ORIGINAL PAPER



High-Order Bound-Preserving Finite Difference Methods for Multispecies and Multireaction Detonations

Jie Du^{1,3} · Yang Yang²

Received: 11 September 2020 / Revised: 18 December 2020 / Accepted: 21 December 2020 © Shanghai University 2021

Abstract

In this paper, we apply high-order finite difference (FD) schemes for multispecies and multireaction detonations (MMD). In MMD, the density and pressure are positive and the mass fraction of the *i*th species in the chemical reaction, say z_i , is between 0 and 1, with $\sum z_i = 1$. Due to the lack of maximum-principle, most of the previous bound-preserving technique cannot be applied directly. To preserve those bounds, we will use the positivity-preserving technique to all the z's and enforce $\sum z_i = 1$ by constructing conservative schemes, thanks to conservative time integrations and consistent numerical fluxes in the system. Moreover, detonation is an extreme singular mode of flame propagation in premixed gas, and the model contains a significant stiff source. It is well known that for hyperbolic equations with stiff source, the transition points in the numerical approximations near the shocks may trigger spurious shock speed, leading to wrong shock position. Intuitively, the high-order weighted essentially non-oscillatory (WENO) scheme, which can suppress oscillations near the discontinuities, would be a good choice for spatial discretization. However, with the nonlinear weights, the numerical fluxes are no longer "consistent", leading to nonconservative numerical schemes and the bound-preserving technique does not work. Numerical experiments demonstrate that, without further numerical techniques such as subcell resolutions, the conservative FD method with linear weights can yield better numerical approximations than the nonconservative WENO scheme.

Keywords Weighted essentially non-oscillatory scheme \cdot Finite difference method \cdot Stiff source \cdot Detonations \cdot Bound-preserving \cdot Conservative

Mathematics Subject Classification 65M06 · 65M12

 Yang Yang yyang7@mtu.edu
 Jie Du jdu@tsinghua.edu.cn

¹ Yau Mathematical Sciences Center, Tsinghua University, Beijing 100084, China

² Department of Mathematical Sciences, Michigan Technological University, Houghton, MI 49931, USA

³ Yanqi Lake Beijing Institute of Mathematical Sciences and Applications, Beijing 101408, China

1 Introduction

In this paper, we construct high-order bound-preserving finite difference (FD) methods for stiff multispecies and multireaction chemical reactive flows. The governing equation in two space dimensions reads

$$\rho_t + m_x + n_y = 0, \tag{1a}$$

$$m_t + (mu + p)_x + (nu)_y = 0,$$
 (1b)

$$n_t + (mv)_x + (nv + p)_y = 0, (1c)$$

$$E_t + ((E+p)u)_x + ((E+p)v)_y = 0,$$
(1d)

$$(r_1)_t + (mz_1)_x + (nz_1)_y = s_1,$$
(1e)

:

$$(r_{M-1})_t + (mz_{M-1})_x + (nz_{M-1})_y = s_{M-1},$$
(1f)

where ρ , u, v, $m = \rho u$, $n = \rho v$, E and p are the total density, velocity in the x direction, velocity in the y direction, momentum in the x direction, momentum in the y direction, the total energy, and pressure, respectively. M is the total number of chemical species. For $1 \le i \le M$, $r_i = \rho z_i$ with z_i being the mass fraction for the *i*th species, and $\sum_{i=1}^{M} z_i = 1$. Therefore, we have

$$\sum_{i=1}^{M} r_i = \rho, \tag{2}$$

and $0 \le z_i \le 1$. We call (2) to be the property of total mass conservation. To close the system, we need one more equation of state given as

$$p = (\gamma - 1) \Big(E - \frac{1}{2} \rho (u^2 + v^2) - \rho z_1 q_1 - \dots - \rho z_M q_M \Big),$$

where q_i is the enthalpy of formation for the *i*th species and the temperature is given as $T = p/\rho$. The source term s_i describes the chemical reactions. We consider *R* reactions of the form:

$$\nu'_{1,r}X_1 + \nu'_{2,r}X_2 + \dots + \nu'_{M,r}X_M \to \nu''_{1,r}X_1 + \nu''_{2,r}X_2 + \dots + \nu''_{M,r}X_M, \quad r = 1, 2, \dots, R,$$

where $v'_{i,r}$ and $v''_{i,r}$ are the stoichiometric coefficients of the reactants and products, respective, of the *i*th species in the *r*th reaction. For non-equilibrium chemistry, the rate of production of the *i*th species can be written as

$$s_i = M_i \sum_{r=1}^{R} (v_{i,r}'' - v_{i,r}') \left[k_r(T) \prod_{j=1}^{M} \left(\frac{r_j}{M_j} \right)^{v_{j,r}'} \right], \quad i = 1, 2, \cdots, M,$$

where M_i is the molar mass of the *i*th species. The reaction rate $k_r(T)$ is given as

Deringer

$$k_r(T) = \begin{cases} B_r T^{\alpha_r}, \ T > T_r, \\ 0, \qquad T \leqslant T_r, \end{cases}$$

where T_r is the ignition temperature for the *r*th reaction, and B_r and α_r are pre-exponential factor and index of temperature, respectively. Moreover, it is easy to check that $\sum_{i=1}^{n} s_i = 0$. We can subtract (1e)–(1f) from (1a) and use the fact $\sum_{i=1}^{M} z_i = 1$ to obtain a fictitious

equation:

$$(r_M)_t + (mz_M)_x + (nz_M)_y = s_M,$$
(3)

which is similar to (1e)-(1f), and this can help us construct the bound-preserving (BP) technique.

Numerical simulation plays an significant role in minimizing hazards in gaseous detonation. In general, single-step models cannot be used to predict correct ignition process of the mixture, and detailed chemical models are commonly used to reproduce results that agree with the experimental data. However, the construction of accurate and efficient numerical methods is not an easy task due to the complexity of chemical kinetics. There are three main difficulties.

Firstly, in gaseous detonations, the density and pressure are positive, the mass fractions are between 0 and 1. It is well known that the exact solution contains shocks, and direct numerical simulations may be highly oscillatory near the shocks and send some variables out of their physical bounds. Physically irrelevant numerical approximations may not yield correct parameters used in the model and the numerical simulations may blow up. Therefore, special BP techniques are necessary in the numerical simulations. There are plenty of works discussing BP techniques for hyperbolic equations in the literature. The idea used in this paper was first introduced in [25], where parametrized maximum principle preserving flux limiters were applied to scalar hyperbolic conservation laws. Later, the extension to problems on unstructured meshes [5] and compressible Euler equations [24] were introduced. The main idea is to combine the high-order and low-order numerical fluxes together to obtain physically relevant numerical approximations. Though the positivity-preserving technique would work for the density and pressure, it cannot be applied to the mass fractions. In fact, the mass fractions do not satisfy maximum-principles and the positivity-preserving technique cannot be used to preserve the upper bound 1. In [6, 10, 26], one of the authors in this paper first introduced the BP discontinuous Galerkin (DG) methods to preserve the two bounds of volume fractions for multiphase flow in oil reservoir simulations, and the FD methods were also discussed in [11]. The basic idea is to use the positivity-preserving technique to each mass fractions and enforce the summation to be 1. The detailed idea is given as follows.

- Apply the positivity-preserving techniques to (1) to obtain positive ρ , p and z_i , i) $i=1,\cdots,M-1.$
- ii) Use the fictitious equation (3) to substitute (2) in the theoretical analysis. Notice that (3) is similar to (1e) and (1f). Therefore, the positivity-preserving technique also works for (3).
- iii) Enforce (2) is satisfied during time evolution.

Therefore, the key point is the construction of "conservative schemes", i.e., namely, (2) is satisfied at time level n + 1, provided it is satisfied at time level n. In [13, 14], the authors applied modified Patankar Runge-Kutta (RK) methods and constructed conservative schemes, extending the ideas in [16, 17]. However, the method cannot preserve the positivity of pressure. In [8, 9], the authors investigated DG methods for gaseous detonation and presented two sufficient conditions for conservative schemes: consistent numerical fluxes and conservative time integrations. We will follow this idea and extend them to FD methods. We emphasize that (3) is only used in the theoretical analysis of the BP technique. By constructing conservative schemes, the numerical solutions also satisfy (2). Hence, we can use (2) directly in the real computations and there is no additional computational cost.

Secondly, due to the rapid reaction rate, the model contains stiff source terms, see, e.g., [7, 18], leading to rather small time steps with explicit strong-stability-preserving (SSP) RK time integrations. In [8, 9], we constructed conservative sign-preserving exponential RK (ERK) time integrations, extending the idea in [12]. Numerical experiments demonstrated that compared with SSPRK, the ERK method yields better numerical approximations. In this paper, we will also use the ERK methods in all the numerical simulations.

Finally, in [18], the authors pointed out that direct numerical simulations on coarse meshes may yield nonphysical shock waves and incorrect shock positions. This is because the transition points near the shocks may trigger the source term, leading to spurious shock speeds. Some strategies to fix this problem can be found in the literature, see, e.g., the level set and front tracking methods [4, 22], subcell resolution [23, 27] and random projections [2, 3] as an incomplete list. The main idea is to remove significant transition points near the shocks. However, this is not the main topic in this paper, and we will only focus on the BP technique. Intuitively, one should apply special numerical techniques, such as weighted essentially non-oscillatory (WENO) schemes [1, 15, 19–21], to suppress oscillations. However, the nonlinear weights used in the WENO algorithm yield "nonconservative numerical schemes", and the bound-preserving techniques fail to work. Moreover, numerical approximations are better than the WENO scheme, indicating that the conservation is more important than oscillations suppression in capturing more accurate shock positions.

The organization of the rest of the paper is as follows. In Sect. 2, we consider the problems in one space dimension and demonstrate the FD methods as well as the WENO algorithm. In Sect. 3, we will consider Euler forward time integration and construct the BP technique for the convection term. In Sect. 4, we extend the formulations to two space dimensions. High-order time integrations will be given in Sect. 5. Numerical experiments in Sect. 6 will be given to compare the WENO algorithm and the proposed BP FD scheme. We will end in Sect. 7 with concluding remarks.

2 Finite Difference Methods in One Space Dimension

In this section, we concentrate on spatial discretizations and will leave the time variable t continuous. We will start with problems in one space dimension. The extension to two-dimensional problems will be discussed later in Sect. 4.

Notice that the exact solution of r_M satisfies both (2) and (3). As we discussed in the introduction, the fictitious equation (3) can help us construct the BP technique. Thus, we combine (1) with (3) and discuss FD methods for solving them. We rewrite the governing equations in one space dimension into the following vector form:

$$\mathbf{w}_t + \mathbf{f}(\mathbf{w})_x = \mathbf{s}(\mathbf{w}),\tag{4}$$

where

$$\mathbf{w} = (\rho, m, E, r_1, \cdots, r_M)^{\mathrm{T}}$$
$$\mathbf{f}(\mathbf{w}) = (m, mu + p, (E + p)u, mz_1, \cdots, mz_M)^{\mathrm{T}},$$
$$\mathbf{s}(\mathbf{w}) = (0, 0, 0, s_1, \cdots, s_M)^{\mathrm{T}}.$$

Notice that the system (4) is hyperbolic. The Jacobian f'(w) has M + 3 real eigenvalues:

$$\lambda_1(\mathbf{w}) \leq \lambda_2(\mathbf{w}) \leq \cdots \leq \lambda_{M+3}(\mathbf{w}),$$

and the corresponding complete set of independent eigenvectors:

$$r_1(\mathbf{w}), r_2(\mathbf{w}), \cdots, r_{M+3}(\mathbf{w}).$$

We denote

$$\Lambda(\mathbf{w}) = \operatorname{diag}(\lambda_1(\mathbf{w}), \cdots, \lambda_{M+3}(\mathbf{w})) \quad \text{and} \quad R(\mathbf{w}) = [r_1(\mathbf{w}), \cdots, r_{M+3}(\mathbf{w})],$$

then

$$R^{-1}(\mathbf{w})\mathbf{f}'(\mathbf{w})R(\mathbf{w}) = \Lambda(\mathbf{w}).$$

We will review the first order FD method and introduce the concept of consistent numerical fluxes in Sect. 2.1. Then, we will review the fifth-order WENO method in Sect. 2.2. The modification of the WENO method to get the high-order linear FD scheme will be discussed in Sect. 2.3.

2.1 First-Order Methods and Consistent Numerical Fluxes

Assume that the computational domain is [a, b]. We give a partition of the computational domain

$$a = x_{\frac{1}{2}} < x_{\frac{3}{2}} < \dots < x_{N+\frac{1}{2}} = b,$$

and denote the *i*th cell as

$$I_i = \left[x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}} \right], \quad i = 1, \cdots, N.$$

The center of the cell I_i is

$$x_i = \frac{1}{2} \left(x_{i-\frac{1}{2}} + x_{i+\frac{1}{2}} \right).$$

We assume that the grid is uniform and denote the cell length as

$$\Delta x = x_{i+\frac{1}{2}} - x_{i-\frac{1}{2}}.$$

We approximate $\mathbf{w}(x, t)$ at grid points $\{x_i, i = 1, \dots, N\}$ and denote the numerical solutions as

2 Springer

$$\mathbf{w}_i(t) = \mathbf{w}(x_i, t).$$

The following first order FD scheme is used to approximate the spatial derivatives in (4):

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbf{w}_{i}(t) = -\frac{1}{\Delta x} \left(\hat{\mathbf{f}}_{i+\frac{1}{2}} - \hat{\mathbf{f}}_{i-\frac{1}{2}} \right) + \mathbf{s}(\mathbf{w}_{i}), \tag{5}$$

where $\hat{\mathbf{f}}_{i+\frac{1}{2}}$ is the commonly used Lax-Friedrichs numerical flux

$$\hat{\mathbf{f}}_{i+\frac{1}{2}} = \frac{1}{2} \Big[\mathbf{f}(\mathbf{w}_i) + \mathbf{f}(\mathbf{w}_{i+1}) - \alpha_{i+\frac{1}{2}} (\mathbf{w}_{i+1} - \mathbf{w}_i) \Big]$$

with

$$\alpha_{i+\frac{1}{2}} = \max_{m=i,i+1} \max_{1 \le k \le M+3} |\lambda_k(\mathbf{w}_m)|$$

Next, we introduce the concept of consistent numerical fluxes, which will be used for constructing conservative schemes. We can check that there exists a constant vector

$$\mathbf{v} = [1, 0, 0, -1, \cdots, -1]^{\mathrm{T}} \in \mathbb{R}^{M+3},\tag{6}$$

such that the total mass conservation property (2) of the exact solution can be written in the following form:

$$\mathbf{w}\cdot\mathbf{v}=0$$

By using this property, we can also check that the exact flux vector in (4) satisfies

$$\mathbf{f} \cdot \mathbf{v} = u \left(\rho - \sum_{k=1}^{M} r_k \right) = u(\mathbf{w} \cdot \mathbf{v}) = 0.$$
(7)

If this property is also satisfied by the numerical flux vector, then we say the numerical flux is consistent.

Definition 1 Considering a semi-discrete scheme for solving the governing system (4), we say the numerical flux $\hat{\mathbf{f}}$ is consistent if $\hat{\mathbf{f}} \cdot \mathbf{v} = 0$, provided $\mathbf{w}_i \cdot \mathbf{v} = 0$ for all $1 \le i \le N$.

Let us consider the numerical flux in (5). By using the property (7) and the assumption that $\mathbf{w}_i \cdot \mathbf{v} = 0$ for all *i*, we can easily check that

$$\hat{\mathbf{f}}_{i+\frac{1}{2}} \cdot \mathbf{v} = \frac{1}{2} \Big[\mathbf{f}(\mathbf{w}_i) \cdot \mathbf{v} + \mathbf{f}(\mathbf{w}_{i+1}) \cdot \mathbf{v} - \alpha_{i+\frac{1}{2}} (\mathbf{w}_{i+1} \cdot \mathbf{v} - \mathbf{w}_i \cdot \mathbf{v}) \Big] = 0.$$
(8)

Hence, we conclude that the first order numerical flux $\hat{\mathbf{f}}_{i+\frac{1}{2}}$ is consistent.

2.2 Fifth Order WENO Methods

Now, we consider high order FD schemes. The only difference between the first order scheme and the high order scheme is the construction of numerical fluxes. We denote the numerical flux computed by the WENO method as $\hat{\mathbf{F}}_{i+\frac{1}{2}}^{\text{WENO}}$, then the fifth order WENO

semi-discrete scheme becomes

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbf{w}_{i}(t) = -\frac{1}{\Delta x} \left(\hat{\mathbf{F}}_{i+\frac{1}{2}}^{\mathrm{WENO}} - \hat{\mathbf{F}}_{i-\frac{1}{2}}^{\mathrm{WENO}} \right) + \mathbf{s}(\mathbf{w}_{i}).$$
(9)

In the following, we will review the FD flux splitting characteristicwise WENO procedure.

We consider each semi-grid point $i + \frac{1}{2}$. For the characteristicwise WENO, we first need to transform variables in the physical field to the local characteristic field. We compute an average state at location $x_{i+\frac{1}{2}}$ using the simple mean

$$\mathbf{w}_{i+\frac{1}{2}} = \frac{1}{2}(\mathbf{w}_i + \mathbf{w}_{i+1}).$$

Then, we "freeze" the matrices $R(\mathbf{w})$ and $R^{-1}(\mathbf{w})$ locally at $x_{i+\frac{1}{2}}$ by evaluating them with the average state at this point. We omit the subscript and still denote these matrices by R and R^{-1} :

$$R = R\left(\mathbf{w}_{i+\frac{1}{2}}\right), \quad R^{-1} = R^{-1}\left(\mathbf{w}_{i+\frac{1}{2}}\right).$$

Then, we can transform the variables in the neighboring points of $x_{i+\frac{1}{2}}$ to the local characteristic field using

$$\mathbf{v}_k = R^{-1}\mathbf{w}_k, \quad \mathbf{h}_k = R^{-1}\mathbf{f}(\mathbf{w}_k), \quad k = i - 2, \cdots, i + 3.$$

Next, we perform the scalar flux splitting WENO procedure component by component on the characteristic variables \mathbf{v}_k and \mathbf{h}_k . In the following, we focus on the ℓ th component with a fixed number ℓ . For simplicity of notations, we omit ℓ and use v_k and \mathbf{h}_k to denote the ℓ th components of \mathbf{v}_k and \mathbf{h}_k , respectively. We first perform the Lax-Friedrichs flux splitting:

$$h_k^{\pm} = \frac{1}{2} (h_k \pm \alpha_\ell v_k), \quad k = i - 2, \cdots, i + 3,$$

where

$$\alpha_{\ell} = \max_{1 \leq i \leq N} |\lambda_{\ell}(\mathbf{w}_i)|.$$

Then, we use $\{h_k^+, k = i - 2, \dots, i + 2\}$ to compute $\hat{h}_{i+\frac{1}{2}}^+$ and use $\{h_k^-, k = i - 1, \dots, i + 3\}$ to compute $\hat{h}_{i+\frac{1}{2}}^-$. The flux $\hat{h}_{i+\frac{1}{2}}^+$ is a nonlinear convex combination of the following third-order approximations:

$$\begin{split} h_{i+\frac{1}{2}}^{(0)} &= \frac{1}{3}h_i^+ + \frac{5}{6}h_{i+1}^+ - \frac{1}{6}h_{i+2}^+, \\ h_{i+\frac{1}{2}}^{(1)} &= -\frac{1}{6}h_{i-1}^+ + \frac{5}{6}h_i^+ + \frac{1}{3}h_{i+1}^+, \\ h_{i+\frac{1}{2}}^{(2)} &= \frac{1}{3}h_{i-2}^+ - \frac{7}{6}h_{i-1}^+ + \frac{11}{6}h_i^+. \end{split}$$

The nonlinear weights are defined as

$$\omega_r = \frac{a_r}{a_0 + a_1 + a_2}, \quad r = 0, 1, 2$$

with

$$a_r = \frac{d_r}{(\epsilon + \beta_r)^2}$$

Here, d_r , r = 0, 1, 2 are linear weights

$$d_0 = \frac{3}{10}, \quad d_1 = \frac{3}{5}, \quad d_2 = \frac{1}{10},$$

 $\epsilon > 0$ is introduced to avoid the denominator to become 0. We take $\epsilon = 10^{-6}$ in the numerical tests. β_r are the smooth indicators

$$\begin{split} \beta_0 &= \frac{13}{12} \left(h_i^+ - 2h_{i+1}^+ + h_{i+2}^+ \right)^2 + \frac{1}{4} \left(3h_{i+1}^+ - 4h_{i+1}^+ + h_{i+2}^+ \right)^2, \\ \beta_1 &= \frac{13}{12} \left(h_{i-1}^+ - 2h_i^+ + h_{i+1}^+ \right)^2 + \frac{1}{4} \left(h_{i-1}^+ - h_{i+1}^+ \right)^2, \\ \beta_2 &= \frac{13}{12} \left(h_{i-2}^+ - 2h_{i-1}^+ + h_i^+ \right)^2 + \frac{1}{4} \left(3h_{i-2}^+ - 4h_{i-1}^+ + h_i^+ \right)^2. \end{split}$$

Then, we can get

$$\hat{h}_{i+\frac{1}{2}}^{+} = \omega_0 h_{i+\frac{1}{2}}^{(0)} + \omega_1 h_{i+\frac{1}{2}}^{(1)} + \omega_2 h_{i+\frac{1}{2}}^{(2)}.$$

The flux $\hat{h}_{i+\frac{1}{2}}^{-}$ can be obtained through a mirror symmetric procedure with respect to $x_{i+\frac{1}{2}}$. For simplicity, we omit the details. Then, the numerical flux $\hat{h}_{i+\frac{1}{2}}^{\ell}$ for the ℓ th component is obtained by

$$\hat{h}^{\ell}_{i+\frac{1}{2}} = \hat{h}^{+}_{i+\frac{1}{2}} + \hat{h}^{-}_{i+\frac{1}{2}}.$$

We repeat the above scalar WENO procedure for each component ($\ell = 1, \dots, M + 3$) of the characteristic variables and can obtain the numerical flux vector in the characteristic field

$$\hat{\mathbf{h}}_{i+\frac{1}{2}} = \left[\hat{h}_{i+\frac{1}{2}}^{1}, \hat{h}_{i+\frac{1}{2}}^{2}, \cdots, \hat{h}_{i+\frac{1}{2}}^{M+3}\right]^{\mathrm{T}}.$$

Finally, we transform back to the physical space by

$$\hat{\mathbf{F}}_{i+\frac{1}{2}}^{\text{WENO}} = R\hat{\mathbf{h}}_{i+\frac{1}{2}}$$

Notice that the weights $\{\omega_0, \omega_1, \omega_2\}$ in the WENO scheme are nonlinear functions of the variables and they can be different for different components of the characteristic vectors. Thus, $\hat{\mathbf{F}}_{i+\frac{1}{2}}^{\text{WENO}}$ is not consistent.

2.3 Fifth Order Linear Method

Since the numerical flux computed by the WENO procedure is not consistent, we modify the WENO method and consider a simple component-wise linear high order FD scheme. Instead of transforming variables to the characteristic field, we perform the flux splitting directly in the physical space:

$$\mathbf{f}_i^{\pm} = \frac{1}{2} \big[\mathbf{f}(\mathbf{w}_i) \pm \alpha \mathbf{w}_i \big], \quad i = -2, \cdots, N+3,$$

where

$$\alpha = \max_{1 \le i \le N} \max_{1 \le k \le M+3} |\lambda_k(\mathbf{w}_i)|.$$
(10)

Similar to the WENO method, we use { \mathbf{f}_{k}^{+} , $k = i - 2, \dots, i + 2$ } to compute $\hat{\mathbf{f}}_{i+\frac{1}{2}}^{+}$ and use { \mathbf{f}_{k}^{-} , $k = i - 1, \dots, i + 3$ } to compute $\hat{\mathbf{f}}_{i+\frac{1}{2}}^{-}$. For computing $\hat{\mathbf{f}}_{i+\frac{1}{2}}^{+}$, we first get third order approximations to the numerical flux

$$\mathbf{f}_{i+\frac{1}{2}}^{(0)} = \frac{1}{3}\mathbf{f}_{i}^{+} + \frac{5}{6}\mathbf{f}_{i+1}^{+} - \frac{1}{6}\mathbf{f}_{i+2}^{+},$$

$$\mathbf{f}_{i+\frac{1}{2}}^{(1)} = -\frac{1}{6}\mathbf{f}_{i-1}^{+} + \frac{5}{6}\mathbf{f}_{i}^{+} + \frac{1}{3}\mathbf{f}_{i+1}^{+},$$

$$\mathbf{f}_{i+\frac{1}{2}}^{(2)} = \frac{1}{3}\mathbf{f}_{i-2}^{+} - \frac{7}{6}\mathbf{f}_{i-1}^{+} + \frac{11}{6}\mathbf{f}_{i}^{+}.$$

Then, we combine them by using linear weights

$$\hat{\mathbf{f}}_{i+\frac{1}{2}}^{+} = d_0 \mathbf{f}_{i+\frac{1}{2}}^{(0)} + d_1 \mathbf{f}_{i+\frac{1}{2}}^{(1)} + d_2 \mathbf{f}_{i+\frac{1}{2}}^{(2)} = \frac{1}{30} \mathbf{f}_{i-2}^{+} - \frac{13}{60} \mathbf{f}_{i-1}^{+} + \frac{47}{60} \mathbf{f}_{i}^{+} + \frac{9}{20} \mathbf{f}_{i+1}^{+} - \frac{1}{20} \mathbf{f}_{i+2}^{+} .$$

By a mirror symmetric procedure with respect to $x_{i+\frac{1}{2}}$, we can compute the flux $\hat{\mathbf{f}}_{i+\frac{1}{2}}^{-}$ through

$$\hat{\mathbf{f}}_{i+\frac{1}{2}}^{-} = \frac{1}{30}\mathbf{f}_{i+3}^{-} - \frac{13}{60}\mathbf{f}_{i+2}^{-} + \frac{47}{60}\mathbf{f}_{i+1}^{-} + \frac{9}{20}\mathbf{f}_{i}^{-} - \frac{1}{20}\mathbf{f}_{i-1}^{-}.$$

Finally, we get the numerical flux at $x_{i+\frac{1}{2}}$ by

$$\hat{\mathbf{F}}_{i+\frac{1}{2}} = \hat{\mathbf{f}}_{i+\frac{1}{2}}^{+} + \hat{\mathbf{f}}_{i+\frac{1}{2}}^{-}.$$
(11)

Then, the fifth order linear FD scheme becomes

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbf{w}_i(t) = -\frac{1}{\Delta x} \left(\hat{\mathbf{F}}_{i+\frac{1}{2}} - \hat{\mathbf{F}}_{i-\frac{1}{2}} \right) + \mathbf{s}(\mathbf{w}_i).$$
(12)

One can easily check that

$$\hat{\mathbf{F}}_{i+\frac{1}{2}} \cdot \mathbf{v} = 0, \tag{13}$$

and hence the fifth order linear flux $\hat{\mathbf{F}}_{i+\frac{1}{2}}$ is consistent.

3 Bound-Preserving Technique for the Convection Term in One Space Dimension

In the previous section, we only considered spatial discretizations. Now we proceed to consider time integrations and get the fully discrete schemes. In the following, we use

$$\mathbf{w}_{i}^{n} = \left(\rho_{i}^{n}, m_{i}^{n}, E_{i}^{n}, (r_{1})_{i}^{n}, \cdots, (r_{M})_{i}^{n}\right)^{1}$$

to denote the numerical approximation to $\mathbf{w}_i(t)$ at time level *n*. We first emphasize the importance of "conservative schemes" and present sufficient conditions for constructing conservative schemes in Sect. 3.1. Then, we will consider the Euler forward time integration and the BP technique for the convection term in Sect. 3.2.

3.1 Conservative Schemes

For the physical solution, the total density ρ and the pressure *p* should be non-negative. In addition, all mass fractions z_k ($k = 1, \dots, M$) should be between 0 and 1. As discussed in the introduction, it is not easy to preserve the upper bound 1 for each z_k directly. Instead, we let $z_k \ge 0$, $k = 1, \dots, M$ and enforce the condition $\sum_{k=1}^{M} z_k = 1$, so that we will have $z_k \le 1, k = 1, \dots, M$. Hence, we define the admissible set of solutions as

$$G = \left\{ \mathbf{w} = \begin{pmatrix} \rho \\ m \\ E \\ r_1 \\ \vdots \\ r_M \end{pmatrix}, \rho \ge 0, p \ge 0, z_k \ge 0, \ k = 1, \cdots, M, \sum_{k=1}^M z_k = 1 \right\}.$$

It is easy to check that *G* is a convex set as *p* is a concave function of **w** [8, 9]. We aim to get numerical solutions that lie in the convex set *G*. We can see that the preservation of $\sum_{k=1}^{M} z_k = 1$ (or equivalently $\sum_{k=1}^{M} r_k = \rho$) numerically in each time level is important. Hence, we introduce the following definition of conservative schemes.

Definition 2 Consider a fully discrete numerical scheme for solving the detonation problem (4), if the mass conservation property can be maintained during the time evolution, namely, the numerical solution satisfies

$$\sum_{k=1}^{M} (r_k)_i^{n+1} = \rho_i^{n+1} \tag{14}$$

as long as $\sum_{k=1}^{M} (r_k)_i^n = \rho_i^n$, then we say the numerical scheme is a conservative scheme.

We rewrite the semi-discrete linear schemes (5) and (12) into the following ODE system:

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbf{w}_i(t) = \mathbf{F}_i + \mathbf{s}(\mathbf{w}_i),\tag{15}$$

where

$$\mathbf{F}_{i} = \begin{cases} -\frac{1}{\Delta x} \left(\hat{\mathbf{f}}_{i+\frac{1}{2}} - \hat{\mathbf{f}}_{i-\frac{1}{2}} \right), & \text{first order FD,} \\ -\frac{1}{\Delta x} \left(\hat{\mathbf{F}}_{i+\frac{1}{2}} - \hat{\mathbf{F}}_{i-\frac{1}{2}} \right), & \text{fifth order FD.} \end{cases}$$

Thanks to the consistent properties (8) and (13) of the numerical fluxes, we have

 $\mathbf{F}_i \cdot \mathbf{v} = 0.$

Also, we have $\mathbf{s} \cdot \mathbf{v} = 0$ for the detonation problem. Taking dot product with \mathbf{v} on both sides of (15), we know that the value of $\mathbf{w}_i(t) \cdot \mathbf{v}$ should remain unchanged. If this property can be preserved numerically, we say the time integration is conservative.

Definition 3 Consider the ODE system (15) with the assumption that there exists a constant vector **v** such that $\mathbf{F}_i \cdot \mathbf{v} = \mathbf{s} \cdot \mathbf{v} = 0$. A time integration for solving this ODE system is conservative if the numerical solutions satisfy

$$\mathbf{w}_i^{n+1} \cdot \mathbf{v} = \mathbf{w}_i^n \cdot \mathbf{v}.$$

A simple example of conservative time integrations is the Euler forward method:

$$\mathbf{w}_i^{n+1} = \mathbf{w}_i^n + \Delta t \mathbf{F}_i + \Delta t \mathbf{s}(\mathbf{w}_i).$$

Taking dot product with **v** on both sides, we can get $\mathbf{w}_i^{n+1} \cdot \mathbf{v} = \mathbf{w}_i^n \cdot \mathbf{v}$. Notice that $\sum_{k=1}^{M} (r_k)_i^n = \rho_i^n$ is just equivalent to $\mathbf{w}_i^n \cdot \mathbf{v} = 0$. Hence, it is straightforward to get the following theorem about the sufficient conditions for constructing conservative schemes.

Theorem 1 Considering the detonation problems, if we apply consistent numerical fluxes in space and conservative time integrations in time, then the fully discrete scheme is conservative.

Recall that the numerical fluxes computed by using the WENO procedure are not consistent. Hence, we can not get conservative numerical schemes, and the bound-preserving techniques fail to work. In the following section, we only consider the linear FD schemes which yield consistent numerical fluxes.

3.2 BP Technique for the Convection Term

Now, we consider the Euler forward time integration which is conservative. We only show the BP technique for the convection term and assume that s = 0. The source term will be discussed in Sect. 5 by using conservative ERK methods.

We apply the Euler forward time integration to the first order semi-discrete scheme (5) and can get the following fully discrete scheme:

$$\mathbf{w}_i^{n+1} = \mathbf{w}_i^n - \frac{\Delta t}{\Delta x} \Big(\hat{\mathbf{f}}_{i+\frac{1}{2}} - \hat{\mathbf{f}}_{i-\frac{1}{2}} \Big),$$

where all numerical fluxes are computed at time level *n*. Assume that $\mathbf{w}_i^n \in G$ for all *i*. Following Remark 2.4 in [28], it is easy to check that $\mathbf{w}_i^{n+1} \in G$ under the CFL condition $\Delta t \leq \Delta \tilde{t}$, where $\Delta \tilde{t}$ satisfies

$$\alpha \frac{\Delta \tilde{t}}{\Delta x} \leqslant 1. \tag{16}$$

Here, α has been defined in (10).

Next, we consider the fifth order linear scheme (12). By applying the Euler forward time integration, we can get

$$\mathbf{w}_i^{n+1} = \mathbf{w}_i^n - \frac{\Delta t}{\Delta x} \Big(\hat{\mathbf{F}}_{i+\frac{1}{2}} - \hat{\mathbf{F}}_{i-\frac{1}{2}} \Big),$$

where all numerical fluxes are again computed by using $\{\mathbf{w}_i^n, i = 1 \cdots, N\}$. Although the scheme provides high order approximations, the numerical solution may be out of the physical bounds. Since the first order fluxes are bound-preserving, we apply the parameterized positivity-preserving flux limiter [25] to modify the high order fluxes $\hat{\mathbf{F}}_{i+\frac{1}{2}}$ towards the first order fluxes $\hat{\mathbf{f}}_{i+\frac{1}{2}}$. The modified scheme becomes

$$\mathbf{w}_{i}^{n+1} = \mathbf{w}_{i}^{n} - \frac{\Delta t}{\Delta x} \Big(\tilde{\mathbf{F}}_{i+\frac{1}{2}} - \tilde{\mathbf{F}}_{i-\frac{1}{2}} \Big), \tag{17}$$

where $\tilde{\mathbf{F}}_{i+\frac{1}{2}}$ is a combination of the high order flux and the first order flux:

$$\tilde{\mathbf{F}}_{i+\frac{1}{2}} = \hat{\mathbf{f}}_{i+\frac{1}{2}} + \theta_{i+\frac{1}{2}} \Big(\hat{\mathbf{F}}_{i+\frac{1}{2}} - \hat{\mathbf{f}}_{i+\frac{1}{2}} \Big).$$

We denote the first order flux, original fifth order flux and the modified fifth order flux as

$$\begin{split} \hat{\mathbf{f}}_{i+\frac{1}{2}} &= \left(\hat{f}_{i+\frac{1}{2}}^{\rho}, \hat{f}_{i+\frac{1}{2}}^{m}, \hat{f}_{i+\frac{1}{2}}^{E}, \hat{f}_{i+\frac{1}{2}}^{1}, \cdots, \hat{f}_{i+\frac{1}{2}}^{M}\right)^{\mathrm{T}}, \\ \hat{\mathbf{F}}_{i+\frac{1}{2}} &= \left(\hat{F}_{i+\frac{1}{2}}^{\rho}, \hat{F}_{i+\frac{1}{2}}^{m}, \hat{F}_{i+\frac{1}{2}}^{E}, \hat{F}_{i+\frac{1}{2}}^{1}, \cdots, \hat{F}_{i+\frac{1}{2}}^{M}\right)^{\mathrm{T}}, \\ \tilde{\mathbf{F}}_{i+\frac{1}{2}} &= \left(\tilde{F}_{i+\frac{1}{2}}^{\rho}, \tilde{F}_{i+\frac{1}{2}}^{m}, \tilde{F}_{i+\frac{1}{2}}^{E}, \tilde{F}_{i+\frac{1}{2}}^{1}, \cdots, \tilde{F}_{i+\frac{1}{2}}^{M}\right)^{\mathrm{T}}, \end{split}$$

respectively. Assuming that $\mathbf{w}_i^n \in G$ for all *i*, we aim to choose suitable parameter $\theta_{i+\frac{1}{2}} \in [0,1]$ such that the solution of (17) satisfies $\mathbf{w}_i^{n+1} \in G$.

We first show that the modified scheme (17) is still conservative. One can check that

$$\tilde{\mathbf{F}}_{i+\frac{1}{2}} \cdot \mathbf{v} = \hat{\mathbf{f}}_{i+\frac{1}{2}} \cdot \mathbf{v} + \theta_{i+\frac{1}{2}} \left(\hat{\mathbf{F}}_{i+\frac{1}{2}} - \hat{\mathbf{f}}_{i+\frac{1}{2}} \right) \cdot \mathbf{v} = 0$$

using the consistent properties (8) and (13), which means that the modified high order numerical flux vector $\tilde{\mathbf{F}}_{i+\frac{1}{2}}$ is also consistent. Moreover, the Euler forward time integration is conservative. By using Theorem 1, we know that the scheme (17) is conservative.

Now, we try to preserve the positivity of each mass fraction z_k ($k = 1, \dots, M$). Equivalently, we need to enforce $\rho_i^{n+1} \ge 0$ and $(r_k)_i^{n+1} \ge 0$, $k = 1, \dots, M$. Notice that our scheme is

conservative and so that $\rho_i^{n+1} = \sum_{i=1}^{M} (r_k)_i^{n+1}$. Hence, we only need to consider $(r_k)_i^{n+1} \ge 0, k = 1, \dots, M$. In the scheme (17), the equation for solving r_k $(k = 1, \dots, M)$ is

$$(r_k)_i^{n+1} = (r_k)_i^n - \frac{\Delta t}{\Delta x} \left(\tilde{F}_{i+\frac{1}{2}}^k - \tilde{F}_{i-\frac{1}{2}}^k \right), \tag{18}$$

where

$$\tilde{F}_{i+\frac{1}{2}}^{k} = \hat{f}_{i+\frac{1}{2}}^{k} + \theta_{i+\frac{1}{2}}F_{i+\frac{1}{2}}^{k}, \quad F_{i+\frac{1}{2}}^{k} := \hat{F}_{i+\frac{1}{2}}^{k} - \hat{f}_{i+\frac{1}{2}}^{k}.$$

We need to find a local pair $\left(\Lambda_{-\frac{1}{2},I_{i}}^{k},\Lambda_{+\frac{1}{2},I_{i}}^{k}\right)$, such that for any pair $\left(\theta_{i-\frac{1}{2}},\theta_{i+\frac{1}{2}}\right) \in \left[0,\Lambda_{-\frac{1}{2},I_{i}}^{k}\right] \times \left[0,\Lambda_{+\frac{1}{2},I_{i}}^{k}\right]$, the solution computed by (18) satisfies $(r_{k})_{i}^{n+1} \ge 0$, i.e.,

$$\lambda \theta_{i-\frac{1}{2}} F_{i-\frac{1}{2}}^{k} - \lambda \theta_{i+\frac{1}{2}} F_{i+\frac{1}{2}}^{k} + \Gamma_{i} \ge 0,$$
(19)

where $\lambda := \frac{\Delta t}{\Delta x}$ and

Then,

$$\Gamma_i = (r_k)_i^n - \lambda \left(\hat{f}_{i+\frac{1}{2}}^k - \hat{f}_{i-\frac{1}{2}}^k \right) \ge 0$$

is the solution computed by using the first order scheme. Following [25], we adopt the following algorithm.

i) If
$$F_{i-\frac{1}{2}}^{k} \ge 0$$
 and $F_{i+\frac{1}{2}}^{k} \le 0$, take $\left(\Lambda_{-\frac{1}{2},I_{i}}^{k}, \Lambda_{+\frac{1}{2},I_{i}}^{k}\right) = (1,1)$.

ii) If
$$F_{i-\frac{1}{2}}^{k} \ge 0$$
 and $F_{i+\frac{1}{2}}^{k} > 0$, take $\left(\Lambda_{-\frac{1}{2},I_{i}}^{k}, \Lambda_{+\frac{1}{2},I_{i}}^{k}\right) = \left(1, \min\left\{1, \frac{\Gamma_{i}}{\lambda F_{i+\frac{1}{2}}^{k} + \epsilon}\right\}\right)$

iii) If
$$F_{i-\frac{1}{2}}^{k} < 0$$
 and $F_{i+\frac{1}{2}}^{k} \leq 0$, take $\left(\Lambda_{-\frac{1}{2},I_{i}}^{k},\Lambda_{+\frac{1}{2},I_{i}}^{k}\right) = \left(\min\left\{1,\frac{-\Gamma_{i}}{\lambda_{F_{i-\frac{1}{2}}}^{k}-\epsilon}\right\},1\right)$
iv) If $F_{i-\frac{1}{2}}^{k} < 0$ and $F_{i+\frac{1}{2}}^{k} > 0$,

(a) when (19) holds with
$$\left(\theta_{i-\frac{1}{2}}, \theta_{i+\frac{1}{2}}\right) = (1,1)$$
, take $\left(\Lambda_{-\frac{1}{2},I_{i}}^{k}, \Lambda_{+\frac{1}{2},I_{i}}^{k}\right) = (1,1)$;
(b) otherwise, take $\Lambda_{-\frac{1}{2},I_{i}}^{k} = \Lambda_{+\frac{1}{2},I_{i}}^{k} = \frac{-\Gamma_{i}}{\lambda F_{i-\frac{1}{2}}^{k} - \lambda F_{i+\frac{1}{2}}^{k} - \epsilon}$.

In this algorithm, $\epsilon = 10^{-13}$ is introduced to avoid the denominator being 0. Based on the above discussion, we let

$$\Lambda_{-\frac{1}{2},I_{i}}^{\rho} := \min_{k=1,\dots,M} \Lambda_{-\frac{1}{2},I_{i}}^{k}, \qquad \Lambda_{+\frac{1}{2},I_{i}}^{\rho} := \min_{k=1,\dots,M} \Lambda_{+\frac{1}{2},I_{i}}^{k}.$$

for any pair $\left(\theta_{i-\frac{1}{2}}, \theta_{i+\frac{1}{2}}\right) \in \left[0, \Lambda_{-\frac{1}{2},I_{i}}^{\rho}\right] \times \left[0, \Lambda_{+\frac{1}{2},I_{i}}^{\rho}\right]$, the solution of (17) satisfies
 $(z_{k})_{i}^{n+1} \ge 0, \ k = 1, \dots, M, \qquad \sum_{k=1}^{M} (z_{k})_{i}^{n+1} = 1, \quad \rho_{i}^{n+1} \ge 0.$

Deringer

Next, we preserve the positivity of the pressure. Recall that the pressure p is a function of the solution \mathbf{w}_{i}^{n+1} :

$$p(\mathbf{w}_i^{n+1}) = (\gamma - 1) \left(E_i^{n+1} - \frac{1}{2} \frac{(m_i^{n+1})^2}{\rho_i^{n+1}} - (r_1)_i^{n+1} q_1 - \dots - (r_M)_i^{n+1} q_M \right).$$

Since \mathbf{w}_i^{n+1} computed by (17) depends on the pair $A := \left(\theta_{i-\frac{1}{2}}, \theta_{i+\frac{1}{2}}\right)$, we also denote

$$p(A) = p\left(\theta_{i-\frac{1}{2}}, \theta_{i+\frac{1}{2}}\right) := p\left(\mathbf{w}_{i}^{n+1}\left(\theta_{i-\frac{1}{2}}, \theta_{i+\frac{1}{2}}\right)\right).$$

Following [24], we perform the following algorithm.

- Denote $A^1 = \left(0, \Lambda_{+\frac{1}{2}, I_i}^{\rho}\right), A^2 = \left(\Lambda_{-\frac{1}{2}, I_i}^{\rho}, 0\right)$ and $A^3 = \left(\Lambda_{-\frac{1}{2}, I_i}^{\rho}, \Lambda_{+\frac{1}{2}, I_i}^{\rho}\right)$. For k = 1, 2, 3, if $p(A^k) \ge 0$, let $B^k = A^k$. Otherwise, find r such that $p(rA^k) \ge 0$ and let $B^k = rA^k$. We i) also denote $B^k = (B_1^k, B_2^k)$ for k = 1, 2, 3. ii) Let $\Lambda_{-\frac{1}{2}, I_i} = \min(B_1^2, B_1^3)$ and $\Lambda_{+\frac{1}{2}, I_i} = \min(B_2^1, B_2^3)$.

Finally, we take $\theta_{i+\frac{1}{2}} = \min\left(\Lambda_{+\frac{1}{2},I_i}, \Lambda_{-\frac{1}{2},I_{i+1}}\right)$ for each *i*. Then, the solution of the scheme (17) satisfies $\mathbf{w}_{i}^{n+1} \in G$.

4 Problems in Two Space Dimensions

We consider the problem in two space dimensions:

$$\mathbf{w}_t + \mathbf{f}(\mathbf{w})_x + \mathbf{g}(\mathbf{w})_y = \mathbf{s}(\mathbf{w}), \tag{20}$$

where

$$\mathbf{w} = \left(\rho, m, n, E, r_1, \cdots, r_M\right)^{\mathrm{T}},$$

$$\mathbf{f}(\mathbf{w}) = \left(m, mu + p, mv, (E + p)u, mz_1, \cdots, mz_M\right)^{\mathrm{T}},$$

$$\mathbf{g}(\mathbf{w}) = \left(n, nu, nv + p, (E + p)v, nz_1, \cdots, nz_M\right)^{\mathrm{T}},$$

$$\mathbf{s}(\mathbf{w}) = \left(0, 0, 0, 0, s_1, \cdots, s_M\right)^{\mathrm{T}}.$$

Recall that we introduced a constant vector \mathbf{v} in the problem in one space dimension. For solving the problem in two space dimensions, v becomes

$$\mathbf{v} = [1, 0, 0, 0, -1, \cdots, -1]^{\mathrm{T}} \in \mathbb{R}^{M+4}.$$
(21)

We can check that the exact solution of (20) satisfies

$$\mathbf{w} \cdot \mathbf{v} = \mathbf{f} \cdot \mathbf{v} = \mathbf{g} \cdot \mathbf{v} = \mathbf{s} \cdot \mathbf{v} = 0.$$

We will first give FD methods in Sect. 4.1 and then show the BP technique in Sect. 4.2.

4.1 Finite Difference Methods

Assume that the computational domain is $[a, b] \times [c, d]$. We will consider Cartesian grids, that is, the domain is covered by rectangular cells:

$$I_{i,j} = \left[x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}} \right] \times \left[y_{j-\frac{1}{2}}, y_{j+\frac{1}{2}} \right], \quad i = 1, \cdots, N_x, \quad j = 1, \cdots, N_y,$$

where

$$a = x_{\frac{1}{2}} < x_{\frac{3}{2}} < \dots < x_{N_x + \frac{1}{2}} = b_x$$

and

$$c = y_{\frac{1}{2}} < y_{\frac{3}{2}} < \dots < y_{N_y + \frac{1}{2}} = d.$$

The center of the cell $I_{i,j}$ is (x_i, y_j) with

$$x_i = \frac{1}{2} \left(x_{i-\frac{1}{2}} + x_{i+\frac{1}{2}} \right), \qquad y_j = \frac{1}{2} \left(y_{j-\frac{1}{2}} + y_{j+\frac{1}{2}} \right).$$

We assume the grid is uniform and denote the cell lengths in x and y directions as

$$\Delta x = x_{i+\frac{1}{2}} - x_{i-\frac{1}{2}}, \qquad \Delta y = y_{j+\frac{1}{2}} - y_{j-\frac{1}{2}},$$

respectively. We approximate the solution **w** at grid points (x_i, y_j) and denote the numerical solutions as

$$\mathbf{w}_{i,j}(t) = \left(\rho_{i,j}(t), m_{i,j}(t), n_{i,j}(t), E_{i,j}(t), (r_1)_{i,j}(t), \cdots, (r_M)_{i,j}(t)\right)^{\mathrm{T}}$$

The following first order finite difference scheme is used to approximate the spatial derivatives:

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbf{w}_{i,j}(t) = -\frac{1}{\Delta x} \left(\hat{\mathbf{f}}_{i+\frac{1}{2},j} - \hat{\mathbf{f}}_{i-\frac{1}{2},j} \right) - \frac{1}{\Delta y} \left(\hat{\mathbf{g}}_{i,j+\frac{1}{2}} - \hat{\mathbf{g}}_{i,j-\frac{1}{2}} \right) + \mathbf{s}(\mathbf{w}_{i,j}).$$

Here $\hat{\mathbf{f}}_{i+\frac{1}{2},j}$ and $\hat{\mathbf{g}}_{i,j+\frac{1}{2}}$ are first order numerical fluxes which can be computed dimension by dimension. For each fixed *j*, the numerical flux $\hat{\mathbf{f}}_{i+\frac{1}{2},j}$ can be obtained in the *x* direction by using the one-dimensional algorithm in Sect. 2. Likewise, the numerical flux $\hat{\mathbf{g}}_{i,j+\frac{1}{2}}$ is obtained in the *y* direction with *i* fixed. We adopt the commonly used Lax-Friedrichs flux:

$$\begin{split} \hat{\mathbf{f}}_{i+\frac{1}{2},j} &= \frac{1}{2} \Big[\mathbf{f}(\mathbf{w}_{i,j}) + \mathbf{f}(\mathbf{w}_{i+1,j}) - \alpha_{i+\frac{1}{2}}(\mathbf{w}_{i+1,j} - \mathbf{w}_{i,j}) \Big], \\ \hat{\mathbf{g}}_{i,j+\frac{1}{2}} &= \frac{1}{2} \Big[\mathbf{g}(\mathbf{w}_{i,j}) + \mathbf{g}(\mathbf{w}_{i,j+1}) - \beta_{j+\frac{1}{2}}(\mathbf{w}_{i,j+1} - \mathbf{w}_{i,j}) \Big], \end{split}$$

where $\alpha_{i+\frac{1}{2}}$ is the maximum absolute values of the eigenvalues of the Jacobian $\mathbf{f}'(\mathbf{w})$ computed over $\mathbf{w}_{i+1,j}$ and $\mathbf{w}_{i,j}$, $\beta_{j+\frac{1}{2}}$ is the maximum absolute values of the eigenvalues of the $\mathbf{g}'(\mathbf{w})$ computed over $\mathbf{w}_{i,j}$ and $\mathbf{w}_{i,j+1}$. Following the same line, we can extend the WENO FD and fifth-order linear FD scheme to the problem in two space dimensions:

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbf{w}_{i,j}(t) = -\frac{1}{\Delta x} \left(\hat{\mathbf{F}}_{i+\frac{1}{2},j}^{\mathrm{WENO}} - \hat{\mathbf{F}}_{i-\frac{1}{2},j}^{\mathrm{WENO}} \right) - \frac{1}{\Delta y} \left(\hat{\mathbf{G}}_{i,j+\frac{1}{2}}^{\mathrm{WENO}} - \hat{\mathbf{G}}_{i,j-\frac{1}{2}}^{\mathrm{WENO}} \right) + \mathbf{s}(\mathbf{w}_{i,j}), \quad (22)$$

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbf{w}_{i,j}(t) = -\frac{1}{\Delta x} \left(\hat{\mathbf{F}}_{i+\frac{1}{2},j} - \hat{\mathbf{F}}_{i-\frac{1}{2},j} \right) - \frac{1}{\Delta y} \left(\hat{\mathbf{G}}_{i,j+\frac{1}{2}} - \hat{\mathbf{G}}_{i,j-\frac{1}{2}} \right) + \mathbf{s}(\mathbf{w}_{i,j}), \tag{23}$$

where all numerical fluxes are computed dimension by dimension. For simplicity, we omit the detailed formulation.

As in the problem in one space dimension, we can define the consistent numerical fluxes:

Definition 4 Considering the semi-discrete scheme for solving the governing system (20), we say the numerical fluxes $\hat{\mathbf{f}}$ and $\hat{\mathbf{g}}$ are consistent if $\hat{\mathbf{f}} \cdot \mathbf{v} = \hat{\mathbf{g}} \cdot \mathbf{v} = 0$, provided $\mathbf{w}_{i,j} \cdot \mathbf{v} = 0$ for all $1 \le i \le N_x$ and $1 \le j \le N_y$.

Since the numerical fluxes are computed dimension by dimension, we can easily extend the analysis in one space dimension to obtain

$$\hat{\mathbf{f}}_{i+\frac{1}{2}j} \cdot \mathbf{v} = \hat{\mathbf{g}}_{i,j+\frac{1}{2}} \cdot \mathbf{v} = \hat{\mathbf{F}}_{i+\frac{1}{2}j} \cdot \mathbf{v} = \hat{\mathbf{G}}_{i,j+\frac{1}{2}} \cdot \mathbf{v} = 0.$$
(24)

Hence, the numerical fluxes in the first-order and fifth-order linear schemes are consistent. Again, the numerical fluxes computed by WENO procedure are not consistent.

4.2 Euler Forward Time Integration and Bound-Preserving Technique

We consider the bound-preserving technique for the convection term and assume that s = 0. We define the admissible set of solutions as

$$G = \left\{ \mathbf{w} = \begin{pmatrix} \rho \\ m \\ n \\ E \\ r_1 \\ \vdots \\ r_M \end{pmatrix}, \rho \ge 0, p \ge 0, z_k \ge 0, \ k = 1, \cdots, M, \sum_{k=1}^M z_k = 1 \right\}.$$

For the problem in two space dimensions, *G* is also a convex set. We apply the Euler forward time integration to the first order scheme and the fifth order linear scheme, respectively, and get the following two fully discrete schemes for the convection terms:

$$\begin{split} \mathbf{w}_{i,j}^{n+1} &= \mathbf{w}_{i,j}^n - \frac{\Delta t}{\Delta x} \Big(\hat{\mathbf{f}}_{i+\frac{1}{2},j} - \hat{\mathbf{f}}_{i-\frac{1}{2},j} \Big) - \frac{\Delta t}{\Delta y} \Big(\hat{\mathbf{g}}_{i,j+\frac{1}{2}} - \hat{\mathbf{g}}_{i,j-\frac{1}{2}} \Big), \\ \mathbf{w}_{i,j}^{n+1} &= \mathbf{w}_{i,j}^n - \frac{\Delta t}{\Delta x} \Big(\hat{\mathbf{F}}_{i+\frac{1}{2},j} - \hat{\mathbf{F}}_{i-\frac{1}{2},j} \Big) - \frac{\Delta t}{\Delta y} \Big(\hat{\mathbf{G}}_{i,j+\frac{1}{2}} - \hat{\mathbf{G}}_{i,j-\frac{1}{2}} \Big), \end{split}$$

where $\mathbf{w}_{i,j}^n = \left(\rho_{i,j}^n, m_{i,j}^n, n_{i,j}^n, E_{i,j}^n, (r_1)_{i,j}^n, \cdots, (r_M)_{i,j}^n\right)^{\mathrm{T}}$ is the numerical approximation at time level *n*, and all numerical fluxes are also computed at time level *n*. Assume $\mathbf{w}_{i,j}^n \in G$ for all *i* and *j*. For the first order scheme, we have $\mathbf{w}_{i,j}^{n+1} \in G$ under the CFL condition $\Delta t \leq \Delta \tilde{t}$ with

$$\alpha \frac{\Delta \tilde{t}}{\Delta x} + \beta \frac{\Delta \tilde{t}}{\Delta y} \leqslant 1, \tag{25}$$

where $\alpha = \|(|u| + c)\|_{\infty}$ and $\beta = \|(|v| + c)\|_{\infty}$ are the global maximum absolute eigenvalues of $\mathbf{f}'(\mathbf{u})$ and $\mathbf{g}'(\mathbf{u})$, respectively. For the high order scheme, the numerical solution may be out of the physical bounds. We also apply the parameterized positivity-preserving flux limiter to get the modified scheme

$$\mathbf{w}_{ij}^{n+1} = \mathbf{w}_{ij}^{n} - \frac{\Delta t}{\Delta x} \Big(\tilde{\mathbf{F}}_{i+\frac{1}{2}j} - \tilde{\mathbf{F}}_{i-\frac{1}{2}j} \Big) - \frac{\Delta t}{\Delta y} \Big(\tilde{\mathbf{G}}_{ij+\frac{1}{2}} - \tilde{\mathbf{G}}_{ij-\frac{1}{2}} \Big), \tag{26}$$

where $\tilde{\mathbf{F}}_{i+\frac{1}{2}j}$ and $\tilde{\mathbf{G}}_{i,j+\frac{1}{2}}$ are combinations of the high-order and the first-order fluxes:

$$\begin{split} \tilde{\mathbf{F}}_{i+\frac{1}{2}j} &= \hat{\mathbf{f}}_{i+\frac{1}{2}j} + \theta_{i+\frac{1}{2}j} \Big(\hat{\mathbf{F}}_{i+\frac{1}{2}j} - \hat{\mathbf{f}}_{i+\frac{1}{2}j} \Big), \\ \tilde{\mathbf{G}}_{i,j+\frac{1}{2}} &= \hat{\mathbf{g}}_{i,j+\frac{1}{2}} + \theta_{i,j+\frac{1}{2}} \Big(\hat{\mathbf{G}}_{i,j+\frac{1}{2}} - \hat{\mathbf{g}}_{i,j+\frac{1}{2}} \Big). \end{split}$$

One can check that the modified numerical fluxes are still consistent using the property (24):

$$\tilde{\mathbf{F}}_{i+\frac{1}{2},j} \cdot \mathbf{v} = \tilde{\mathbf{G}}_{i,j+\frac{1}{2}} \cdot \mathbf{v} = 0.$$

Assuming that $\mathbf{w}_{i,j}^n \in G$ for all *i* and *j*, we aim to choose suitable parameters $\theta_{i+\frac{1}{2},j} \in [0, 1]$ and $\theta_{i,j+\frac{1}{2}} \in [0, 1]$, such that the solution of (26) satisfies $\mathbf{w}_{i,j}^{n+1} \in G$.

Notice that the conservative property of the scheme is important for constructing the BP technique in one space dimension. Now we extend the definition of conservative schemes to the problem in two space dimensions.

Definition 5 Consider a fully discrete numerical scheme for solving the detonation problem (20), if the mass conservation property can be maintained during the time evolution, namely, the numerical solution satisfies

$$\sum_{k=1}^{M} (r_k)_{i,j}^{n+1} = \rho_{i,j}^{n+1}$$

as long as $\sum_{k=1}^{M} (r_k)_{i,j}^n = \rho_{i,j}^n$, then we say the numerical scheme is a conservative scheme.

We first state the conservative property of the fully discrete scheme (26). Recall that Theorem 1 gives two sufficient conditions for constructing conservative schemes: consistent numerical fluxes and conservative time integrations. Although this theorem is derived in one space dimension, it also applies to the problem in two space dimensions. As the numerical fluxes in (26) are consistent and the Euler forward time integration is conservative, the scheme (26) is conservative.

Next, we enforce $(r_k)_{i,j}^{n+1} \ge 0$, $k = 1, \dots, M$. In the scheme (26), we denote the equation for solving the density of the *k*th $(k = 1, \dots, M)$ species as

$$(r_k)_{ij}^{n+1} = (r_k)_{ij}^n - \frac{\Delta t}{\Delta x} \left(\tilde{F}_{i+\frac{1}{2}j}^k - \tilde{F}_{i-\frac{1}{2}j}^k \right) - \frac{\Delta t}{\Delta y} \left(\tilde{G}_{ij+\frac{1}{2}}^k - \tilde{G}_{ij-\frac{1}{2}}^k \right).$$

We need to find local pairs $\left(\Lambda_{L,I_{ij}}^{k},\Lambda_{R,I_{ij}}^{k}\right)$ and $\left(\Lambda_{D,I_{ij}}^{k},\Lambda_{U,I_{ij}}^{k}\right)$, such that for any pairs $\left(\theta_{i-\frac{1}{2},j},\theta_{i+\frac{1}{2},j}\right) \in \left[0,\Lambda_{L,I_{ij}}^{k}\right] \times \left[0,\Lambda_{R,I_{ij}}^{k}\right]$ and $\left(\theta_{i,j-\frac{1}{2}},\theta_{i,j+\frac{1}{2}}\right) \in \left[0,\Lambda_{D,I_{ij}}^{k}\right] \times \left[0,\Lambda_{U,I_{ij}}^{k}\right]$, we have $(r_{k})_{i,j}^{n+1} \ge 0$. We divide this problem into two parts:

$$\frac{1}{2}(r_k)_{i,j}^n - \frac{\Delta t}{\Delta x} \left(\tilde{F}_{i+\frac{1}{2},j}^k - \tilde{F}_{i-\frac{1}{2},j}^k \right) \ge 0,$$
(27)

$$\frac{1}{2}(r_k)_{i,j}^n - \frac{\Delta t}{\Delta y} \left(\tilde{G}_{i,j+\frac{1}{2}}^k - \tilde{G}_{i,j-\frac{1}{2}}^k \right) \ge 0.$$
(28)

For each fixed *j*, (27) is just a one-dimensional problem along the *x* direction. Following the algorithm in Sect. 3.2, we can obtain the local pair $\left(\Lambda_{L,I_{ij}}^{k}, \Lambda_{R,I_{ij}}^{k}\right)$. Similarly, we can get the local pair $\left(\Lambda_{D,I_{ij}}^{k}, \Lambda_{U,I_{ij}}^{k}\right)$ by solving (28). Then, we let

$$\begin{split} \Lambda^{\rho}_{L,I_{ij}} &:= \min_{k=1,\cdots,M} \Lambda^{k}_{L,I_{ij}}, \quad \Lambda^{\rho}_{R,I_{ij}} &:= \min_{k=1,\cdots,M} \Lambda^{k}_{R,I_{ij}}, \\ \Lambda^{\rho}_{D,I_{ij}} &:= \min_{k=1,\cdots,M} \Lambda^{k}_{D,I_{ij}}, \quad \Lambda^{\rho}_{U,I_{ij}} &:= \min_{k=1,\cdots,M} \Lambda^{k}_{U,I_{ij}}. \end{split}$$

For any pairs $\left(\theta_{i-\frac{1}{2},j},\theta_{i+\frac{1}{2},j}\right) \in \left[0,\Lambda_{L,I_{i,j}}^{\rho}\right] \times \left[0,\Lambda_{R,I_{i,j}}^{\rho}\right] \quad \text{and}\left(\theta_{i,j-\frac{1}{2}},\theta_{i,j+\frac{1}{2}}\right) \in \left[0,\Lambda_{D,I_{i,j}}^{\rho}\right] \times \left[0,\Lambda_{U,I_{i,j}}^{\rho}\right]$ the solution of (26) satisfies

$$(z_k)_{i,j}^{n+1} \ge 0, \ k = 1, \cdots, M, \qquad \sum_{k=1}^M (z_k)_{i,j}^{n+1} = 1, \quad \rho_{i,j}^{n+1} \ge 0.$$

Finally, we preserve the positivity of the pressure. As in the one-dimensional case, for $A = \left(\theta_{i-\frac{1}{2}j}, \theta_{i+\frac{1}{2}j}, \theta_{i,j-\frac{1}{2}}, \theta_{i,j+\frac{1}{2}}\right)$, we denote

$$p(A) := p\Big(\mathbf{w}_{i,j}^{n+1}\Big(\theta_{i-\frac{1}{2},j},\theta_{i+\frac{1}{2},j},\theta_{i,j-\frac{1}{2}},\theta_{i,j+\frac{1}{2}}\Big)\Big).$$

Following [24], we perform the following algorithm.

- i) Denote $A^{k_1,k_2,k_3,k_4} = \left(k_1 \Lambda_{L,I_{ij}}^{\rho}, k_2 \Lambda_{R,I_{ij}}^{\rho}, k_3 \Lambda_{D,I_{ij}}^{\rho}, k_4 \Lambda_{U,I_{ij}}^{\rho}\right)$, where each $k_l \ (l = 1, 2, 3, 4)$ can be 0 or 1.
- ii) For each $(k_1, k_2, k_3, k_4) \neq (0, 0, 0, 0)$, if $p(A^{k_1, k_2, k_3, k_4}) \ge 0$, let $B^{k_1, k_2, k_3, k_4} = A^{k_1, k_2, k_3, k_4}$. Otherwise, find *r* such that $p(rA^{k_1, k_2, k_3, k_4}) \ge 0$ and let $B^{k_1, k_2, k_3, k_4} = rA^{k_1, k_2, k_3, k_4}$. We also denote

$$B^{k_1,k_2,k_3,k_4} = \left(B_1^{k_1,k_2,k_3,k_4}, B_2^{k_1,k_2,k_3,k_4}, B_3^{k_1,k_2,k_3,k_4}, B_4^{k_1,k_2,k_3,k_4}\right).$$

iii) Take

$$\begin{split} \Lambda_{L,I_{ij}} &= \min\left(B_1^{1,1,1,0}, B_1^{1,1,0,1}, B_1^{1,0,1}\right), \quad \Lambda_{R,I_{ij}} &= \min\left(B_2^{1,1,1,0}, B_2^{1,1,0,1}, B_2^{0,1,1,1}\right), \\ \Lambda_{D,I_{ij}} &= \min\left(B_3^{1,1,1,0}, B_3^{1,0,1,1}, B_3^{0,1,1,1}\right), \quad \Lambda_{U,I_{ij}} &= \min\left(B_4^{1,1,0,1}, B_4^{1,0,1,1}, B_4^{0,1,1,1}\right). \end{split}$$
iv) Let $\theta_{i+\frac{1}{2},j} &= \min\left(\Lambda_{R,I_{ij}}, \Lambda_{L,I_{i+1j}}\right)$ and $\theta_{i,j+\frac{1}{2}} &= \min\left(\Lambda_{U,I_{ij}}, \Lambda_{D,I_{ij+1}}\right).$

5 High-Order Time Integrations

In the previous sections, we have assumed that $\mathbf{s} = \mathbf{0}$ and considered the Euler forward time integration. In this section, we proceed to consider high-order time integrations for the general case with $\mathbf{s} \neq \mathbf{0}$. For simplicity, we rewrite different semi-discrete schemes in the following uniform form:

$$\mathbf{w}_t = \mathbf{F}(\mathbf{w}) + \mathbf{s}(\mathbf{w}),\tag{29}$$

where **F** represents the spatial discretization of the flux and **s** is the source term. For the one-dimensional fifth order linear scheme (12), we can take

$$\mathbf{F} = -\frac{1}{\Delta x} \Big(\hat{\mathbf{F}}_{i+\frac{1}{2}} - \hat{\mathbf{F}}_{i-\frac{1}{2}} \Big).$$

For the two-dimensional scheme (23), we have

$$\mathbf{F} = -\frac{1}{\Delta x} \Big(\hat{\mathbf{F}}_{i+\frac{1}{2}j} - \hat{\mathbf{F}}_{i-\frac{1}{2}j} \Big) - \frac{1}{\Delta y} \Big(\hat{\mathbf{G}}_{i,j+\frac{1}{2}} - \hat{\mathbf{G}}_{i,j-\frac{1}{2}} \Big).$$

For both cases, we summarize the common properties of (29) as follows.

• There exists a constant vector \mathbf{v} such that $\mathbf{F} \cdot \mathbf{v} = \mathbf{s} \cdot \mathbf{v} = 0$. For the one-dimensional detonation problem (4), we can adopt

$$\mathbf{v} = [1, 0, 0, -1, \cdots, -1]^{\mathrm{T}} \in \mathbb{R}^{M+3}.$$

For the two-dimensional problem (20), v becomes

$$\mathbf{v} = [1, 0, 0, 0, -1, \cdots, -1]^{\mathrm{T}} \in \mathbb{R}^{M+4}.$$

It is easy to check that $\mathbf{s} \cdot \mathbf{v} = 0$. In addition, using the consistent properties (13) and (24) of the linear numerical fluxes, we can further get $\mathbf{F} \cdot \mathbf{v} = 0$.

• Assume that $\mathbf{w} \in G$. Then, we can apply the flux limiter and replace \mathbf{F} with $\tilde{\mathbf{F}}$, such that

$$\tilde{\mathbf{F}} \cdot \mathbf{v} = 0, \quad \mathbf{w} + \Delta t \tilde{\mathbf{F}} \in G.$$

For the one-dimensional case as in (17), we have

$$\tilde{\mathbf{F}} = -\frac{1}{\Delta x} \Big(\tilde{\mathbf{F}}_{i+\frac{1}{2}} - \tilde{\mathbf{F}}_{i-\frac{1}{2}} \Big).$$

For the two-dimensional case (26), we have

$$\tilde{\mathbf{F}} = -\frac{1}{\Delta x} \left(\tilde{\mathbf{F}}_{i+\frac{1}{2}} - \tilde{\mathbf{F}}_{i-\frac{1}{2}} \right) - \frac{1}{\Delta y} \left(\tilde{\mathbf{G}}_{i,j+\frac{1}{2}} - \tilde{\mathbf{G}}_{i,j-\frac{1}{2}} \right).$$

Now, we solve for (29) by using high order time integrations. Since the source term may be stiff, we adopt the third order conservative modified ERK methods [9]:

$$\mathbf{w}^{(1)} = \left[\alpha_{10}\mathbf{w}^n + \beta_{10}\Delta t\mathbf{F}(\mathbf{w}^n) + \beta_{10}\Delta t(\mathbf{s}(\mathbf{w}^n) + \mu\mathbf{w}^n)\right]/A_1,\tag{30}$$

$$\mathbf{w}^{(2)} = \left[\alpha_{20}\mathbf{w}^{n} + \beta_{20}\Delta t\mathbf{F}(\mathbf{w}^{n}) + \beta_{20}\Delta t(\mathbf{s}(\mathbf{w}^{n}) + \mu\mathbf{w}^{n})\right]/A_{2} + e^{\beta_{10}\mu\Delta t} \left[\alpha_{21}\mathbf{w}^{(1)} + \beta_{21}\Delta t\mathbf{F}(\mathbf{w}^{(1)}) + \beta_{21}\Delta t(\mathbf{s}(\mathbf{w}^{(1)}) + \mu\mathbf{w}^{(1)})\right]/A_{2},$$
(31)

$$\mathbf{w}^{n+1} = \left[\alpha_{30}\mathbf{w}^{n} + \beta_{30}\Delta t \mathbf{F}(\mathbf{w}^{n}) + \beta_{30}\Delta t(\mathbf{s}(\mathbf{w}^{n}) + \mu \mathbf{w}^{n})\right] / A_{3} + e^{\beta_{10}\mu\Delta t} \left[\alpha_{31}\mathbf{w}^{(1)} + \beta_{31}\Delta t \mathbf{F}(\mathbf{w}^{(1)}) + \beta_{31}\Delta t(\mathbf{s}(\mathbf{w}^{(1)}) + \mu \mathbf{w}^{(1)})\right] / A_{3} + e^{A\mu\Delta t} \left[\alpha_{32}\mathbf{w}^{(2)} + \beta_{32}\Delta t \mathbf{F}(\mathbf{w}^{(2)}) + \beta_{32}\Delta t(\mathbf{s}(\mathbf{w}^{(2)}) + \mu \mathbf{w}^{(2)})\right] / A_{3},$$
(32)

where

$$A_{1} = \alpha_{10} + \beta_{10}\mu\Delta t, \quad A_{2} = \alpha_{20} + \beta_{20}\mu\Delta t + e^{\beta_{10}\mu\Delta t} (\alpha_{21} + \beta_{21}\mu\Delta t), A_{3} = \alpha_{30} + \beta_{30}\mu\Delta t + e^{\beta_{10}\mu\Delta t} (\alpha_{31} + \beta_{31}\mu\Delta t) + e^{A\mu\Delta t} (\alpha_{32} + \beta_{32}\mu\Delta t),$$

where μ is a parameter to be determined later. To achieve the third-order accuracy, the conditions on all α and β was given in [9]. In this paper, we adopt the following choice:

$$\begin{aligned} \alpha_{10} &= 1, \quad \beta_{10} = \frac{2}{3}, \quad \alpha_{20} = \frac{7}{8}, \quad \beta_{20} = \frac{1}{12}, \quad \alpha_{21} = \frac{1}{8}, \quad \beta_{21} = \frac{1}{2}, \\ \alpha_{30} &= \frac{1}{2}, \quad \beta_{30} = \frac{1}{12}, \quad \alpha_{31} = \frac{1}{6}, \quad \beta_{31} = \frac{1}{12}, \quad \alpha_{32} = \frac{1}{3}, \quad \beta_{32} = \frac{1}{2} \end{aligned}$$

Remark 1 The time integration given above separates the convection and source terms. We can rewrite the source as, e.g.,

$$\mathbf{s}(\mathbf{w}^n) + \mu \mathbf{w}^n = \mu \left(\mathbf{w}^n + \frac{1}{\mu} \mathbf{s}(\mathbf{w}^n) \right),$$

then the parameter μ can be regarded as the reciprocal of the time step in the Euler forward time integration. Following [23], we can apply the technique of subcell resolution to the source term.

The conservative property of the above time integration was proved in [9]. Since the numerical fluxes are consistent, the fully discrete scheme is conservative. Hence, we only need to preserve the positivity of each mass fraction and the upper bound 1 will be preserved automatically. Finally, we can state the following theorem for bound-preserving technique.

Theorem 2 Consider the fifth order linear FD scheme (12) (or (23) for the problem in two space dimensions) coupled with the three-stage ERK method (30)–(32), where we apply flux limiter on each stage and μ satisfies

$$\mu \ge \max_{0 \le i \le M} \left\{ -\frac{\sum_{j=1}^{M} s_j q_j}{r_i, \frac{p}{p}, 0} \right\}$$
(33)

for $\mathbf{w} = \mathbf{w}^n, \mathbf{w}^{(1)}, \mathbf{w}^{(2)}$. If $\mathbf{w}^n \in G$, then we have $\mathbf{w}^{n+1} \in G$ under the condition

$$\Delta t \leqslant \min\left\{\frac{\alpha_{10}}{\beta_{10}}, \frac{\alpha_{20}}{\beta_{20}}, \frac{\alpha_{21}}{\beta_{21}}, \frac{\alpha_{30}}{\beta_{30}}, \frac{\alpha_{31}}{\beta_{31}}, \frac{\alpha_{32}}{\beta_{32}}\right\} \Delta \tilde{t},$$

where $\Delta \tilde{t}$ satisfies (16) (or (25) for the problem in two space dimensions).

Proof For simplicity, we consider the first stage of the ERK method only and prove $\mathbf{w}^{(1)} \in G$. After simple computations, we get

$$\mathbf{w}^{(1)} = \left(\alpha_{10}\mathbf{R}_1 + \beta_{10}\mu\Delta t\mathbf{R}_2\right)/A_1,$$

where

$$\mathbf{R}_1 = \mathbf{w}^n + \frac{\beta_{10}}{\alpha_{10}} \Delta t \mathbf{F}(\mathbf{w}^n)$$
 and $\mathbf{R}_2 = \frac{1}{\mu} (\mathbf{s}(\mathbf{w}^n) + \mu \mathbf{w}^n).$

For the convection term, we can apply the flux limiter and replace **F** with $\tilde{\mathbf{F}}$, such that $\mathbf{R}_1 \in G$ under the condition $\Delta t \leq \frac{\alpha_{10}}{\beta_{10}} \Delta \tilde{t}$. For the source term, Lemma 4.1 in [8] shows that $\mathbf{R}_2 \in G$ under the condition (33). Recall that $A_1 = \alpha_{10} + \beta_{10}\mu\Delta t$ and $\mathbf{w}^{(1)}$ is a convex combination of \mathbf{R}_1 and \mathbf{R}_2 . Since *G* is a convex set, we have $\mathbf{w}^{(1)} \in G$. Following the same analysis above, we can also prove that $\mathbf{w}^{(2)}, \mathbf{w}^{\mathbf{n}+1} \in G$.

6 Numerical Examples

In this section, we show some numerical experiments to demonstrate the good performance of the high order linear FD scheme. We will show that, without further numerical techniques such as subcell resolutions, the conservative FD method with linear weights can yield better numerical approximations than the nonconservative WENO scheme. The reference solutions are computed by the regular 5th order characteristic-wise WENO-LF method coupled with SSP-RK3 with $N = 10\ 000$ grids and CFL = 0.5. In all figures, we use "new RK" to represent the conservative ERK methods (30)–(32).

Example 1 (Accuracy test for one dimensional PDE system)

In this example, we test the one dimensional convection reaction system problem. We consider periodic boundary condition and take u = 1 and p = 0 in the exact solution. We choose M = 2 and the source term is given as $s_1 = -cr_1^7$. Hence, we need to solve the following system:

$$\begin{cases} \rho_t + \rho_x = 0, \\ (r_1)_t + (r_1)_x = -c(r_1)^7, \quad x \in [0, 2\pi]. \end{cases}$$

The parameter *c* can be used to adjust the stiffness of the equation. The initial conditions are given as $r_1(x, 0) = 0.1(1 + \sin(x))$ and $\rho(x, 0) = 0.1(2 + \sin(x) + \cos(x))$. The final time is taken as t = 0.5. For this problem, the total density ρ should be non-negative and the mass fraction r_1/ρ should be between 0 and 1.

Both non-stiff (c = 100) and stiff ($c = 10\ 000$) cases are calculated. We first adopt the 5th order conservative linear scheme and the numerical errors are listed in the left part

N	Without limiter				With limiter					
	L^2 norm	Order	L^{∞} norm	Order	L^2 norm	Order	L^{∞} norm	Order	Percent- age/%	
c = 1	00									
10	9.49E-05	_	1.89E-04	_	9.49E-05	_	1.89E-04	_	0	
20	3.26E-06	4.86	7.52E-06	4.65	3.26E-06	4.86	7.52E-06	4.65	0	
40	1.07E-07	4.93	2.69E-07	4.80	1.07E-07	4.93	2.70E-07	4.80	1.48E-03	
80	3.39E-09	4.98	8.65E-09	4.96	3.43E-09	4.96	1.07E-08	4.66	3.70E-03	
160	1.06E-10	4.99	2.72E-10	4.99	1.56E-10	4.46	9.94E-10	3.43	2.08E-03	
c = 1	0 000									
10	1.51E-04	-	3.66E-04	-	4.43E-04	-	8.25E-04	-	0	
20	9.45E-06	4.00	2.42E-05	3.92	3.98E-05	3.47	9.86E-05	3.06	0	
40	3.71E-07	4.67	9.79E-07	4.63	1.94E-07	4.36	5.46E-06	4.17	1.48E-02	
80	1.22E-08	4.92	3.25E-08	4.91	6.85E-08	4.83	2.03E-07	4.75	3.70E-03	
160	3.88E-10	4.98	1.03E-09	4.98	2.21E-10	4.96	7.23E-09	4.81	2.08E-03	

 Table 1
 Accuracy test for the one-dimensional system

of Table 1. We can observe the optimal convergence rate as expected. Next, we combine the 5th order flux with the 1st order LF flux by using the flux limiter and aim to preserve the bounds of the total density and mass fraction. It is reported in [25] that the time step size needs to be small for recovering the original high order of accuracy. Here we take $\Delta t = 0.02\Delta x^2$. The results are shown in the right part of Table 1 and the last column shows the percentage of grids that have been modified by the flux limiter. We can see that the flux limiter will lead to a small order decreasing, especially when *N* is large. The 5th order accuracy can be fully recovered if we use an even small time step size. This phenomenon has been demonstrated and analyzed in [25]. Some alternatives to improve the accuracy were also given in [25]. However, those techniques can hardly be applied to our time integration. Therefore, we will discuss this issue in the future.

Example 2 (A 1D detonation wave with 3 species and 1 reaction)

In this example, we solve a reacting model with three species and one reaction:

$$2H_2 + O_2 \rightarrow 2H_2O.$$

The parameters are taken as $T_1 = 2.0, B_1 = 500, \alpha_1 = 1, q_1 = 1000, q_2 = 0, q_3 = 0, M_1 = 2, M_2 = 32, M_3 = 18$. The computational domain is [0, 50]. Initially there is a mixture of hydrogen and oxygen on the right-hand side. On the left-hand side, the hydrogen and oxygen generate water. The initial condition is given as piecewise constants:

$$(\rho, u, p, z_1, z_2, z_3)(x, 0) = \begin{cases} (2.0, 10.0, 40.0, 0.325, 0, 0.675), \ x \leq 2.5, \\ (1.0, 0, 1.0, 0.4, 0.6, 0), \qquad x > 2.5. \end{cases}$$

We take the final time to be t = 3. This is a simple one-step chemical model for hydrogenoxygen mixtures. The exact solution consists of a detonation wave, followed by a contact discontinuity and a shock, all moving to the right.

We first compare different temporal and spacial discretizations in Fig. 1 by taking N = 1070 and CFL = 0.03. As shown by the orange solid lines with circle symbols,



Fig. 1 Example 2. N = 1070. Comparison of different temporal and spacial discretizations

the WENO scheme coupled with the classical SSP-RK method can not capture the correct shock location and produces spurious numerical results. Based on this combination, if we replace the temporal discretization with our conservative ERK method (blue dashed lines with square symbols) or replace the spacial discretization with the linear



Fig. 2 Example 2. N = 1000. Comparison of linear scheme with WENO

FD method (green dashed lines), the correct propagation speed of the detonation wave can be captured. Of course if we combine the linear method in space and new ERK method in time (long red dashed lines), no spurious numerical results will be observed. Notice that for the high order linear scheme in space, there are oscillations. We will add the flux limiters to eliminate the oscillations later.

Based on the above results, we know that both the new conservative ERK method in time and the linear scheme in space help to capture the correct solutions. From now on,



Fig. 3 Example 2. N = 4000. Comparison of linear scheme with WENO

we fixed the temporal discretization as our new RK method and further compare different spacial discretizations.

We first take a relatively coarser mesh N = 1000 with CFL = 0.03 and compare our linear scheme (conservative) with the WENO scheme (not conservative). In Fig. 2, red lines with symbols represent the linear scheme without any limiter. We can see that the linear method is able to capture the correct propagation speed on this mesh. Since the scheme is linear, there are oscillations. Also, we can observe some negative values of pressure and the mass fraction z_3 . After we add the flux limiter, we can obtain the green dashed lines. We see that all oscillations together with the non-physical negative values are eliminated. As shown by the blue dash dot lines, the complex WENO scheme produces spurious numerical results even when we use the new RK method. The results for the component-wise WENO-LF scheme are similar and hence we do not show the plots.

Next, we take a denser grid with N = 4000. In this case, CFL = 0.1 is enough for our method to capture the correction shock position. Figure 3a, c show the results of the linear scheme. We can observe a few oscillations in this case. In addition, the linear method is able to resolve the thin reaction zone near the shock (Fig. 3a, near x = 40). As shown by the green dashed lines, the flux limiter will eliminate non-physical values as well as oscillations, while maintaining the thin reaction zone. Figure 3b, d show the results of the WENO scheme. We can see that WENO can control oscillations. However, it is not able to resolve the thin reaction zone on this mesh. Notice that although we use the WENO scheme on a very dense mesh to get the reference solutions, there still exists unreasonable values near the contact discontinuity. We zoom in this area and compare both methods in Fig. 3e. We can see that the results of linear scheme are reasonable, even when we do not add the limiter.

Based on the above observations, we know that the conservative property among components is important. Now we try to modify the componentwise WENO-LF method. We still use nonlinear weights to control oscillations. But we use the same set of weights for all components and hence the scheme is still conservative. Figure 4 shows the plots of density by using the nonlinear weights of the mass species ρ_1 , ρ_2 , ρ_3 , the momentum and the energy, respectively. Here we take N = 1000 and CFL = 0.03. No limiter is added. We can see that the choice of nonlinear weights does impact the numerical results. Hence, we will still use the simple linear method coupled with flux limiter in the remaining part of the paper.

Example 3 (A 1D detonation wave with 4 species and 1 reaction) In this example, we consider a reacting model with four species and one reaction. A prototype reaction for this model is

$$CH_4 + 2O_2 \rightarrow CO_2 + 2H_2O_2$$

The parameters are $T_1 = 2.0, B_1 = 10^6, \alpha_1 = 0, q_1 = 500, q_2 = 0, q_3 = 0, q_4 = 0, M_1 = 16, M_2 = 32, M_3 = 44, M_4 = 18$. The initial data are as follows:

$$(\rho, u, p, z_1, z_2, z_3, z_4)(x, 0) = \begin{cases} (2, 10, 40, 0, 0.2, 0.475, 0.325), & x \leq 2.5, \\ (1, 0, 1, 0.1, 0.6, 0.2, 0.1), & x > 2.5. \end{cases}$$

The computational domain is [0, 50] and final time is t = 3. The exact solution consists of a detonation wave, followed by a contact discontinuity and a shock, all moving to the right.



Fig. 4 Example 2. Conservative component-wise WENO-LF with different nonlinear weights. N = 1000

We take N = 760 and CFL = 0.1 and compare different methods in Fig. 5. We first use the regular WENO method in space which is not conservative, and then combine it with the traditional RK method (orange dash dot lines) and the new RK method (blue dash dot lines) in time, respectively. We can see that the numerical results are similar and both



Fig. 5 Example 3 at t = 3. N = 760

methods will lead to spurious numerical results. Next, we replace the WENO method with the 5th order conservative linear method and use RK in time (rose red dash dot lines). The numerical results are much better and we see that the conservative property is very important. But there are still some spurious numerical results due to the RK method. Next, we



Fig. 6 Example 3 at t = 3. N = 1500

further replace the temporal discretization with the new RK method. We can capture the correct shock positions with no spurious numerical results (red solid lines with circle symbols). Since this method is only linear, we observe some oscillations. After applying the flux limiter, all these oscillations will disappear.

On the above mesh, when we fix the spacial discretization as the WENO scheme, the results of the RK method and the new RK method are almost the same. Now we take a fine mesh with N = 1500 and CFL = 0.3 and further compare different RK methods in Fig. 6. We can see that the new RK method is better for this example.

Example 4 (Accuracy test for 2D system) From now on, we consider the two dimensional problems. In this example, we consider periodic boundary condition and take u = v = 1 and p = 0 in the exact solution. We choose M = 2 and the source is given as $s_1 = -cr_1^7$. Hence, we need to solve the following system:

$$\begin{cases} \rho_t + \rho_x + \rho_y = 0, \\ (r_1)_t + (r_1)_x + (r_1)_y = -c(r_1)^7, \\ \end{cases} (x, y) \in [0, 2\pi]^2.$$

The initial conditions are given as $\rho(x, y, 0) = 0.1(2 + \sin(x + y) + \cos(x + y))$ and $r_1(x, y, 0) = 0.1(1 + \sin(x + y))$, respectively. For this problem, the total density ρ should be non-negative and the mass fraction r_1/ρ should be between 0 and 1.

Numerical errors at the final time t = 0.5 are listed in Table 2. The left part of the table shows the results for the 5th order conservative linear scheme. We can again observe the expected optimal convergence rate. We further add the flux limiter to preserve the lower bound of ρ and the two bounds of r_1/ρ , and show the results in the right part of the error table. Here we take $\Delta t = 0.01\Delta x^2$. We can see that the flux limiter will lead to a small order deficiency, especially when N is large. The 5th order accuracy can be fully recovered if we use an even smaller time step size.

Example 5 (A 2D detonation wave with 4 species and 1 reaction) In this example, we test a 2D reacting model with four species and one reaction. A prototype reaction for this model is

$$CH_4 + 2O_2 \rightarrow CO_2 + 2H_2O.$$

N	Without limiter				With limiter					
	L^2 norm	Order	L^{∞} norm	Order	$\overline{L^2 \text{ norm}}$	Order	L^{∞} norm	Order	Percent- age/%	
c = 1	00									
10	1.89E-04	_	3.59E-04	_	1.89E-04	_	3.59E-04	_	0	
20	6.50E-06	4.86	1.48E-05	4.60	6.50E-06	4.86	1.48E-05	4.60	0	
40	2.14E-07	4.93	5.38E-07	4.78	2.14E-07	4.93	5.38E-07	4.78	0	
80	6.78E-09	4.98	1.73E-08	4.96	6.79E-09	4.98	1.73E-08	4.96	4.54E-03	
160	2.13E-10	4.99	5.44E-10	4.99	2.30E-10	4.88	1.04E-09	4.05	1.83E-03	
c = 1	0 000									
10	8.22E-04	-	1.81E-03	-	8.22E-04	-	1.81E-03	-	0	
20	7.40E-05	3.47	1.78E-04	3.35	7.40E-05	3.47	1.78E-04	3.35	0	
40	3.84E-06	4.27	1.06E-05	4.07	3.84E-06	4.27	1.06E-05	4.07	0	
80	1.37E-07	4.81	3.95E-07	4.74	1.37E-07	4.81	3.95E-07	4.74	4.54E-03	
160	4.41E-09	4.96	1.31E-08	4.91	4.41E-09	4.96	1.36E-08	4.86	1.83E-03	

 Table 2
 Accuracy test for the two dimensional problem



Fig. 7 Numerical solutions of Example 5 along the line x = y at t = 2

The parameters are $T_1 = 2$, $B_1 = 10^6$, $\alpha_1 = 0$, $q_1 = 200$, $q_2 = 0$, $q_3 = 0$, $q_4 = 0$, $M_1 = 16$, $M_2 = 32$, $M_3 = 44$, $M_4 = 18$. The initial values consist of totally burnt gas inside of a circle with radius 10 and totally unburnt gas everywhere outside this circle. The setup is as follows:

$$(\rho, u, v, p, z_1, z_2, z_3, z_4)(x, y, 0) = \begin{cases} (2, 10x/r, 10y/r, 40, 0, 0.2, 0.475, 0.325), & r \leq 10, \\ (1, 0, 0, 1, 0.1, 0.6, 0.2, 0.1), & r > 10, \end{cases}$$

where $r = \sqrt{x^2 + y^2}$. The computational domain is $[0, 50] \times [0, 50]$. This is a radially symmetric problem and the detonation front is circular. The boundary conditions are solid-wall boundary conditions on the left and lower boundaries, and outflow boundary conditions on the right and upper boundaries.

We take $N_x = N_y = 600$ and CFL = 0.1. Figure 7 shows the one dimensional cuts of pressure, density and mass fractions along the line x = y at t = 2. We can see that our scheme preserves the positivity of the density and pressure, and the two bounds 0 and 1 of each mass fraction. In addition, our linear scheme with flux limiter can capture the detonations well.

7 Conclusion

In this paper, we constructed high-order BP FD methods for MMD. The key point is to construct consistent numerical fluxes and apply conservative time integrations. We compared the proposed work with the well-developed WENO algorithm, where the numerical fluxes are not consistent. We find that without special numerical techniques such as subcell resolutions, our methods yield better numerical approximations than the WENO algorithm.

Acknowledgements Jie Du is supported by the National Natural Science Foundation of China under Grant Number NSFC 11801302 and Tsinghua University Initiative Scientific Research Program. Yang Yang is supported by the NSF Grant DMS-1818467.

Compliance with Ethical Standards

Conflict of interest On behalf of all authors, the corresponding author states that there is no conflict of interest.

References

- Balsara, D.S., Shu, C.-W.: Monotonicity preserving weighted essentially non-oscillatory schemes with increasingly high order of accuracy. J. Comput. Phys. 160, 405–452 (2000)
- Bao, W., Jin, S.: The random projection method for stiff detonation capturing. SIAM J. Sci. Comput. 23, 1000–1025 (2001)
- Bao, W., Jin, S.: The random projection method for stiff multispecies detonation capturing. J. Comput. Phys. 178, 37–57 (2002)
- Bihari, B., Schwendeman, D.: Multiresolution schemes for the reactive Euler equations. J. Comput. Phys. 154, 197–230 (1999)
- Christlieb, A., Liu, Y., Tang, Q., Xu, Z.: High order parametrized maximum-principle-preserving and positivity-preserving WENO schemes on unstructured meshes. J. Comput. Phys. 281, 334–351 (2015)
- Chuenjarern, N., Xu, Z., Yang, Y.: High-order bound-preserving discontinuous Galerkin methods for compressible miscible displacements in porous media on triangular meshes. J. Comput. Phys. 378, 110–128 (2019)
- Clarke, J.F., Karni, S., Quirk, J.J., Roe, P.L., Simmonds, L.G., Toro, E.F.: Numerical computation of two-dimensional unsteady detonation waves in high energy solids. J. Comput. Phys. 106, 215–233 (1993)
- Du, J., Wang, C., Qian, C., Yang, Y.: High-order bound-preserving discontinuous Galerkin methods for stiff multispecies detonation. SIAM J. Sci. Comput. 41, B250–B273 (2019)
- Du, J., Yang, Y.: Third-order conservative sign-preserving and steady-state-preserving time integrations and applications in stiff multispecies and multireaction detonations. J. Comput. Phys. 395, 489–510 (2019)
- Guo, H., Yang, Y.: Bound-preserving discontinuous Galerkin method for compressible miscible displacement problem in porous media. SIAM J. Sci. Comput. 39, A1969–A1990 (2017)
- 11. Guo, H., Liu, X., Yang, Y.: High-order bound-preserving finite difference methods for miscible displacements in porous media. J. Comput. Phys. **406**(24), 109219 (2020)
- Huang, J., Shu, C.-W.: Bound-preserving modified exponential Runge-Kutta discontinuous Galerkin methods for scalar hyperbolic equations with stiff source terms. J. Comput. Phys. 361, 111–135 (2018)
- 13. Huang, J., Shu, C.-W.: Positivity-preserving time discretizations for production-destruction equations with applications to non-equilibrium flows. J. Sci. Comput. **78**, 1811–1839 (2019)
- Huang, J., Zhao, W., Shu, C.-W.: A third-order unconditionally positivity-preserving scheme for production-destruction equations with applications to non-equilibrium flows. J. Sci. Comput. 79, 1015–1056 (2019)
- Jiang, G., Shu, C.-W.: Efficient implementation of weighted ENO schemes. J. Comput. Phys. 126, 202–228 (1996)
- Kopecz, S., Meister, A.: On order conditions for modified Patankar-Runge-Kutta schemes. Appl. Numer. Math. 123, 159–179 (2018)
- Kopecz, S., Meister, A.: Unconditionally positive and conservative third order modified Patankar-Runge-Kutta discretizations of production-destruction systems. BIT Numer. Math. 58, 691–728 (2018)

- LeVeque, R.J., Yee, H.C.: A study of numerical methods for hyperbolic conservation laws with stiff source terms. J. Comput. Phys. 86, 187–210 (1990)
- Liu, X.-D., Osher, S., Chan, T.: Weighted essentially non-oscillatory schemes. J. Comput. Phys. 115, 200–212 (1994)
- Shu, C.-W.: Essentially non-oscillatory and weighted essentially non-oscillatory schemes for hyperbolic conservation laws. In: Quarteroni, A. (ed.) Advanced Numerical Approximation of Nonlinear Hyperbolic Equations. Lecture Notes in Mathematics, vol. 1697. Springer, Berlin, Heidelberg (1998)
- Shu, C.-W., Osher, S.: Efficient implementation of essentially non-oscillatory shock-capturing schemes. J. Comput. Phys. 77, 439–471 (1988)
- 22. Sun, Y., Engquist, B.: Heterogeneous multiscale methods for interface tracking of combustion fronts. Multiscale Model. Simul. 5, 532–563 (2006)
- Wang, W., Shu, C.-W., Yee, H.C., Kotov, D.V., Sjögreen, B.: High order finite difference methods with subcell resolution for stiff multispecies detonation capturing. Commun. Comput. Phys. 17, 317–336 (2015)
- Xiong, T., Qiu, J.-M., Xu, Z.: Parametrized positivity preserving flux limiters for the high order finite difference WENO scheme solving compressible Euler equations. J. Sci. Comput. 67, 1066– 1088 (2016)
- Xu, Z.: Parametrized maximum principle preserving flux limiters for high order schemes solving hyperbolic conservation laws: one dimensional scalar problem. Math. Comput. 83, 2213–2238 (2014)
- Xu, Z., Yang, Y., Guo, H.: High-order bound-preserving discontinuous Galerkin methods for wormhole propagation on triangular meshes. J. Comput. Phys. 390, 323–341 (2019)
- Yee, H.C., Kotov, D.V., Wang, W., Shu, C.-W.: Spurious behavior of shock-capturing methods by the fractional step approach: problems containing stiff source terms and discontinuities. J. Comput. Phys. 241, 266–291 (2013)
- Zhang, X., Shu, C.-W.: On positivity preserving high order discontinuous Galerkin schemes for compressible Euler equations on rectangular meshes. J. Comput. Phys. 229, 8918–8934 (2010)