*  
664 *independent schemes and convex limiting for hyperbolic systems*, Comput. Methods Appl.  
665 Mech. Engrg., 347 (2019), pp. 143–175.
- 666 [22] H. HAJDUK, *Monolithic convex limiting in discontinuous Galerkin discretizations of hyperbolic*  
667 *conservation laws*, Comput. Math. Appl., 87 (2021), pp. 120–138.
- 668 [23] D. HÄNEL, R. SCHWANE, AND G. SEIDER, *On the accuracy of upwind schemes for the solution*  
669 *of the Navier-Stokes equations*, Fluid Dynamics and Co-located Conferences, American  
670 Institute of Aeronautics and Astronautics, June 1987.
- 671 [24] A. HARTEN, P. D. LAX, AND B. V. LEER, *On upstream differencing and Godunov-type*  
672 *schemes for hyperbolic conservation laws*, SIAM Rev., 25 (1983), pp. 35–61.
- 673 [25] C. HELZEL, D. KERKMANN, AND L. SCANDURRA, *A new ADER method inspired by the active*

- 674 *flux method*, J. Sci. Comput., 80 (2019), pp. 1463–1497.
- 675 [26] X. Y. HU, N. A. ADAMS, AND C.-W. SHU, *Positivity-preserving method for high-order con-*  
 676 *servative schemes solving compressible Euler equations*, J. Comput. Phys., 242 (2013),  
 677 pp. 169–180.
- 678 [27] G. S. JIANG AND C. W. SHU, *Efficient implementation of weighted ENO schemes*, J. Comput.  
 679 Phys., 126 (1996), pp. 202–228.
- 680 [28] J. R. KAMM AND F. X. TIMMES, *On efficient generation of numerically robust Sedov solutions*,  
 681 Tech. Report LA-UR-07-2849, 2007.
- 682 [29] D. KUZMIN, *Monolithic convex limiting for continuous finite element discretizations of hyper-*  
 683 *bolic conservation laws*, Comput. Methods Appl. Mech. Engrg., 361 (2020), p. 112804.
- 684 [30] D. KUZMIN, R. LÖHNER, AND S. TUREK, eds., *Flux-Corrected Transport: Principles, Algo-*  
 685 *ritms, and Applications*, Scientific Computation, Springer Netherlands, Dordrecht, 2012.
- 686 [31] X.-D. LIU AND S. OSHER, *Nonoscillatory high order accurate self-similar maximum principle*  
 687 *satisfying shock capturing schemes I*, SIAM J. Numer. Anal., 33 (1996), pp. 760–779.
- 688 [32] C. LOHMANN, D. KUZMIN, J. N. SHADID, AND S. MABUZA, *Flux-corrected transport algo-*  
 689 *ritms for continuous Galerkin methods based on high order Bernstein finite elements*, J.  
 690 Comput. Phys., 344 (2017), pp. 151–186.
- 691 [33] B. PERTHAME AND C.-W. SHU, *On positivity preserving finite volume schemes for Euler equa-*  
 692 *tions*, Numer. Math., 73 (1996), pp. 119–130.
- 693 [34] P. ROE, *Is discontinuous reconstruction really a good idea?*, J. Sci. Comput., 73 (2017),  
 694 pp. 1094–1114.
- 695 [35] L. I. SEDOV, *Similarity and Dimensional Methods in Mechanics*, Academic Press, New York,  
 696 1959.
- 697 [36] J. L. STEGER AND R. F. WARMING, *Flux vector splitting of the inviscid gasdynamic equations*  
 698 *with application to finite-difference methods*, J. Comput. Phys., 40 (1981), pp. 263–293.
- 699 [37] H. TANG, *On the sonic point glitch*, J. Comput. Phys., 202 (2005), pp. 507–532.
- 700 [38] E. F. TORO, *Riemann Solvers and Numerical Methods for Fluid Dynamics*, Springer Berlin  
 701 Heidelberg, 2009.
- 702 [39] B. VAN LEER, *Towards the ultimate conservative difference scheme. IV. A new approach to*  
 703 *numerical convection*, J. Comput. Phys., 23 (1977), pp. 276–299.
- 704 [40] B. VAN LEER, *Flux-vector splitting for the Euler equations*, in Eighth International Conference  
 705 on Numerical Methods in Fluid Dynamics, E. Krause, ed., Lecture Notes in Physics, Berlin,  
 706 Heidelberg, 1982, Springer, pp. 507–512.
- 707 [41] P. WOODWARD AND P. COLELLA, *The numerical simulation of two-dimensional fluid flow*  
 708 *with strong shocks*, J. Comput. Phys., 54 (1984), pp. 115–173.
- 709 [42] K. WU AND C.-W. SHU, *Geometric quasilinearization framework for analysis and design of*  
 710 *bound-preserving schemes*, SIAM Rev., 65 (2023), pp. 1031–1073.
- 711 [43] Z. XU, *Parametrized maximum principle preserving flux limiters for high order schemes solving*  
 712 *hyperbolic conservation laws: one-dimensional scalar problem*, Math. Comput., 83 (2014),  
 713 pp. 2213–2238.
- 714 [44] X. ZHANG AND C.-W. SHU, *On positivity-preserving high order discontinuous Galerkin*  
 715 *schemes for compressible Euler equations on rectangular meshes*, J. Comput. Phys., 229  
 716 (2010), pp. 8918–8934.
- 717 [45] X. ZHANG AND C.-W. SHU, *Maximum-principle-satisfying and positivity-preserving high-order*  
 718 *schemes for conservation laws: survey and new developments*, Proceedings of the Royal  
 719 Society A: Mathematical, Physical and Engineering Sciences, 467 (2011), pp. 2752–2776.
- 720 [46] X. ZHANG AND C.-W. SHU, *Positivity-preserving high order discontinuous Galerkin schemes*  
 721 *for compressible Euler equations with source terms*, J. Comput. Phys., 230 (2011),  
 722 pp. 1238–1248.