



Entropy dissipative higher order accurate positivity preserving time-implicit discretizations for nonlinear degenerate parabolic equations

Fengna Yan ^{a,b}, J.J.W. Van der Vegt ^b, Yinhua Xia ^c, Yan Xu ^{c,*}

^a School of Mathematics, Hefei University of Technology, Hefei, Anhui, 230000, PR China

^b Department of Applied Mathematics, Mathematics of Computational Science Group, University of Twente, Enschede, 7500 AE, The Netherlands

^c School of Mathematics, University of Science and Technology of China, Hefei, Anhui, 230026, PR China

ARTICLE INFO

Keywords:

Local discontinuous Galerkin discretizations
DIRK methods
Nonlinear degenerate parabolic equations
Entropy dissipation
KKT limiter

ABSTRACT

We develop entropy dissipative higher order accurate local discontinuous Galerkin (LDG) discretizations coupled with Diagonally Implicit Runge–Kutta (DIRK) methods for nonlinear degenerate parabolic equations with a gradient flow structure. Using the simple alternating numerical flux, we construct DIRK-LDG discretizations that combine the advantages of higher order accuracy, entropy dissipation and proper long-time behavior. We theoretically prove the entropy dissipation of the implicit Euler-LDG discretization without any time-step restrictions when no positivity constraint is imposed. Next, in order to ensure the positivity of the numerical solution, we use the Karush–Kuhn–Tucker (KKT) limiter, which achieves a positive solution by coupling the positivity preserving KKT conditions with higher order accurate DIRK-LDG discretizations using Lagrange multipliers. In addition, mass conservation of the positivity-limited solution is ensured by imposing a mass conservation equality constraint to the KKT equations. Under a time step restriction, the unique solvability and entropy dissipation for implicit first order accurate in time, but higher order accurate in space, positivity-preserving LDG discretizations with periodic boundary conditions are proved, which provide a first theoretical analysis of the KKT limiter. Finally, numerical results demonstrate the higher order accuracy and entropy dissipation of the positivity-preserving DIRK-LDG discretizations for problems requiring a positivity limiter. In addition, we can observe from the numerical results that the implicit time-discrete methods alleviate the time-step restrictions needed for the stability of the numerical discretizations, which improves computational efficiency.

1. Introduction

Consider the following degenerate parabolic equation [1]

$$\begin{cases} u_t = \nabla \cdot (f(u)\nabla(\Psi(\mathbf{x}) + H'(u))), & \text{in } \Omega \times (0, T], \\ u(\mathbf{x}, 0) = u_0(\mathbf{x}), & \text{in } \Omega, \end{cases} \quad (1.1)$$

with zero-flux boundary condition

$$\nabla(\Psi(\mathbf{x}) + H'(u)) \cdot \mathbf{v} = 0, \quad \text{on } \partial\Omega \times (0, T], \quad (1.2)$$

* Corresponding author.

E-mail addresses: fnyan@hfut.edu.cn (F. Yan), j.j.w.vandervegt@utwente.nl (J.J.W. Van der Vegt), yhxia@ustc.edu.cn (Y. Xia), yxu@ustc.edu.cn (Y. Xu).

where Ω is an open bounded polygonally connected domain in $\mathbb{R}^d, d = 1, 2$, with unit outward normal vector \mathbf{v} at the boundary $\partial\Omega, u(\mathbf{x}, t) \geq 0$ represents a nonnegative density with time derivative denoted as u_t , and initial solution $u_0(\mathbf{x}) \geq 0, \Psi(\mathbf{x})$ is a given potential function for $\mathbf{x} \in \mathbb{R}^d$, and f, H are given functions such that

$$f : \mathbb{R}^+ \rightarrow \mathbb{R}^+, \quad H : \mathbb{R}^+ \rightarrow \mathbb{R}, \quad f(u)H''(u) \geq 0, \tag{1.3}$$

where \mathbb{R}^+ is the nonnegative real space. Here $f(u)H''(u)$ can vanish for certain values of u , resulting in degenerate cases. The entropy corresponding to (1.1) is defined by

$$E(u) = \int_{\Omega} (u\Psi(\mathbf{x}) + H(u))d\Omega. \tag{1.4}$$

Multiplying (1.1) with $\Psi(\mathbf{x}) + H'(u)$ and integrating over Ω , with the zero-flux boundary condition (1.2), together with (1.4), we obtain that the time derivative of the entropy satisfies

$$\frac{d}{dt} E(u) = - \int_{\Omega} f(u)|\nabla(\Psi(\mathbf{x}) + H'(u))|^2 d\Omega \leq 0. \tag{1.5}$$

System (1.1) can represent different physical problems, such as the porous media equation [2,3], the nonlinear nonlocal equation with a double-well potential [4], the nonlinear Fokker–Planck model for fermion and boson gases [5–7].

Recently, many numerical discretizations have been proposed for (1.1); e.g. mixed finite element methods [8], finite volume methods [1,4], DG methods [9–11] and LDG methods [3]. Regarding positivity preserving discretizations, Liu and Yu developed in [10,11], respectively, for the linear Fokker–Planck equation a maximum preserving DG scheme and an entropy dissipative DG scheme, but these discretizations cannot be directly applied to the general case given by (1.1). Liu and Wang subsequently developed in [9] an explicit Runge–Kutta (RK) time-discrete method for (1.1) in one dimension together with a positivity preserving high order accurate DG scheme under some Courant–Friedrichs–Lewy (CFL) constraints. For the porous media equation, an LDG discretization coupled with an explicit RK method was considered in [3], which is similar to the DG method in [9]. Still, it uses a special numerical flux to ensure the non-negativity of the numerical solution. Cheng and Shen in [12] propose a Lagrange multiplier approach to construct positivity preserving schemes for a class of parabolic equations, which is different from (1.1), but contains the porous media equation.

For the time-step τ and mesh size h , the condition $\tau = O(h^2)$ is needed for stability in [3,9]. Therefore, these explicit time discretizations suffer from severe time step restrictions, but there are currently no feasible positivity preserving time-implicit LDG discretizations for (1.1). In this paper, we present higher order accurate Diagonally Implicit Runge–Kutta (DIRK) LDG discretizations, which ensure positivity and mass conservation of the numerical solution without the severe time step restrictions of explicit methods.

The LDG method proposed by Cockburn and Shu in [13] has many advantages, including high parallelizability, high order accuracy, a simple choice of trial and test spaces and easy handling of complicated geometries. We refer to [14–17] for examples of applications of the LDG method.

For many physical problems, it is crucial that the numerical discretization preserves the positivity properties of the partial differential equations (PDEs). Not only is this necessary to obtain physically meaningful solutions, but also negative values may result in ill-posedness of the problem and divergence of the numerical discretization. Positivity preserving DG methods have been extensively studied by many mathematicians. However, most positivity preserving DG methods are combined with explicit time-discretizations [9,18–20], for which numerical stability frequently imposes severe time step restrictions. These severe time-step constraints make explicit methods impractical for parabolic PDEs, such as (1.1).

Recently, Qin and Shu extended in [21] the general framework for establishing positivity-preserving schemes, proposed in [19,20], from explicit to implicit time discretizations. They developed a positivity preserving DG method with high-order spatial accuracy combined with the first-order backward Euler implicit temporal discretization for one-dimensional conservation laws. This approach requires, however, a detailed analysis of the numerical discretization to ensure positivity and it is not straightforward to extend this approach to higher order accurate time-implicit methods. Huang and Shen in [22] constructed higher order linear bound preserving implicit discretizations for the Keller–Segel and Poisson–Nernst–Planck equations. Van der Vegt, Xia and Xu proposed in [23] the KKT limiter concept to construct positivity preserving time-implicit discretizations. The KKT limiter in [23] is obtained by coupling the inequality and equality constraints imposed by the physical problem with higher order accurate DIRK-DG discretizations using Lagrange multipliers. The resulting semi-smooth nonlinear equations are solved by an efficient active set semi-smooth Newton method.

In this paper, we consider a general class of nonlinear degenerate parabolic equations given by (1.1) and aim at developing higher order accurate entropy dissipative and positivity preserving time-implicit LDG discretizations. For the spatial discretization, we use an LDG method with simple alternating numerical fluxes, which results in entropy dissipation of the semi-discrete LDG discretization. For the temporal discretization, we consider DIRK methods, which significantly enlarge the time step for stability. Without any time-step restrictions, the entropy dissipation of the LDG discretization combined with an implicit Euler time integration method is proved theoretically. We construct positivity preserving discretizations using the KKT limiter by imposing the positivity constraint on the numerical discretization using Lagrange multipliers. Under a time-step restriction, the unique solvability of the resulting positivity preserving KKT system is proved. In addition, we also prove the entropy dissipation of the positivity preserving LDG discretization when it is combined with the backward Euler time integration method. Numerical results demonstrate the accuracy and entropy dissipation of the higher order accurate positivity preserving DIRK-LDG discretizations.

This paper is organized as follows. In Section 2, we present the semi-discrete LDG discretization with simple alternating numerical fluxes for the nonlinear degenerate parabolic equation stated in (1.1) and prove that the numerical approximation is entropy

dissipative. Higher order accurate DIRK-LDG discretizations, which enlarge the stable time step to a great extent, are discussed in Section 3. Without any time-step restrictions, the entropy dissipation of the implicit Euler LDG discretizations is proved in Section 3.1. In order to ensure the positivity of the numerical solution and mass conservation of the positivity limited numerical discretizations, we introduce in Section 4.1 the KKT system. The higher order DIRK-LDG discretizations with positivity and mass conservation constraints are formulated in Section 4.2 as a KKT mixed complementarity problem. Under a time-step restriction, the unique solvability and entropy dissipation of the algebraic system resulting from a time implicit Euler-LDG discretization with positivity constraint are proved in Section 4.3. In Section 5, numerical results demonstrate the higher order accuracy, positivity and entropy dissipation of the positivity preserving DIRK-LDG discretizations. Concluding remarks are given in Section 6.

2. Semi-discrete LDG schemes

2.1. Definitions, notations

Let \mathcal{T}_h be a shape-regular tessellation of $\Omega \subset \mathbb{R}^d$, $d = 1, 2$, with line or convex quadrilateral elements K . Given the reference element $\hat{K} = [-1, 1]^d$. Let $\mathcal{Q}_k(\hat{K})$ denote the space composed of the tensor product of Legendre polynomials $\mathcal{P}_k(\hat{K})$ on $[-1, 1]$ of degree at most $k \geq 0$. The space $\mathcal{Q}_k(K)$ is obtained by using an isoparametric transformation from element K to the reference element \hat{K} . The finite element spaces V_h^k and \mathbf{W}_h^k are defined by

$$V_h^k = \{v \in L^2(\Omega) : v|_K \in \mathcal{Q}_k(K), \forall K \in \mathcal{T}_h\},$$

$$\mathbf{W}_h^k = \{\mathbf{w} \in [L^2(\Omega)]^d : \mathbf{w}|_K \in [\mathcal{Q}_k(K)]^d, \forall K \in \mathcal{T}_h\},$$

and are allowed to have discontinuities across element interfaces. Let e be an interior edge connected to the “left” and “right” elements denoted, respectively, by K_L and K_R . If u is a function on K_L and K_R , we set $u^L := (u|_{K_L})|_e$ and $u^R := (u|_{K_R})|_e$ for the left and right trace of u at e .

Note that $L^1(\Omega)$, $L^2(\Omega)$ and $L^\infty(\Omega)$ are Lebesgue spaces, $\|u\|_{L^2(\Omega)}$ is the $L^2(\Omega)$ -norm and $(\cdot, \cdot)_\Omega$ is the $L^2(\Omega)$ inner product. For simplicity, we denote the inner product as $(u, v) := (u, v)_\Omega$.

2.2. LDG discretization in space

For the LDG discretization of (1.1), we first rewrite this equation as a first order system

$$u_t = \nabla \cdot \mathbf{q},$$

$$\mathbf{q} = f(u)\mathbf{s},$$

$$\mathbf{s} = \nabla p,$$

$$p = \Psi(\mathbf{x}) + H'(u).$$

Then, the LDG discretization can be readily obtained by multiplying the above equations with arbitrary test functions, integrating by parts over each element $K \in \mathcal{T}_h$, and finally a summation of element and face contributions. The LDG discretization can be stated as: find $u_h, p_h \in V_h^k$, $\mathbf{q}_h, \mathbf{s}_h \in \mathbf{W}_h^k$, such that for all $\rho, \varphi \in V_h^k$ and $\boldsymbol{\theta}, \boldsymbol{\eta} \in \mathbf{W}_h^k$, we have

$$(u_h, \rho) + L_h^1(\mathbf{q}_h; \rho) = 0, \tag{2.1a}$$

$$(\mathbf{q}_h, \boldsymbol{\theta}) + L_h^2(u_h, \mathbf{s}_h; \boldsymbol{\theta}) = 0, \tag{2.1b}$$

$$(\mathbf{s}_h, \boldsymbol{\eta}) + L_h^3(p_h; \boldsymbol{\eta}) = 0, \tag{2.1c}$$

$$(p_h, \varphi) + L_h^4(u_h; \varphi) = 0, \tag{2.1d}$$

where

$$L_h^1(\mathbf{q}_h; \rho) := (\mathbf{q}_h, \nabla \rho) - \sum_{K \in \mathcal{T}_h} (\hat{\mathbf{q}}_h \cdot \mathbf{v}, \rho)_{\partial K}, \tag{2.2a}$$

$$L_h^2(u_h, \mathbf{s}_h; \boldsymbol{\theta}) := -(f(u_h)\mathbf{s}_h, \boldsymbol{\theta}), \tag{2.2b}$$

$$L_h^3(p_h; \boldsymbol{\eta}) := (p_h, \nabla \cdot \boldsymbol{\eta}) - \sum_{K \in \mathcal{T}_h} (\hat{p}_h, \mathbf{v} \cdot \boldsymbol{\eta})_{\partial K}, \tag{2.2c}$$

$$L_h^4(u_h; \varphi) := -(\Psi(\mathbf{x}) + H'(u_h), \varphi). \tag{2.2d}$$

Note that \mathbf{v} is the unit outward normal vector of an element K at its boundary ∂K . The “hat” terms in L_h^1 and L_h^3 are the so-called “numerical fluxes”, whose choices play an important role in ensuring stability. We remark that the choices for the numerical fluxes are not unique. Here, we use the alternating numerical fluxes

$$\hat{\mathbf{q}}_h = \mathbf{q}_h^R, \quad \hat{p}_h = p_h^L, \tag{2.3}$$

or

$$\hat{\mathbf{q}}_h = \mathbf{q}_h^L, \quad \hat{p}_h = p_h^R. \tag{2.4}$$

Considering the zero-flux boundary condition $\nabla(\Psi(\mathbf{x}) + H'(u)) \cdot \mathbf{v} = 0$, we take

$$\hat{\mathbf{q}}_h \cdot \mathbf{v} = 0, \quad p_h = (p_h)^{in} \tag{2.5}$$

at $\partial\Omega$, where ‘‘in’’ refers to the value obtained by taking the boundary trace from the inside of the domain Ω .

2.3. Entropy dissipation

Theorem 2.1. For $u_h \in V_h^k, \mathbf{s}_h \in \mathbf{W}_h^k$, the LDG scheme (2.1)–(2.5) with f satisfying (1.3) is entropy dissipative and satisfies

$$\frac{d}{dt} E(u_h) = -(f(u_h)\mathbf{s}_h, \mathbf{s}_h) \leq 0,$$

which is consistent with the entropy dissipation property (1.5) of the PDE (1.1).

Proof. By taking

$$\rho = p_h, \quad \boldsymbol{\theta} = -\mathbf{s}_h, \quad \boldsymbol{\eta} = \mathbf{q}_h, \quad \varphi = -u_{ht},$$

in (2.1a)–(2.1d), respectively, and after integration by parts, we have

$$\begin{aligned} & (\Psi(\mathbf{x}) + H'(u_h), u_{ht}) \\ &= -(f(u_h)\mathbf{s}_h, \mathbf{s}_h) - (\mathbf{q}_h, \nabla p_h) + \sum_{K \in \mathcal{T}_h} (\hat{\mathbf{q}}_h \cdot \mathbf{v}, p_h)_{\partial K} - (p_h, \nabla \cdot \mathbf{q}_h) + \sum_{K \in \mathcal{T}_h} (\hat{p}_h, \mathbf{v} \cdot \mathbf{q}_h)_{\partial K} \\ &= -(f(u_h)\mathbf{s}_h, \mathbf{s}_h) - \sum_{K \in \mathcal{T}_h} (\mathbf{q}_h \cdot \mathbf{v}, p_h)_{\partial K} + \sum_{K \in \mathcal{T}_h} (\hat{\mathbf{q}}_h \cdot \mathbf{v}, p_h)_{\partial K} + \sum_{K \in \mathcal{T}_h} (\hat{p}_h, \mathbf{v} \cdot \mathbf{q}_h)_{\partial K}. \end{aligned} \tag{2.6}$$

Assume that e is an interior edge shared by elements K_L and K_R , then $\mathbf{v}^R = -\mathbf{v}^L$, and together with the numerical fluxes (2.3), we obtain

$$\begin{aligned} & - \sum_{K_L \cup K_R} (\mathbf{q}_h \cdot \mathbf{v}, p_h)_e + \sum_{K_L \cup K_R} (\hat{\mathbf{q}}_h \cdot \mathbf{v}, p_h)_e + \sum_{K_L \cup K_R} (\hat{p}_h, \mathbf{v} \cdot \mathbf{q}_h)_e \\ &= -(\mathbf{q}_h^L \cdot \mathbf{v}^L, p_h^L)_e + (\mathbf{q}_h^R \cdot \mathbf{v}^L, p_h^R)_e + (\mathbf{q}_h^R \cdot \mathbf{v}^L, p_h^L)_e - (\mathbf{q}_h^R \cdot \mathbf{v}^L, p_h^R)_e \\ & \quad + (\mathbf{q}_h^L \cdot \mathbf{v}^L, p_h^L)_e - (\mathbf{q}_h^R \cdot \mathbf{v}^L, p_h^L)_e = 0. \end{aligned} \tag{2.7}$$

Combining (2.6)–(2.7), using (1.4), boundary conditions (2.5) and the condition on f (1.3), we get

$$\frac{d}{dt} E(u_h) = (\Psi(\mathbf{x}) + H'(u_h), u_{ht}) = -(f(u_h)\mathbf{s}_h, \mathbf{s}_h) \leq 0. \quad \square$$

Remark 2.1. For brevity, we will only consider in the remaining article the numerical fluxes (2.3) and omit the discussion of the numerical fluxes (2.4), but all results also apply to the numerical fluxes (2.4).

Remark 2.2. Compared to the spatial discretizations in [3,9], we choose the simpler alternating numerical fluxes (2.3) and (2.4), which significantly simplifies the theoretical analysis of the entropy dissipation property of the LDG discretization.

3. Time-implicit LDG schemes

The numerical discretization of the nonlinear parabolic Eqs. (1.1) using explicit time discretization methods suffers from the rather severe time-step constraint $\tau = O(h_{\mathcal{T}}^2)$, with $h_{\mathcal{T}}$ the mesh size for the tessellation \mathcal{T}_h . In this section, we will discuss implicit time discretizations, which will be coupled with positivity constraints in Section 4.

We divide the time interval $[0, T]$ into N parts $0 = t_0 < t_1 < \dots < t_N = T$, with $\tau^n = t_n - t_{n-1}$ ($n = 1, 2, \dots, N$). For $n = 0, 1, \dots, N$, let $u_n = u(\cdot, t_n)$ and u_h^n , respectively, denote the exact and approximate values of u at time t_n .

3.1. Backward Euler LDG discretization

Discretizing (2.1) in time with the implicit Euler method gives the following discrete system

$$\left(\frac{u_h^{n+1} - u_h^n}{\tau^{n+1}}, \rho \right) + L_h^1(\mathbf{q}_h^{n+1}; \rho) = 0, \tag{3.1a}$$

$$(\mathbf{q}_h^{n+1}, \boldsymbol{\theta}) + L_h^2(u_h^{n+1}, \mathbf{s}_h^{n+1}; \boldsymbol{\theta}) = 0, \tag{3.1b}$$

$$(\mathbf{s}_h^{n+1}, \boldsymbol{\eta}) + L_h^3(p_h^{n+1}; \boldsymbol{\eta}) = 0, \tag{3.1c}$$

$$(p_h^{n+1}, \varphi) + L_h^4(u_h^{n+1}; \varphi) = 0. \tag{3.1d}$$

Define the discrete entropy as

$$E_h(u_h^n) = \int_{\Omega} (u_h^n \Psi(\mathbf{x}) + H(u_h^n)) dx. \tag{3.2}$$

We have the following relation for the discrete entropy dissipation.

Theorem 3.1. For all time levels n , the numerical solutions $u_h^n, u_h^{n+1} \in V_h^k$ of the LDG discretization (3.1), with boundary condition (2.5) and conditions on f, H stated in (1.3), satisfy the following entropy dissipation relation

$$E_h(u_h^{n+1}) \leq E_h(u_h^n), \tag{3.3}$$

which implies that the LDG discretization is entropy dissipative without any time-step restrictions.

Proof. By choosing, respectively, in (3.1a)–(3.1d) the following test functions

$$\rho = p_h^{n+1}, \quad \theta = -s_h^{n+1}, \quad \eta = \mathbf{q}_h^{n+1}, \quad \varphi = -\frac{u_h^{n+1} - u_h^n}{\tau^{n+1}},$$

we get

$$\begin{aligned} & \left(\Psi(\mathbf{x}), \frac{u_h^{n+1} - u_h^n}{\tau^{n+1}} \right) + \left(H'(u_h^{n+1}), \frac{u_h^{n+1} - u_h^n}{\tau^{n+1}} \right) \\ &= - (f(u_h^{n+1}) \mathbf{s}_h^{n+1}, \mathbf{s}_h^{n+1}) - (\mathbf{q}_h^{n+1}, \nabla p_h^{n+1}) + \sum_{K \in \mathcal{T}_h} (\hat{\mathbf{q}}_h^{n+1} \cdot \mathbf{v}, p_h^{n+1})_{\partial K} \\ & \quad - (p_h^{n+1}, \nabla \cdot \mathbf{q}_h^{n+1}) + \sum_{K \in \mathcal{T}_h} (\hat{p}_h^{n+1}, \mathbf{v} \cdot \mathbf{q}_h^{n+1})_{\partial K} \\ &= - (f(u_h^{n+1}) \mathbf{s}_h^{n+1}, \mathbf{s}_h^{n+1}) - \sum_{K \in \mathcal{T}_h} (\mathbf{q}_h^{n+1} \cdot \mathbf{v}, p_h^{n+1})_{\partial K} + \sum_{K \in \mathcal{T}_h} (\hat{\mathbf{q}}_h^{n+1} \cdot \mathbf{v}, p_h^{n+1})_{\partial K} \\ & \quad + \sum_{K \in \mathcal{T}_h} (\hat{p}_h^{n+1}, \mathbf{v} \cdot \mathbf{q}_h^{n+1})_{\partial K}. \end{aligned}$$

Together with (2.7), the numerical fluxes (2.3) and the boundary condition (2.5), we obtain then

$$\left(\Psi(\mathbf{x}), \frac{u_h^{n+1} - u_h^n}{\tau^{n+1}} \right) + \left(H'(u_h^{n+1}), \frac{u_h^{n+1} - u_h^n}{\tau^{n+1}} \right) = - (f(u_h^{n+1}) \mathbf{s}_h^{n+1}, \mathbf{s}_h^{n+1}).$$

Because of the following Taylor expansion

$$H(u_h^n) = H(u_h^{n+1}) + H'(u_h^{n+1})(u_h^n - u_h^{n+1}) + \frac{1}{2} H''(\xi^{n+1})(u_h^{n+1} - u_h^n)^2, \quad \xi^{n+1} \in (u_h^n, u_h^{n+1}),$$

we have, using the conditions on f, H stated in (1.3) and the definition of E_h in (3.2),

$$\begin{aligned} E_h(u_h^{n+1}) - E_h(u_h^n) &= (\Psi(\mathbf{x}), u_h^{n+1} - u_h^n) + (H(u_h^{n+1}) - H(u_h^n), 1) \\ &= -\tau^{n+1} (f(u_h^{n+1}) \mathbf{s}_h^{n+1}, \mathbf{s}_h^{n+1}) - \frac{1}{2} (H''(\xi^{n+1}), (u_h^{n+1} - u_h^n)^2) \\ &\leq 0. \quad \square \end{aligned}$$

3.2. Higher order DIRK-LDG discretizations

For higher order accurate implicit in time discretizations of the system (2.1), we use a Diagonally Implicit Runge–Kutta (DIRK) method [24]. Assuming we know the numerical solution at time level n , we obtain the solution at time level $n + 1$ with a DIRK method by solving for each DIRK stage $i, i = 1, 2, \dots, s$ the following equations.

$$\left(\frac{u_h^{n+1,i} - u_h^n}{\tau^{n+1}}, \rho \right) + \sum_{j=1}^i a_{ij} L_h^1(\mathbf{q}_h^{n+1,j}; \rho) = 0, \tag{3.4a}$$

$$(\mathbf{q}_h^{n+1,i}, \theta) + L_h^2(u_h^{n+1,i}, \mathbf{s}_h^{n+1,i}; \theta) = 0, \tag{3.4b}$$

$$(\mathbf{s}_h^{n+1,i}, \eta) + L_h^3(p_h^{n+1,i}; \eta) = 0, \tag{3.4c}$$

$$(p_h^{n+1,i}, \varphi) + L_h^4(u_h^{n+1,i}; \varphi) = 0. \tag{3.4d}$$

Then the solution at time t_{n+1} is

$$(u_h^{n+1}, \rho) = (u_h^n, \rho) - \tau^{n+1} \sum_{i=1}^s b_i L_h^1(\mathbf{q}_h^{n+1,i}; \rho). \tag{3.5}$$

The coefficient matrices (a_{ij}) in (3.4a) and (b_i) in (3.5) for the second, third and fourth order accurate DIRK methods used in this paper are

- Second order DIRK method [25, Theorem 5]

$$(a_{ij}) = \begin{pmatrix} \alpha & 0 \\ 1-\alpha & \alpha \end{pmatrix}, (b_i) = (1-\alpha \quad \alpha), (c_i) = (\alpha \quad 1), \tag{3.6}$$

where $\alpha = 1 - \frac{\sqrt{2}}{2}$.

- Third order DIRK method [26, top of page 2117]

$$(a_{ij}) = \begin{pmatrix} \gamma & 0 & 0 \\ 1/2 - \gamma/2 & \gamma & 0 \\ 1 - \delta - \gamma & \delta & \gamma \end{pmatrix}, (b_i) = (1 - \delta - \gamma \quad \delta \quad \gamma),$$

$$(c_i) = (\gamma \quad 1/2 + \gamma/2 \quad 1), \tag{3.7}$$

where $\gamma = 0.435866521508$, $\delta = 0.25(5 - 20\gamma + 6\gamma^2)$.

- Fourth order DIRK method [26, top of page 2118]

$$(a_{ij}) = \begin{pmatrix} 1/4 & 0 & 0 & 0 & 0 \\ -1/4 & 1/4 & 0 & 0 & 0 \\ 1/8 & 1/8 & 1/4 & 0 & 0 \\ -3/2 & 3/4 & 3/2 & 1/4 & 0 \\ 0 & 1/6 & 2/3 & -1/12 & 1/4 \end{pmatrix},$$

$$(b_j) = (0 \quad 1/6 \quad 2/3 \quad -1/12 \quad 1/4),$$

$$(c_i) = (1/4 \quad 0 \quad 1/2 \quad 1 \quad 1). \tag{3.8}$$

For DG discretizations using polynomial basis functions of order k we use a DIRK method of order $k + 1$.

The above DIRK methods satisfy $a_{si} = b_i$, $i = 1, 2, \dots, s$, which implies $u_h^{n+1} = u_h^{n+1,s}$. The above time discretization methods are easy to implement since the matrix (a_{ij}) in the DIRK methods has a lower triangular structure, which means that we can compute the DIRK stages one after another, starting from $i = 1$ up to $i = s$. For detailed information about the DIRK time integration method, we refer to [24].

4. Higher order accurate positivity preserving DIRK-LDG discretizations

The positivity constraints on the LDG solution will be enforced by transforming the DIRK-LDG equations with positivity constraints into a mixed complementarity problem using the Karush–Kuhn–Tucker (KKT) equations [27]. In the following sections, we will first define the positivity preserving DIRK-LDG discretization. Next, we will consider the unique solvability and entropy dissipation of the discrete KKT system.

4.1. KKT-system

For the KKT equations [27], we define the set

$$\mathbb{K} := \{\tilde{U} \in \mathbb{R}^{dof} \mid h(\tilde{U}) = 0, g(\tilde{U}) \leq 0\}, \tag{4.1}$$

with equality constraints $h : \mathbb{R}^{dof} \rightarrow \mathbb{R}^l$ and inequality constraints $g : \mathbb{R}^{dof} \rightarrow \mathbb{R}^m$ being vector-valued continuously differentiable functions. In the following $dof = N_k \cdot N_e$, where N_k is the number of basis functions in one element and N_e the number of elements in the tessellation \mathcal{T}_h of the domain Ω . The inequality constraints g in (4.8) are used to ensure positivity. The equality constraint h in (4.9) ensures that the limited DIRK-LDG discretization is mass conservative. Mass conservation is a property of the unlimited DIRK-LDG discretization, but one has to ensure that this property also holds after applying the positivity preserving limiter.

We assume that L is a continuously differentiable function from \mathbb{K} to \mathbb{R}^{dof} . The function L representing the LDG discretization for each DIRK stage $i = 1, 2, \dots, s$ is given by (4.5)–(4.6). The corresponding KKT-system [27] then is

$$L(\tilde{U}) + \nabla_{\tilde{U}} h(\tilde{U})^T \mu + \nabla_{\tilde{U}} g(\tilde{U})^T \lambda = 0, \tag{4.2a}$$

$$h(\tilde{U}) = 0, \tag{4.2b}$$

$$0 \geq g(\tilde{U}) \perp \lambda \geq 0, \tag{4.2c}$$

where $\mu \in \mathbb{R}^l$ and $\lambda \in \mathbb{R}^m$ are the Lagrange multipliers used to ensure $h(\tilde{U}) = 0$ and $g(\tilde{U}) \leq 0$, respectively, $\tilde{U} \in \mathbb{R}^{dof}$ are the LDG coefficients in the positivity-preserving DIRK-LDG discretization, and $\nabla_{\tilde{U}}$ denotes the gradient with respect to \tilde{U} . The compatibility condition (4.2c) implies that $g(\tilde{U}) \leq 0$, $\lambda \geq 0$ and $g(\tilde{U})^T \lambda = 0$, which can be expressed as

$$\min(-g_j(\tilde{U}), \lambda_j) = 0, \quad j = 1, 2, \dots, m.$$

The KKT-system then can be formulated as

$$F(z) = \begin{pmatrix} L(\tilde{U}) + \nabla_{\tilde{U}} h(\tilde{U})^T \mu + \nabla_{\tilde{U}} g(\tilde{U})^T \lambda \\ h(\tilde{U}) \\ \min(-g_1(\tilde{U}), \lambda_1) \\ \vdots \\ \min(-g_m(\tilde{U}), \lambda_m) \end{pmatrix} = 0. \tag{4.3}$$

Here $z = (\tilde{U}, \mu, \lambda) \in \mathbb{R}^{dof+l+m}$, and $F : \mathbb{R}^{dof+l+m} \rightarrow \mathbb{R}^{dof+l+m}$ represents the DIRK-LDG discretization combined with the positivity and mass conservation constraints. Note, the KKT system (4.3) is nonlinear and $F(z)$ is not continuously differentiable, as is necessary for standard Newton methods, but semi-smooth. We will therefore solve (4.3) with the active set semi-smooth Newton method presented in [23].

4.2. Positivity preserving LDG discretizations

In this section, we will provide the details of the higher order accurate positivity preserving DIRK-LDG discretizations (3.4) coupled with the positivity and mass conservation constraints using Lagrange multipliers as stated in (4.2).

We introduce the following notation for the element-wise positivity preserving LDG solution

$$U_h|_K := \sum_{j=1}^{N_k} \tilde{U}_j^K \phi_j^K, \quad \mathcal{Q}_h|_K := \sum_{j=1}^{N_k} \tilde{\mathcal{Q}}_j^K \phi_j^K$$

with $K \in \mathcal{T}_h$, ϕ_j^K the tensor product Legendre basis functions in $\mathcal{Q}_k(K)$, and LDG coefficients $\tilde{U}_j^K \in \mathbb{R}$, $\tilde{\mathcal{Q}}_j^K \in \mathbb{R}^d$. We denote \tilde{U} and $\tilde{\mathcal{Q}}$ by

$$\tilde{U} = \begin{pmatrix} \tilde{U}_1^{K_1} \\ \vdots \\ \tilde{U}_{N_k}^{K_1} \\ \vdots \\ \tilde{U}_1^{K_{N_e}} \\ \vdots \\ \tilde{U}_{N_k}^{K_{N_e}} \end{pmatrix} \in \mathbb{R}^{dof}, \quad \tilde{\mathcal{Q}} = \begin{pmatrix} \tilde{\mathcal{Q}}_1^{K_1} \\ \vdots \\ \tilde{\mathcal{Q}}_{N_k}^{K_1} \\ \vdots \\ \tilde{\mathcal{Q}}_1^{K_{N_e}} \\ \vdots \\ \tilde{\mathcal{Q}}_{N_k}^{K_{N_e}} \end{pmatrix} \in \mathbb{R}^{d \cdot dof}.$$

with $K_l \in \mathcal{T}_h$ for all $l = 1, \dots, N_e$.

Taking in each element $K \in \mathcal{T}_h$ the test function $\rho = \phi_j^K$, $j = 1, 2, \dots, N_k$ in the operator $L_h^1(\mathcal{Q}_h; \rho)$, stated in (2.2a), we can define

$$\mathbb{L}_h^1(\tilde{\mathcal{Q}}) := \begin{pmatrix} L_h^1(\mathcal{Q}_h; \phi_1^{K_1}) \\ \vdots \\ L_h^1(\mathcal{Q}_h; \phi_{N_k}^{K_1}) \\ \vdots \\ L_h^1(\mathcal{Q}_h; \phi_1^{K_{N_e}}) \\ \vdots \\ L_h^1(\mathcal{Q}_h; \phi_{N_k}^{K_{N_e}}) \end{pmatrix} \in \mathbb{R}^{N_k N_e}, \tag{4.4}$$

with $K_l \in \mathcal{T}_h$ for all $l = 1, \dots, N_e$. We use similar definitions of \mathbb{L}_h^k for L_h^k , $k = 2, 3, 4$ stated in (2.2b)–(2.2d).

Representing the block-diagonal mass matrices for the scalar and vector variables as $M \in \mathbb{R}^{N_k N_e \times N_k N_e}$ and $\mathbf{M} \in \mathbb{R}^{d N_k N_e \times d N_k N_e}$, respectively, the operator L for DIRK stage i ($i = 1, 2, \dots, s$), as stated in (3.4a), can be expressed as

$$L(\tilde{U}^{n+1,i}) := M(\tilde{U}^{n+1,i} - \tilde{U}^n) + \tau^{n+1} \sum_{j=1}^i a_{ij} \mathbb{L}_h^1(\tilde{\mathcal{Q}}^{n+1,j}), \tag{4.5}$$

with LDG coefficients $\tilde{U}^{n+1,i} \in \mathbb{R}^{N_k N_e}$. The function L in (4.5) is constructed using $\tilde{\mathcal{Q}}^{n+1,i}$, $\tilde{\mathcal{S}}^{n+1,i}$ and $\tilde{P}^{n+1,i}$, which are obtained using (3.4b), (3.4c) and (3.4d) and are defined as

$$\tilde{\mathcal{Q}}^{n+1,i} = -\mathbf{M}^{-1} \mathbb{L}_h^2(\tilde{U}^{n+1,i}, \tilde{\mathcal{S}}^{n+1,i}), \tag{4.6a}$$

$$\tilde{\mathcal{S}}^{n+1,i} = -\mathbf{M}^{-1} \mathbb{L}_h^3(\tilde{P}^{n+1,i}), \tag{4.6b}$$

$$\tilde{P}^{n+1,i} = -M^{-1} \mathbb{L}_h^4(\tilde{U}^{n+1,i}), \tag{4.6c}$$

with LDG coefficients $\tilde{\mathcal{Q}}^{n+1,i} \in \mathbb{R}^{d N_k N_e}$, $\tilde{\mathcal{S}}^{n+1,i} \in \mathbb{R}^{d N_k N_e}$, $\tilde{P}^{n+1,i} \in \mathbb{R}^{N_k N_e}$.

The constraints on the DIRK-LDG discretization can be directly imposed on the DG coefficients for each DIRK stage using the equality and inequality constraints in the KKT-system (4.3). We obtain for each DIRK stage i , with $i = 1, 2, \dots, s$, the LDG coefficients $\tilde{U}^{n+1,i}$ by solving the following KKT system for $\tilde{U}^{n+1,i}$,

$$\begin{pmatrix} L(\tilde{U}^{n+1,i}) + \nabla_{\tilde{U}} h(\tilde{U}^{n+1,i})^T \mu + \nabla_{\tilde{U}} g(\tilde{U}^{n+1,i})^T \lambda \\ h(\tilde{U}^{n+1,i}) \\ \min(-g_1(\tilde{U}^{n+1,i}), \lambda_1) \\ \vdots \\ \min(-g_m(\tilde{U}^{n+1,i}), \lambda_m) \end{pmatrix} = 0, \tag{4.7}$$

where the positivity preserving inequality constraint $g(\tilde{U}^{n+1,i})$ and the mass conservation equality constraint $h(\tilde{U}^{n+1,i})$ are defined as follows.

1. Positivity preserving inequality constraint

In each element $K \in \mathcal{T}_h$, we define the function g stated in (4.7) as

$$g_p^K(\tilde{U}^{n+1,i}) = u_{\min} - \sum_{j=1}^{N_k} \tilde{U}_j^{K,(n+1,i)} \phi_j^K(\mathbf{x}_p), \quad p = 1, \dots, N_p, \tag{4.8}$$

with N_p the number of Gauss–Lobatto quadrature points, and \mathbf{x}_p the Gauss–Lobatto quadrature points where the inequality constraints $U_h(\mathbf{x}_p) \geq u_{\min}$ are imposed. The use of Gauss–Lobatto quadrature rules ensures that the positivity constraint is also imposed in the computation of the numerical fluxes at the element edges where Gauss–Lobatto rules have, next to the element itself, also quadrature points. Note, the Gauss–Lobatto quadrature points \mathbf{x}_p are the only points used in the LDG discretization and the positivity constraint u_{\min} therefore only needs to be enforced at these points. In practice, one only needs to impose the inequality constraints on quadrature points where the solution is close to the bounds.

2. Mass conservation equality constraint

In order to ensure mass conservation of the LDG discretization when the positivity constraint is enforced, we impose the following equality constraint, which is obtained by setting $\rho = 1$ in (3.4a) and using the numerical flux (2.3) or (2.4).

$$\begin{aligned} h(\tilde{U}^{n+1,i}) &= \sum_{K \in \mathcal{T}_h} \int_K U_h^n dK + \tau^{n+1} \sum_{j=1}^i a_{ij} \sum_{\substack{K \in \mathcal{T}_h \\ \partial K \cap \partial \Omega \neq \emptyset}} (\hat{\mathcal{Q}}_h^{n+1,j} \cdot \mathbf{v}, 1)_{\partial K} \\ &\quad - \sum_{K \in \mathcal{T}_h} \sum_{j=1}^{N_k} \tilde{U}_j^{K,(n+1,i)} \int_K \phi_j^K(\mathbf{x}) dK \\ &= \mathbb{1}^T M(\tilde{U}^n - \tilde{U}^{n+1,i}) + \tau^{n+1} \sum_{j=1}^i a_{ij} \sum_{\substack{K \in \mathcal{T}_h \\ \partial K \cap \partial \Omega \neq \emptyset}} (\hat{\mathcal{Q}}_h^{n+1,j} \cdot \mathbf{v}, 1)_{\partial K}, \end{aligned} \tag{4.9}$$

with \tilde{U}^n the DG coefficients of U_h^n , the positivity-preserving DIRK-LDG solution at time t_n , $\mathbb{1} = (\mathbb{1}_{N_k}, \dots, \mathbb{1}_{N_k})^T \in \mathbb{R}^{dof}$ with $\mathbb{1}_{N_k} = (1, \underbrace{0, \dots, 0}_{N_k-1})$.

Remark 4.1. We compute $\nabla_{\tilde{U}} h(\tilde{U}^{n+1,i})^T$ for the term

$$\tau^{n+1} \sum_{j=1}^i a_{ij} \sum_{\substack{K \in \mathcal{T}_h \\ \partial K \cap \partial \Omega \neq \emptyset}} (\hat{\mathcal{Q}}_h^{n+1,j} \cdot \mathbf{v}, 1)_{\partial K} \neq 0,$$

in (4.9) using the chain rule, the relation between $\hat{\mathcal{Q}}_h^{n+1,i}$ and $\tilde{U}^{n+1,i}$ given by (4.6), and the numerical flux (2.3) or (2.4).

For each DIRK stage i , the KKT-system (4.7) for the higher order accurate positivity preserving LDG discretization is now defined. After solving the KKT Eqs. (4.7) for $i = 1, \dots, s$, the numerical solution at time t^{n+1} is directly obtained from the last DIRK stage, $U_h^{n+1} = U_h^{n+1,s}$ since we use DIRK methods with $a_{si} = b_i$.

Remark 4.2. In order to ensure the positivity of the discrete initial solution U_h^0 , we use the L^2 -projection coupled with the positivity constraint (4.8), which is obtained by replacing $\tilde{U}^{n+1,i}$ with \tilde{U}^0 . The equality constraint ensures mass conservation of the positivity limited initial solution

$$h(\tilde{U}^0) = \sum_{K \in \mathcal{T}_h} \int_K u_0(\mathbf{x}) dK - \sum_{K \in \mathcal{T}_h} \sum_{j=1}^{N_k} \tilde{U}_j^{K,0} \int_K \phi_j^K(\mathbf{x}) dK.$$

The constraints on the L^2 -projection are imposed using KKT equations similar to (4.3). To prevent pathological cases, we assume that the limited initial solution satisfies

$$\frac{1}{|\Omega|} \sum_{K \in \mathcal{T}_h} \int_K u_0(\mathbf{x}) dK \geq u_{\min}.$$

Remark 4.3. We emphasize that u_{\min} must be chosen strictly positive to ensure that errors do not violate the positivity of the numerical solution due to the finite precision of the computer arithmetic. For more details, we refer to the test cases in Section 5.

4.3. Unique solvability and stability of the positivity preserving LDG discretization

In Section 4.2, we have presented the positivity preserving LDG discretization for (1.1). In this section, we will consider the unique solvability of the algebraic equations resulting from the positivity-preserving backward Euler LDG discretization. In the theoretical analysis, we will also consider the entropy dissipation of the positivity preserving backward Euler LDG discretization and use periodic boundary conditions.

With (4.5)–(4.9), the positivity preserving backward Euler LDG discretization results now in the following KKT system,

$$L(\tilde{U}^{n+1}) + \nabla_{\tilde{U}} h(\tilde{U}^{n+1})^T \mu^{n+1} + \nabla_{\tilde{U}} g(\tilde{U}^{n+1})^T \lambda^{n+1} = 0, \tag{4.10a}$$

$$-h(\tilde{U}^{n+1}) = 0, \tag{4.10b}$$

$$\min(-g(\tilde{U}^{n+1}), \lambda^{n+1}) = 0. \tag{4.10c}$$

Here $L : \mathbb{R}^{N_k N_e} \rightarrow \mathbb{R}^{N_k N_e}$ and

$$L(\tilde{U}^{n+1}) := M(\tilde{U}^{n+1} - \tilde{U}^n) + \tau^{n+1} B\tilde{Q}^{n+1}, \tag{4.11}$$

$$M\tilde{Q}^{n+1} = C_d(\tilde{U}^{n+1})\tilde{S}^{n+1}, \tag{4.12}$$

$$M\tilde{S}^{n+1} = A\tilde{P}^{n+1}, \tag{4.13}$$

$$M\tilde{P}^{n+1} = D(\tilde{U}^{n+1}). \tag{4.14}$$

From (4.4)–(4.6), we obtain that

$$B\tilde{Q}^{n+1} = \mathbb{L}_h^1(\tilde{Q}^{n+1}) \in \mathbb{R}^{N_k N_e}, \tag{4.15}$$

$$C_d(\tilde{U}^{n+1})\tilde{S}^{n+1} = -\mathbb{L}_h^2(\tilde{U}^{n+1}, \tilde{S}^{n+1}) \in \mathbb{R}^{d N_k N_e}, \tag{4.16}$$

$$A\tilde{P}^{n+1} = -\mathbb{L}_h^3(\tilde{P}^{n+1}) \in \mathbb{R}^{d N_k N_e}, \tag{4.17}$$

$$D(\tilde{U}^{n+1}) = -\mathbb{L}_h^4(\tilde{U}^{n+1}) \in \mathbb{R}^{N_k N_e}, \tag{4.18}$$

where

$$C_d(\tilde{U}^{n+1}) = \begin{pmatrix} C(\tilde{U}^{n+1}) & & \\ & \ddots & \\ & & C(\tilde{U}^{n+1}) \end{pmatrix} \in \mathbb{R}^{d N_k N_e \times d N_k N_e}, \quad C(\tilde{U}^{n+1}) \in \mathbb{R}^{N_k N_e}. \tag{4.19}$$

The constraints $h : \mathbb{R}^{N_k N_e} \rightarrow \mathbb{R}$, $g : \mathbb{R}^{N_k N_e} \rightarrow \mathbb{R}^{N_p N_e}$ are defined by

$$h(\tilde{U}^{n+1}) = \sum_{K \in \mathcal{T}_h} \int_K U_h^0 dK - \sum_{K \in \mathcal{T}_h} \sum_{j=1}^{N_k} \tilde{U}_j^{K,(n+1)} \int_K \phi_j^K(\mathbf{x}) dK = \mathbb{1}^T M(\tilde{U}^0 - \tilde{U}^{n+1}), \tag{4.20}$$

$$g(\tilde{U}^{n+1}) = (g_1^{K_1}(\tilde{U}^{n+1}), \dots, g_{N_p}^{K_1}(\tilde{U}^{n+1}), \dots, g_1^{K_{N_e}}(\tilde{U}^{n+1}), \dots, g_{N_p}^{K_{N_e}}(\tilde{U}^{n+1})), \tag{4.21}$$

with the definition of the constraints $g_p^{K_j}$, $1 \leq p \leq N_p$, $1 \leq j \leq N_e$, given in (4.8), and $\mathbb{1} = (\mathbb{1}_{N_k}, \dots, \mathbb{1}_{N_k})^T \in \mathbb{R}^{dof}$ with $\mathbb{1}_{N_k} = \underbrace{(1, 0, \dots, 0)}_{N_k-1}$.

4.3.1. Auxiliary results used to prove the solvability of the KKT-system

In this section, we will introduce some auxiliary results, which will be used in Section 4.3.2 to prove the unique solvability of the KKT-system (4.10).

Definition 4.4 ([27, Sections 1.1, 3.2]). Let \mathbb{K} be given by (4.1). Given a map $L : \mathbb{K} \rightarrow \mathbb{R}^{dof}$, the Variational Inequality (VI(\mathbb{K} , L)) is to find $\tilde{U} \in \mathbb{K}$ such that

$$(y - \tilde{U})^T L(\tilde{U}) \geq 0, \quad y \in \mathbb{K}. \tag{4.22}$$

The set of solutions for VI(\mathbb{K} , L)(4.22) is denoted by SOL(\mathbb{K} , L).

Lemma 4.1 ([27, Proposition 1.3.4]). Let $\tilde{U} \in \text{SOL}(\mathbb{K}, L)$ solve (4.22) with \mathbb{K} given by (4.1). If Abadie’s Constraint Qualification holds at \tilde{U} , then there exist vectors $\mu \in \mathbb{R}^l$ and $\lambda \in \mathbb{R}^m$ satisfying the KKT system (4.10).

Lemma 4.2 ([27, Proposition 1.3.4]). If $(\tilde{U}, \mu, \lambda)$ satisfies (4.10), and if each function h_j ($1 \leq j \leq l$) is affine and each function g_i ($1 \leq i \leq m$) is convex, then \tilde{U} solves VI(\mathbb{K} , L) given by (4.22) with \mathbb{K} given by (4.1).

For the inequality constraints (4.8), the set \mathbb{K} given by (4.1), with equality constraint $h \equiv 0$, reduces to the following box constraint problem

$$\mathbb{K}_b := \{\tilde{U} \in \mathbb{R}^{dof} \mid \tilde{U}_i^{\min} \leq \tilde{U}_i \leq \tilde{U}_i^{\max}, i \in \{1, \dots, dof\}\}. \tag{4.23}$$

The values of the LDG coefficients \tilde{U}_i^{\min} follow from the inequality constraints g in (4.8), where u_{\min} depends on the physical constraint to be satisfied. A similar constraint can also be imposed for \tilde{U}_i^{\max} using u_{\max} for the maximum allowed physical constraints. Depending on the problem considered, \tilde{U}_i^{\max} can be infinity.

Remark 4.5. Abadie’s Constraint Qualification states that the tangent cone at $\tilde{U} \in \mathbb{K}$ must be equal to the linearization cone of \mathbb{K} at \tilde{U} . This is true for \mathbb{K}_b since the domain \mathbb{K}_b is a box (or polyhedral) domain, see [27, Section 1.3.1].

We write \mathbb{K}_b as

$$\mathbb{K}_b = \prod_{\vartheta=1}^N \mathbb{K}_{n_\vartheta}, \tag{4.24}$$

where \mathbb{K}_{n_ϑ} is a subset of \mathbb{R}^{n_ϑ} with $\sum_{\vartheta=1}^N n_\vartheta = dof$. Thus for a vector $\tilde{U} \in \mathbb{K}_b$, we write $\tilde{U} = (\tilde{U}_\vartheta)$, where each \tilde{U}_ϑ belongs to \mathbb{K}_{n_ϑ} .

Definition 4.6 ([27, Section 3.5.2]). Let \mathbb{K}_b be given by (4.23). A map $L : \mathbb{K}_b \rightarrow \mathbb{R}^{dof}$ is said to be

(a) a P-function on \mathbb{K}_b if for all pairs of distinct vectors \tilde{U} and \tilde{U}' in \mathbb{K}_b ,

$$\max_{1 \leq \vartheta \leq N} (\tilde{U}_\vartheta - \tilde{U}'_\vartheta)^T (L_\vartheta(\tilde{U}) - L_\vartheta(\tilde{U}')) > 0,$$

(b) a uniformly P-function on \mathbb{K}_b if there exists a constant $\varpi > 0$ such that for all pairs of distinct vectors \tilde{U} and \tilde{U}' in \mathbb{K}_b ,

$$\max_{1 \leq \vartheta \leq N} (\tilde{U}_\vartheta - \tilde{U}'_\vartheta)^T (L_\vartheta(\tilde{U}) - L_\vartheta(\tilde{U}')) \geq \varpi \|\tilde{U} - \tilde{U}'\|^2.$$

Lemma 4.3 ([27, Proposition 3.5.10]). Let \mathbb{K}_b be given by (4.23).

(a) If L is a P-function on \mathbb{K}_b , then $\text{VI}(\mathbb{K}_b, L)$ has at most one solution.

(b) If each \mathbb{K}_{n_ϑ} is a closed convex set and L is a continuous uniformly P-function on \mathbb{K}_b , then the $\text{VI}(\mathbb{K}_b, L)$ has a unique solution.

4.3.2. Existence and uniqueness of LDG discretization with positivity and mass conservation constraints

In this section, we will prove the existence and uniqueness of the KKT system (4.10)–(4.21) using the unique solvability conditions discussed in Section 4.3.1.

Lemma 4.4. For periodic boundary conditions, the matrices B in (4.15) and A in (4.17) satisfy $B^T = A$.

Proof. In order to prove the symmetry of B in (4.15) and A in (4.17), we define the bilinear function $a : (V_h^k \times \mathbf{W}_h^k) \times (V_h^k \times \mathbf{W}_h^k) \rightarrow \mathbb{R}$ by

$$\begin{aligned} a(P_h^{n+1}, \mathbf{Q}_h^{n+1}; \rho, \boldsymbol{\theta}) &= (\mathbf{Q}_h^{n+1}, \nabla \rho) - \sum_{K \in \mathcal{T}_h} (\hat{\mathbf{Q}}_h^{n+1} \cdot \mathbf{v}, \rho)_{\partial K} \\ &\quad - (P_h^{n+1}, \nabla \cdot \boldsymbol{\theta}) + \sum_{K \in \mathcal{T}_h} (\hat{P}_h^{n+1}, \mathbf{v} \cdot \boldsymbol{\theta})_{\partial K}. \end{aligned}$$

Based on the definition of B in (4.15) using (2.2a), A in (4.17) using (2.2c), we rewrite the above bilinear function a as follows:

$$a(P_h^{n+1}, \mathbf{Q}_h^{n+1}; \rho, \boldsymbol{\theta}) = (\rho, \Theta) \begin{pmatrix} 0 & B \\ A & 0 \end{pmatrix} (\tilde{P}^{n+1}, \tilde{\mathbf{Q}}^{n+1})^T,$$

with ρ, Θ the LDG coefficients of $\rho, \boldsymbol{\theta}$ and $\tilde{P}^{n+1}, \tilde{\mathbf{Q}}^{n+1}$ the LDG coefficients of $P_h^{n+1}, \mathbf{Q}_h^{n+1}$, respectively.

Interchanging the arguments of a , we get

$$\begin{aligned} a(\rho, \boldsymbol{\theta}; P_h^{n+1}, \mathbf{Q}_h^{n+1}) &= (\boldsymbol{\theta}, \nabla P_h^{n+1}) - \sum_{K \in \mathcal{T}_h} (\hat{\boldsymbol{\theta}} \cdot \mathbf{v}, P_h^{n+1})_{\partial K} \\ &\quad - (\rho, \nabla \cdot \mathbf{Q}_h^{n+1}) + \sum_{K \in \mathcal{T}_h} (\hat{\rho}, \mathbf{v} \cdot \mathbf{Q}_h^{n+1})_{\partial K} \\ &= - (P_h^{n+1}, \nabla \cdot \boldsymbol{\theta}) + \sum_{K \in \mathcal{T}_h} (\boldsymbol{\theta} \cdot \mathbf{v}, P_h^{n+1})_{\partial K} - \sum_{K \in \mathcal{T}_h} (\hat{\boldsymbol{\theta}} \cdot \mathbf{v}, P_h^{n+1})_{\partial K} \\ &\quad + (\mathbf{Q}_h^{n+1}, \nabla \rho) - \sum_{K \in \mathcal{T}_h} (\rho, \mathbf{v} \cdot \mathbf{Q}_h^{n+1})_{\partial K} + \sum_{K \in \mathcal{T}_h} (\hat{\rho}, \mathbf{v} \cdot \mathbf{Q}_h^{n+1})_{\partial K}. \end{aligned}$$

Using equality (2.7), the alternating numerical fluxes for $\hat{\theta}$ and $\hat{\rho}$ in (2.3) or (2.4), and the periodic boundary conditions, we obtain

$$\begin{aligned} & \sum_{K \in \mathcal{T}_h} (\theta \cdot \mathbf{v}, P_h^{n+1})_{\partial K} - \sum_{K \in \mathcal{T}_h} (\hat{\theta} \cdot \mathbf{v}, P_h^{n+1})_{\partial K} = \sum_{K \in \mathcal{T}_h} (\hat{P}_h^{n+1}, \mathbf{v} \cdot \theta)_{\partial K}, \\ & - \sum_{K \in \mathcal{T}_h} (\rho, \mathbf{v} \cdot \mathbf{Q}_h^{n+1})_{\partial K} + \sum_{K \in \mathcal{T}_h} (\hat{\rho}, \mathbf{v} \cdot \mathbf{Q}_h^{n+1})_{\partial K} = - \sum_{K \in \mathcal{T}_h} (\hat{\mathbf{Q}}_h^{n+1} \cdot \mathbf{v}, \rho)_{\partial K}. \end{aligned}$$

Hence,

$$a(P_h^{n+1}, \mathbf{Q}_h^{n+1}; \rho, \theta) = a(\rho, \theta; P_h^{n+1}, \mathbf{Q}_h^{n+1}),$$

which implies

$$\begin{aligned} (\varrho, \Theta) \begin{pmatrix} 0 & B \\ A & 0 \end{pmatrix} (\tilde{P}^{n+1}, \tilde{\mathbf{Q}}^{n+1})^T &= (\tilde{P}^{n+1}, \tilde{\mathbf{Q}}^{n+1}) \begin{pmatrix} 0 & B \\ A & 0 \end{pmatrix} (\varrho, \Theta)^T \\ &= (\varrho, \Theta) \begin{pmatrix} 0 & A^T \\ B^T & 0 \end{pmatrix} (\tilde{P}^{n+1}, \tilde{\mathbf{Q}}^{n+1})^T. \end{aligned} \tag{4.25}$$

Since $(P_h^{n+1}, \mathbf{Q}_h^{n+1}) \in V_h^k \times \mathbf{W}_h^k$ and $(\rho, \theta) \in V_h^k \times \mathbf{W}_h^k$ are arbitrary functions, relation (4.25) implies that $A = B^T$. \square

Using (4.12)–(4.14) and Lemma 4.4, the operator $L(\tilde{U}^{n+1})$ in (4.11) can be written as

$$L(\tilde{U}^{n+1}) = M(\tilde{U}^{n+1} - \tilde{U}^n) + \tau^{n+1} \mathbf{B} \mathbf{M}^{-1} C_d(\tilde{U}^{n+1}) \mathbf{M}^{-1} \mathbf{B}^T M^{-1} D(\tilde{U}^{n+1}). \tag{4.26}$$

Lemma 4.5. Given \tilde{U}^n , the operator L in (4.26) is a uniformly P -function on \mathbb{K}_b for $0 < \tau^{n+1} \leq \frac{\sigma}{4c}$, with σ the smallest eigenvalue of the symmetric positive definite mass matrix M , and c a positive constant independent of \tilde{U} .

Proof. Using relation (4.26) for L , for arbitrary $\tilde{U}_I^{n+1}, \tilde{U}_{II}^{n+1} \in \mathbb{K}_b$, there holds

$$\begin{aligned} L(\tilde{U}_I^{n+1}) - L(\tilde{U}_{II}^{n+1}) &= M(\tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1}) + \tau^{n+1} \mathbf{B} \mathbf{M}^{-1} C_d(\tilde{U}_I^{n+1}) \mathbf{M}^{-1} \mathbf{B}^T M^{-1} D(\tilde{U}_I^{n+1}) \\ &\quad - \tau^{n+1} \mathbf{B} \mathbf{M}^{-1} C_d(\tilde{U}_{II}^{n+1}) \mathbf{M}^{-1} \mathbf{B}^T M^{-1} D(\tilde{U}_{II}^{n+1}). \end{aligned} \tag{4.27}$$

After subtracting and adding $\tau^{n+1} \mathbf{B} \mathbf{M}^{-1} C_d(\tilde{U}_I^{n+1}) \mathbf{M}^{-1} \mathbf{B}^T M^{-1} D(\tilde{U}_{II}^{n+1})$ in (4.27), we obtain

$$\begin{aligned} & L(\tilde{U}_I^{n+1}) - L(\tilde{U}_{II}^{n+1}) \\ &= M(\tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1}) + \tau^{n+1} \mathbf{B} \mathbf{M}^{-1} C_d(\tilde{U}_I^{n+1}) \mathbf{M}^{-1} \mathbf{B}^T M^{-1} (D(\tilde{U}_I^{n+1}) - D(\tilde{U}_{II}^{n+1})) \\ &\quad + \tau^{n+1} \mathbf{B} \mathbf{M}^{-1} (C_d(\tilde{U}_I^{n+1}) - C_d(\tilde{U}_{II}^{n+1})) \mathbf{M}^{-1} \mathbf{B}^T M^{-1} D(\tilde{U}_{II}^{n+1}). \end{aligned} \tag{4.28}$$

With the definition of D in (4.18) using (2.2d), we obtain that

$$\begin{aligned} (D(\tilde{U}_I^{n+1}) - D(\tilde{U}_{II}^{n+1}))_i &= \int_{\Omega} \left(H' \left(\sum_{j=1}^{N_k N_e} \tilde{U}_{I,j}^{n+1} \phi_j \right) - H' \left(\sum_{j=1}^{N_k N_e} \tilde{U}_{II,j}^{n+1} \phi_j \right) \right) \phi_i d\Omega \\ &= \sum_{j=1}^{N_k N_e} (\tilde{U}_{I,j}^{n+1} - \tilde{U}_{II,j}^{n+1}) \int_{\Omega} H''(\xi_1^{n+1}) \phi_j \phi_i d\Omega, \quad i \in \{1, \dots, N_k N_e\}, \xi_1^{n+1} \in (U_{h,I}^{n+1}, U_{h,II}^{n+1}), \end{aligned}$$

and write

$$D(\tilde{U}_I^{n+1}) - D(\tilde{U}_{II}^{n+1}) := D_{\tilde{U}}(\xi_1^{n+1})(\tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1}). \tag{4.29}$$

Similarly, from the definition of C_d in (4.16), (4.19) using (2.2b), we obtain that

$$\begin{aligned} C_d(\tilde{U}_I^{n+1}) - C_d(\tilde{U}_{II}^{n+1}) &= \begin{pmatrix} C(\tilde{U}_I^{n+1}) - C(\tilde{U}_{II}^{n+1}) & & \\ & \ddots & \\ & & C(\tilde{U}_I^{n+1}) - C(\tilde{U}_{II}^{n+1}) \end{pmatrix}, \\ (C(\tilde{U}_I^{n+1}) - C(\tilde{U}_{II}^{n+1}))_{ij} &= \int_{\Omega} \left(f \left(\sum_{k=1}^{N_k N_e} \tilde{U}_{I,k}^{n+1} \phi_k \right) - f \left(\sum_{k=1}^{N_k N_e} \tilde{U}_{II,k}^{n+1} \phi_k \right) \right) \phi_j \phi_i d\Omega \\ &= \sum_{k=1}^{N_k N_e} (\tilde{U}_{I,k}^{n+1} - \tilde{U}_{II,k}^{n+1}) \int_{\Omega} f'(\xi_2^{n+1}) \phi_k \phi_j \phi_i d\Omega, \quad i, j, k \in \{1, \dots, N_k N_e\}, \xi_2^{n+1} \in (U_{h,I}^{n+1}, U_{h,II}^{n+1}), \end{aligned}$$

and write

$$C(\tilde{U}_I^{n+1}) - C(\tilde{U}_{II}^{n+1}) := \sum_{k=1}^{N_k N_e} [C_d \tilde{U}(\xi_2^{n+1})]_k (\tilde{U}_{I,k}^{n+1} - \tilde{U}_{II,k}^{n+1}). \tag{4.30}$$

Assume for arbitrary $\tilde{U} \in \mathbb{K}_b$ in (4.23), that

$$\begin{aligned} |C(\tilde{U})_{ij}| &\leq c, \quad |D(\tilde{U})_i| \leq c, \\ \|[C_{\tilde{U}}(\tilde{U})_{ij}]_k\| &\leq c, \quad |D_{\tilde{U}}(\tilde{U})_{ij}| \leq c, \quad i, j, k \in \{1, \dots, N_k N_e\}, \end{aligned} \tag{4.31}$$

with c a positive constant, independent of \tilde{U} . In the remainder of this section, c is a positive constant, but not necessarily the same.

Using (4.29)–(4.30) and assumption (4.31), we obtain the following two estimates

$$\begin{aligned} &(\tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1})^T \mathbf{B} \mathbf{M}^{-1} C_d(\tilde{U}_I^{n+1}) \mathbf{M}^{-1} \mathbf{B}^T M^{-1} (D(\tilde{U}_I^{n+1}) - D(\tilde{U}_{II}^{n+1})) \\ &\leq \| \mathbf{B} \| \| \mathbf{M}^{-1} \| \| C_d(\tilde{U}_I^{n+1}) \| \| \mathbf{M}^{-1} \| \| \mathbf{B}^T \| \| M^{-1} \| \| D_{\tilde{U}}(\xi_1^{n+1}) \| \| \tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1} \|^2 \\ &\leq c \| \tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1} \|^2, \end{aligned}$$

and

$$\begin{aligned} &(\tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1})^T \mathbf{B} \mathbf{M}^{-1} (C_d(\tilde{U}_I^{n+1}) - C_d(\tilde{U}_{II}^{n+1})) \mathbf{M}^{-1} \mathbf{B}^T M^{-1} D(\tilde{U}_{II}^{n+1}) \\ &\leq \| \mathbf{B} \| \| \mathbf{M}^{-1} \| \sum_{k=1}^{N_k N_e} \| [C_{d\tilde{U}}(\xi_2^{n+1})]_k \| \| \mathbf{M}^{-1} \| \| \mathbf{B}^T \| \| M^{-1} \| \| D(\tilde{U}_{II}^{n+1}) \| \| \tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1} \|^2 \\ &\leq c \| \tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1} \|^2. \end{aligned}$$

Then multiplying (4.28) with $(\tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1})^T$ gives

$$\begin{aligned} &(\tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1})^T (L(\tilde{U}_I^{n+1}) - L(\tilde{U}_{II}^{n+1})) = (\tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1})^T M (\tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1}) \\ &+ \tau^{n+1} (\tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1})^T \mathbf{B} \mathbf{M}^{-1} C_d(\tilde{U}_I^{n+1}) \mathbf{M}^{-1} \mathbf{B}^T M^{-1} (D(\tilde{U}_I^{n+1}) - D(\tilde{U}_{II}^{n+1})) \\ &+ \tau^{n+1} (\tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1})^T \mathbf{B} \mathbf{M}^{-1} (C_d(\tilde{U}_I^{n+1}) - C_d(\tilde{U}_{II}^{n+1})) \mathbf{M}^{-1} \mathbf{B}^T M^{-1} D(\tilde{U}_{II}^{n+1}) \\ &\geq \sigma \| \tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1} \|^2 - 2c\tau^{n+1} \| \tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1} \|^2, \end{aligned} \tag{4.32}$$

where $\sigma > 0$ is the smallest eigenvalue of the symmetric positive mass matrix M .

Choosing $0 < \tau^{n+1} \leq \frac{\sigma}{4c}$, we obtain that

$$(\tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1})^T (L(\tilde{U}_I^{n+1}) - L(\tilde{U}_{II}^{n+1})) \geq \frac{\sigma}{2} \| \tilde{U}_I^{n+1} - \tilde{U}_{II}^{n+1} \|^2, \quad \forall \tilde{U}_I^{n+1}, \tilde{U}_{II}^{n+1} \in \mathbb{K}_b, \tag{4.33}$$

which implies that for τ^{n+1} sufficiently small $L(\tilde{U}^{n+1})$ is a uniformly function of \mathbb{K}_b , \square

From Lemmas 4.1, 4.2, 4.3 and 4.5, we obtain the main result of this section.

Theorem 4.6. *Given the DG coefficients \tilde{U}^n and the positivity preserving backward Euler LDG discretization (4.10)–(4.21) with equality constraint $h \equiv 0$. Given a time step τ^{n+1} satisfying $0 < \tau^{n+1} \leq \frac{\sigma}{4c}$, with σ the smallest eigenvalue of the symmetric positive definite mass matrix M . If assumption (4.31) is satisfied, then the KKT system (4.10)–(4.21) has only one solution.*

Proof. From Lemmas 4.3 and 4.5, we have that the $\text{VI}(\mathbb{K}_b, L)$ has a unique solution denoted by \tilde{U}^{n+1} with \mathbb{K}_b given by (4.23) and L given by (4.11). Then from Lemma 4.1, there exists a solution \tilde{U}^{n+1} , and vectors $\mu \in \mathbb{R}^l$ and $\lambda \in \mathbb{R}^m$ satisfying the KKT system (4.10)–(4.21) with equality constraint $h \equiv 0$.

Since in (4.20)–(4.21), each function h_j ($1 \leq j \leq l$) is affine and each function g_i ($1 \leq i \leq m$) is convex (linear), the KKT system (4.10)–(4.21) has only one solution, which follows from Lemma 4.2 and the uniqueness of the solution for $\text{VI}(\mathbb{K}_b, L)$. \square

Corollary 4.7. *Given the DG coefficients \tilde{U}^n . Given a time step τ^{n+1} satisfying $0 < \tau^{n+1} \leq \frac{\sigma}{4c}$, with σ the smallest eigenvalue of the symmetric positive definite mass matrix M . If assumption (4.31) is satisfied, then for the degenerate parabolic Eq. (1.1) with periodic boundary conditions, there exists only one solution satisfying the higher order accurate in time, positivity preserving DIRK-LDG discretizations (4.7) with equality constraint $h \equiv 0$.*

Proof. Since the DIRK coefficient matrix (a_{ij}) introduced in Section 3.2 is a lower triangular matrix, the structure of the DIRK-LDG discretizations is similar to the form obtained for the backward Euler LDG discretization. The analysis therefore is completely analogous to Theorem 4.6. \square

4.3.3. Stability of the positivity preserving DIRK-LDG discretization

Theorem 4.8. *Given the numerical solution $U_h^n \in V_h^k$ of the positivity preserving backward Euler LDG discretization (4.10)–(4.21). Given a time step τ^{n+1} satisfying $0 < \tau^{n+1} \leq \frac{\sigma}{4c}$, with σ the smallest eigenvalue of the symmetric positive definite mass matrix M , and c a strictly positive constant. If assumption (4.31) is satisfied, then the discrete entropy E_h stated in (3.2) satisfies for $n = 0, 1, \dots$,*

$$E_h(U_h^{n+1}) \leq E_h(U_h^n), \tag{4.34}$$

which implies that the positivity preserving backward Euler LDG discretization is entropy dissipative.

Proof. From Lemma 4.2, we obtain that the LDG coefficients \tilde{U}^{n+1} of the positivity preserving solution U_h^{n+1} solve

$$(y - \tilde{U}^{n+1})^T L(\tilde{U}^{n+1}) \geq 0, \quad \forall y \in \mathbb{K}, \tag{4.35}$$

with L given by (4.26) and \mathbb{K} given by (4.1).

From assumption (4.31), we have that there exists a positive constant $c \geq c_0 > 0$ such that

$$\tilde{U}^{n+1} - cM^{-1}D(\tilde{U}^{n+1}) \in \mathbb{K}. \tag{4.36}$$

Next, we choose $y = \tilde{U}^{n+1} - cM^{-1}D(\tilde{U}^{n+1})$ in (4.35), which implies

$$-c(M^{-1}D(\tilde{U}^{n+1}))^T L(\tilde{U}^{n+1}) \geq 0. \tag{4.37}$$

Using (4.26) and the fact that $c > 0$, we obtain that (4.37) implies the inequality

$$D(\tilde{U}^{n+1})^T (\tilde{U}^{n+1} - \tilde{U}^n) + \tau^{n+1} D(\tilde{U}^{n+1})^T M^{-1} B M^{-1} C_d(\tilde{U}^{n+1}) M^{-1} B^T M^{-1} D(\tilde{U}^{n+1}) \leq 0. \tag{4.38}$$

From the definition of C_d in (4.16), (4.19) using (2.2b) and the conditions on f stated in (1.3), we obtain that $C_d(\tilde{U}^{n+1})$ is symmetric positive definite. Hence using $\tau^{n+1} > 0$, we have

$$\tau^{n+1} D(\tilde{U}^{n+1})^T M^{-1} B M^{-1} C_d(\tilde{U}^{n+1}) M^{-1} B^T M^{-1} D(\tilde{U}^{n+1}) \geq 0,$$

which with (4.38) yields

$$D(\tilde{U}^{n+1})^T (\tilde{U}^{n+1} - \tilde{U}^n) \leq 0. \tag{4.39}$$

From the definition of D in (4.18) using (2.2d) and (4.39), we obtain the bound

$$(\Psi(\mathbf{x}), U_h^{n+1} - U_h^n) + (H'(U_h^{n+1}), U_h^{n+1} - U_h^n) \leq 0. \tag{4.40}$$

Using the following Taylor expansion

$$H(U_h^n) = H(U_h^{n+1}) + H'(U_h^{n+1})(U_h^n - U_h^{n+1}) + \frac{1}{2} H''(\xi_3^{n+1})(U_h^{n+1} - U_h^n)^2, \quad \xi_3^{n+1} \in (U_h^n, U_h^{n+1}),$$

we obtain that (4.40) gives

$$(\Psi(\mathbf{x}), U_h^{n+1} - U_h^n) + (H(U_h^{n+1}) - H(U_h^n), 1) + \frac{1}{2} (H''(\xi_3^{n+1}), (U_h^{n+1} - U_h^n)^2) \leq 0,$$

which implies, using the definition of E_h in (3.2), that

$$E_h(U_h^{n+1}) - E_h(U_h^n) = (\Psi(\mathbf{x}), U_h^{n+1} - U_h^n) + (H(U_h^{n+1}) - H(U_h^n), 1) \leq 0,$$

since (1.3) gives $H''(\xi_3^{n+1}) \geq 0$. This proves (4.34). \square

5. Numerical tests

In this section, we will discuss several numerical experiments to demonstrate the performance of the positivity preserving DIRK-LDG algorithm for the degenerate parabolic Eq. (1.1). In the computations, we will consider the porous medium equation, the nonlinear diffusion equation with a double-well potential and the nonlinear Fokker–Planck equation for fermion and boson gases. Firstly, we will present in Section 5.1 the order of accuracy of the DIRK-LDG discretizations with and without positivity preserving limiter to investigate if the limiter negatively affects the accuracy of the discretizations. Next, we will present in Sections 5.2–5.5 test cases for which the positivity preserving limiter is essential. Without the positivity constraint, obtaining a numerical solution is not possible or only for extremely small time steps.

In the computations, we take $\tau = \alpha \cdot h_{\mathcal{T}}$ with $h_{\mathcal{T}}$ the mesh size for the tessellation \mathcal{T}_h . For numerical efficiency, it is important to have a good balance between the number of Newton iterations and the time step size. If the Newton method during strongly nonlinear stages requires a large number of iterations, it is generally more efficient to reduce the time step to $\frac{1}{2}\tau$ and restart the Newton iterations. When the Newton method converges well, then τ is increased each time step to 1.2τ , till the maximum predefined time step is obtained.

In order to avoid round-off effects, a positivity bound $u_{\min} = 10^{-10}$ is used in the numerical simulations, except for Section 5.1 where $u_{\min} = 10^{-14}$. If it is not stated otherwise, the numerical results for 1D problems are obtained on a mesh containing 100 elements and Legendre polynomials of order 2. For 2D problems, a mesh consisting of 30×30 square elements and tensor product Legendre polynomial basis functions of order 2 are used.

Table 5.1

Error in L^∞ - and L^1 - norms for Example 5.1 at time $T = 1$ without positivity preserving limiter.

\mathcal{P}_k	\mathbb{M}	$\ u_n - u_h^n\ _{L^\infty(\Omega)}$	Order	$\ u_n - u_h^n\ _{L^1(\Omega)}$	Order	$\min u_h^n$
1	40	7.33E-003	-	1.03E-003	-	-8.87e-005
	80	1.24e-003	2.56	2.27e-004	2.18	-1.08e-005
	160	2.63e-004	2.24	5.44e-005	2.06	-4.41e-007
	320	6.05e-005	2.12	1.35e-005	2.01	-1.57e-008
2	40	1.70E-003	-	8.73E-005	-	-1.60e-005
	80	1.43e-004	3.57	8.07e-006	3.44	-1.79e-007
	160	1.36e-005	3.39	9.40e-007	3.10	-6.24e-009
	320	1.34e-006	3.34	1.16e-007	3.02	-2.07e-010
3	40	1.45e-004	-	6.00e-006	-	-2.14e-006
	80	9.87e-006	3.88	3.11e-007	4.27	-9.56e-008
	160	5.51e-007	4.16	1.76e-008	4.14	-3.51e-009
	320	3.50e-008	3.98	1.11e-009	3.99	-1.19e-010

Table 5.2

Error in L^∞ - and L^1 - norms for Example 5.1 at time $T = 1$ with positivity preserving limiter.

\mathcal{P}_k	\mathbb{M}	$\ u_n - U_h^n\ _{L^\infty(\Omega)}$	Order	$\ u_n - U_h^n\ _{L^1(\Omega)}$	Order	$\min U_h^n$
1	40	7.33E-003	-	1.05E-003	-	2.05e-005
	80	1.24e-003	2.56	2.27e-004	2.21	8.15e-007
	160	2.63e-004	2.24	5.44e-005	2.06	2.77e-008
	320	6.05e-005	2.12	1.35e-005	2.01	8.55e-010
2	40	1.70E-003	-	8.73E-005	-	6.15e-008
	80	1.43e-004	3.57	8.08e-006	3.43	3.03e-007
	160	1.36e-005	3.39	9.40e-007	3.10	1.08e-008
	320	1.34e-006	3.34	1.16e-007	3.02	4.55e-010
3	40	1.45e-004	-	6.02e-006	-	1.00e-014
	80	9.87e-006	3.88	3.13e-007	4.27	4.45e-008
	160	5.51e-007	4.16	1.77e-008	4.14	1.21e-009
	320	3.50e-008	3.98	1.11e-009	4.00	2.55e-011

5.1. Accuracy tests

For the accuracy test, we use a uniform mesh with \mathbb{M} elements and positivity bound $u_{\min} = 10^{-14}$.

Example 5.1. We consider (1.1) on the domain $\Omega = (-1, 1)$ with Dirichlet boundary conditions based on the exact solution and select the following parameters

$$f(u) = u, \quad H'(u) = u^2, \quad \Psi(x) = 0, \quad x \in \Omega.$$

Then (1.1) with a properly chosen source term has the nonnegative solution

$$u(x, t) = \exp(-t)(1 - x^4)^5, \quad x \in \Omega.$$

We take α in the definition of the time step as $\alpha = 1$. Tables 5.1–5.2 show that the DIRK-LDG discretizations with and without positivity preserving limiter are convergent at the rate $O(h_T^{k+1})$ for basis functions with polynomial order ranging from 1 to 3. Note, for polynomials \mathcal{P}_k of order k the DIRK methods, see Section 3.2, have order of accuracy $k + 1$. The errors and orders of accuracy presented in Tables 5.1–5.2 indicate that the positivity preserving limiter is necessary and does not negatively affect accuracy.

5.2. Porous media equation

For the porous media equation, $f(u)H''(u)$ can locally vanish, resulting in degenerate cases [1]. We test the asymptotic behavior of the numerical solution and will show that the KKT limiter is necessary. The entropy defined in (1.4), which should be non-increasing, is also computed.

Example 5.2. In order to test degenerate cases, we choose the following parameters in (1.1) on the domain $\Omega = (0, 1)$ with zero-flux boundary condition (1.2)

$$f(u) = u, \quad H'(u) = \frac{4}{3} \left(u - \frac{1}{2}\right)^3 \max\left(u, \frac{1}{2}\right), \quad \Psi(x) = 0, \quad x \in \Omega,$$

and initial data

$$u(x, 0) = \frac{1}{2} - \frac{1}{2} \cos(2\pi x), \quad x \in \Omega.$$

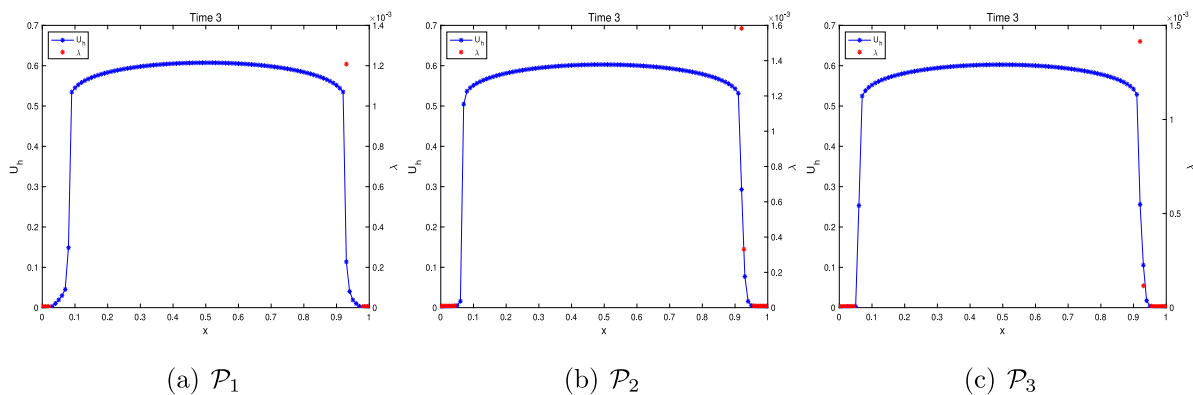


Fig. 5.1. (Example 5.2) Numerical solution U_h for different orders of polynomial basis functions \mathcal{P}_1 - \mathcal{P}_3 with the KKT limiter enforced and Lagrange multiplier λ (red dots).

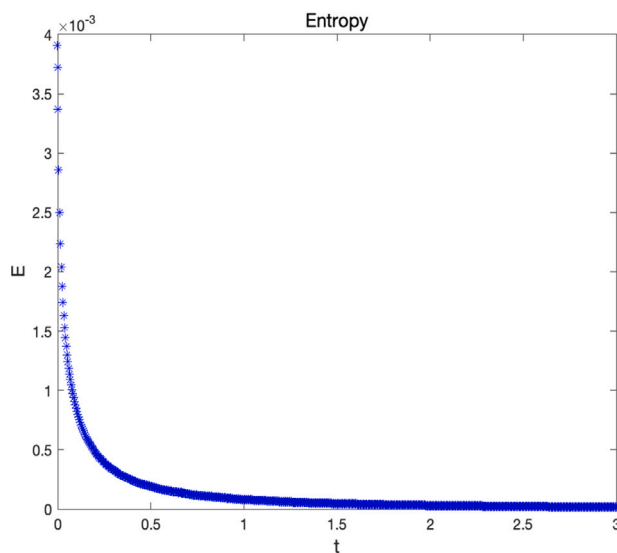


Fig. 5.2. (Example 5.2) Entropy E_h for \mathcal{P}_3 basis functions with the KKT limiter enforced.

During the computations, the value of α for optimal convergence of the semi-smooth Newton algorithm is usually close to 0.1. We present the numerical solution in Fig. 5.1 for basis functions with polynomial order ranging from 1 to 3 and with the KKT limiter enforced. Values of the Lagrange multiplier λ larger than 10^{-10} are shown in Fig. 5.1, which indicate that the positivity constraint works well since it is only active at locations where the solution is close to the minimum value. The entropy decay using the KKT limiter and polynomial basis functions of order 3 is presented in Fig. 5.2, which the result is consistent with the entropy analysis. In Fig. 5.3, the numerical solution without KKT limiter and for polynomial basis functions with order 3 is plotted. This computation breaks down due to unphysical oscillations.

Example 5.3. We consider a 2D test case on the domain $\Omega = (-6, 6)^2$ with zero-flux boundary condition (1.2) by choosing in (1.1) the following parameters

$$f(u) = u, \quad H'(u) = 2u, \quad \Psi(\mathbf{x}) = 0, \quad \mathbf{x} \in \Omega,$$

and initial data

$$u(\mathbf{x}, 0) = \exp\left(-\frac{1}{2}|\mathbf{x}|^2\right), \quad \mathbf{x} \in \Omega.$$

In this case, the value of α in the definition of the time step ranges between 0.1 and 1. Fig. 5.4 presents the numerical solution with the KKT limiter active and also the Lagrange multiplier λ . Considering the position of the non-zero Lagrange multipliers, we can see that the limiter also works well in the two-dimensional case since it is only active in areas where positivity must be enforced. The entropy decay is plotted in Fig. 5.5, which is consistent with the stability result of the numerical solution. Without the KKT limiter, there will be unphysical oscillations, and the computation will break down at some point in the computations.

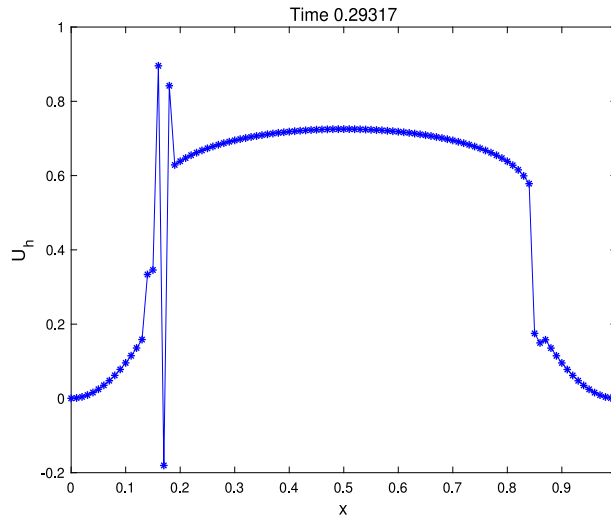


Fig. 5.3. (Example 5.2) Numerical solution U_h for P_3 basis functions without KKT limiter just before blow up.

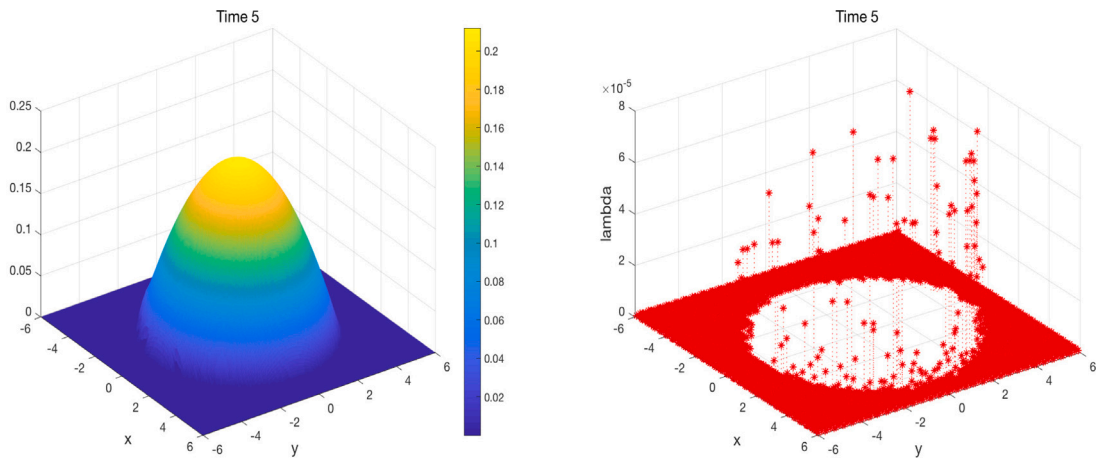


Fig. 5.4. (Example 5.3) Numerical solution U_h for P_2 basis functions with KKT limiter enforced (Left) and Lagrange multiplier λ (Right).

5.3. Nonlinear diffusion with a double-well potential

Consider the nonlinear diffusion equation with double-well potential [28] on the domain $\Omega = (-1.4, 1.4)$, which is obtained by choosing in (1.1) zero-flux boundary condition (1.2) and the following parameters

$$f(u) = u, \quad H'(u) = u, \quad \Psi(x) = \frac{1}{4}x^4 - \frac{1}{2}x^2, \quad x \in \Omega. \tag{5.1}$$

This model is taken from [4]. We will test the evolution of the numerical solution with and without KKT limiter, and also the decay of the entropy (1.4). The value of α to compute the time step ranges between 0.01 to 0.1.

Example 5.4. We consider (1.1) with (5.1) and the initial data

$$u(x, 0) = \frac{0.2}{\sqrt{0.4\pi}} \exp\left(-\frac{x^2}{0.4}\right), \quad x \in \Omega.$$

The numerical solution with the KKT limiter enforced and the values of the Lagrange multiplier λ larger than 10^{-10} are shown in Fig. 5.6. These results indicate that the numerical solution tends to a steady state and that the KKT limiter is only active at places where the positivity constraint needs to be imposed. The entropy dissipation is presented in Fig. 5.7, in which uniform decay coincides with our theoretical analysis. For the numerical solution without the KKT limiter, we observe that violating the positivity constraint will result in discontinuities in the solution and a computation breakdown, even for a very small CFL number.

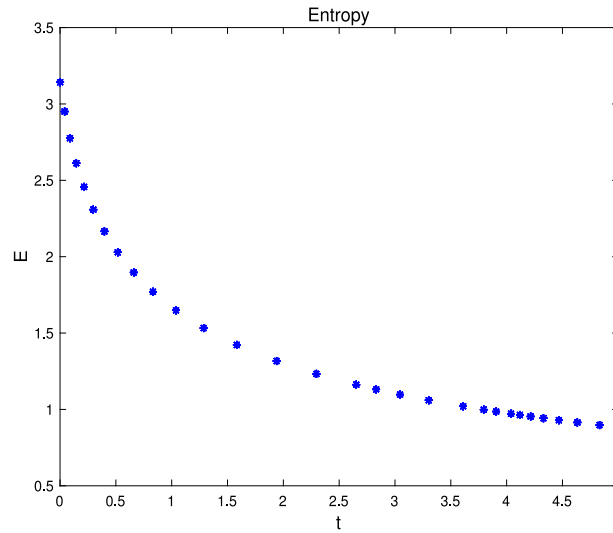


Fig. 5.5. (Example 5.3) Entropy E_h for \mathcal{P}_2 basis functions with KKT limiter enforced.

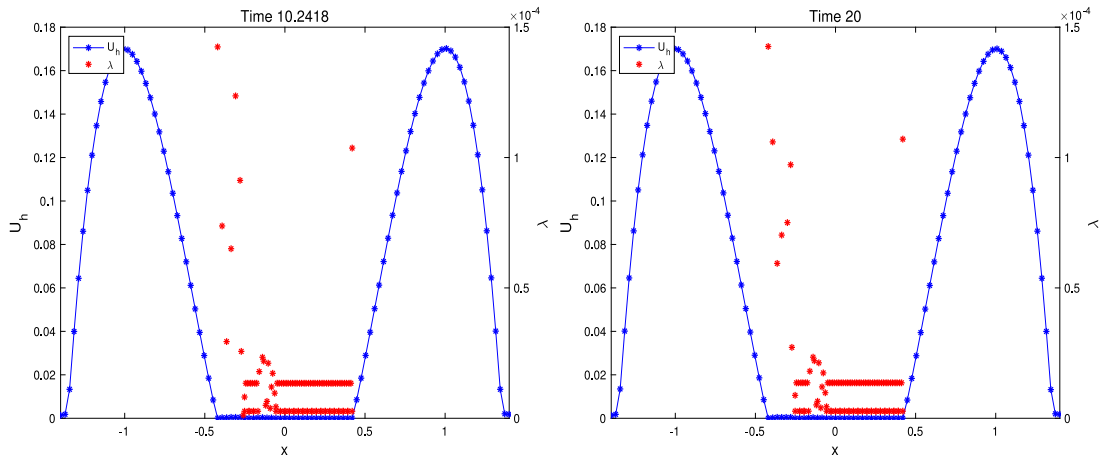


Fig. 5.6. (Example 5.4) Numerical solution U_h for \mathcal{P}_2 basis functions with KKT limiter enforced and Lagrange multiplier λ (red dots).

5.4. Nonlinear fokker–planck equation for fermion gases

Example 5.5. We consider the nonlinear Fokker–Planck equation for fermion gases [1] on the domain $\Omega = (-10, 10)^2$, for which we select the following parameters in (1.1)

$$f(u) = u(1 - u), \quad H'(u) = \log \frac{u}{1 - u}, \quad \Psi(\mathbf{x}) = \frac{1}{2}|\mathbf{x}|^2, \quad \mathbf{x} \in \Omega, \tag{5.2}$$

together with zero-flux boundary condition (1.2). The initial data is given by

$$u(\mathbf{x}, 0) = \frac{1}{2\sqrt{2\pi}} \left(\exp\left(-\frac{1}{2}|\mathbf{x} - (2, 2)|^2\right) + \exp\left(-\frac{1}{2}|\mathbf{x} - (2, -2)|^2\right) + \exp\left(-\frac{1}{2}|\mathbf{x} - (-2, 2)|^2\right) + \exp\left(-\frac{1}{2}|\mathbf{x} - (-2, -2)|^2\right) \right), \quad \mathbf{x} \in \Omega.$$

During the computations, the value of α in the definition of the time step ranges between 0.1 and 1, but for most time steps $\alpha = 1$. The numerical solutions at several time levels with the KKT limiter enforced and the entropy dissipation are presented in Figs. 5.8 and 5.9, respectively, showing the time-asymptotic convergence of the numerical solution towards a steady state. Without the KKT limiter, the computations break down, even for very small CFL numbers.

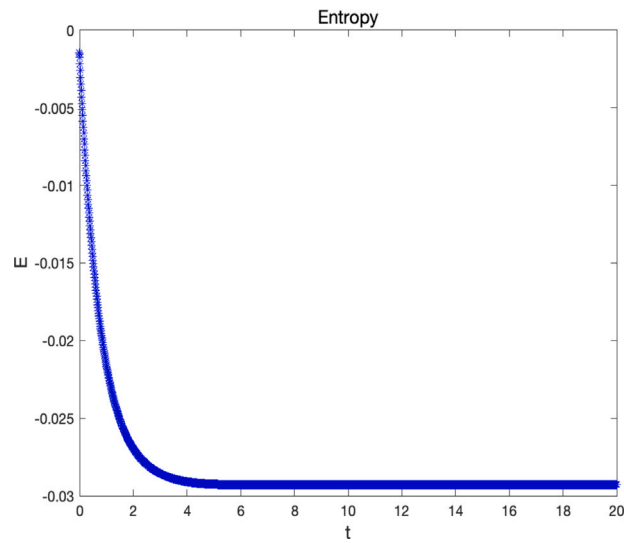


Fig. 5.7. (Example 5.4) Entropy E_h for \mathcal{P}_2 basis functions with KKT limiter enforced.

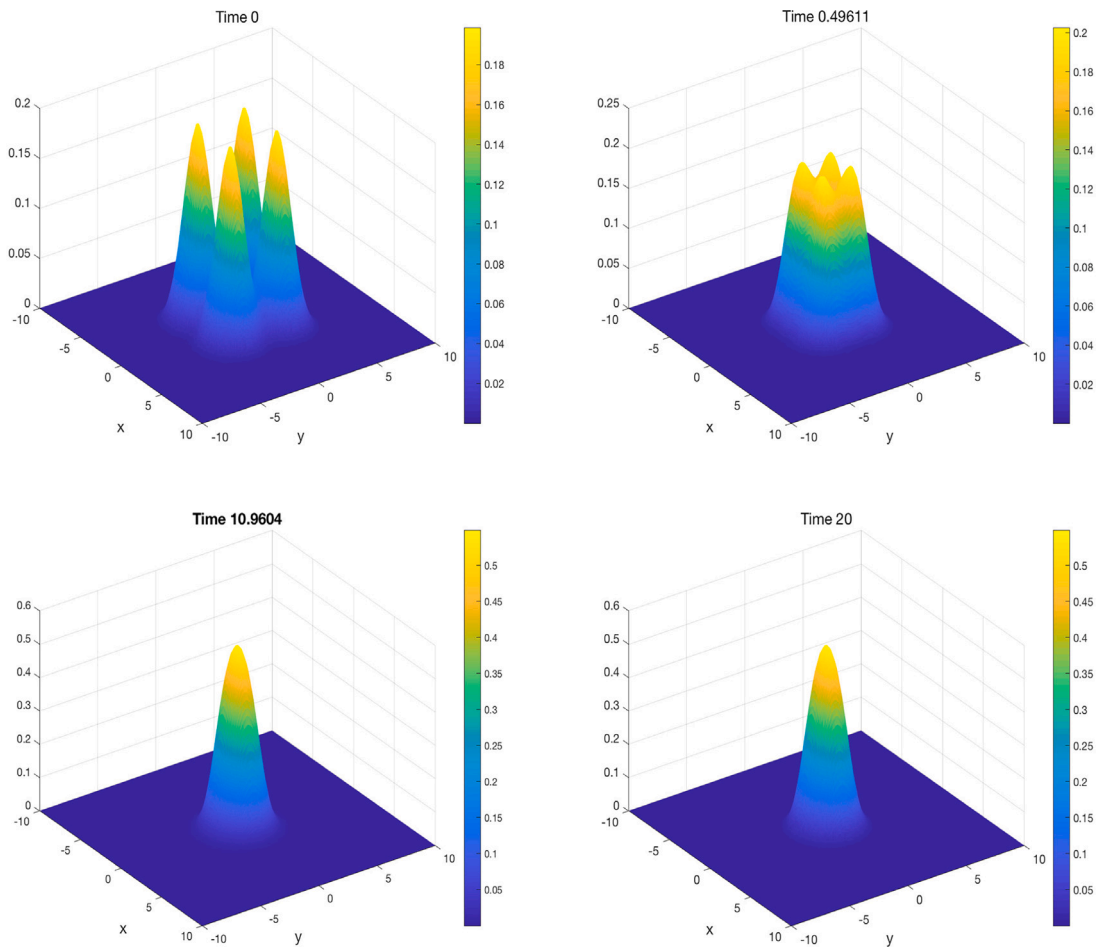


Fig. 5.8. (Example 5.5) Numerical solution U_h for \mathcal{P}_2 basis functions with KKT limiter enforced.

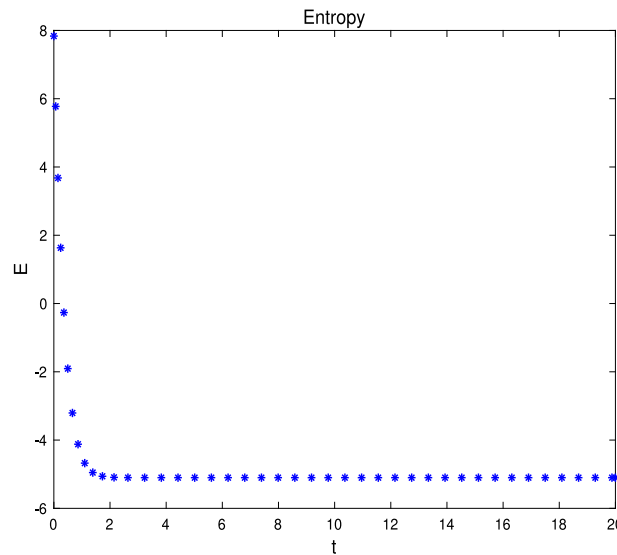


Fig. 5.9. (Example 5.5) Entropy E_h for \mathcal{P}_2 basis functions with KKT limiter enforced.

5.5. Nonlinear fokker-planck equation for boson gases

Example 5.6. We consider a nonlinear Fokker-Planck equation for boson gases with zero-flux boundary condition (1.2) on a domain $\Omega = (-10, 10)$, which requires the following parameters in (1.1)

$$f(u) = u(1 + u^3), \quad H'(u) = \log \frac{u}{(1 + u^3)^{\frac{1}{3}}}, \quad \Psi(x) = \frac{x^2}{2}, \quad x \in \Omega.$$

The initial data is [1,9]

$$u(x, 0) = \frac{M}{2\sqrt{2\pi}} \left(\exp\left(-\frac{(x-2)^2}{2}\right) + \exp\left(-\frac{(x+2)^2}{2}\right) \right), \quad x \in \Omega,$$

where $M \geq 0$ is the mass of $u(x, 0)$.

For most time steps, the value of α in the definition of the time step is 1. For the case $M = 1$, Fig. 5.10 displays the numerical solution at various times. Also, the locations and values of the Lagrange multiplier λ and the entropy with the KKT limiter enforced are shown. The results in Figs. 5.10 and 5.11 indicate that the numerical solution tends to a steady state, and that the Lagrange multiplier λ is needed to ensure that the positivity constraint is satisfied. Without the KKT limiter, the computations break down, even for very small CFL numbers.

We also compute the entropy for the mass $M = 1$ using the second order non-algebraic stable DIRK method

$$(a_{ij}) = \begin{pmatrix} 0 & 0 \\ 1/2 & 1/2 \end{pmatrix}, (b_i) = (1/2 \quad 1/2), (c_i) = (0 \quad 1), \tag{5.3}$$

the second order algebraic stable DIRK method

$$(a_{ij}) = \begin{pmatrix} 1/2 & 0 \\ 1/2 & 1/2 \end{pmatrix}, (b_i) = (1/2 \quad 1/2), (c_i) = (1/2 \quad 1), \tag{5.4}$$

and the second order DIRK method (3.6). Fig. 5.12 shows that the entropy for all three DIRK methods is dissipative.

For this model equation, there is a critical mass phenomenon [5], which states that solutions with a large initial mass blow-up in a finite time, while solutions with a small mass at an initial time will not. The numerical solutions with sub-critical mass $M = 1$ at times $t = 5$ and $t = 10$ and with super-critical mass $M = 10$ at times $t = 0.2$ and $t = 1$ are shown in Fig. 5.13 and Fig. 5.14, respectively, and agree with the results shown in [5] and the numerical observation in [1,9].

6. Conclusions

The main topic of this paper is the formulation of higher order accurate positivity preserving DIRK-LDG discretizations for the nonlinear degenerate parabolic Eq. (1.1). The presented numerical discretizations allow the combination of a positivity preserving limiter and time-implicit numerical discretizations for PDEs and alleviate the time step restrictions of currently available positivity preserving DG discretizations, which generally require the use of explicit time integration methods. For the spatial discretization an

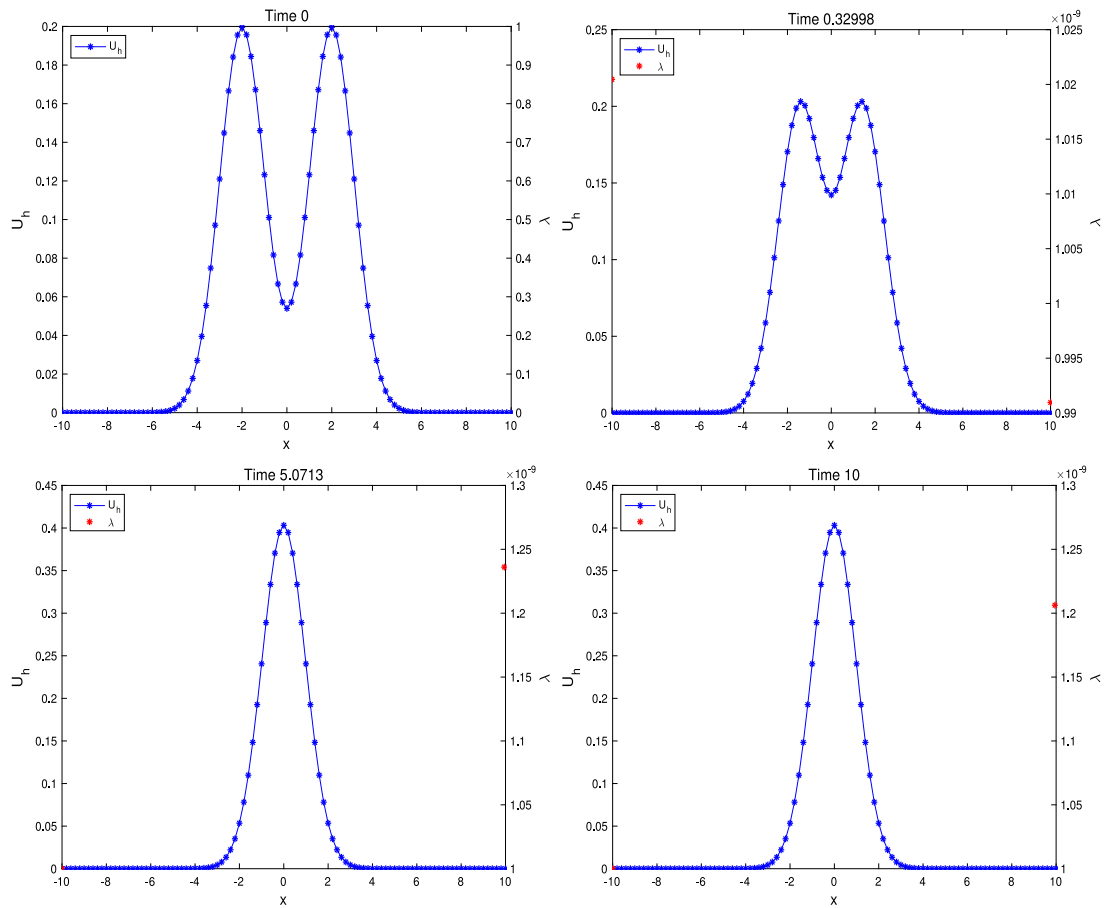


Fig. 5.10. (Example 5.6: mass $M=1$): Numerical solution U_h for \mathcal{P}_2 basis functions with KKT limiter enforced.

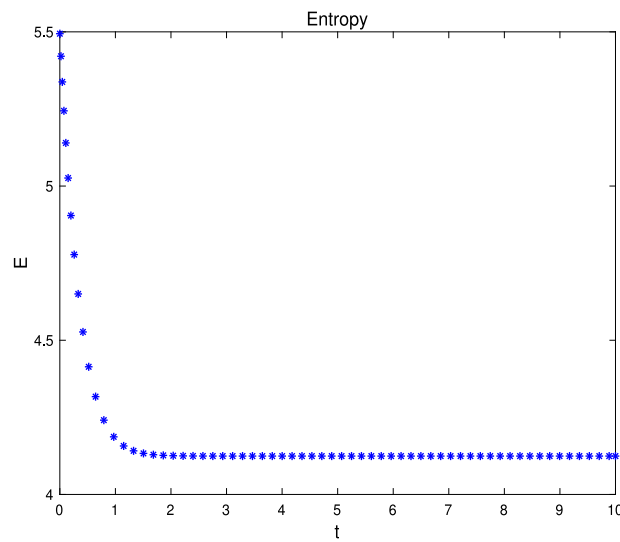


Fig. 5.11. (Example 5.6: mass $M=1$): Entropy E_h for \mathcal{P}_2 basis functions with KKT limiter enforced.

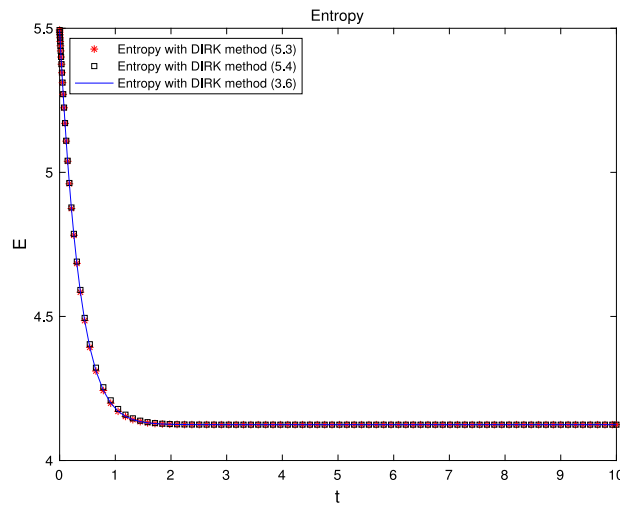


Fig. 5.12. (Example 5.6: mass $M=1$): Entropy E_h for \mathcal{P}_1 basis functions with KKT limiter enforced.

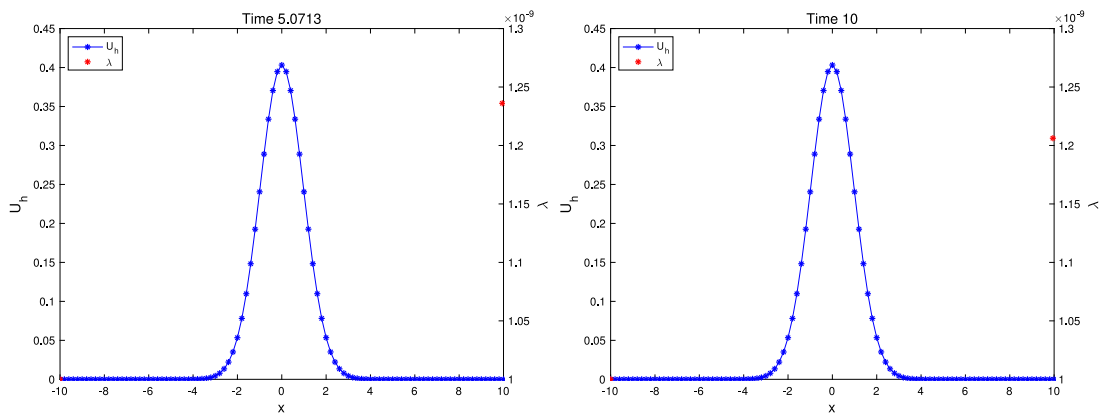


Fig. 5.13. (Example 5.6: mass $M = 1$): Numerical solution U_h for \mathcal{P}_2 basis functions with KKT limiter enforced.

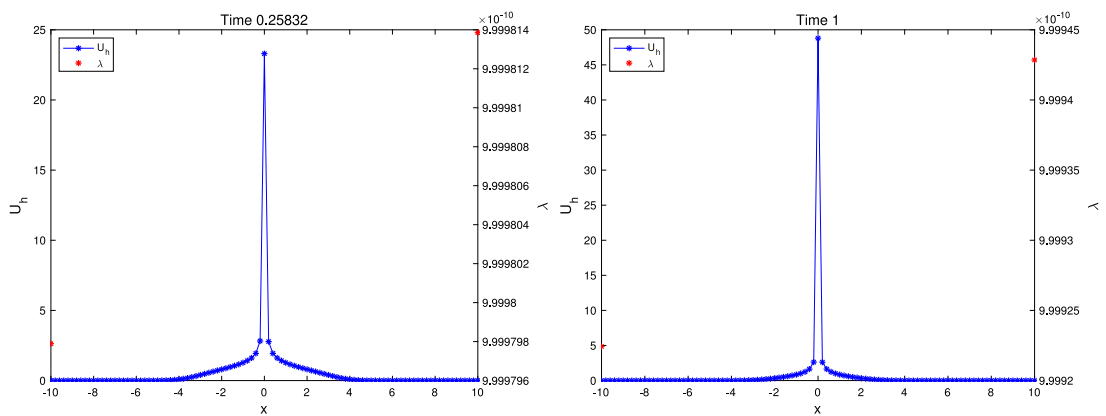


Fig. 5.14. (Example 5.6: mass $M = 10$) Numerical solution U_h for \mathcal{P}_2 basis functions with KKT limiter enforced.

LDG method combined with a simple alternating numerical flux is used, which simplifies the theoretical analysis for the entropy dissipation. For the temporal discretization, the numerical results show that implicit DIRK methods enlarge the time-step required for stability of the numerical discretization. Under a time-step restriction, we prove existence, uniqueness and entropy dissipation

of the positivity preserving high order accurate LDG discretization combined with an implicit Euler time discretization. Numerical results are presented to demonstrate the accuracy of the higher order accurate positivity preserving DIRK-LDG discretizations, which are of optimal order and not affected by the positivity preserving KKT limiter. The numerical solutions satisfy the entropy decay condition.

Data availability

Data will be made available on request.

Acknowledgments

The research of Fengna Yan was supported by a fellowship from the China Scholarship Council (No. 201806340058), and the Fundamental Research Funds for the Central Universities, China JZ2021HGTA0179, JZ2022HGQA0157. The research of J.J.W. van der Vegt was partially supported by the University of Science and Technology of China (USTC), Hefei, Anhui, China, while the author was in residence at USTC. The research of Yinhua Xia was partially supported by National Key R&D Program of China No. 2022YFA1005202/2022YFA1005200 and National Natural Science Foundation of China grant No. 12271498. The research of Yan Xu was partially supported by National Natural Science Foundation of China grant No. 12071455.

References

- [1] M. Bessemoulin-Chatard, F. Filbet, A finite volume scheme for nonlinear degenerate parabolic equations, *SIAM J. Sci. Comput.* 34 (2012) B559–B583.
- [2] J.L. Vázquez, *The Porous Medium Equation: Mathematical Theory*, Oxford University Press, 2007.
- [3] Q. Zhang, Z.-L. Wu, Numerical simulation for porous medium equation by local discontinuous Galerkin finite element method, *J. Sci. Comput.* 38 (2009) 127–148.
- [4] J.A. Carrillo, A. Chertock, Y. Huang, A finite-volume method for nonlinear nonlocal equations with a gradient flow structure, *Commun. Comput. Phys.* 17 (2015) 233–258.
- [5] N.B. Abdallah, I.M. Gamba, G. Toscani, On the minimization problem of sub-linear convex functionals, *Kinet. Relat. Models* 4 (2011) 857–871.
- [6] J.A. Carrillo, P. Laurençot, J. Rosado, Fermi-Dirac-Fokker-Planck equation: Well-posedness and long-time asymptotics, *J. Differential Equations* 247 (2009) 2209–2234.
- [7] G. Toscani, Finite time blow up in Kaniadakis-Quarati model of Bose-Einstein particles, *Comm. Partial Differential Equations* 37 (2012) 77–87.
- [8] M. Burger, J.A. Carrillo, M.-T. Wolfram, A mixed finite element method for nonlinear diffusion equations, *Kinet. Relat. Models* 3 (2010) 59–83.
- [9] H. Liu, Z. Wang, An entropy satisfying discontinuous Galerkin method for nonlinear Fokker-Planck equations, *J. Sci. Comput.* 68 (2016) 1217–1240.
- [10] H. Liu, H. Yu, Maximum-principle-satisfying third order discontinuous Galerkin schemes for Fokker-Planck equations, *SIAM J. Sci. Comput.* 36 (2014) A2296–A2325.
- [11] H. Liu, H. Yu, The entropy satisfying discontinuous Galerkin method for Fokker-Planck equations, *J. Sci. Comput.* 62 (2015) 803–830.
- [12] X. Cheng, J. Shen, A new Lagrange multiplier approach for constructing structure preserving schemes, I. Positivity preserving, *Comput. Methods Appl. Mech. Engrg.* 391 (2022) 114585.
- [13] B. Cockburn, C.-W. Shu, The local discontinuous Galerkin method for time-dependent convection–diffusion systems, *SIAM J. Numer. Anal.* 35 (1998) 2440–2463.
- [14] B. Cockburn, G. Kanschat, I. Perugia, D. Schötzau, Superconvergence of the local discontinuous Galerkin method for elliptic problems on Cartesian grids, *SIAM J. Numer. Anal.* 39 (2001) 264–285.
- [15] R. Guo, Y. Xu, A high order adaptive time-stepping strategy and local discontinuous Galerkin method for the modified phase field crystal equation, *Commun. Comput. Phys.* 24 (2018) 123–151.
- [16] L. Tian, Y. Xu, J.G. Kuerten, J.J.W. van der Vegt, An h-adaptive local discontinuous Galerkin method for the Navier-Stokes-Korteweg equations, *J. Comput. Phys.* 319 (2016) 242–265.
- [17] L. Zhou, Y. Xu, Stability analysis and error estimates of semi-implicit spectral deferred correction coupled with local discontinuous Galerkin method for linear convection–diffusion equations, *J. Sci. Comput.* 77 (2018) 1001–1029.
- [18] Y. Yang, D. Wei, C.-W. Shu, Discontinuous Galerkin method for Krause's consensus models and pressureless Euler equations, *J. Comput. Phys.* 252 (2013) 109–127.
- [19] X. Zhang, C.-W. Shu, On maximum-principle-satisfying high order schemes for scalar conservation laws, *J. Comput. Phys.* 229 (2010) 3091–3120.
- [20] X. Zhang, C.-W. Shu, On positivity-preserving high order discontinuous Galerkin schemes for compressible Euler equations on rectangular meshes, *J. Comput. Phys.* 229 (2010) 8918–8934.
- [21] T. Qin, C.-W. Shu, Implicit positivity-preserving high-order discontinuous Galerkin methods for conservation laws, *SIAM J. Sci. Comput.* 40 (2018) A81–A107.
- [22] F. Huang, J. Shen, Bound/Positivity preserving and energy stable scalar auxiliary variable schemes for dissipative systems: applications to Keller–Segel and Poisson-Nernst–Planck equations, *SIAM J. Sci. Comput.* 43 (2021) A1832–A1857.
- [23] J.J.W. van der Vegt, Y. Xia, Y. Xu, Positivity preserving limiters for time-implicit higher order accurate discontinuous Galerkin discretizations, *SIAM J. Sci. Comput.* 41 (2019) A2037–A2063.
- [24] E. Hairer, G. Wanner, *Solving Ordinary Differential Equations II. Stiff and Differential–Algebraic Problem*, Springer Science & Business Media, 2010.
- [25] R. Alexander, Diagonally implicit Runge-Kutta methods for stiff ODE's, *SIAM J. Numer. Anal.* 14 (1977) 1006–1021.
- [26] L. Skvortsov, Diagonally implicit Runge-Kutta methods for stiff problems, *Comput. Math. Math. Phys.* 46 (2006) 2110–2123.
- [27] F. Facchinei, J.-S. Pang, *Finite-Dimensional Variational Inequalities and Complementarity Problems, Volume I*, Springer-Verlag, New York, 2003.
- [28] Y.P. Kalmykov, W. Coffey, S. Titov, On the Brownian motion in a double-well potential in the overdamped limit, *Physica A* 377 (2007) 412–420.