

1 **IMPLICIT POSITIVITY-PRESERVING HIGH ORDER**
2 **DISCONTINUOUS GALERKIN METHODS**
3 **FOR CONSERVATION LAWS***

4 TONG QIN[†] AND CHI-WANG SHU[‡]

5 **Abstract.** Positivity-preserving discontinuous Galerkin (DG) methods for solving hyperbolic
6 conservation laws have been extensively studied in the last several years. But nearly all the devel-
7 oped schemes are coupled with explicit time discretizations. Explicit discretizations suffer from the
8 constraint for the Courant-Friedrichs-Levis (CFL) number. This makes explicit methods impractical
9 for problems involving unstructured and extremely varying meshes or long-time simulations. Instead,
10 implicit DG schemes are often popular in practice, especially in the computational fluid dynamics
11 (CFD) community. In this paper we develop a high-order positivity-preserving DG method with
12 the backward Euler time discretization for conservation laws. We focus on one spatial dimension,
13 however the result easily generalizes to multidimensional tensor product meshes and polynomial
14 spaces. This work is based on a generalization of the positivity-preserving limiters in (X. Zhang and
15 C.-W. Shu, Journal of Computational Physics, 229 (2010), pp. 3091–3120) and (X. Zhang and C.-W.
16 Shu, Journal of Computational Physics, 229 (2010), pp. 8918–8934) to implicit time discretizations.
17 Both the analysis and numerical experiments indicate that a lower bound for the CFL number is
18 required to obtain the positivity-preserving property. The proposed scheme not only preserves the
19 positivity of the numerical approximation without compromising the designed high-order accuracy,
20 but also helps accelerate the convergence towards the steady-state solution and add robustness to
21 the nonlinear solver. Numerical experiments are provided to support these conclusions.

22 **Key words.** Positivity-preserving; Discontinuous Galerkin method; Backward Euler

23 **AMS subject classifications.** 65M60, 65M12

24 **1. Introduction.** In this paper, we consider the conservation law

$$\begin{aligned} u_t + f(u)_x &= 0, & (x, t) \in [0, 2\pi] \times [0, +\infty), \\ u(x, 0) &= u_0(x), & x \in [0, 2\pi], \end{aligned} \tag{1.1}$$

25 and its system version with appropriate boundary conditions. We focus on this one-
26 dimensional case, even though the result can be easily generalized to multidimensional
27 tensor product meshes and polynomial spaces.

28 For scalar conservation laws, it is well known that the entropy solution satisfies
29 the following maximum principle

$$\min_{x \in [0, 2\pi]} u_0(x) \leq u(x, t) \leq \max_{x \in [0, 2\pi]} u_0(x), \quad \forall t \geq 0.$$

30 In particular, if the initial condition is positive, then the entropy solution **must** satisfy
31 the following positivity-preserving property

$$u_0(x) \geq 0 \implies u(x, t) \geq 0, \quad \forall t \geq 0. \tag{1.2}$$

32 For systems, even though the entropy solution does not satisfy the maximum principle
33 in general, the physically relevant solution, for example the density and pressure in

*Submitted to the editors on April 3, 2017.

Funding: Research supported by ARO grant W911NF-15-1-0226 and NSF grants DMS-1418750 and DMS-1719410.

[†]Division of Applied Mathematics, Brown University, Providence, RI 02912, U.S.A. E-mail: tong_qin@brown.edu

[‡]Division of Applied Mathematics, Brown University, Providence, RI 02912, U.S.A. E-mail: shu@dam.brown.edu.

34 the compressible Euler system, is always positive. In this paper, the words “positive”
 35 and “positivity” are used actually to mean “nonnegative” and “nonnegativity”. We
 36 shall use “strictly positive” to mean the usual “positive”.

37 When designing numerical methods, we would like our numerical approximations
 38 to respect this positivity-preserving property (1.2), not only because it makes the num-
 39 erical approximation physically meaningful, but also it makes the numerical scheme
 40 more robust, since negative values sometimes cause ill-posedness of the problem and
 41 blow-ups of the numerical algorithm [14]. In recent years, the positivity-preserving
 42 DG schemes have been actively designed and applied for solving hyperbolic conser-
 43 vation laws [41, 42, 37, 38, 31, 7, 40]. All these methods are coupled with explicit
 44 temporal discretizations, such as strong stability preserving (SSP) Runge-Kutta (RK)
 45 methods [34, 13] and multi-step methods [33]. Explicit temporal discretizations en-
 46 joy many advantages, for example, the easiness in handling the nonlinear terms and
 47 boundary conditions, high-order accuracy with SSP properties [13], low storage re-
 48 quirement and so on. However, they suffer from the CFL constraint. For DG methods,
 49 to obtain the linear stability [1] or the maximum-principle stability [41], the CFL con-
 50 straint becomes more and more severe as we increase the polynomial degree in the
 51 approximation space. Such stringent time stepping restriction makes explicit methods
 52 impractical in computations involving unstructured and extremely varying meshes,
 53 viscous effect [29], low Mach numbers [3] or long-time simulations for steady-state
 54 calculation [15].

55 To circumvent the severe CFL constraint of explicit methods, implicit time dis-
 56 cretizations, which allow larger CFL numbers especially for stiff problems, are widely
 57 used in practice, especially in the CFD community to solve compressible flow prob-
 58 lems [15, 16, 6, 29, 28, 27, 25, 2] and see also the book [12]. Although most of the
 59 effort has been made for increasing accuracy of the time discretization and for in-
 60 creasing the efficiency of the nonlinear solver, only a few works exist in the literature
 61 concerning the positivity-preserving property of implicit methods. For compressible
 62 turbulent flow problems, Batten et al. [4] have proposed a positive finite difference
 63 scheme by splitting the fluxes into “implicit” and “correction” parts and the source
 64 term into positive and negative parts. The “implicit” and negative terms are treated
 65 implicitly via the Patankar trick [26]. In [22, 23, 24], Moryossef and Levy have de-
 66 veloped implicit unconditional positive finite volume schemes for unsteady turbulent
 67 flows. Their main idea to preserve the positivity is to make the Jacobian matrix in
 68 each implicit time step an M -matrix. All these methods mentioned are low-order
 69 accurate and are complicated to generalize to high order. For DG methods, in [21],
 70 Meister and Ortleb have constructed an unconditional implicit positive DG scheme
 71 for solving shallow water equations. The positivity of the numerical approximation is
 72 preserved via a modified Patankar trick [26]. The method is shown to be conservative
 73 and unconditional positivity-preserving, but only third-order accuracy is proved by
 74 a truncation error argument with no rigorous proof for arbitrary high-order [spatial](#)
 75 accuracy. In [39], Yuan, Cheng and Shu have developed a high-order unconditionally
 76 positive implicit DG method for radiative transfer equations. The positivity is pre-
 77 served by utilizing the particular boundary conditions of the problem and by designing
 78 a novel rotational limiter.

79 In this paper, we extend the general framework for constructing positivity-pres-
 80 erving schemes proposed by Zhang and Shu in [41, 42] to implicit temporal discretiza-
 81 tions and develop a positivity-preserving DG method with high-order [spatial](#) accuracy
 82 for one-dimensional conservation laws. The DG methods were first introduced by Reed

83 and Hill [32] for solving neutron transport equations and were further developed by
 84 Cockburn et al. in [10, 9, 8, 11] for solving the hyperbolic conservation laws. The
 85 DG method enjoys mathematically provable high-order accuracy and stability. The
 86 discontinuous feature of its approximation space makes it a good fit for parallel im-
 87 plementation and for handling unstructured meshes. Moreover, for a class of implicit
 88 temporal discretizations, it has been shown in [17] via the cell entropy inequality
 89 that the fully discrete scheme for the nonlinear conservation law is unconditionally
 90 L^2 -stable, which works for arbitrary triangulation and any **spatial** order of accuracy.
 91 We adopt the backward Euler temporal discretization in this paper. Our focus is
 92 on constructing a spatially high-order positivity-preserving DG scheme. The main
 93 conclusion is that in order to generalize the Zhang-Shu positivity-preserving limiter
 94 [41, 42] to the backward Euler DG scheme, a lower bound for the CFL number is
 95 required. This is proved theoretically for linear scalar equations and numerically veri-
 96 fied for nonlinear equations. The proposed positivity-preserving limiter is inexpensive
 97 and easy to implement. It not only preserves the positivity and high-order **spatial**
 98 accuracy but also makes the numerical scheme more robust, in the sense that it accel-
 99 erates the convergence towards the steady-state solution and adds robustness to the
 100 nonlinear solver for extreme test cases.

101 The organization of the paper is as follows. In Section 2 we describe the DG
 102 scheme. Then in Section 3 the positivity-preserving technique is introduced for scalar
 103 equations. In particular, a CFL condition is derived for linear equations to ensure the
 104 positiveness of the scheme. The positivity-preserving DG scheme for the compressible
 105 Euler system follows in Section 4. Numerical experiments are presented in Section 5
 106 and concluding remarks are given in Section 6.

107 **2. Implicit DG scheme.**

2.1. The DG discretization. In this section, we define the DG scheme for
 (1.1). First, let us fix some notations. We decompose the domain $\Omega = [0, 2\pi]$ into
 N subintervals, $I_j = [x_{j-\frac{1}{2}}, x_{j+\frac{1}{2}}]$, for $j = 1, 2, \dots, N$. The size of each subinterval
 is denoted by h_j . Define $\hat{I} = [-1, 1]$ to be the reference cell and define $T_j(x) =$
 $2(x - x_j)/h_j$ to be the affine mapping between the intervals I_j and \hat{I} , where $x_j =$
 $(x_{j-\frac{1}{2}} + x_{j+\frac{1}{2}})/2$ is the midpoint of I_j . Moreover, let $(\cdot, \cdot)_j$ denote the usual L^2 inner
 product on I_j and $(\cdot, \cdot)_{\hat{I}}$ the one on \hat{I} . Then we define the approximation space

$$V_h = \{v \in L^2(\Omega) : v|_{I_j} \in P_k(I_j), \forall j = 1, \dots, N\}$$

108 where $P_k(I_j)$ denotes the polynomial space on I_j with degree up to k .

109 The semi-discrete DG scheme is to seek the approximation $u_h(t) \in V_h$, such that
 110 in each subinterval I_j ,

$$\frac{d}{dt}(u_h(t), v)_j - (f(u_h(t)), v_x)_j + \hat{f}_{j+\frac{1}{2}}(u_h(t))v(x_{j+\frac{1}{2}}^-) - \hat{f}_{j-\frac{1}{2}}(u_h(t))v(x_{j-\frac{1}{2}}^+) = 0 \quad (2.1)$$

111 holds for any $v \in V_h$, where $v(x_{j+\frac{1}{2}}^+)$ and $v(x_{j+\frac{1}{2}}^-)$ denote the right and the left limits
 112 of the function v at $x_{j+\frac{1}{2}}$. The single valued function $\hat{f}_{j+\frac{1}{2}}(u) = \hat{f}(u(x_{j+\frac{1}{2}}^-), u(x_{j+\frac{1}{2}}^+))$
 113 is the numerical flux, which depends on both the left and right limits of u at $x_{j+\frac{1}{2}}$.
 114 In this paper, we consider the global Lax-Friedrichs flux

$$\hat{f}(a, b) = \frac{1}{2}[f(a) + f(b) - \alpha(b - a)], \quad (2.2)$$

115 where $\alpha = \max_{x \in \Omega} |f'(u_0(x))|$.

116 **2.2. Time discretization.** With shorthand notation, the semidiscrete scheme
 117 (2.1) can be rewritten as below

$$\frac{d}{dt}(u_h(t), v)_j = L_j(u_h(t), v), \quad \forall v \in V_h, \quad \forall j = 1, \dots, N \quad (2.3)$$

118 where $L_j(u_h(t), v) = (f(u_h(t)), v_x)_j - \left[\hat{f}_{j+\frac{1}{2}}(u_h(t))v(x_{j+\frac{1}{2}}^-) - \hat{f}_{j-\frac{1}{2}}(u_h(t))v(x_{j-\frac{1}{2}}^+) \right]$.
 119 We use the backward Euler method to further discretize this ODE system. Then the
 120 fully discrete scheme is defined by seeking the approximation at time t^{n+1} , which is
 121 denoted by $u_h^{n+1} \in V_h$, such that in each cell I_j , we have

$$(u_h^{n+1}, v)_j - \Delta t L_j(u_h^{n+1}, v) = (u_h^n, v)_j \quad (2.4)$$

122 for all $v \in V_h$. In the following, we use u_j^n to denote $u_h^n|_{I_j}$, $(u_{j+\frac{1}{2}}^n)^\pm$ to denote $u_j^n(x_{j+\frac{1}{2}}^\pm)$
 123 and use \bar{u}_j^n to denote the cell average of u_j^n in the interval I_j .

124 To further solve the nonlinear system (2.4), there have been many works on how
 125 to build efficient solvers, such as the work in [29, 28, 27]. But since our main focus is
 126 on the positivity preserving property rather than the efficiency of the nonlinear solver,
 127 we use the Newton method [12] for the nonlinear system up to accuracy 10^{-13} . For
 128 the robustness and accuracy reasons, in each Newton iteration, the Jacobian matrix
 129 is solved with the direct solver.

130 **3. Positivity-preserving DG scheme for scalar equations.** In this section,
 131 we introduce how to add the positivity-preserving property to the scheme (2.4). First,
 132 let us give the definition of the positivity-preserving DG scheme for the scalar equation
 133 as that in [41].

134 **DEFINITION 3.1.** *A DG scheme is defined to be positivity preserving if given $u_h^n(x) \geq$
 135 0 , for any $x \in \Omega$, then we have $u_h^{n+1}(x) \geq 0$, $\forall x \in \Omega$.*

136 This definition will be slightly modified later (requiring positivity on specified quadra-
 137 ture points rather than on all points) in order to obtain a more efficient implementa-
 138 tion. Generally, the original high-order DG method is not positivity-preserving. We
 139 follow the general approach in [41] and construct high-order positivity-preserving DG
 140 methods in the following two steps.

141 **Step 1** First, given $u_j^n(x) \geq 0$ for any $x \in I_j$ and any j , find a sufficient condition
 142 such that we have the cell average \bar{u}_j^{n+1} positive for all j .

143 **Step 2** Next, we make the whole polynomial $u_j^{n+1}(x) \geq 0$ by invoking the scaling
 144 limiter in [20, 41].

145 The main difficulty lies in the first step. The implicit DG approximation u_j^{n+1}
 146 depends on the approximation at the previous time step u_h^n in a global and implicit
 147 way. Effort is needed to represent the cell average \bar{u}_j^{n+1} in terms of u_h^n . In this section,
 148 we would first show how to overcome this difficulty for scalar linear equations and
 149 then we derive a CFL condition, under which, the step 1 is fulfilled. Then we will
 150 introduce the scaling limiter and summarize the algorithm.

151 **3.1. Preliminaries.** Let us first recall some definitions and results that will be
 152 useful in the following analysis. The first useful tool is the so-called M -matrix. For a
 153 thorough introduction, one can refer to [5]. To define it let us first set

$$\mathcal{Z}^{n \times n} = \{A = (a_{ij}) \in \mathbb{R}^{n \times n} : a_{ij} \leq 0, i \neq j\}$$

154 which is the set of all the $n \times n$ real matrices with nonpositive off-diagonal entries.
 155 In [30], the author listed forty equivalent characterizations for M -matrices. For our
 156 purpose, we adopt the following one as the definition.

157 **DEFINITION 3.2.** *A matrix $A \in \mathcal{Z}^{n \times n}$ is called an M -matrix if A is inverse-positive,
 158 that is, A^{-1} exists and each entry of A^{-1} is nonnegative.*

159 M -matrices have the following equivalent characterization [30].

160 **THEOREM 3.3.** *A matrix $A = (a_{ij}) \in \mathcal{Z}^{n \times n}$ is an M -matrix if and only if $a_{ii} > 0$,
 161 $1 \leq i \leq n$, and there exists a positive diagonal matrix $D = \text{diag}\{d_1, \dots, d_n\}$ such that
 162 AD is strictly diagonally dominant, that is, $a_{ii}d_i > \sum_{j \neq i} |a_{ij}|d_j$ for $1 \leq i \leq n$.*

163 In particular, if D is the identity matrix, we have the following corollary.

164 **COROLLARY 3.4.** *A matrix $A = (a_{ij}) \in \mathcal{Z}^{n \times n}$ is an M -matrix if $a_{ii} > 0$ and it is
 165 strictly diagonally dominant.*

166 In the following, we also utilize properties of Legendre polynomials. We consider
 167 the standard Legendre polynomials $\{p_n(x)\}$ defined on \hat{I} by the following recursive
 168 relationship

$$(n+1)p_{n+1}(x) = (2n+1)xp_n(x) - np_{n-1}(x), \quad p_0(x) = 1, \quad p_1(x) = x, \quad x \in \hat{I}. \quad (3.1)$$

169 In the following lemma, we collect some properties of Legendre polynomials that will
 170 be useful in the following analysis. For the proof, one can refer to [35, Sections 3.2,
 171 4.1, 4.3, 4.7, 7.2].

172 **LEMMA 3.5.** *Legendre polynomials defined in (3.1) have the following properties*

- 173 (i) $p_n(1) = 1$, $p_n(-x) = (-1)^n p_n(x)$, $\forall x \in \hat{I}$ and $|p_n(x)| < 1$, $\forall x \in (-1, 1)$.
- 174 (ii) $\int_{\hat{I}} p_n(x) p_m(x) dx = \frac{2}{2n+1} \delta_{nm}$, where δ_{nm} is the Kronecker delta.
- 175 (iii) $(2n+1)p_n(x) = \frac{d}{dx}[p_{n+1}(x) - p_{n-1}(x)]$.
- 176 (iv) Rodrigues' formula $p_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} [(x^2 - 1)^n]$.
- 177 (v) *Christoffel-Darboux formula*

$$\sum_{n=0}^k \alpha_n p_n(x) p_n(y) = \frac{\alpha_k(k+1)}{2k+1} \frac{p_{k+1}(x)p_k(y) - p_{k+1}(y)p_k(x)}{x-y} \quad (3.2)$$

178 where $\alpha_n > 0$.

179 **3.2. CFL condition for linear equations.** We consider the linear equation

$$\begin{aligned} u_t + u_x &= 0, \quad x \in \Omega, \\ u(x, 0) &= u_0(x), \end{aligned} \quad (3.3)$$

180 with periodic boundary condition. Then the scheme (2.4) becomes

$$(u_h^{n+1}, v)_j - \Delta t (u_h^{n+1}, v_x)_j + \Delta t [(u_{j+\frac{1}{2}}^{n+1})^- v_{j+\frac{1}{2}}^- - (u_{j-\frac{1}{2}}^{n+1})^- v_{j-\frac{1}{2}}^+] = (u_h^n, v)_j \quad (3.4)$$

181 for all $v \in V_h$. Given $u_h^n(x) \geq 0, \forall x \in \Omega$, we want to derive a CFL condition, under
 182 which the cell average $\bar{u}_j^{n+1} \geq 0, \forall j$. In order to do that, we first express \bar{u}_j^{n+1} in
 183 terms of u_h^n . The idea is to take $k+1$ different test functions as probes to extract the
 184 information out from $u_h^{n+1}(x)$ in terms of $u_h^n(x)$.

185 First, let us take $v = 1$ in the scheme (3.4), we have

$$\bar{u}_j^{n+1} + \lambda_j [(u_{j+\frac{1}{2}}^{n+1})^- - (u_{j-\frac{1}{2}}^{n+1})^-] = \bar{u}_j^n, \quad j = 1, \dots, N,$$

186 where $\lambda_j = \frac{\Delta t}{h_j}$. We can rewrite the system above in the matrix form as below

$$\Lambda^{-1} \bar{\mathbf{u}}^{n+1} + A(\mathbf{u}^{n+1})^- = \Lambda^{-1} \bar{\mathbf{u}}^n \quad (3.5)$$

187 where $\bar{\mathbf{u}}^n = (\bar{u}_1^n, \dots, \bar{u}_N^n)^T$ and $(\mathbf{u}^{n+1})^- = ((u_{\frac{3}{2}}^{n+1})^-, \dots, (u_{N+\frac{1}{2}}^{n+1})^-)$. The $N \times N$
188 matrices Λ and A take the following form

$$\Lambda = \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \lambda_N \end{pmatrix}, \quad A = \begin{pmatrix} 1 & 0 & \cdots & -1 \\ -1 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & -1 & 1 \end{pmatrix}.$$

189 Next, we need to express the cell boundary value $(u_{j+\frac{1}{2}}^{n+1})^-$ in terms of \bar{u}_j^{n+1} and u_h^n .
190 To this end, we need to take other special test functions. Recall that the Dirac delta
191 distribution can be approximated by the following series in the distribution sense [19]
192

$$\delta(x - y) = \frac{1}{2} \sum_{l=0}^{\infty} (2l+1) p_l(x) p_l(y), \quad x, y \in \hat{I}, \quad (3.6)$$

193 where $p_l(x)$ is the standard Legendre polynomial defined in (3.1). Then we set $y = 1$,
194 truncate the series (3.6) at the $(k+1)$ th term and define

$$\hat{\delta}^k(x) = \frac{1}{2} \sum_{l=0}^k (2l+1) p_l(x) = \frac{k+1}{2} \frac{p_{k+1}(x) - p_k(x)}{x-1}, \quad x \in \hat{I}. \quad (3.7)$$

195 We have employed (3.2) with $\alpha_n = (2n+1)/2$ in the last equality.

196 The following lemma says that this polynomial is an analogue to the Dirac delta
197 distribution in $P^k(\hat{I})$ at the point $y = 1$.

198 LEMMA 3.6. *The polynomial $\hat{\delta}^k$ has the following properties*

- 199 (i) $\hat{\delta}^k(x) \in P^k(\hat{I})$ and for any $w \in P^k(\hat{I})$, we have $(w, \hat{\delta}^k)_{\hat{I}} = w(1)$.
200 (ii) In the cell I_j , define

$$\delta_j^k(x) = \frac{2}{h_j} \hat{\delta}^k(T_j(x)), \quad x \in I_j, \quad (3.8)$$

201 then for any $w \in P^k(I_j)$ we have $(w, \delta_j^k)_j = w(x_{j+\frac{1}{2}})$.

- 202 (iii) The mass is concentrated at $x = 1$, in the sense that for $k \geq 1$ and $j = 0, \dots, k-1$,
203 we have $(\hat{\delta}^k)^{(j)}(1) - (\hat{\delta}^k)^{(j)}(x) > 0$ for any $x \in [-1, 1)$.

204 *Proof.* It is obvious that $\hat{\delta}^k \in P_k(\hat{I})$. For any polynomial $w \in P_k(\hat{I})$, we can write
205 it as a linear combination of Legendre polynomials as below

$$w(x) = \sum_{l=0}^k c_l p_l(x), \quad x \in \hat{I}.$$

206 Then by the definition of $\hat{\delta}^k$ and Lemma 3.5, we have

$$(w, \hat{\delta}^k)_{\hat{I}} = \sum_{l=0}^k c_l (p_l, \hat{\delta}^k)_{\hat{I}} = \sum_{l=0}^k c_l = \sum_{l=0}^k c_l p_l(1) = w(1).$$

207 The property (ii) can be shown by a simple change of variable.

208 For property (iii), by (3.7), we have

$$\begin{aligned} (\hat{\delta}^k)^{(j)}(1) - (\hat{\delta}^k)^{(j)}(x) &= \sum_{l=0}^k \frac{2l+1}{2} [p_l^{(j)}(1) - p_l^{(j)}(x)] \\ &= \sum_{l=j+1}^k \frac{2l+1}{2} [p_l^{(j)}(1) - p_l^{(j)}(x)]. \end{aligned} \quad (3.9)$$

209 When $x = -1$. By Lemma 3.5 (i), we have

$$(\hat{\delta}^k)^{(j)}(1) - (\hat{\delta}^k)^{(j)}(-1) = \sum_{l=j+1}^k \frac{2l+1}{2} [1 - (-1)^{l+j}] p_l^{(j)}(1). \quad (3.10)$$

210 We claim that $p_l^{(j)}(1) > 0$ for any $l = 0, \dots, k$ and any $j = 0, \dots, l$. By Lemma 3.5
 211 (i), this holds for $j = 0$. For $j \geq 1$, it can be checked that the claim holds for $l = 0, 1$.
 212 And for $l \geq 2$, we can use Lemma 3.5 (iii) and show the claim by induction. Then we
 213 can conclude that the summation (3.10) is positive, since for fixed $j = 0, \dots, k-1$,
 214 there has to be at least one $l = j+1, \dots, k$ such that $[1 - (-1)^{l+j}] = 2$.

215 Next, when $x \in (-1, 1)$, by (3.9), it suffices to show

$$p_l^{(j)}(1) - p_l^{(j)}(x) > 0 \quad (3.11)$$

216 for $j = 0, \dots, k-1$ and $l = j+1, \dots, k$ or equivalently, for $l = 1, \dots, k$ and $j =$
 217 $0, \dots, l-1$. First, by Lemma 3.5(i), (3.11) holds for any l with $j = 0$. Therefore, we
 218 only need to consider $l = 2, \dots, k$ and $j = 1, \dots, l-1$. It is straightforward to verify
 219 that (3.11) holds for $l = 2$. If (3.11) holds for $l \leq m-1$ with $m \geq 3$, then we have
 220 when $l = m$, for fixed $j = 1, \dots, l-1$ and any $x \in (-1, 1)$, by Lemma 3.5 (iii)

$$p_m^{(j)}(1) = (2m-1)p_{m-1}^{(j-1)}(1) + p_{m-2}^{(j)}(1) > (2m-1)p_{m-1}^{(j-1)}(x) + p_{m-2}^{(j)}(x) = p_m^{(j)}(x).$$

221 Then by induction, we have proved (3.11). \square

222 With the help of the **delta approximation** (3.8), we have the following represen-
 223 tation of $(u_{j+\frac{1}{2}}^{n+1})^-$.

224 **LEMMA 3.7.** *For linear scalar conservation law with $f(u) = u$ discretized by the*
 225 *scheme (2.4), we have*

$$\sigma_j^k (u_{j+\frac{1}{2}}^{n+1})^- = \xi_j^k \bar{u}_j^{n+1} + (\hat{u}_j^n, g_j^k)_{\hat{I}} \quad (3.12)$$

226 where

$$\sigma_j^k = 1 + \sum_{i=0}^{k-1} (2\lambda_j)^{i+1} (\alpha_i^k - \beta_i^k), \quad \xi_j^k = 2 \sum_{i=0}^k (2\lambda_j)^i \beta_i^k,$$

$$g_j^k(x) = \sum_{i=0}^{k-1} (2\lambda_j)^i \left[(\hat{\delta}^k)^{(i)}(x) - \beta_i^k \right]$$

227 and

$$\alpha_i^k = (\hat{\delta}^k)^{(i)}(1), \quad \beta_i^k = (\hat{\delta}^k)^{(i)}(-1). \quad (3.13)$$

228 *Proof.* For fixed $l = 0, 1, \dots, k-1$, take the test function to be $(\delta_j^k)^{(l)}(x) -$
 229 $(\delta_j^k)^{(l)}(x_{j-\frac{1}{2}})$ in the scheme (3.4). By the definition of $\delta_j^k(x)$ and Lemma 3.6, we have

$$\begin{aligned} & (u_j^{n+1}, (\hat{\delta}^k)^{(l)}(T_j(x)))_j - h_j \beta_l^k \bar{u}_j^{n+1} - 2\lambda_j (u_j^{n+1}, (\hat{\delta}^k)^{(l+1)}(T_j(x)))_j \\ & + \Delta t (u_{j+\frac{1}{2}}^{n+1})^- (\alpha_l^k - \beta_l^k) = (u_j^n, (\hat{\delta}^k)^{(l)}(T_j(x)))_j - h_j \beta_l^k \bar{u}_j^n. \end{aligned}$$

230 If we expand u_j^{n+1} in terms of the basis $\phi_j^l(x) = p_l(T_j(x))$, $l = 0, \dots, k$, as $u_j^{n+1} =$
 231 $\sum_{l=0}^k (c_j^l)^{n+1} \phi_j^l$ and by a change of variable, we obtain

$$\begin{aligned} & (\hat{u}_j^{n+1}, (\hat{\delta}^k)^{(l)})_{\hat{I}} - 2\beta_l^k \bar{u}_j^{n+1} - 2\lambda_j (\hat{u}_j^{n+1}, (\hat{\delta}^k)^{(l+1)})_{\hat{I}} + 2\lambda_j (u_{j+\frac{1}{2}}^{n+1})^- (\alpha_l^k - \beta_l^k) = \\ & (\hat{u}_j^n, (\hat{\delta}^k)^{(l)})_{\hat{I}} - 2\beta_l^k \bar{u}_j^n \end{aligned}$$

232 where $\hat{u}_j^n = \sum_{l=0}^k (c_j^l)^n p_l(x)$. Or equivalently,

$$\begin{aligned} & (\hat{u}_j^{n+1}, (\hat{\delta}^k)^{(l)})_{\hat{I}} - 2\lambda_j (\hat{u}_j^{n+1}, (\hat{\delta}^k)^{(l+1)})_{\hat{I}} = \\ & 2\beta_l^k \bar{u}_j^{n+1} - 2\lambda_j (u_{j+\frac{1}{2}}^{n+1})^- (\alpha_l^k - \beta_l^k) + (\hat{u}_j^n, (\hat{\delta}^k)^{(l)})_{\hat{I}} - 2\beta_l^k \bar{u}_j^n. \end{aligned}$$

233 If we set

$$D_l = (\hat{u}_j^{n+1}, (\hat{\delta}^k)^{(l)})_{\hat{I}}, \quad C_l = 2\beta_l^k \bar{u}_j^{n+1} - 2\lambda_j (u_{j+\frac{1}{2}}^{n+1})^- (\alpha_l^k - \beta_l^k) + (\hat{u}_j^n, (\hat{\delta}^k)^{(l)})_{\hat{I}} - 2\beta_l^k \bar{u}_j^n,$$

234 then we have

$$D_l - 2\lambda_j D_{l+1} = C_l, \quad l = 0, \dots, k-1 \quad (3.14)$$

235 and in particular, when $l = k-1$, we have

$$D_{k-1} = 2\lambda_j D_k + C_{k-1} = 2\lambda_j (\hat{u}_j^{n+1}, (\hat{\delta}^k)^{(k)})_{\hat{I}} + C_{k-1} = 4\lambda_j \beta_k^k \bar{u}_j^{n+1} + C_{k-1}. \quad (3.15)$$

236 By Lemma 3.6 and by using (3.14), we have the following representation of $(u_{j+\frac{1}{2}}^{n+1})^-$

$$(u_{j+\frac{1}{2}}^{n+1})^- = (u_j^{n+1}, \delta_j^k)_j = (\hat{u}_j^{n+1}, \hat{\delta}^k)_{\hat{I}} = D_0 = 2\lambda_j D_1 + C_0.$$

237 If we continue using (3.14) for another $k-2$ times and by using (3.15), we arrive at

$$\begin{aligned} (u_{j+\frac{1}{2}}^{n+1})^- &= (2\lambda_j)^{k-1} D_{k-1} + \sum_{i=0}^{k-2} (2\lambda_j)^i C_i \\ &= (2\lambda_j)^{k-1} [4\lambda_j \beta_k^k \bar{u}_j^{n+1} + C_{k-1}] + \sum_{i=0}^{k-2} (2\lambda_j)^i C_i \end{aligned}$$

$$= 2(2\lambda_j)^k \beta_k^k \bar{u}_j^{n+1} + \sum_{i=0}^{k-1} (2\lambda_j)^i C_i.$$

238 Then after plugging in the definition of C_i and some manipulations, we obtain

$$\left[1 + \sum_{i=0}^{k-1} (2\lambda_j)^{i+1} (\alpha_i^k - \beta_i^k) \right] (u_{j+\frac{1}{2}}^{n+1})^- = \left[2 \sum_{i=0}^k \beta_i^k (2\lambda_j)^i \right] \bar{u}_j^{n+1} + \sum_{i=0}^{k-1} (2\lambda_j)^i (\hat{u}_j^n, (\hat{\delta}^k)^{(i)} - \beta_i^k)_{\hat{I}}.$$

239 \square

240 For the parameters $\{\beta_j^k, \alpha_j^k\}_{j=1}^N$, we have the following lemma, which will be useful
241 in the proof of Proposition 3.11.

242 LEMMA 3.8. For any $k \geq 0$, the following results hold

243 (i) $\alpha_i^k > \beta_i^k$ for $i = 0, \dots, k-1$ and hence $\sigma_j^k > 0$ for any j .

244 (ii) $\beta_{k-2i}^k > 0$, $\beta_{k-2i-1}^k < 0$, for $i = 0, \dots, \lfloor k/2 \rfloor$.

245 (iii) $\beta_{k-2i}^k + \beta_{k-2i-1}^k > 0$, for $i = 0, \dots, \lfloor k/2 \rfloor$.

246 *Proof.* The statement (i) is a direct conclusion of Lemma 3.6 (iii).

247 For the statement (ii), let us first derive an explicit formula for β_i^k . By the
248 Rodrigues' formula in Lemma 3.5 (iv) we have

$$p_l(x) = \frac{1}{2^l l!} \frac{d^l}{dx^l} (x^2 - 1)^l = \frac{1}{2^l l!} \frac{d^l}{dx^l} [(x-1)^l (x+1)^l].$$

249 Then by the Leibnitz's rule, we obtain for $i \leq l$,

$$\begin{aligned} p_l^{(i)}(-1) &= \frac{1}{2^l l!} \binom{l+i}{l} l! [(x-1)^l]^{(i)} \Big|_{x=-1} = \frac{1}{2^l} \frac{(i+l)!}{i! l!} \frac{l!}{(l-i)!} (-2)^{l-i} \\ &= \frac{(-2)^{-i} (i+l)!}{i! (l-i)!} (-1)^l \end{aligned}$$

250 and hence we have

$$\beta_i^k = \frac{1}{2} \sum_{l=i}^k (2l+1) \frac{(-2)^{-i} (l+i)!}{i! (l-i)!} (-1)^l = \frac{1}{2^{i+1} i!} \sum_{l=i}^k (2l+1) \frac{(l+i)!}{(l-i)!} (-1)^{l-i} = \frac{1}{C_i} \sum_{l=i}^k \gamma_i^l$$

251 where $C_i = 2^{i+1} i!$ and $\gamma_i^l = (2l+1) \frac{(l+i)!}{(l-i)!} (-1)^{l-i}$.

252 For γ_i^l , we have

$$|\gamma_i^{l+1}| = (2l+3) \frac{(l+i+1)!}{(l+1-i)!} = \frac{2l+3}{2l+1} \frac{l+i+1}{l-i+1} |\gamma_i^l| > |\gamma_i^l|.$$

253 Next for $j = 0, \dots, \lfloor k/2 \rfloor$, let us consider β_{k-2j}^k . If we replace i with $k-2j$, we obtain

254

$$C_{k-2j} \beta_{k-2j}^k = (|\gamma_{k-2j}^k| - |\gamma_{k-2j-1}^k|) + \dots + (|\gamma_{k-2j}^{k-2j+2}| - |\gamma_{k-2j}^{k-2j+1}|) + |\gamma_{k-2j}^{k-2j}| > 0. \quad (3.16)$$

255 For β_{k-2j-1}^k , we have

$$C_{k-2j-1} \beta_{k-2j-1}^k = -[(|\gamma_{k-2j-1}^k| - |\gamma_{k-2j-1}^{k-1}|) + \dots + (|\gamma_{k-2j-1}^{k-2j}| - |\gamma_{k-2j-1}^{k-2j-1}|)] < 0. \quad (3.17)$$

256 Therefore, we can conclude that $\beta_{k-2j}^k > 0$ and $\beta_{k-2j-1}^k < 0$.

257 For the last statement, let us consider the sum $\beta_{k-2j}^k + \beta_{k-2j-1}^k$. To this end, first
 258 for general γ_j^l , let us consider the following expression

$$\tau_j^l := \frac{1}{2j} (|\gamma_j^l| - |\gamma_j^{l-1}|) - |\gamma_{j-1}^l| + |\gamma_{j-1}^{l-1}|.$$

259 If we plug the definition of γ_j^l in and after direct calculation, we obtain

$$\tau_j^l = \frac{l(l+j-2)!}{j(l-j+1)!} [(l^2 - j^2)(2j+1) + 2j - 1] \geq 0.$$

260 If we combine (3.16) and (3.17) together we would obtain,

$$C_{k-2j-1}(\beta_{k-2j}^k + \beta_{k-2j-1}^k) = \sum_{i=0}^{j+1} \tau_{k-2j}^{k-2i} + |\gamma_{k-2j}^{k-2j}| > 0,$$

261 which implies the desired conclusion. \square

262 With Lemma 3.7 and the equation (3.5), we can obtain the following cell average
 263 equation

$$T\bar{\mathbf{u}}^{n+1} = \mathcal{L}(\mathbf{u}^n) \quad (3.18)$$

264 where

$$T = \begin{pmatrix} \frac{\xi_1^k}{\sigma_1^k} + \frac{1}{\lambda_1} & 0 & \cdots & -\frac{\xi_N^k}{\sigma_N^k} \\ -\frac{\xi_1^k}{\sigma_1^k} & \frac{\xi_2^k}{\sigma_2^k} + \frac{1}{\lambda_2} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & -\frac{\xi_{N-1}^k}{\sigma_{N-1}^k} & \frac{\xi_N^k}{\sigma_N^k} + \frac{1}{\lambda_N} \end{pmatrix}$$

265 and

$$(\mathcal{L}(\mathbf{u}^n))_j = \left(\hat{u}_j^n, \frac{1}{2\lambda_j} - \frac{g_j^k}{\sigma_j^k} \right)_{\hat{I}} + \left(\hat{u}_{j-1}^n, \frac{g_{j-1}^k}{\sigma_{j-1}^k} \right)_{\hat{I}}, \quad j = 1, \dots, N.$$

266 A set of sufficient conditions to make the cell average \bar{u}_j^{n+1} positive for any $j =$
 267 $1, \dots, N$ are the following

268 **Condition I** T is an M -matrix, or by Corollary 3.4 and Lemma 3.8,

$$\xi_j^k \geq 0, \quad \forall j = 1, \dots, N. \quad (3.19)$$

269 **Condition II** $\mathcal{L}(\mathbf{u}^n)$ is positive, or

$$(\hat{u}_j^n, 1/(2\lambda_j) - g_j^k/\sigma_j^k)_{\hat{I}} + (\hat{u}_{j-1}^n, g_{j-1}^k/\sigma_{j-1}^k)_{\hat{I}} \geq 0, \quad j = 1, \dots, N. \quad (3.20)$$

270 By Lemma 3.6 (iii), we obtain $g_j^k(x) < \frac{\sigma_j^k}{2\lambda_j}, \forall x \in \hat{I}$, for any k and j and hence the first
 271 term in (3.20) is always positive. Since \hat{u}_j^n and \hat{u}_{j-1}^n can be any independent positive
 272 polynomials, the condition (3.20) is further reduced to

$$(v, g_j^k)_{\hat{I}} \geq 0, \quad \forall v \in P^k(\hat{I}) \text{ and } v \geq 0. \quad (3.21)$$

273 In summary, if we set $F_k(\lambda, x) = \sum_{i=0}^k (2\lambda)^i (\hat{\delta}^k)^{(i)}(x)$ and then $\xi_j^k = 2F_k(\lambda_j, -1)$,
 274 $g_j^k = F_k(\lambda_j, x) - F_k(\lambda_j, -1)$, sufficient conditions (3.19) and (3.20) actually require
 275 that the CFL numbers λ_j satisfy

$$F_k(\lambda_j, -1) \geq 0, \quad (3.22)$$

$$(F_k(\lambda_j, \cdot) - F_k(\lambda_j, -1), v(\cdot))_{\hat{I}} \geq 0, \quad \forall v \in P^k(\hat{I}) \text{ and } v \geq 0. \quad (3.23)$$

276 The following theorem says that in order to make these two conditions hold simulta-
 277 neously, the CFL numbers λ_j can not be arbitrarily small.

278 **THEOREM 3.9.** *When λ_j is small, conditions (3.22) and (3.23) can not hold at the*
 279 *same time. More specifically, we have the following two cases:*

- 280 1. *When the polynomial degree k is odd, there exists $\eta_1^k > 0$ such that when*
 281 *$\lambda_j < \eta_1^k$, the condition (3.22) does not hold.*
- 282 2. *When $k \geq 2$ is even, there exists $\eta_2^k > 0$ such that when $\lambda_j < \eta_2^k$, the second*
 283 *condition (3.23) does not hold.*

284 *Proof.* When k is odd, $F_k(\lambda_j, -1) = \sum_{i=0}^k \beta_i^k (2\lambda_j)^i$ is a polynomial of odd degree.
 285 By Lemma 3.8, we have the leading coefficient $\beta_k^k > 0$ and hence $F_k(\lambda_j, -1) > 0$ when
 286 λ_j is large. On the other hand, when $\lambda_j = 0$, $F_k(0, -1) = \beta_0^k < 0$, again by Lemma 3.8.
 287 Therefore, the polynomial $F_k(\lambda_j, -1)$ must have at least one and at most k positive
 288 roots. If we take η_1^k to be the smallest one, we would have the first statement.

289 When $k \geq 2$ is even, in (3.23) take $v = 1$ and we obtain

$$(F_k(\lambda_j, x) - F_k(\lambda_j, -1), 1)_{\hat{I}} = \sum_{i=0}^{k-1} (2\lambda_j)^i \left[(\hat{\delta}^k)^{(i-1)}(1) - (\hat{\delta}^k)^{(i-1)}(-1) - 2(\hat{\delta}^k)^{(i)}(-1) \right]$$

290 where $(\hat{\delta}^k)^{(-1)}(1) - (\hat{\delta}^k)^{(-1)}(-1) := \int_{\hat{I}} \hat{\delta}^k(x) dx$. To simplify the notation, let us set
 291 $y = 2\lambda_j$ and set

$$G_k(y) = \sum_{i=0}^{k-1} y^i \left[(\hat{\delta}^k)^{(i-1)}(1) - (\hat{\delta}^k)^{(i-1)}(-1) - 2(\hat{\delta}^k)^{(i)}(-1) \right].$$

292 Since $k-1$ is odd, again, we would like to show it has positive real roots. First, when
 293 $y = 0$, by (3.7) and Lemma 3.5 (i) we have

$$G_k(0) = \int_{\hat{I}} \hat{\delta}^k(x) dx - 2\hat{\delta}^k(-1) = 1 - (k+1) = -k \leq -2.$$

294 Next, let us check the leading coefficient of $G_k(y)$. Since $(\hat{\delta}^k)^{(k)} = \beta_k^k$, we have
 295 $(\hat{\delta}^k)^{(k-2)} = \frac{1}{2}\beta_k^k x^2 + C_1 x + C_2$, where C_1 and C_2 are constants. As a consequence, we
 296 have

$$\alpha_{k-2}^k = \frac{1}{2}\beta_k^k + C_1 + C_2, \quad \beta_{k-2}^k = \frac{1}{2}\beta_k^k - C_1 + C_2, \quad \beta_{k-1}^k = -\beta_k^k + C_1,$$

297 and hence the leading coefficient of $G_k(y)$ satisfies

$$(\hat{\delta}^k)^{(k-2)}(1) - (\hat{\delta}^k)^{(k-2)}(-1) - 2(\hat{\delta}^k)^{(k-1)}(-1) = \alpha_{k-2}^k - \beta_{k-2}^k - 2\beta_{k-1}^k = 2\beta_k^k > 0.$$

298 Therefore, the odd-degree polynomial $G_k(y)$ must have at least one and at most $k-1$
 299 positive real roots. If we take η_2^k to be the smallest one, we can conclude the second
 300 statement. \square

301 *Remark 3.10.* This theorem shows that a lower bound for the CFL number is necessary for conditions (3.22) and (3.23), which are sufficient conditions for the positivity of the cell averages at the next time step. It does not imply the necessity of the lower bounds to guarantee the cell averages' positivity. The latter necessity will be confirmed by the numerical evidence in Table 3.2 and Table 3.3.

306 Theorem 3.9 indicates that unlike the situation for the DG method with Euler forward time discretization in [41], where an upper bound for the CFL number is sufficient for the cell average at the next time level to be positive, for the DG scheme with the Euler backward time discretization, an lower bound may be required. The following analysis and numerical experiments confirm this statement.

311 If we re-examine the condition (3.23) and note that for fixed λ_j , $F_k(\lambda_j, x) \in P^k(\hat{I})$, the inner product in (3.23) can actually be approximated exactly by certain quadrature rules, say $\{(x^\alpha, \omega^\alpha)\}_{\alpha=1}^{N_q}$, where $\{x^\alpha\}$ are the abscissas in \hat{I} , $\{\omega^\alpha\}$ the weights and N_q is large enough such that the quadrature rule is exact for polynomials of degree $2k$. We denote $\{(x_j^\alpha, \omega_j^\alpha)\}_{\alpha=1}^{N_q}$ to be the transformed quadrature rule in I_j . Then the condition (3.23) can be reduced to require

$$J_\alpha^k(\lambda_j) := F_k(\lambda_j, x^\alpha) - F_k(\lambda_j, -1) \geq 0, \quad \alpha = 1, \dots, N_q. \quad (3.24)$$

317 If we define $J_0^k(\lambda_j) = F_k(\lambda_j, -1)$, together with condition (3.22), we require the CFL number to make $N_q + 1$ polynomials $\{J_\alpha^k(\lambda_j)\}_{\alpha=0}^{N_q}$ positive. The following result states that it suffices to require $\lambda_j \geq \frac{1}{2}$.

320 **PROPOSITION 3.11.** *When $\min_j \lambda_j \geq \frac{1}{2}$, we have*

$$F_k(\lambda_j, -1) > 0, \quad (3.25)$$

$$F_k(\lambda_j, x) - F_k(\lambda_j, -1) > 0, \quad \forall x \in (-1, 1]. \quad (3.26)$$

321 *Proof.* First, let us consider (3.25). By the definition we have $F_k(\lambda_j, -1) = \sum_{i=0}^k (2\lambda_j)^k \beta_i^k$, which can be rewritten as

$$F_k(\lambda_j, -1) = (2\lambda_j)^{k-1} (2\lambda_j \beta_k^k + \beta_{k-1}^k) + \dots + \begin{cases} (2\lambda_j \beta_2^k + \beta_1^k) + \beta_0^k, & \text{if } k \text{ is even} \\ (2\lambda_j \beta_1^k + \beta_0^k), & \text{if } k \text{ is odd} \end{cases}. \quad (3.27)$$

323 When $\lambda_j > 1/2$, i.e., $2\lambda_j > 1$, by Lemma 3.8 (ii) and (iii) we can show that each term in (3.27) is strictly positive and hence $F_k(\lambda_j, -1) > 0$.

325 For the condition (3.26), consider the i th derivative of F_k with respect to x

$$\frac{\partial^i}{\partial x^i} F_k(\lambda_j, x) = \sum_{l=i}^k (2\lambda_j)^{l-i} (\hat{\delta}^k)^{(l)}(x), \quad i = 0, \dots, k.$$

326 When $i = k$, by Lemma 3.8 (ii), we have $\frac{\partial^k}{\partial x^k} F_k(\lambda_j, x) = \beta_k^k > 0$, and hence

$$\frac{\partial^{k-1}}{\partial x^{k-1}} F_k(\lambda_j, x) > \frac{\partial^{k-1}}{\partial x^{k-1}} F_k(\lambda_j, -1) = \beta_{k-1}^k + (2\lambda_j) \beta_k^k > 0, \quad \forall x \in (-1, 1].$$

327 We have used the fact $2\lambda_j \geq 1$ and Lemma 3.8 (ii) and (iii) to derive the last inequality.

328 Then we continue the same procedure till $i = 0$ and obtain

$$F_k(\lambda_j, x) > F_k(\lambda_j, -1), \quad \forall x \in (-1, 1].$$

329 Therefore, in summary, when $\min_j \lambda_j \geq \frac{1}{2}$, we have conditions (3.25) and (3.26) hold
 330 and hence we can draw the conclusion. \square

331 Consequently, when $\lambda_j \geq \frac{1}{2}$ all the polynomials $\{J_\alpha^k\}_{\alpha=0}^k$ are strictly positive. On
 332 the other hand by the proof of Theorem 3.9, at $\lambda_j = 0$ these polynomials can not be
 333 all nonnegative. This implies that if we denote \mathcal{R}_k to be the set of all the positive
 334 roots of each polynomial J_α^k , i.e.,

$$\mathcal{R}_k = \{r \in \mathbb{R}^+ : J_\alpha^k(r) = 0 \text{ for some } \alpha = 0, \dots, N_q\},$$

335 then we must have $\mathcal{R}_k \neq \emptyset$, $\mathcal{R}_k \in (0, \frac{1}{2})$ and \mathcal{R}_k is finite. Then if we set

$$r_k = \max_{r \in \mathcal{R}_k} r, \tag{3.28}$$

336 we have the following theorem.

337 **THEOREM 3.12.** *Let $\{x_j^\alpha\}_{\alpha=1}^{N_q} \in I_j$ be a set of quadrature points which is exact for*
 338 *polynomials of degree $2k$ and let $u_j^n \in P^k(I_j)$ be the backward Euler DG approximation*
 339 *for the linear equation (3.3) in the cell I_j at time t^n . Then given $u_j^n(x_j^\alpha) \geq 0$ for*
 340 *$\alpha = 1, \dots, N_q$ and $j = 1, \dots, N$, we have $\bar{u}_j^{n+1} \geq 0$ for any j under the following CFL*
 341 *condition*

$$\min_j \lambda_j \geq r_k \tag{3.29}$$

342 where $r_k \in (0, \frac{1}{2})$ is defined as in (3.28).

343 *Proof.* When $\lambda_j \geq r_k$ by the definition of \mathcal{R}_k , we have (3.24) and (3.22) hold.
 344 Therefore T in (3.18) is an M -matrix and $\mathcal{L} \geq 0$. By the definition of M -matrix, we
 345 can conclude the result. \square

346 *Remark 3.13.* We only require u_h^n to be positive on quadrature points $\{x_j^\alpha\}_{\alpha=1}^{N_q}$, which
 347 is weaker than the condition $u_h^n(x) \geq 0$, for any $x \in \Omega$ in the Definition 3.1.

348 *Remark 3.14.* Even though the Theorem 3.12 is only proved for linear equations,
 349 numerical experiments suggest that for nonlinear equations, a lower bound for the
 350 CFL number is still necessary to make the cell average at the next time level positive.

351 *Remark 3.15.* The results also hold for problems with positive source terms and posi-
 352 tive inflow boundary conditions.

353 *Remark 3.16.* The lower bound depends on the polynomial degree k as well as the
 354 quadrature rule we choose. For each k and fixed quadrature rule we can actually
 355 obtain the lower bound r_k by solving the positive roots of each J_α^k . In Table 3.1, we
 356 record the lower bounds for Legendre-Gauss-Lobatto (LGL) quadrature rule and the
 357 Legendre-Gauss (LG) rule respectively. We see that the LGL rule gives smaller lower
 358 bound. In practice we will limit the polynomial u_j^n to make it positive at least on the
 359 LGL points $\{x_j^\alpha\}_{\alpha=1}^{N_q}$ in each cell I_j .

360 *Remark 3.17.* The lower bounds in Table 3.1 are sharp for odd k and sufficient for
 361 even k . If we start from the following initial condition,

$$u_0(x) = \begin{cases} 1, & \text{if } x \in I_M \\ 0, & \text{otherwise} \end{cases},$$

TABLE 3.1

Values of r_k for $k = 1, \dots, 5$ and for Legendre-Gauss-Lobatto (LGL) and Legendre-Gauss quadrature rules respectively.

k	LGL rule		LG rule	
	N_q	r^k	N_q	r^k
1	3	0.333	2	0.333
2	4	0.262	3	0.344
3	5	0.177	4	0.177
4	6	0.177	5	0.212
5	7	0.121	6	0.121

where $M = \arg \max_j h_j$. After one step we record the minimum cell averages for different odd k and λ in Table 3.2. We see that when λ is slightly smaller than the lower bound r_k , after one time step, at least one of the cell averages will become negative. If λ is larger than r_k , the average will be uniformly positive.

TABLE 3.2

Minimum cell average after one time step for odd k .

k	$\lambda = r^k - 0.001$	$\min_j \bar{u}_j^1$	$\lambda = r^k + 0.001$	$\min_j \bar{u}_j^1$
1	0.332	-3.498 E-04	0.334	5.284 E-35
3	0.176	-4.177 E-05	0.178	7.941 E-46
5	0.120	-2.135 E-06	0.122	1.980 E-59

365
366

For even k , we consider a different initial condition

$$u_0(x) = \begin{cases} \left(\frac{2(x-x_M)}{h_M} - 0.72 \right)^k, & \text{if } x \in I_M \\ 0, & \text{otherwise} \end{cases}.$$

The Table 3.3 shows the minimum cell averages for different k and λ . We see that the lower bound for the CFL number is still necessary for the positivity of \bar{u}_j^1 and r_k listed in Table 3.3 is sufficient.

TABLE 3.3

Minimum cell average after one time step for even k .

k	λ	$\min_j \bar{u}_j^1$	$\lambda = r^k$	$\min_j \bar{u}_j^1$
2	0.170	-1.022 E-03	0.262	1.221 E-28
4	0.120	-1.343 E-03	0.177	5.021 E-46

369

3.3. Scaling limiter. Once in each cell I_j , the cell average \bar{u}_j^{n+1} is positive, we limit the whole polynomial $u_j^{n+1}(x)$ towards its cell average by utilizing the following scaling limiter [20, 41].

$$\tilde{u}_j^{n+1} = \theta_j [u_j^{n+1} - \bar{u}_j^{n+1}] + \bar{u}_j^{n+1} \quad (3.30)$$

where

$$\theta_j = \begin{cases} \frac{\bar{u}_j^{n+1}}{\bar{u}_j^{n+1} - \min_{x \in I_j} u_j^{n+1}(x)}, & \text{if } \min_{x \in I_j} u_j^{n+1}(x) < 0 \\ 1, & \text{otherwise} \end{cases}. \quad (3.31)$$

374 This procedure preserves the original high-order accuracy [41].

375 LEMMA 3.18. For the modified polynomial $\tilde{u}_j^{n+1}(x)$, we have

$$|\tilde{u}_j^{n+1}(x) - u_j^{n+1}(x)| \leq C_k \max_{x \in I_j} |u_j^{n+1}(x) - u(x)|$$

376 where u is the smooth solution.

377 Basically, what this lemma says is that the error we commit in the limiting procedure
378 is bounded by the error of the original approximation up to a constant depending on
379 k . The proof can be found in [40].

380 *Remark 3.19.* In order to calculate the scaling parameter θ_j in (3.31) we need to
381 calculate the minimum value of u_j^{n+1} in each cell I_j , which can be done efficiently
382 up to $k = 5$ via the root formulas. For larger k , this calculation becomes expensive.
383 But recall that in Theorem 3.12, we only require each polynomial u_j^{n+1} to be positive
384 at the LGL quadrature points $\{x_j^\alpha\}_{\alpha=1}^{N_q}$ and hence we can instead use the following
385 scaling parameter

$$\tilde{\theta}_j = \begin{cases} \frac{\bar{u}_j^{n+1}}{\bar{u}_j^{n+1} - \min_\alpha u_j^{n+1}(x_j^\alpha)}, & \text{if } \min_\alpha u_j^{n+1}(x_j^\alpha) < 0 \\ 1, & \text{otherwise} \end{cases} \quad (3.32)$$

386 in the limiter. This is similar with the scaling limiter in [41] and the same argument
387 can be conducted here to prove that the modified limiter does not kill the original
388 high-order accuracy either.

389 **3.4. Algorithm for scalar equations.** Now, we can summarize the positivity-
390 preserving algorithm for scalar equations as below.

- 391 1. At time level t^n , given $u_j^n(x)$ being positive at least on the LGL quadrature
392 points $\{x_j^\alpha\}_{\alpha=1}^{N_q}$.
- 393 2. Choose a sufficiently large CFL number, a priori or adaptively enlarge it in
394 each time step until $\bar{u}_j^{n+1} \geq 0, \forall j$.
- 395 3. Apply the scaling limiter (3.30) with θ_j defined in (3.31) or (3.32) to u_j^{n+1}
396 such that $u_j^{n+1}(x_j^\alpha) \geq \epsilon$, for any $\alpha = 1, \dots, N_q$ and j , where ϵ is a small
397 number to help get rid of the round-off effect. In the numerical examples, we
398 take $\epsilon = 10^{-13}$.

399 **4. Positivity-preserving DG scheme for compressible Euler systems.**

400 We consider the following compressible Euler system for ideal gas

$$\mathbf{u}_t + \mathbf{f}(\mathbf{u})_x = \mathbf{0}, \quad t \geq 0, \quad x \in [0, 1], \quad (4.1)$$

401 with

$$\mathbf{u} = (\rho, m, E)^T, \quad \mathbf{f} = (m, \rho v^2 + p, (E + p)v)^T$$

402 and $m = \rho v$, $E = \frac{1}{2}\rho v^2 + \rho e$, $p = (\gamma - 1)\rho e$, where ρ is the density, v the velocity,
403 m the momentum, p the pressure, E the total energy and e is the internal energy.
404 The constant $\gamma > 1$ is the ratio of specific heats. The Jacobian matrix $\frac{\partial \mathbf{f}}{\partial \mathbf{u}}$ has the
405 following three eigenvalues $\zeta_1 = v - c$, $\zeta_2 = v$, $\zeta_3 = v + c$, where $c = \sqrt{\gamma p / \rho}$ is the
406 sound speed.

407 The physical solution lies in the following admissible set

$$G = \left\{ \mathbf{u} = (\rho, m, E)^T : \rho \geq 0, p = (\gamma - 1) \left(E - \frac{1}{2} \frac{m^2}{\rho} \right) \geq 0 \right\}. \quad (4.2)$$

408 It can be verified that G is a convex set [42].

409 The DG scheme for (4.1) is to seek an approximation vector $\mathbf{u}_h(t) \in \mathbf{V}_h = [V_h]^3$
410 such that in each cell I_j we have

$$\frac{d}{dt}(\mathbf{u}_h(t), \mathbf{v})_j = L_j(\mathbf{u}_h(t), \mathbf{v})_j, \quad \forall \mathbf{v} \in \mathbf{V}_h, \quad (4.3)$$

411 where $L_j(\mathbf{u}_h(t), \mathbf{v}) = (\mathbf{f}(\mathbf{u}_h(t)), \mathbf{v}_x)_j - \left[\hat{\mathbf{f}}_{j+\frac{1}{2}}(\mathbf{u}_h(t)) \mathbf{v}(x_{j+\frac{1}{2}}^-) - \hat{\mathbf{f}}_{j-\frac{1}{2}}(\mathbf{u}_h(t)) \mathbf{v}(x_{j-\frac{1}{2}}^+) \right]$.

412 The numerical flux $\hat{\mathbf{f}}(\cdot, \cdot)$ is taken to be the global Lax-Friedrichs flux

$$\hat{\mathbf{f}}(\mathbf{a}, \mathbf{b}) = \frac{1}{2}[\mathbf{f}(\mathbf{a}) + \mathbf{f}(\mathbf{b}) - \alpha(\mathbf{b} - \mathbf{a})]$$

413 with $\alpha = \|c + |v|\|_\infty$.

414 As for the scalar case, the ODE system (4.3) is solved by the backward Euler
415 method. If we use \mathbf{u}_h^n to denote the DG approximation at time level n , then the
416 positivity-preserving DG scheme for the compressible Euler system is defined as below.

417

418 **DEFINITION 4.1.** *A DG scheme for the compressible Euler system is positivity-preserving if at time level n given $\mathbf{u}_h^n(x) \in G$ for all $x \in \Omega$, then at the next time level $(n+1)$, we have $\mathbf{u}_h^{n+1}(x) \in G$ for all $x \in \Omega$.*

421 We design the positivity-preserving DG scheme by extending the positivity-preserving
422 limiter in [42] and the more robust version in [36] for the explicit time stepping
423 to the backward Euler time stepping. We stress on the applicability of the proposed
424 method rather than its theoretical justification. The analysis for the linear scalar
425 equation suggests that starting from $\mathbf{u}_j^n(x) \in G$, in order to have the cell average
426 $\bar{\mathbf{u}}_j^{n+1} \in G$, a lower bound for the CFL number may be required. Therefore, we
427 formulate the algorithm as follows.

- 428 1. At time level t^n , in each cell I_j , given $\mathbf{u}_j^n(x_j^\alpha) \in G$ at the LGL quadrature
429 points $\{x_j^\alpha\}_{\alpha=1}^{N_q}$.
- 430 2. Choose a large enough CFL number, a priori or adaptively enlarge it in each
431 time step until $\bar{\mathbf{u}}_j^{n+1} \in G$ for any j .
- 432 3. In each cell I_j , apply the following scaling limiter to the first component of
433 \mathbf{u}_j^{n+1} to obtain $\rho_j^{n+1}(x) \geq 0$ at the LGL quadrature points $\{x_j^\alpha\}_{\alpha=1}^{N_q}$

$$\tilde{\rho}_j^{n+1} = \theta_1(\rho_j^{n+1} - \bar{\rho}_j^{n+1}) + \bar{\rho}_j^{n+1}$$

434 where

$$\theta_1 = \begin{cases} \frac{\bar{\rho}_j^{n+1}}{\bar{\rho}_j^{n+1} - \min_\alpha \rho_j^{n+1}(x_j^\alpha)}, & \text{if } \min_\alpha \rho_j^{n+1}(x_j^\alpha) < 0 \\ 1, & \text{otherwise} \end{cases}.$$

435 Denote the modified polynomial by $\tilde{\mathbf{u}}_j^{n+1}$.

436 4. In each cell I_j , apply the scaling limiter again to the whole modified polyno-
 437 mial $\tilde{\mathbf{u}}_j^{n+1}$ such that $p(x_j^\alpha) \geq 0$ for each α as below

$$\tilde{\mathbf{u}}_j^{n+1} = \theta_2(\tilde{\mathbf{u}}_j^{n+1} - \bar{\mathbf{u}}_j^{n+1}) + \bar{\mathbf{u}}_j^{n+1},$$

438 where

$$\theta_2 = \begin{cases} \frac{\bar{p}_j^{n+1}}{\bar{p}_j^{n+1} - \min_\alpha \tilde{p}_j^{n+1}(x_j^\alpha)}, & \text{if } \min_\alpha \tilde{p}_j^{n+1}(x_j^\alpha) < 0 \\ 1, & \text{otherwise} \end{cases}$$

439 and the pressure average is defined by $\bar{p}_j^{n+1} = p(\bar{\mathbf{u}}_j^{n+1})$.

440 **5. Numerical experiments.** In this section, we present numerical examples.
 441 First, we verify the high-order **spatial** accuracy of the proposed method by testing it on
 442 both linear and nonlinear steady-state problems. An acceleration of the convergence
 443 towards the steady state solution is also observed. Next, we test the methods on
 444 moving-shock problems. At last, examples for **the** compressible Euler system will
 445 be presented. **In all the examples, the domain is first uniformly decomposed with**
 446 **meshsize h and then each node $x_{j+\frac{1}{2}}$ is randomly perturbed in the range $[x_{j+\frac{1}{2}} -$**
 447 **$\frac{h}{5}, x_{j+\frac{1}{2}} + \frac{h}{5}]$.**

448 **5.1. Accuracy tests.** First, let us test the accuracy of the proposed method.
 449 We check the **spatial** accuracy with the steady-state solution to both of the linear
 450 equation and the **Burgers'** equation. We take $\Delta t = 10 \max_j h_j$ and march in time
 451 until $\|u_h^{n+1} - u_h^n\|_2 \leq 10^{-12}$.

452 *Example 5.1 (Steady-state solution to linear problem).* For the linear equation, we co-
 453 nsider the following problem

$$u_t + u_x = \sin^4(x), \quad u(x, 0) = \sin^2(x), \quad u(0, t) = 0, \quad (5.1)$$

454 with the outflow boundary condition at $x = 2\pi$. The exact solution $u(x, t)$ can be
 455 derived by the characteristic theory and can be shown to be positive for all $t > 0$.
 456 In Table 5.1 and Table 5.2 we record the errors, numerical orders of accuracy and
 457 the minimum value of the numerical approximation, without and with the positivity-
 458 preserving limiter respectively. We see that without the positivity-preserving limiter
 459 the minimum value of the steady-state approximation is negative. When the limiter
 460 is put on, the minimum value becomes positive and the high-order accuracy is not
 461 destroyed.

462 *Example 5.2 (Steady-state solution to Burgers' problems).* For the Burgers' equation,
 463 we consider the steady state solution to the following problem

$$u_t + \left(\frac{u^2}{2}\right)_x = \sin\left(\frac{x}{4}\right), \quad u(x, 0) = x, \quad u(0, t) = 0, \quad (5.2)$$

464 with the outflow boundary condition at $x = 2\pi$. Again, by the characteristic theory,
 465 one can show that the solution to this problem is always positive. In Tables 5.3 and
 466 5.4, we present numerical results for the cases where the positivity-preserving limiter
 467 is off and on respectively. This example shows the effectiveness of the positivity
 468 preserving limiter for the nonlinear scalar problem. With the limiter on, the solution
 469 stays positive and the high-order accuracy is preserved.

TABLE 5.1

Error table for Example 5.1, approximation of the steady state solution to the linear problem (5.1), without the positivity preserving limiter.

k	N	L^2 error	order	L^∞ error	order	$\min u_h$
1	20	4.555 E-2	–	6.376 E-2	–	-6.037 E-2
	40	1.177 E-2	1.95	1.704 E-2	1.90	-5.075 E-3
	80	2.967 E-3	1.99	4.347 E-3	1.97	-2.808 E-4
	160	7.434 E-4	2.00	1.092 E-3	1.99	-1.170 E-5
	320	1.859 E-4	2.00	2.733 E-4	2.00	-3.901 E-7
2	20	5.749 E-3	–	9.561 E-3	–	-1.667 E-3
	40	7.482 E-4	2.94	1.121 E-3	3.09	-6.550 E-5
	80	9.449 E-5	2.99	1.543 E-4	2.86	-2.163 E-6
	160	1.184 E-5	3.00	1.975 E-5	2.97	-6.853 E-8
	320	1.481 E-6	3.00	2.484 E-6	2.99	-2.149 E-9
3	20	6.987 E-4	–	6.013 E-4	–	-8.652 E-4
	40	4.564 E-5	3.94	4.743 E-5	3.66	-3.909 E-5
	80	2.885 E-6	3.98	2.986 E-6	3.99	-1.329 E-6
	160	1.808 E-7	4.00	1.887 E-7	3.98	-4.240 E-8
	320	1.131 E-8	4.00	1.178 E-8	4.00	-1.332 E-9
4	20	6.094 E-5	–	4.622 E-5	–	-6.545 E-6
	40	1.955 E-5	4.96	1.335 E-6	5.11	-1.329 E-6
	80	6.149 E-8	4.99	4.597 E-8	4.86	-5.211 E-8
	160	1.925 E-9	5.00	1.471 E-9	4.97	-1.715 E-9
	320	6.018 E-11	5.00	4.623 E-10	4.99	-5.426 E-11

Next, let us turn to another steady-state Burgers' problem

$$u_t + \left(\frac{u^2}{2}\right)_x = \sin^3\left(\frac{x}{4}\right), \quad u(x, 0) = \sin^2\left(\frac{x}{4}\right), \quad u(0, t) = 0, \quad (5.3)$$

with the outflow boundary condition at $x = 2\pi$. This problem is more difficult than the previous one, since the source is much closer to zero around $x = 0$. In Table 5.5, we present the error, the numerical convergence rate as well as the number of time steps taken to reach the steady state solution, which is denoted by N_T . Both cases where the positivity preserving limiter is on and off are presented respectively. We use this example to illustrate that the stability added to the scheme by the positivity-preserving limiter helps accelerate the convergence towards the steady-state solution.

With this example, we also show the advantage of implicit methods over explicit ones in steady-state simulations. In Table 5.6, we record the CPU time for the backward Euler and the TVD-RK3 [34] time discretizations. The space is discretized with P^2 -DG element with positivity-preserving limiter on. The error tables for both time discretizations are exactly the same, as shown in Table 5.5. However, since the backward Euler method allows large CFL number, $\lambda = 10$ in this example, it is 20 times faster than the TVD-RK3 method to reach the steady state.

5.2. Moving Shocks. Next, we test the proposed scheme on problems involving moving shocks. In all the numerical experiments below, quadratic P^2 elements and non-uniform meshes are employed.

TABLE 5.2

Error table for Example 5.1, approximation of the steady state solution to the linear problem (5.1) with the positivity preserving limiter.

k	N	L^2 error	order	L^∞ error	order	$\min u_h$
1	20	4.253 E-2	–	6.376 E-2	–	1.000 E-13
	40	1.173 E-2	1.86	1.704 E-2	1.90	1.000 E-13
	80	2.966 E-3	1.98	4.347 E-3	1.97	1.000 E-13
	160	7.434 E-4	2.00	1.092 E-3	1.99	1.000 E-13
	320	1.859 E-4	2.00	2.733 E-4	2.00	1.000 E-13
2	20	5.762 E-3	–	9.561 E-3	–	1.000 E-13
	40	7.482 E-4	2.95	1.121 E-3	3.09	1.000 E-13
	80	9.449 E-5	2.99	1.543 E-4	2.86	1.000 E-13
	160	1.184 E-5	3.00	1.975 E-5	2.97	1.000 E-13
	320	1.481 E-6	3.00	2.484 E-6	2.99	1.000 E-13
3	20	1.015 E-3	–	2.240 E-3	–	1.000 E-13
	40	5.077 E-5	4.32	9.673 E-5	4.53	1.000 E-13
	80	2.932 E-6	4.11	3.253 E-6	4.89	1.000 E-13
	160	1.812 E-7	4.02	1.887 E-7	4.11	1.000 E-13
	320	1.131 E-8	4.00	1.178 E-8	4.00	1.000 E-13
4	20	6.141 E-5	–	4.622 E-5	–	1.000 E-13
	40	2.230 E-5	4.78	4.356 E-6	3.41	1.000 E-13
	80	6.830 E-8	5.03	1.700 E-7	4.68	1.000 E-13
	160	2.045 E-9	5.06	5.591 E-9	4.93	1.000 E-13
	320	6.214 E-10	5.04	1.773 E-10	4.98	1.000 E-13

TABLE 5.3

Error table for Example 5.2, approximation of the steady state solution to the Burgers' equation (5.2), without the positivity preserving limiter.

k	N	L^2 error	order	L^∞ error	order	$\min u_h$
2	20	2.915 E-6	–	5.085 E-6	–	-8.073 E-6
	40	3.508 E-7	3.05	6.358 E-7	3.00	-1.009 E-6
	80	4.293 E-8	3.03	7.948 E-8	3.00	-1.261 E-7
	160	5.304 E-9	3.02	9.934 E-9	3.00	-1.577 E-8
3	20	2.336 E-9	–	1.648 E-9	–	-3.855 E-10
	40	1.445 E-10	4.02	1.030 E-10	4.00	-1.204 E-11
	80	9.010 E-11	4.00	6.525 E-11	3.98	-3.763 E-13
	160	6.031 E-12	3.90	4.978 E-12	3.71	-1.175 E-14
4	20	3.403 E-10	–	6.312 E-10	–	-1.497 E-9
	40	1.010 E-11	5.07	1.970 E-11	5.00	-4.678 E-11
	80	3.116 E-13	5.02	5.398 E-13	5.19	-1.462 E-12

489 Example 5.3 (Linear Problem). The first example is the linear equation with the ini-
 490 tial data

$$u_0(x) = \begin{cases} 1, & \text{if } x \in [3, 4] \\ 0, & \text{otherwise} \end{cases}.$$

491 In Figure 5.1, we present the numerical results at $T = 0.2$, with and without the

TABLE 5.4

Error table for Example 5.2, approximation of the steady state solution to the Burgers' equation (5.2), with the positivity preserving limiter.

k	N	L^2 error	order	L^∞ error	order	$\min u_h$
2	20	3.207 E-6	–	4.408 E-6	–	1.000 E-13
	40	3.696 E-7	3.12	6.010 E-7	2.87	1.000 E-13
	80	4.435 E-8	3.06	8.171 E-8	2.88	1.000 E-13
	160	5.414 E-9	3.03	1.083 E-8	2.92	1.000 E-13
3	20	2.338 E-9	–	1.648 E-9	–	1.000 E-13
	40	1.445 E-10	4.02	1.030 E-10	4.00	1.000 E-13
	80	9.010 E-11	4.00	6.511 E-11	3.98	1.000 E-13
	160	6.031 E-12	3.90	4.987 E-12	3.71	1.051 E-14
4	20	4.792 E-10	–	9.061 E-10	–	1.000 E-13
	40	1.253 E-11	5.26	3.102 E-11	4.87	1.000 E-13
	80	3.653 E-13	5.10	1.086 E-12	4.84	1.000 E-13

TABLE 5.5

Error table for Example 5.2, approximation of the steady-state solution to the Burgers' equation (5.3), with and without the positivity preserving limiter.

k	N	without limiter			with limiter		
		L^2 error	order	N_T	L^2 error	order	N_T
2	20	1.71 E-5	–	670	1.71 E-5	–	158
	40	2.11 E-6	3.02	1962	2.11 E-6	3.02	500
	80	2.62 E-7	3.01	4973	2.62 E-7	3.01	1560
	160	3.29 E-8	2.99	8352	3.27 E-8	3.01	4560
3	20	9.10 E-8	–	1162	9.10 E-8	–	970
	40	5.36 E-9	4.08	3440	5.36 E-9	4.08	2211
	80	3.44 E-10	3.96	9007	3.35 E-10	4.00	2653

TABLE 5.6

CPU time for Example 5.2, problem (5.3) solved by P^2 -DG with backward Euler and TVD-RK3 temporal discretizations. The positivity-preserving limiter is on.

N	20	40	80	160
Backward Euler	0.31	1.57	9.73	58.12
TVD-RK3	5.16	32.16	199.75	1209.09

positivity-preserving limiter respectively. For the linear equation, we take $\min_j \lambda_j = r_2 = 0.262$ as listed in Table 3.1. As the zoom-in plots show, the limiter helps the solution to stay positive. Without the limiter a negative undershoot appears.

Example 5.4 (Burgers' problem). Next, let us consider the Burgers' equation with the initial condition $u_0(x) = 1 + \sin(x)$ and periodic boundary conditions. The initial condition is positive and takes value zero at $x = \pi$. At $T = 1.5$, a shock is developed. In Figure 5.2, we show the numerical approximation with and without the limiter respectively. We see that even though the profiles are smeared due to the first-order accuracy of the backward Euler time discretization, the effectiveness of the positivity-preserving limiter can still be observed in the zoom-in plots around the shock.

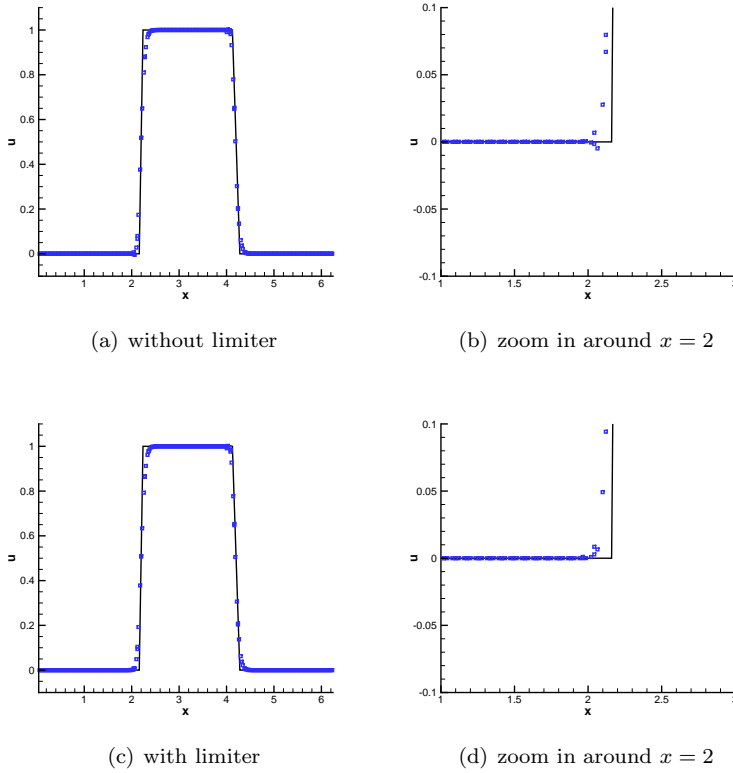


FIG. 5.1. *Example 5.3: at $T = 0.2$, with $\Delta t = 0.266 \max_j h_j$ and $N = 120$. Solid line is the exact solution. Blue squares are point values at the Legendre-Gauss-Lobatto points.*

502 *Example 5.5 (Buckley-Leverett problem).* In this example, the Buckley-Leverett prob-
 503 lem in $[-1, 1]$ is considered, with $f(u) = \frac{4u^2}{4u^2 + (1-u)^2}$. For this flux function, we have
 504 $f'(0) = f'(1) = 0$. If we start with the usual step function $u_0(x) = I_{[-0.5, 1]}$, numerical
 505 experiments indicate that no matter how large Δt is, we can not make the cell average
 506 \bar{u}_j^1 positive for all j . The same phenomenon is also observed for the Burgers' problem
 507 with step function as the initial condition. This indicates that the lower bound for the
 508 CFL number is also necessary for the positivity-preserving limiter to work for nonlin-
 509 ear problems, and in general, the limiter may not be an effective positivity-preserving
 510 tool when applied to problems with sonic points. In this example, we change the
 511 initial condition to

$$u_0(x) = \begin{cases} 0.9, & \text{if } x \in [-0.5, 1] \\ 10^{-3}, & \text{otherwise} \end{cases}.$$

512 In Figure 5.3, we present the plots with and without limiter. Even though the pro-
 513 file is smeared around the shocks and the rarefaction wave by the first-order time-
 514 discretization, the positivity-preserving property is observed in the zoom-in plots.
 515

516 **5.3. Compressible Euler System.** At last, we turn to the examination of the
 517 applicability of the proposed scheme to the compressible Euler system (4.1). The

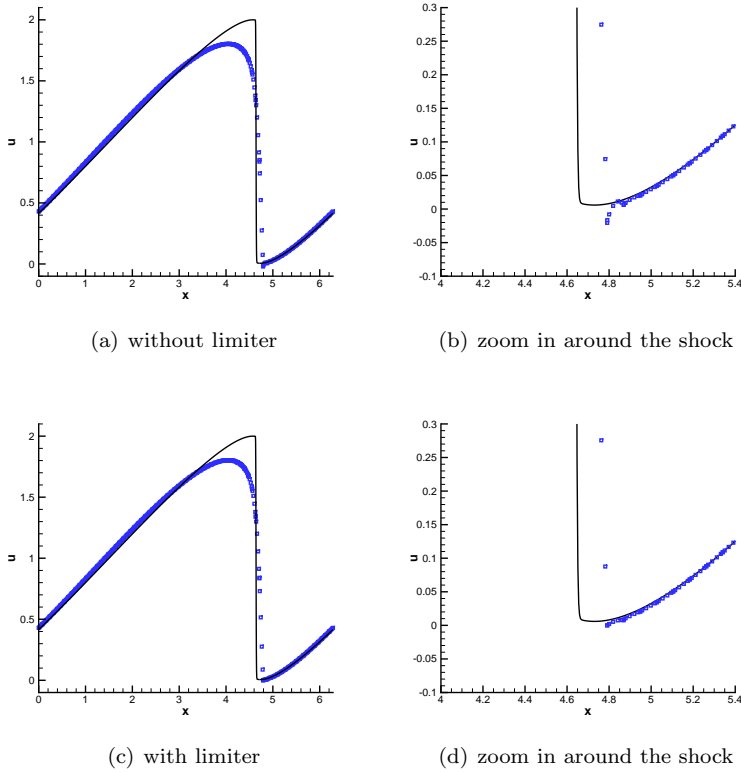


FIG. 5.2. Example 5.4 at $T = 1.5$. Solid line is the exact solution. Blue squares are point values of the numerical approximation on the Legendre-Gauss-Lobatto points. Mesh with $N = 120$ cells and $CFL = 2$.

518 computational domain $\Omega = [0, 1]$ is decomposed into a nonuniform mesh. The generic
 519 ratio of specific heats is taken to be $\gamma = 1.4$. For all the examples, quadratic element
 520 is employed with the positivity-preserving limiter on. And the time stepping size is
 521 set to be $\Delta t = \frac{CFL}{\|c+|v|\|_\infty} \min_j h_j$, where c is the sound speed and v is the velocity. We
 522 take $CFL = 2$ for all the examples.

523 *Example 5.6 (Shock tube)*. First, let us consider the following shock tube problem

$$\begin{cases} \rho = 1, & \text{if } x < 0.5 \\ v = 0, & \text{if } x < 0.5, \\ p = 1000, & \text{if } x < 0.5 \end{cases}, \quad \begin{cases} \rho = 1, & \text{if } x \geq 0.5 \\ v = 0, & \text{if } x \geq 0.5. \\ p = 0.01, & \text{if } x \geq 0.5 \end{cases}.$$

524 The numerical approximation is presented in Figure 5.4. The solutions consists of a
 525 strong shock wave, a contact discontinuity and a rarefaction wave. Due to the first-
 526 order temporal discretization and the large CFL number, the contact discontinuity
 527 and the shock wave in the density are not very well captured. However, with the
 528 positivity-preserving limiter on, both the pressure and the density stay positive all
 529 the time during the simulation.

530 *Example 5.7 (Double rarefaction)*. The double rarefaction problem starts from the

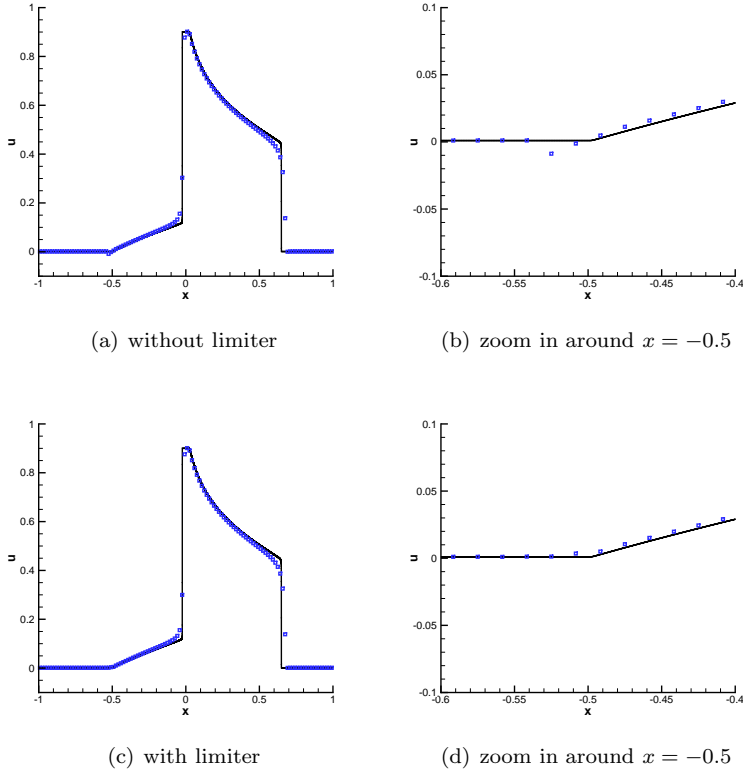


FIG. 5.3. Example 5.5 at $T = 0.4$. Solid line is the exact solution. Blue squares are point values of the numerical approximation on the Legendre-Gauss-Lobatto points. Mesh with $N = 120$ cells and $CFL = 3.5$.

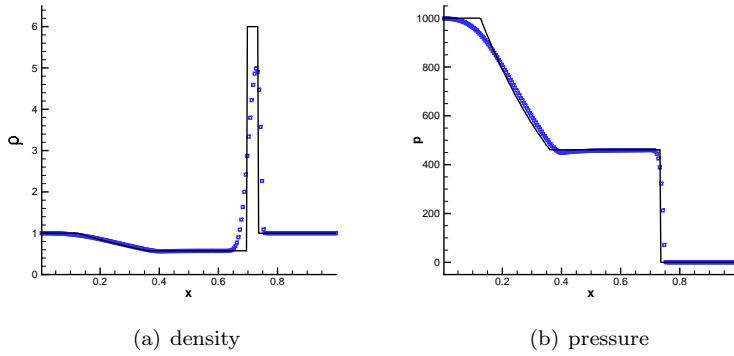


FIG. 5.4. Example 5.6 at $T = 0.01$ with $N = 200$ and $CFL = 2$. With the positivity-preserving limiter on. Solid line is the exact solution and blue squares are numerical approximations.

531 following initial condition

$$\begin{cases} \rho = 1, & \text{if } x < 0.5 \\ v = -2, & \text{if } x < 0.5, \\ p = 0.4, & \text{if } x < 0.5 \end{cases} \quad \begin{cases} \rho = 1, & \text{if } x \geq 0.5 \\ v = 2, & \text{if } x \geq 0.5. \\ p = 0.4, & \text{if } x \geq 0.5 \end{cases}$$

532 This problem has a solution consisting of two symmetric rarefaction waves and a
 533 trivial contact wave of zero speed. The region between the nonlinear waves around
 534 $x = 0.5$ is close to vacuum, which brings difficulty to the simulation. Without the
 535 limiter, the nonlinear solver will experience a hard time to converge. In Figure 5.5,
 536 we show the plots for the pressure and the density and we see that the near-vacuum
 region is well resolved with the help of the positivity-preserving limiter.

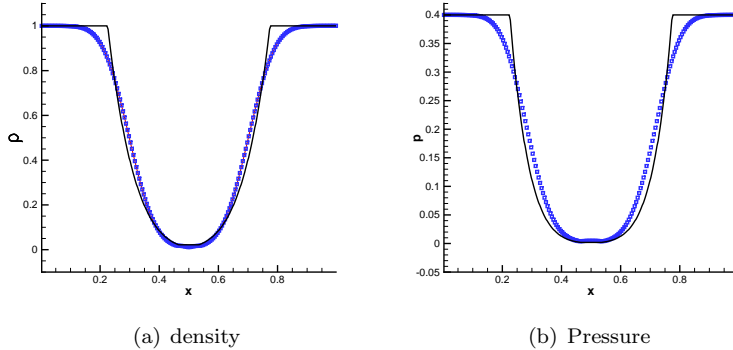


FIG. 5.5. Example 5.7 at $T = 0.1$, with $N = 200$ and $CFL = 2$. With the positivity-preserving limiter on. Solid line is the exact solution and blue squares are the numerical approximations.

537

538 *Example 5.8 (Blast wave).* The last example is the Sedov point-blast wave [18]. Ini-
 539 tially the gas is steady with uniform density one in the whole domain. The pressure is
 540 set to be $p = 10^{-9}$, except in the central cell, where the pressure is as high as $p = 10^4$.
 541 Then a blast-wave starts to propagate from the central cell with a shock front. This
 542 problem is difficult, since it involves a low density region and strong shocks. In Fig-
 543 ure 5.6 we present the numerical approximation and the exact solution [18] for the
 544 pressure and density respectively. The positivity-preserving limiter not only helps
 545 keep the low-density region positive, but also add robustness to the nonlinear solver.
 546 Without it, the code breaks down due to the failure of convergence of the nonlinear
 547 solver.

548 **6. Concluding remarks.** In this paper, we develop an implicit positivity-pres-
 549 erving DG method with high-order **spatial** accuracy for one-dimensional conservation
 550 laws. This work is an extension of the positivity-preserving limiter in [41, 42] for
 551 explicit schemes to implicit ones with backward Euler time discretization. To make
 552 the scheme positive, a lower bound for the CFL number is necessary. This conclusion
 553 is verified via both theoretical analysis and numerical experiments. The positivity-
 554 preserving limiter not only makes the numerical approximation physically meaningful
 555 but also brings robustness to the scheme and accelerates convergence towards the
 556 steady-state solution. The scheme also sees its success on the compressible Euler sys-
 557 tem. In the future, we have the following three directions to further explore. First,
 558 even though the result in this paper easily generalizes to multidimensional tensor
 559 product meshes and polynomial spaces, positivity-preserving DG schemes for multi-
 560 dimensional domain with unstructured mesh needs to be further developed. Second,
 561 we plan to generalize the proposed implicit positivity-preserving DG scheme to other
 562 types of equations such as the convection-diffusion equations. Thirdly, high-order im-
 563 plicit temporal discretizations need to be considered. Since there are no implicit SSP

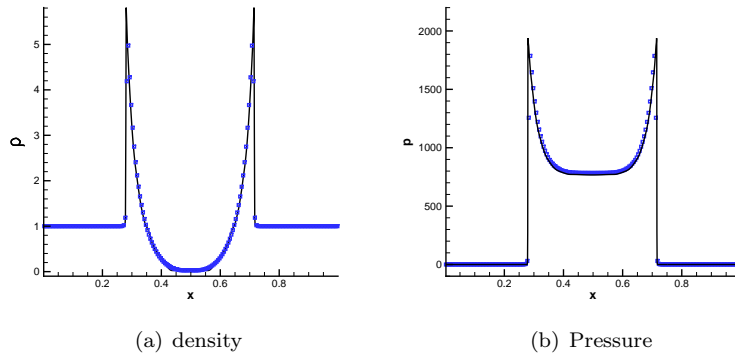


FIG. 5.6. Example 5.8 at $T = 0.003$, with $N = 200$ and $CFL = 2$. With the positivity-preserving limiter on. Solid line is the exact solution and blue squares are the numerical approximations.

564 methods with order greater than one [13], we need to turn to other types of implicit
 565 time discretizations, such as BDF methods and fully implicit RK methods.

REFERENCES

566

567 [1] H. L. ATKINS AND C.-W. SHU, *Quadrature-free implementation of discontinuous*
 568 *Galerkin method for hyperbolic equations*, AIAA J., 36 (1998), pp. 775–782,
 569 <https://doi.org/10.2514/2.436>.

570 [2] F. BASSI, L. BOTTI, A. COLOMBO, A. GHIDONI, AND F. MASSA, *Linearly implicit Rosenbrock-*
 571 *type Runge–Kutta schemes applied to the discontinuous Galerkin solution of compressible*
 572 *and incompressible unsteady flows*, Comput. Fluids, 118 (2015), pp. 305–320,
 573 <https://doi.org/10.1016/j.compfluid.2015.06.007>.

574 [3] F. BASSI, C. DE BARTOLO, R. HARTMANN, AND A. NIGRO, *A discontinuous Galerkin method*
 575 *for inviscid low Mach number flows*, J. Comput. Phys., 228 (2009), pp. 3996–4011,
 576 <https://doi.org/10.1016/j.jcp.2009.02.021>.

577 [4] P. BATTEN, M. A. LESCHZNER, AND U. C. GOLDBERG, *Average-state Jacobians and implicit*
 578 *methods for compressible viscous and turbulent flows*, J. Comput. Phys., 137 (1997), pp. 38–
 579 78, <https://doi.org/10.1006/jcph.1997.5793>.

580 [5] A. BERMAN AND R. PLEMMONS, *Nonnegative Matrices in the Mathematical Sciences*, Clas-
 581 *sics in Applied Mathematics*, Society for Industrial and Applied Mathematics, Jan. 1994,
 582 <https://doi.org/10.1137/1.9781611971262>.

583 [6] M. H. CARPENTER, C. A. KENNEDY, H. BIJL, S. A. VIKEN, AND V. N. VATSA, *Fourth-*
 584 *order Runge–Kutta schemes for fluid mechanics applications*, J. Sci. Comput., 25 (2005),
 585 pp. 157–194, <https://doi.org/10.1007/s10915-004-4637-3>.

586 [7] Z. CHEN, H. HUANG, AND J. YAN, *Third order maximum-principle-satisfying direct*
 587 *discontinuous Galerkin methods for time dependent convection diffusion equations*
 588 *on unstructured triangular meshes*, J. Comput. Phys., 308 (2016), pp. 198–217,
 589 <https://doi.org/10.1016/j.jcp.2015.12.039>.

590 [8] B. COCKBURN, S. HOU, AND C.-W. SHU, *The Runge-Kutta local projection*
 591 *discontinuous Galerkin finite element method for conservation laws.*
 592 *IV. The multidimensional case*, Math. Comp., 54 (1990), pp. 545–581,
 593 <https://doi.org/10.1090/S0025-5718-1990-1010597-0>.

594 [9] B. COCKBURN, S.-Y. LIN, AND C.-W. SHU, *TVB Runge-Kutta local projection discontinuous*
 595 *Galerkin finite element method for conservation laws III: One-dimensional systems*, J.
 596 *Comput. Phys.*, 84 (1989), pp. 90–113, [https://doi.org/10.1016/0021-9991\(89\)90183-6](https://doi.org/10.1016/0021-9991(89)90183-6).

597 [10] B. COCKBURN AND C.-W. SHU, *TVB Runge-Kutta local projection discontinuous Galerkin finite*
 598 *element method for conservation laws. II. General framework*, Math. Comp., 52 (1989),
 599 pp. 411–435, <https://doi.org/10.1090/S0025-5718-1989-0983311-4>.

600 [11] B. COCKBURN AND C.-W. SHU, *The Runge–Kutta discontinuous Galerkin method for conser-*

- 601 *vation laws V: Multidimensional systems*, J. Comput. Phys., 141 (1998), pp. 199–224,
 602 <https://doi.org/10.1006/jcph.1998.5892>.
- [12] V. DOLEJŠÍ AND M. FEISTAUER, *Discontinuous Galerkin Method - Analysis and Applications to Compressible Flow*, vol. 48 of Springer Series in Computational Mathematics, Springer International Publishing, 2015.
- [13] S. GOTTLIEB, C. SHU, AND E. TADMOR, *Strong stability-preserving high-order time discretization methods*, SIAM Rev., 43 (2001), pp. 89–112,
 607 <https://doi.org/10.1137/S003614450036757X>.
- [14] Y. HA, C. L. GARDNER, A. GELB, AND C.-W. SHU, *Numerical simulation of high Mach number astrophysical jets with radiative cooling*, J. Sci. Comput., 24 (2005), pp. 29–44,
 610 <https://doi.org/10.1007/s10915-004-4786-4>.
- [15] A. JAMESON, *Time dependent calculations using multigrid, with applications to unsteady flows past airfoils and wings*, in 10th Computational Fluid Dynamics Conference, Fluid Dynamics and Co-located Conferences, American Institute of Aeronautics and Astronautics, June 1991, <https://doi.org/10.2514/6.1991-1596>.
- [16] A. JAMESON, *Application of dual time stepping to fully implicit Runge Kutta schemes for unsteady flow calculations*, in 22nd AIAA Computational Fluid Dynamics Conference, AIAA AVIATION Forum, American Institute of Aeronautics and Astronautics, June 2015, <https://doi.org/10.2514/6.2015-2753>.
- [17] G. S. JIANG AND C.-W. SHU, *On a cell entropy inequality for discontinuous Galerkin methods*, Math. Comp., 62 (1994), pp. 531–538,
 622 <https://doi.org/10.1090/S0025-5718-1994-1223232-7>.
- [18] V. KOROBENIKOV, *Problems of Point Blast Theory*, AIP-Press, 1991.
- [19] Y.-T. LI AND R. WONG, *Integral and series representations of the Dirac delta function*, *Comm. Pure Appl. Math.*, 7 (2008), pp. 229–247, <https://doi.org/10.3934/cpaa.2008.7.229>.
- [20] X. LIU AND S. OSHER, *Nonoscillatory high order accurate self-similar maximum principle satisfying shock capturing schemes I*, SIAM J. Numer. Anal., 33 (1996), pp. 760–779,
 628 <https://doi.org/10.1137/0733038>.
- [21] A. MEISTER AND S. ORTLIEB, *On unconditionally positive implicit time integration for the DG scheme applied to shallow water flows*, Int. J. Numer. Meth. Fluids, 76 (2014), pp. 69–94,
 630 <https://doi.org/10.1002/fld.3921>.
- [22] Y. MORYOSSEF AND Y. LEVY, *Unconditionally positive implicit procedure for two-equation turbulence models: Application to $k-\omega$ turbulence models*, J. Comput. Phys., 220 (2006), pp. 88–108, <https://doi.org/10.1016/j.jcp.2006.05.001>.
- [23] Y. MORYOSSEF AND Y. LEVY, *Designing a positive second-order implicit time integration procedure for unsteady turbulent flows*, *Comput. Methods in Appl. Mech. Eng.*, 196 (2007), pp. 4196–4206, <https://doi.org/10.1016/j.cma.2007.04.001>.
- [24] Y. MORYOSSEF AND Y. LEVY, *The unconditionally positive-convergent implicit time integration scheme for two-equation turbulence models: Revisited*, *Comput. Fluids*, 38 (2009), pp. 1984–1994, <https://doi.org/10.1016/j.compfluid.2009.06.005>.
- [25] A. NIGRO, A. GHIDONI, S. REBAY, AND F. BASSI, *Modified extended BDF scheme for the discontinuous Galerkin solution of unsteady compressible flows*, Int. J. Numer. Meth. Fluids, 76 (2014), pp. 549–574, <https://doi.org/10.1002/fld.3944>.
- [26] S. V. PATANKAR, *Numerical Heat Transfer and Fluid Flow*, McGraw Hill, 1980.
- [27] W. PAZNER AND P.-O. PERSSON, *Stage-parallel fully implicit Runge–Kutta solvers for discontinuous Galerkin fluid simulations*, J. Comput. Phys., 335 (2017), pp. 700–717,
 647 <https://doi.org/10.1016/j.jcp.2017.01.050>.
- [28] P.-O. PERSSON, *Scalable parallel Newton-Krylov solvers for discontinuous Galerkin discretizations*, in 47th AIAA Aerospace Sciences Meeting Including The New Horizons Forum and Aerospace Exposition, Aerospace Sciences Meetings, American Institute of Aeronautics and Astronautics, Jan. 2009, <https://doi.org/10.2514/6.2009-606>.
- [29] P.-O. PERSSON AND J. PERAIRE, *An efficient low memory implicit DG algorithm for time dependent problems*, in 44th AIAA Aerospace Sciences Meeting and Exhibit, Aerospace Sciences Meetings, American Institute of Aeronautics and Astronautics, Jan. 2006, <https://doi.org/10.2514/6.2006-113>.
- [30] R. J. PLEMMONS, *M-matrix characterizations. I—nonsingular M-matrices*, *Linear Algebra and its Applications*, 18 (1977), pp. 175–188,
 658 [https://doi.org/10.1016/0024-3795\(77\)90073-8](https://doi.org/10.1016/0024-3795(77)90073-8).
- [31] T. QIN, C.-W. SHU, AND Y. YANG, *Bound-preserving discontinuous Galerkin methods for relativistic hydrodynamics*, J. Comput. Phys., 315 (2016), pp. 323–347,
 661 <https://doi.org/10.1016/j.jcp.2016.02.079>.
- [32] W. H. REED AND T. R. HILL, *Triangular mesh methods for the neutron transport equation*,

663 Los Alamos Scientific Laboratory Report LA-UR-73-479, Los Alamos, NM, 1973.

664 [33] C.-W. SHU, *Total-variation-diminishing time discretizations*, SIAM J. Sci. and Stat. Comput.,
 665 9 (1988), pp. 1073–1084, <https://doi.org/10.1137/0909073>.

666 [34] C.-W. SHU AND S. OSHER, *Efficient implementation of essentially non-oscillatory*
 667 *shock-capturing schemes*, J. Comput. Phys., 77 (1988), pp. 439–471,
 668 [https://doi.org/10.1016/0021-9991\(88\)90177-5](https://doi.org/10.1016/0021-9991(88)90177-5).

669 [35] G. SZEGÖ, *Orthogonal polynomials*, American Mathematical Society, New York, 1939.

670 [36] C. WANG, X. ZHANG, C.-W. SHU, AND J. NING, *Robust high order discontinuous Galerkin*
 671 *schemes for two-dimensional gaseous detonations*, J. Comput. Phys., 231 (2012), pp. 653–
 672 665, <https://doi.org/10.1016/j.jcp.2011.10.002>.

673 [37] Y. XING, X. ZHANG, AND C.-W. SHU, *Positivity-preserving high order well-balanced discontin-*
 674 *uous Galerkin methods for the shallow water equations*, Advances in Water Resources, 33
 675 (2010), pp. 1476–1493, <https://doi.org/10.1016/j.advwatres.2010.08.005>.

676 [38] Y. YANG, D. WEI, AND C.-W. SHU, *Discontinuous Galerkin method for Krause’s consen-*
 677 *sus models and pressureless Euler equations*, J. Comput. Phys., 252 (2013), pp. 109–127,
 678 <https://doi.org/10.1016/j.jcp.2013.06.015>.

679 [39] D. YUAN, J. CHENG, AND C.-W. SHU, *High order positivity-preserving discontinuous Galerkin*
 680 *methods for radiative transfer equations*, SIAM J. Sci. Comput., 38 (2016), pp. A2987–
 681 A3019, <https://doi.org/10.1137/16M1061072>.

682 [40] X. ZHANG, *On positivity-preserving high order discontinuous Galerkin schemes for com-*
 683 *pressible Navier–Stokes equations*, J. Comput. Phys., 328 (2017), pp. 301–343,
 684 <https://doi.org/10.1016/j.jcp.2016.10.002>.

685 [41] X. ZHANG AND C.-W. SHU, *On maximum-principle-satisfying high order schemes*
 686 *for scalar conservation laws*, J. Comput. Phys., 229 (2010), pp. 3091–3120,
 687 <https://doi.org/10.1016/j.jcp.2009.12.030>.

688 [42] X. ZHANG AND C.-W. SHU, *On positivity-preserving high order discontinuous Galerkin schemes*
 689 *for compressible Euler equations on rectangular meshes*, J. Comput. Phys., 229 (2010),
 690 pp. 8918–8934, <https://doi.org/10.1016/j.jcp.2010.08.016>.