# AN ENTROPY SATISFYING DISCONTINUOUS GALERKIN METHOD FOR NONLINEAR FOKKER-PLANCK EQUATIONS

HAILIANG LIU† AND ZHONGMING WANG‡

ABSTRACT. We propose a high order discontinuous Galerkin (DG) method for solving nonlinear Fokker-Planck equations with a gradient flow structure. For some of these models it is known that the transient solutions converge to steady-states when time tends to infinity. The scheme is shown to satisfy a discrete version of the entropy dissipation law and preserve steady-states, therefore providing numerical solutions with satisfying long-time behavior. The positivity of numerical solutions is enforced through a reconstruction algorithm, based on positive cell averages. For the model with trivial potential, a parameter range sufficient for positivity preservation is rigorously established. For other cases, cell averages can be made positive at each time step by tuning the numerical flux parameters. A selected set of numerical examples is presented to confirm both the high-order accuracy and the efficiency to capture the large-time asymptotic.

## 1. INTRODUCTION

In this paper, we propose a high order accurate discontinuous Galerkin (DG) method for solving the following problem

$$\partial_t u = \nabla_x \cdot (f(u) \nabla_x (\Phi(x) + H'(u))), \quad x \in \Omega, \ t > 0, \tag{1a}$$

$$u(x, 0) = u_0(x), \tag{1b}$$

subject to appropriate boundary conditions. Here $u(t, x) \geq 0$ is the unknown, $\Omega$ is a bounded domain in $\mathbb{R}^d$, $H : \mathbb{R}^+ \to \mathbb{R}$ and $f : \mathbb{R}^+ \to \mathbb{R}^+$ are given functions, and $\Phi(x)$ is a given potential function.

This equation has a gradient flow structure corresponding to the entropy functional

$$E = \int_\Omega (H(u) + u\Phi(x)) dx.$$

A simple calculation shows that the time derivative of this entropy along the equation (1a) with zero flux boundary condition is

$$\frac{d}{dt} E(t) = -\int_\Omega f(u) |\nabla_x (\Phi + H'(u))|^2 dx \leq 0, \tag{2}$$

which reveals the entropy dissipation property of the underlying system. Certain entropy dissipation inequalities are recognized to characterize the fine details of the convergence to steady states, see e.g., [7, 9, 11, 24].

Equations such as (1a) appear in a wide range of applications. In the case $f(u) = u$, the equation becomes

$$\partial_t u = \nabla_x \cdot (u \nabla_x (\Phi(x) + H'(u))). \tag{3}$$

If $H'(u) = u^m (m > 1)$ and $\Phi = 0$, it is the porous medium equation [11, 24], and for $H'(u) = \nu u^{m-1}$ and $\Phi = x^4/4 - x^2/2$, it is the nonlinear diffusion equation confined by a double-well

potential [6]. A particular example with nonlinear $f(u)$ is

(4) $$\partial_t u = \nabla_x \cdot (xu(1 + ku) + \nabla_x u),$$

which is known as a model for fermion ($k = -1$) and boson ($k = 1$) gases [8, 10, 28]. A more general class of the form

(5) $$\partial_t u = \nabla_x \cdot (xu(1 + u^N) + \nabla_x u), \quad N > 2,$$

is known to develop finite time concentration beyond some critical mass [1].

In order to capture the rich dynamics of solutions to (1), it is highly desirable to develop high order schemes which can preserve the entropy dissipation law (2) at the discrete level. In this work, we propose such a scheme for (1) using the discontinuous Galerkin discretization.

A related finite volume method was already proposed in [5] for (1), and further generalized to cover the nonlocal terms and general dimension in [6]. For (1) with $f(u) = u$ and an additional nonlocal interaction term, a mixed finite element method was studied in [4] based on their interpretation as gradient flows in optimal transportation metrics, following the so called JKO formulation, which is a variational scheme proposed by Jordan, Kinderlehrer and Otto [13] for linear Fokker-Planck equations. Regarding the use of relative entropy functionals we refer to [2] for the study of the large time behavior of a fully implicit semi-discretization applied to linear parabolic Fokker-Planck type equations in the form of (1) with $f(u) = u$, $H = u\log u$. A free energy satisfying finite difference method was proposed in [18] for the Poisson-Nernst-Planck (PNP) equations, which correspond to (1) with $f = u$, $H = u\log u$, further coupled with a Poisson equation for governing the potential $\Phi$. However, these existing schemes are only up to second-order.

An entropy satisfying DG method has been recently developed in [22] for the linear Fokker-Planck equation

(6) $$\partial_t u = \nabla_x \cdot (\nabla_x u + u\nabla_x \Phi),$$

which corresponds to (3) with $H = u\log u$. The obtained DG method generalizes and improves upon the finite volume method introduced in [21]. The idea in [22] is to apply the DG discretization to the non-logarithmic Landau formulation of (6),

$$\partial_t u = \nabla_x \cdot \left( M \nabla_x \left( \frac{u}{M} \right) \right), \quad M = e^{-\Phi(x)},$$

so that the quadratic entropy dissipation law is satisfied. Again based on this formulation, a third order DG scheme was further developed in [23] to numerically preserve the maximum principle: if $c_1 \leq u_0(x)/M \leq c_2$, then $c_1 \leq u(x,t)/M \leq c_2$ for all $t > 0$. However, the non-logarithmic Landau formulation does not apply directly to the more general class of equations (1a).

In this work, we construct an arbitrary high order entropy satisfying DG scheme for solving (1). The main idea behind the scheme construction is to apply the DG discretization to the following reformulation

(7) $$\partial_t u = \partial_x(f(u)\partial_x q), \quad q = \Phi(x) + H'(u),$$

by using a special numerical flux for $\partial_x q$. The resulting scheme is shown to feature several nice properties: (i) the entropy dissipation law (2) is satisfied at the discrete level; (ii) the steady states are shown to be preserved; (iii) for the third order scheme applied to the model with a trivial potential, a sufficient condition on the range of flux parameters is rigorously established so that cell averages remain positive at each time step, as long as each cell polynomial is positive at three test points. For the numerical positivity a reconstruction algorithm based on positive cell averages is introduced so that the positivity of cell polynomials is enforced, without destroying the accuracy, at least for smooth solutions. This reconstruction also serves as a limiter imposed

upon the numerical solution to suppress spurious oscillations at the solution singularity near zero. For the general case the positivity of cell averages can be achieved by carefully tuning the parameters in the numerical flux, as illustrated in the numerical experiments.

The discontinuous Galerkin (DG) method we discuss in this paper is a class of finite element methods, using a completely discontinuous piecewise polynomial space for the numerical solution and the test functions. One main advantage of the DG method was the flexibility afforded by local approximation spaces combined with the suitable design of numerical fluxes crossing cell interfaces. More general information about DG methods for elliptic, parabolic, and hyperbolic PDEs can be found in the recent books and lecture notes [12, 14, 26, 27]. Following the methodology of the direct discontinuous Galerkin (DDG) method proposed in [19, 20], we adopt a similar numerical flux formula for $\partial_x q$ in (7). The main feature in the DDG schemes proposed in [19, 20] lies in numerical flux choices for the solution gradient, which involve higher order derivatives evaluated crossing cell interfaces.

The plan of the paper is as follows. In Section 2, we present our DG scheme in one dimensional setting. In Section 3 we prove several important properties of the scheme, including the semi-discrete entropy dissipation law in Theorem 3.1, the fully-discrete entropy dissipation law in Theorem 3.3, the preservation of positive cell averages for the model with trivial potential in Theorem 3.4, and the preservation of steady states in Theorem 3.5. In Section 4, we elaborate various details in numerical implementation, including the reconstruction algorithm, the time discretization, and the spatial Numerical results are in Section 5, where we verify experimentally the high order spatial accuracy of our scheme and simulate the long-time behavior of numerical solutions. The proposed scheme is applied to several physical models including the porous medium equation, the nonlinear diffusion with a double-well potential, and the general Fokker–Planck equation. The numerical results confirm both the high order of accuracy and the numerical efficiency to capture the large-time asymptotic. Concluding remarks are given in Section 6.

## 2. DG DISCRETIZATION IN SPACE

In this section, we present our DG scheme for (1). For clarity of presentaiton , we restrict ourselves to the problem in one spatial dimension. It is straightforward to generalize this construction for Cartesian meshes in multidimensional case.

In one-dimensional setting, let $\Omega = [a, b]$ be a bounded interval. We divide $\Omega$ with a mesh

$$a = x_{1/2} < x_1 < \cdots < x_{N-1/2} < x_N < x_{N+1/2} = b,$$

and the mesh size $\Delta x_j = x_{j+1/2} - x_{j-1/2}$, and a family of $N$ control cells $I_j = (x_{j-1/2}, x_{j+1/2})$ with cell center $x_j = (x_{j-1/2} + x_{j+1/2})/2$. We denote by $v^+$ and $v^-$ the right and left limits of function $v$, and define

$$[v] = v^+ - v^-, \quad \{v\} = \frac{v^+ + v^-}{2}.$$

Define an $k-$degree discontinuous finite element space

$$V_h = \left\{ v \in L^2(\Omega), \quad v|_{I_j} \in P^k(I_j), j \in \mathbb{Z}_N \right\},$$

where $P^k(I_j)$ denotes the set of all polynomials of degree at most $k$ on $I_j$, and $\mathbb{Z}_r = \{1, \cdots, r\}$ for any positive integer $r$.

We rewrite the equation (1) as follows

(8a) $$\partial_t u = \partial_x(f(u)\partial_x q),$$

(8b) $$q = \Phi(x) + H'(u).$$

The DG scheme is to find $(u_h, q_h) \in V_h \times V_h$ such that for all $v, r \in V_h$ and $j \in \mathbb{Z}_N$,

(9a) $\quad \int_{I_j} \partial_t u_h v dx = -\int_{I_j} f(u_h) \partial_x q_h \partial_x v dx + \{f(u_h)\}\widehat{\partial_x q_h} v|_{\partial I_j} + \{f(u_h)\}\partial_x v(q_h - \{q_h\})|_{\partial I_j},$

(9b) $\quad \int_{I_j} q_h r dx = \int_{I_j} (\Phi(x) + H'(u_h)) r dx.$

Here

$$v|_{\partial I_j} = v(x^-_{j+1/2}) - v(x^+_{j-1/2}),$$

and $\widehat{\partial_x q_h}$ is the numerical flux, following [20], taken as

(10) $\quad\quad\quad\quad\quad\quad \widehat{\partial_x q_h} = \beta_0 \frac{[q_h]}{h} + \{\partial_x q_h\} + \beta_1 h [\partial_x^2 q_h],$

where $h = \Delta x$ for uniform meshes and $h = (\Delta x_j + \Delta x_{j+1})/2$ at $x_{j+1/2}$ for non-uniform meshes. Here $\beta_i, i = 0, 1$ are parameters satisfying a condition of the form

$$\beta_0 > \Gamma(\beta_1),$$

where $\Gamma(\beta_1)$ is chosen to ensure certain stability property of the underlying PDE.

Note that if zero-flux boundary conditions of the form $\partial_x(\Phi(x) + H'(u)) = 0$ are specified, we simply set $q$-related terms on the domain boundary to be zero. If a Dirichlet boundary condition for $u$ is given at $\partial\Omega$, we define the boundary numerical flux (10) in the following way:

(11a) $\quad\quad \{f(u_h)\} = \dfrac{f(u(a,t)) + f(u_h^+)}{2}$ if $x = a; \quad \dfrac{f(u_h^-) + f(u(b,t))}{2}$ if $x = b,$

(11b) $\quad\quad [q_h] = \begin{cases} q_h^+ - (\Phi(a) + H'(u(a,t))) & \text{for } x = a, \\ (\Phi(b) + H'(u(b,t))) - q_h^- & \text{for } x = b, \end{cases}$

(11c) $\quad\quad \{\partial_x q_h\} = \partial_x q_h^+$ if $x = a; \quad \partial_x q_h^-$ if $x = b,$

(11d) $\quad\quad [\partial_x^2 q_h] = 0.$

Here the boundary conditions are built into the scheme in such a way that the boundary data are used when available, otherwise the value of the numerical solution in corresponding end cells will be used.

## 3. Properties of the DG scheme

In this section, we investigate several desired properties of the semi-discrete DG scheme (9), and its time discretization.

3.1. **Entropy dissipation.** We first state the entropy satisfying property of DG scheme (9), using the following notation:

(12) $\quad\quad\quad ||q_h||_E^2 := \left[ \sum_{j=1}^N \int_{I_j} f(u_h)|\partial_x q_h|^2 dx + \sum_{j=1}^{N-1} \{f(u_h)\} \left( \frac{\beta_0}{h}[q_h]^2 \right) \Big|_{x_{j+\frac{1}{2}}} \right].$

**Theorem 3.1.** *Consider the DG scheme (9)-(10), subject to zero-flux boundary condition. If $f(u_h) \geq 0$, then the semi-discrete entropy*

$$E(t) = \sum_{j=1}^N \int_{I_j} (\Phi u_h + H(u_h)) dx$$

*satisfies*

$$(13) \qquad \frac{d}{dt} E(t) \leq -\gamma \|q_h\|_E^2$$

*for* $\gamma = 1 - \sqrt{\frac{\Gamma}{\beta_0}} \in (0, 1)$, *provided*

$$(14) \qquad \beta_0 > \Gamma(\beta_1) := \max_{1 \leq j \leq N-1} \frac{\{f(u_h)\} \left(\{\partial_x q_h\} + \frac{\beta_1}{2} h [\partial_x^2 q_h]\right)^2 \Big|_{x_{j+1/2}}}{\frac{1}{2h} \left(\int_{I_j} + \int_{I_{j+1}}\right) f(u_h) |\partial_x q_h|^2 dx}.$$

*Proof.* Summing (9)-(10) over all index $j$ we obtain a global formulation:

$$(15) \qquad \int_\Omega \partial_t u_h v \, dx = -\sum_{j=1}^N \int_{I_j} f(u_h) \partial_x q_h \partial_x v \, dx - \sum_{j=1}^{N-1} \{f(u_h)\} \left(\widehat{\partial_x q_h}[v] + \{\partial_x v\}[q_h]\right)_{j+1/2},$$

$$(16) \qquad \int_\Omega q_h r \, dx = \int_\Omega (\Phi + H'(u_h)) r \, dx.$$

Take $r = \partial_t u_h$ in (16) to obtain

$$\int_\Omega \partial_t u_h q_h \, dx = \int_\Omega (\Phi(x) + H'(u_h)) \partial_t u_h \, dx = \frac{d}{dt} \int_\Omega (\Phi u_h + H(u_h)) dx = \frac{d}{dt} E(t).$$

The right hand side from taking $v = q_h$ in (15) becomes

$$\frac{d}{dt} E(t) = -\sum_{j=1}^N \int_{I_j} f(u_h) |\partial_x q_h|^2 dx - \sum_{j=1}^{N-1} \{f(u_h)\} \left(\widehat{\partial_x q_h}[q_h] + \{\partial_x q_h\}[q_h]\right)_{j+1/2}$$

$$= -\sum_{j=1}^N \int_{I_j} f(u_h) |\partial_x q_h|^2 dx - \sum_{j=1}^{N-1} \{f(u_h)\} \left(\beta_0 [q_h]^2/h + [q_h](2\{\partial_x q_h\} + \beta_1 h [\partial_x^2 q_h])\right)_{j+1/2}.$$

Using Young's inequality we obtain

$$-(2\{\partial_x q_h\} + \beta_1 h [\partial_x^2 q_h])[q_h] \leq \beta_0 (1 - \gamma)[q_h]^2/h + \frac{h}{4\beta_0(1-\gamma)} \left(2\{\partial_x q_h\} + \beta_1 h [\partial_x^2 q_h]\right)^2$$

for some $0 < \gamma < 1$. Hence

$$(17) \qquad \frac{d}{dt} E(t) \leq -\gamma \left[\sum_{j=1}^N \int_{I_j} f(u_h) |\partial_x q_h|^2 dx + \sum_{j=1}^{N-1} \left(\frac{\{f(u_h)\}\beta_0}{h}[q_h]^2\right)_{j+1/2}\right]$$

$$-\left[(1-\gamma) \sum_{j=1}^N \int_{I_j} f(u_h) |\partial_x q_h|^2 dx - \sum_{j=1}^{N-1} \frac{h\{f(u_h)\}}{4\beta_0(1-\gamma)} \left(2\{\partial_x q_h\} + \beta_1 h [\partial_x^2 q_h]\right)^2\right]$$

$$\leq -\gamma \left[\sum_{j=1}^N \int_{I_j} f(u_h) |\partial_x q_h|^2 dx + \sum_{j=1}^{N-1} \left(\frac{\{f(u_h)\}\beta_0}{h}[q_h]^2\right)_{j+1/2}\right]$$

$$-\frac{1-\gamma}{2} \int_{I_1 \cup I_N} f(u_h) |\partial_x q_h|^2 dx,$$

since $\beta_0$ satisfies (14), hence

$$\beta_0(1-\gamma)^2 = \Gamma \geq \frac{\sum_{j=1}^{N-1} h\{f(u_h)\}\left(\{\partial_x q_h\} + \frac{\beta_1}{2}h[\partial_x^2 q_h]\right)_{j+1/2}^2}{\left(\sum_{j=2}^{N-1} \int_{I_j} + \frac{1}{2}\int_{I_1 \cup I_N}\right) f(u_h)|\partial_x q_h|^2 dx}.$$

This finishes the proof of (13). $\qquad\square$

*Remark* 3.1. We remark that a larger, yet simpler, $\Gamma(\beta_1)$ can be found for sufficiently small $h$ since the variation of ratio $\frac{\{f\}}{f}$ is also small. Assume that this ratio is bounded by a factor 2, i.e., $2 \geq \frac{f}{\{f\}} \geq \frac{1}{2}$, then

$$\Gamma(\beta_1) \leq 2 \max_{1 \leq j \leq N-1} \frac{\left(\{\partial_x q_h\} + \frac{\beta_1}{2}h[\partial_x^2 q_h]\right)^2\Big|_{x_{j+1/2}}}{\frac{1}{2h}\left(\int_{I_j} + \int_{I_{j+1}}\right)|\partial_x q_h|^2 dx}$$

$$\leq 2 \max_{1 \leq j \leq N-1} \frac{\left(\frac{\partial_x q_h^- - \beta_1 h\partial_x^2 q_h^-}{2}\right)_{x_{j+1/2}}^2 + \left(\frac{\partial_x q_h^+ + \beta_1 h\partial_x^2 q_h^+}{2}\right)_{x_{j+1/2}}^2}{\frac{1}{2h}\left(\int_{I_j}|\partial_x q_h|^2 dx + \int_{I_{j+1}}|\partial_x q_h|^2 dx\right)}$$

It is clear that this inequality is implied by

$$(18) \qquad \Gamma(\beta_1) \leq 2 \max_{1 \leq j \leq N-1} \left\{\frac{(\partial_x q_h^- - \beta_1 h\partial_x^2 q_h^-)^2}{\frac{1}{2h}\int_{I_j}|\partial_x q_h|^2}, \frac{(\partial_x q_h^+ + \beta_1 h\partial_x^2 q_h^+)^2}{\frac{1}{2h}\int_{I_{j+1}}|\partial_x q_h|^2}\right\}.$$

By setting $v(\xi) = \partial_x q_h\left(x_j + \frac{h}{2}\xi\right)$ for $q_h(x)|_{I_j}$, and $v(\xi) = \partial_x q_h\left(x_{j+1} - \frac{h}{2}\xi\right)$ for $q_h|_{I_{j+1}}$, we have

$$\Gamma(\beta_1) \leq 2 \sup_{v \in P^{k-1}} \frac{(v(1) - 2\beta_1 \partial_\xi v(1))^2}{\frac{1}{2}\int_{-1}^1 |v|^2 d\xi} = 2k^2\left(1 - \beta_1(k^2-1) + \frac{\beta_1^2}{3}(k^2-1)^2\right),$$

here we have used the exact formula in [15, Lemma 3.1]. Hence it suffices to choose $\beta_0$ such that

$$(19) \qquad \beta_0 > 2k^2\left(1 - \beta_1(k^2-1) + \frac{\beta_1^2}{3}(k^2-1)^2\right).$$

*Remark* 3.2. The positivity of numerical solutions are realized through a reconstruction algorithm at each time step, based on positive cell averages, as detailed in Section 4.1. It is shown in Theorem 3.4 that the use of non-zero $\beta_1$ is crucial in the sense that the positivity of cell averages can be ensured. Indeed, this is proved for the third order DG scheme in solving (1) with zero potential. For the model with non-trivial potential, our numerical experiments again confirm the special role of $\beta_1$ in the preservation of positivity of numerical cell averages.

3.2. **The fully-discrete DG scheme.** In order to preserve the entropy dissipation law for $u_h^n$ at each time step, the time step restriction is needed when using an explicit time discretization. We now discuss this issue by taking the Euler first order time discretization of (9): find $u_h^{n+1}(x) \in V_h$ such that for any $r(x), v(x) \in V_h$,

$$(20a) \qquad \int_{I_j} q_h^n r \, dx = \int_{I_j} \left(\Phi(x) + H'(u_h^n)\right) r \, dx,$$

$$(20b) \qquad \int_{I_j} D_t u_h^n v \, dx = -\int_{I_j} f(u_h^n)\partial_x q_h^n \partial_x v \, dx + \{f(u_h^n)\} \left[\widehat{\partial_x q_h^n} v + \partial_x v(q_h^n - \{q_h^n\})\right]\Big|_{\partial I_j}.$$

Here and in what follows, we use the notation for any function $w^n(x)$ as

$$D_t w^n = \frac{w^{n+1} - w^n}{\Delta t},$$

and $\mu = \frac{\Delta t}{h^2}$ as the mesh ratio.

**Lemma 3.2.** *The following inverse inequalities hold for any $v \in V_h$:*

(21a)
$$\sum_{j=1}^N \int_{I_j} v_x^2 dx \leq \frac{k(k+1)^2(k+2)}{h^2} \sum_{j=1}^N \int_{I_j} v^2 dx,$$

(21b)
$$\sum_{j=1}^{N-1} [v]_{j+1/2} \leq \frac{4(k+1)^2}{h} \sum_{j=1}^N \int_{I_j} v^2 dx,$$

(21c)
$$\sum_{j=1}^{N-1} \{v_x\}_{j+1/2}^2 \leq \frac{k^3(k+1)^2(k+2)}{h^2} \sum_{j=1}^N \int_{I_j} v^2 dx.$$

*Proof.* These follow from the repeated use of the two inverse inequalities:

(22a)
$$\max\{|w(a)|, |w(b)|\} \leq (m+1)|I|^{-1/2} \|w\|_{L^2(I)},$$

(22b)
$$\|\partial_x w\|_{L^2(I)} \leq (m+1)\sqrt{m(m+2)}|I|^{-1} \|w\|_{L^2(I)},$$

provided $w \in P^m(I)$ with $I = (a, b)$ and $|I| = b - a$. The first bound is well known, see e.g. [29]. The second inequality may be found in [17, Lemma 3.1] $\square$

**Theorem 3.3.** *Let the fully discrete entropy be defined as*

$$E^n = \sum_{j=1}^N \int_{I_j} (\Phi(x) u_h^n(x) + H(u_h^n(x))) \, dx.$$

*The DG scheme (20), subject to zero-flux boundary condition, satisfies*

(23)
$$D_t E^n \leq -\frac{\gamma}{2} \|q_h^n\|_E^2$$

*for some $\gamma \in (0, 1)$, provided $u_h^n(x)$ remains positive, $\beta_0 > \Gamma(\beta_1)$, and*

(24)
$$\mu \leq \frac{\gamma}{C(k, \beta_0, \beta_1)\| \max\{0, H''(u_h^n(\cdot))\}\|_\infty \|f(u_h^n(\cdot))\|_\infty},$$

*where $C(k, \beta_0, \beta_1)$ is given in (29) below.*

*Proof.* Summing (20) over all index $j$'s we obtain

(25)
$$\sum_{j=1}^N \int_{I_j} q_h^n r \, dx = \sum_{j=1}^N \int_{I_j} (\Phi(x) + H'(u_h^n)) \, r \, dx,$$

(26)
$$\sum_{j=1}^N \int_{I_j} D_t u_h^n v \, dx = -\sum_{j=1}^N \int_{I_j} f(u_h^n)\partial_x q_h^n \partial_x v \, dx - \sum_{j=1}^{N-1} \{f(u_h^n)\} \left( \widehat{\partial_x q_h^n}[v] + \{\partial_x v\}[q_h^n] \right)\Big|_{x_{j+\frac{1}{2}}}.$$

Take $r = D_t u_h^n$ in (25) to obtain

$$\int_\Omega D_t u_h^n q_h^n dx = \int_\Omega \left(\Phi(x) + H'(u_h^n(x))\right) D_t u_h^n \, dx$$

$$= D_t E^n - \frac{1}{\Delta t} \int_\Omega (H(u_h^{n+1}) - H(u_h^n) - H'(u_h^n)(u_h^{n+1} - u_h^n))dx$$

$$= D_t E^n - \frac{\Delta t}{2} \int_\Omega H''(\cdot)(D_t u_h^n)^2 dx.$$

Here $(\cdot)$ denotes the intermediate value between $u_h^n$ and $u_h^{n+1}$. Taking $v = q_h^n$, (26) becomes

$$\int_\Omega D_t u_h^n q_h^n dx = -\sum_{j=1}^N \int_{I_j} f(u_h^n)|\partial_x q_h^n|^2 \, dx - \sum_{j=1}^{N-1} \{f(u_h^n)\}[q_h^n]\left(\widehat{\partial_x q_h^n} + \{\partial_x q_h^n\}\right)\Big|_{x_{j+\frac{1}{2}}}$$

$$\leq -\gamma \|q_h^n\|_E^2,$$

for $\beta_0$ satisfying (14) at each interface $x_{j+\frac{1}{2}}$, $j = 1, \dots, N-1$. Hence

$$D_t E^n \leq -\gamma \|q_h^n\|_E^2 + \frac{\Delta t}{2} \int_\Omega H''(\cdot)(D_t u_h^n)^2 \, dx.$$

The claimed estimate follows if

(27) $$\Delta t \leq \frac{\gamma \|q_h^n\|_E^2}{\int_\Omega \max\{0, H''(\cdot)\}(D_t u_h^n)^2 \, dx}.$$

For convex $H$, this indeed imposes a time restriction.

It remains to show that the bound in (24) is smaller than the right side of (27). In (26), we take $v = D_t u_h^n$ and use the Young inequality $ab \leq \frac{1}{4\epsilon}a^2 + \epsilon b^2$ to obtain

$$\sum_{j=1}^N \int_{I_j} v^2 \, dx = -\sum_{j=1}^N \int_{I_j} f(u_h^n)\partial_x q_h^n \partial_x v \, dx - \sum_{j=1}^{N-1} \{f(u_h^n)\}\left(\widehat{\partial_x q_h^n}[v] + \{\partial_x v\}[q_h^n]\right)\Big|_{x_{j+\frac{1}{2}}}$$

$$\leq \frac{1}{4\epsilon_1 h^2}\sum_{j=1}^N \int_{I_j} f^2(u_h^n)|\partial_x q_h^n|^2 \, dx + \epsilon_1 h^2 \sum_{j=1}^N \int_{I_j} |\partial_x v|^2 \, dx$$

$$+ \frac{1}{4\epsilon_2 h}\sum_{j=1}^{N-1} \{f(u_h^n)\}^2|\widehat{\partial_x q_h^n}|^2\Big|_{x_{j+\frac{1}{2}}} + \epsilon_2 h \sum_{j=1}^{N-1} [v]^2\Big|_{x_{j+\frac{1}{2}}}$$

$$+ \frac{1}{4\epsilon_3 h^3}\sum_{j=1}^{N-1} \{f(u_h^n)\}^2[q_h^n]^2\Big|_{x_{j+\frac{1}{2}}} + \epsilon_3 h^3 \sum_{j=1}^{N-1} \{\partial_x v\}^2\Big|_{x_{j+\frac{1}{2}}}.$$

The use of inequalities in (21) leads to

$$\epsilon_1 h^2 \sum_{j=1}^N \int_{I_j} |\partial_x v|^2 \, dx + \epsilon_2 h \sum_{j=1}^{N-1} [v]^2\Big|_{x_{j+\frac{1}{2}}} + \epsilon_3 h^3 \sum_{j=1}^{N-1} [\partial_x v]^2\Big|_{x_{j+\frac{1}{2}}}$$

$$\leq (k+1)^2(k(k+2)\epsilon_1 + 4\epsilon_2 + k^3(k+2)\epsilon_3) \sum_{j=1}^N \int_{I_j} v^2 \, dx$$

$$= \frac{3}{4}\sum_{j=1}^N \int_{I_j} v^2 \, dx,$$

provided
$$(4\epsilon_1)^{-1} = k(k+1)^2(k+2), \ (4\epsilon_2)^{-1} = 4(k+1)^2, \quad (4\epsilon_3)^{-1} = k^3(k+1)^2(k+2).$$

This gives

$$(28) \quad \frac{1}{4}\sum_{j=1}^{N}\int_{I_j} v^2 \, dx \leq \frac{k(k+1)^2(k+2)}{h^2}\sum_{j=1}^{N}\int_{I_j} f^2(u_h^n)|\partial_x q_h^n|^2 \, dx$$

$$+ \frac{k^3(k+1)^2(k+2)}{h^3}\sum_{j=1}^{N-1}\{f(u_h^n)\}^2[q_h^n]^2\big|_{x_{j+\frac{1}{2}}} + \frac{4(k+1)^2}{h}\sum_{j=1}^{N-1}\{f(u_h^n)\}^2|\widehat{\partial_x q_h^n}|^2\big|_{x_{j+\frac{1}{2}}}.$$

It is clear that the first two terms are bounded by $\|f(u_h^n(\cdot)\|_\infty \|q_h^n\|_E^2$. We now show that the last term is also bounded by $\|f(u_h^n(\cdot)\|_\infty \|q_h^n\|_E^2$, up to constant multiplication factors.

$$\sum_{j=1}^{N-1}\{f(u_h^n)\}|\widehat{\partial_x q_h^n}|^2\big|_{x_{j+\frac{1}{2}}} = \sum_{j=1}^{N-1}\{f(u_h^n)\}\left|\{\partial_x q_h^n\} + \beta_0\frac{[q_h^n]}{h} + \beta_1 h[\partial_x^2 q_h^n]\right|^2\Bigg|_{x_{j+\frac{1}{2}}}$$

$$\leq 2\sum_{j=1}^{N-1}\{f(u_h^n)\}\left(\beta_0^2\frac{[q_h^n]^2}{h^2} + \left(\{\partial_x q_h^n\} + \beta_1 h[\partial_x^2 q_h^n]\right)^2\right)\Bigg|_{x_{j+\frac{1}{2}}}.$$

From (14) it follows that

$$\{f(u_h^n)\}\left(\{\partial_x q_h^n\} + \beta_1 h[\partial_x^2 q_h^n]\right)^2\big|_{x_{j+\frac{1}{2}}} \leq \frac{\Gamma(2\beta_1)}{2h}\left(\int_{I_j} + \int_{I_{j+1}}\right)f(u_h)|\partial_x q_h|^2 dx.$$

Hence

$$\sum_{j=1}^{N-1}\{f(u_h^n)\}\left(\{\partial_x q_h^n\} + \beta_1 h[\partial_x^2 q_h^n]\right)^2\big|_{x_{j+\frac{1}{2}}} \leq \frac{\Gamma(2\beta_1)}{h}\sum_{j=1}^{N}\int_{I_j} f(u_h)|\partial_x q_h|^2 dx.$$

These together yield

$$\sum_{j=1}^{N-1}\{f(u_h^n)\}|\widehat{\partial_x q_h^n}|^2\big|_{x_{j+\frac{1}{2}}} \leq \frac{2}{h}\max\{\beta_0, \Gamma(2\beta_1)\}\|q_h^n\|_E^2.$$

Upon insertion into (28) we obtain

$$\sum_{j=1}^{N}\int_{I_j} v^2 \, dx \leq \frac{C(k,\beta_0,\beta_1)\|f(u_h^n(\cdot))\|_\infty}{h^2}\|q_h^n\|_E^2,$$

where

$$(29) \qquad C(k,\beta_0,\beta_1) := 4(k+1)^2\left(k(k+2)\max\{1, k^2/\beta_0\} + 8\max\{\beta_0, \Gamma(2\beta_1)\}\right).$$

Hence (27) is implied by (24).

This ends the proof. □

3.3. **Preservation of positive cell averages.** It is known to be difficult, if not impossible, to preserve point-wise solution bounds for high order numerical approximations. A popular strategy after the work [30] is to combine an accuracy preserving reconstruction with the bound preserving property of cell averages. For the DG scheme applied to (1) with $\Phi = 0$, following [23], we are able to identify a range of $\beta_1$ so that positive averages are ensured for at least the third order scheme. We have not been able to prove this property for the general case.

By taking the test function $v = 1$ on $I_j$ in (20b), we obtain the evolutionary equation for the cell average,

$$(30) \qquad \bar{u}_j^{n+1} = \bar{u}_j^n + \mu h \left. \{f(u_h^n)\} \widehat{\partial_x q_h^n} \right|_{\partial I_j}.$$

For the case that $H$ is convex and $\Phi(x) = 0$, we reformulate (8) as

$$\partial_t u = \partial_x (f H'' \partial_x q), \quad q = u.$$

At the discrete level, we simply set $q_h = u_h$ and replace $f$ by $f H''$ in (20b). Assuming that $\bar{u}_j^n \in [c_1, c_2]$ for all $j$'s, we can derive some sufficient conditions such that $\bar{u}_j^{n+1} \in [c_1, c_2]$ under certain CFL condition on $\mu$.

For piecewise quadratic polynomials, we have the following result.

**Theorem 3.4.** $(k = 2)$ *The scheme (30) with $q_h = u_h$, and*

$$(31) \qquad \frac{1}{8} < \beta_1 < \frac{1}{4} \quad and \quad \beta_0 \geq 1$$

*is bound preserving, namely, $\bar{u}_j^{n+1} \in [c_1, c_2]$ if $u_h^n(x) \in [c_1, c_2]$ on the set $S_j$'s where*

$$S_j = x_j + \frac{h}{2} \{-1, 0, 1\},$$

*under the CFL condition*

$$(32) \qquad \mu \leq \mu_0 = \frac{1}{12 \max_{1 \leq j \leq N} |f(u_{j-1/2}^n)|} \min \left\{ \frac{1}{\beta_0 + 8\beta_1 - 2}, \frac{1}{1 - 4\beta_1} \right\}.$$

*Proof.* Let

$$p(\xi) = u_h \left( x_j + \frac{h}{2} \xi \right) \text{ for } \xi \in [-1, 1], \quad \text{i.e.,} \quad p = u_h|_{I_j},$$

we have

$$(33) \qquad \bar{u}_j = \frac{1}{6} p(-1) + \frac{2}{3} p(0) + \frac{1}{6} p(1).$$

In what follows we denote $p_- = u_h|_{I_{j-1}}$ and $p_+ = u_h|_{I_{j+1}}$.

We represent the diffusion flux in terms of solution values over the set $S_j$; see [23].

$$(34) \qquad \left. h \widehat{\partial_x u_h} \right|_{x_{j+\frac{1}{2}}} = \alpha_3 p_+(-1) + \alpha_2 p_+(0) + \alpha_1 p_+(1) - (\alpha_1 p(-1) + \alpha_2 p(0) + \alpha_3 p(1)),$$

where

$$(35) \qquad \alpha_1 = \frac{8\beta_1 - 1}{2}, \quad \alpha_2 = 2(1 - 4\beta_1), \quad \alpha_3 = \beta_0 + \frac{8\beta_1 - 3}{2}.$$

It is easy to verify that (31) ensures $\alpha_i \geq 0$ for $i = 1, 2, 3$.

Upon substitution into (30) we obtain

$$
\begin{aligned}
(36) \qquad \bar{u}_j^{n+1} =& \bar{u}_j + 2\mu \left( h\{f(u_h)\}\widehat{\partial_x u_h}\Big|_{x_{j+\frac{1}{2}}} - h\{f(u_h)\}\widehat{\partial_x u_h}\Big|_{x_{j-\frac{1}{2}}} \right) \\
=& \left[ \frac{1}{6} - 2\mu \left( \alpha_3 f_{j-\frac{1}{2}} + \alpha_1 f_{j+\frac{1}{2}} \right) \right] p(-1) \\
&+ \left[ \frac{2}{3} - 2\mu \left( \alpha_2 f_{j-\frac{1}{2}} + \alpha_2 f_{j+\frac{1}{2}} \right) \right] p(0) \\
&+ \left[ \frac{1}{6} - 2\mu \left( \alpha_1 f_{j-\frac{1}{2}} + \alpha_3 f_{j+\frac{1}{2}} \right) \right] p(1) \\
&+ 2\mu f_{j+\frac{1}{2}} \left[ \alpha_3 p_+(-1) + \alpha_2 p_+(0) + \alpha_1 p_+(1) \right] \\
&+ 2\mu f_{j-\frac{1}{2}} \left[ \alpha_1 p_-(-1) + \alpha_2 p_-(0) + \alpha_3 p_-(1) \right].
\end{aligned}
$$

Here we have used the notation

$$
f_{j+\frac{1}{2}} := \{f(u_h)\}|_{x_{j+\frac{1}{2}}} = \left. \frac{f(u_h^-) + f(u_h^+)}{2} \right|_{x_{j+\frac{1}{2}}}.
$$

Note that the sum of all coefficients of above polynomial values is one. Hence $\bar{u}_j^{n+1} \in [c_1, c_2]$ as long as $u_h^n \in [c_1, c_2]$ on $S_j$ and all coefficients are nonnegative. The nonnegativity imposes a CFL condition $\mu \le \mu_0$ with $\mu_0$ being

$$
\frac{1}{12} \min_{1 \le j \le N} \left\{ \frac{1}{\alpha_3 f_{j-\frac{1}{2}} + \alpha_1 f_{j+\frac{1}{2}}}, \frac{4}{\alpha_2 f_{j-\frac{1}{2}} + \alpha_2 f_{j+\frac{1}{2}}}, \frac{1}{\alpha_1 f_{j-\frac{1}{2}} + \alpha_3 f_{j+\frac{1}{2}}} \right\}.
$$

Here we assume that $f_{N+1/2} = 0$ so that $j = N$ can be included in the above expression. It suffices to take smaller

$$
\mu_0 = \frac{1}{12 \max |f(u_{j-1/2}^n)|} \min \left\{ \frac{1}{\alpha_3 + \alpha_1}, \frac{2}{\alpha_2} \right\}.
$$

That is (32), as claimed. $\qquad \square$

*Remark* 3.3. The CFL condition (32) is sufficient conditions rather than necessary to preserve the bound of solutions. Therefore, in practice, these CFL conditions are strictly enforced only in the case the bound preserving property is violated.

*Remark* 3.4. For general case, we expect there is still a proper set of parameters $(\beta_0, \beta_1)$ with which the scheme can preserve positivity of cell averages. Our numerical simulations in Example 2 confirms this expectation.

3.4. **Preservation of steady states.** If we start with an initial data $u_h^0$, already at steady states, i.e., $\Phi(x) + H'(u_h^0(x)) = C$, it follows from (20a) that $q_h^0 = C$. Furthermore, (20b) implies that $u_h^1 = u_h^0 \in V_h$. By induction we have

$$
\Phi(x) + H'(u_h^n(x)) = C \quad \forall n \in \mathbb{N}.
$$

This says that the DG scheme (20a) preserves the steady states. Moreover, we can show that in some cases the numerical solution tends asymptotically toward a steady state, independent of initial data. More precisely, we have the following result.

**Theorem 3.5.** *Let the assumptions in Theorem 3.3 be met, and $(u_h^n, q_h^n)$ be the numerical solution to the fully discrete DG scheme (20), then the limits of $(u_h^n, q_h^n)$ as $n \to \infty$ satisfy*

$$
q_h^* = C, \quad \Phi(x) + H'(u_h^*) \in C + V_h^\perp,
$$

*where $C$ is a constant. For quadratic $H(u)$, $C$ can be determined explicitly by*

$$C = \frac{1}{|\Omega|} \int_{\Omega} (\Phi(x) + H'(u_0)(x))dx.$$

*In addition, if $\Phi(x) \in P^m (m \leq k)$, then we must have $\Phi(x) + H'(u_h^*(x)) \equiv C$.*

*Proof.* Since $E^n$ is non-increasing and bounded from below, we have

$$\lim_{n \to \infty} E^n = \inf\{E^n\}.$$

Observe from (23) that

$$E^{n+1} - E^n \leq -\frac{\gamma \Delta t}{2} \|q_h^n\|_E^2 \leq 0.$$

When passing the limit $n \to \infty$ we have $\lim_{n \to \infty} \|q_h^n\|_E^2 = 0$. This implies that each term in this energy norm must have zero as its limit, that is

$$(37) \qquad \lim_{n \to \infty} \sum_{j=1}^{N} \int_{I_j} f(u_h^n)|\partial_x q_h^n|^2 dx = 0, \quad \lim_{n \to \infty} \sum_{j=1}^{N-1} \frac{\beta_0}{h}\{f(u_h^n)\}[q_h^n]^2\Big|_{j+\frac{1}{2}} = 0.$$

The first relation in (37) tells that the limit of $q_h^n$, denoted by $q_h^*$, must be constant in each computational cell. The second relation in (37) infers that $q_h^*$ must be a constant in the whole domain. These when inserted into (20a) gives the desired result. For quadratic $H(u)$, we use the mass conservation $\int_{\Omega} H'(u_h^*(x))dx = \int_{\Omega} H'(u_0(x))dx$ to determine the constant $C$. The proof is complete. $\qquad \square$

*Remark* 3.5. The above result shows that for quadratic $H(u)$ and potential $\Phi(x)$ being polynomials of degree up to $k$, the steady states are approached by numerical solutions. For other cases, such asymptotic convergence holds only in the projection sense.

## 4. NUMERICAL IMPLEMENTATION

In this section, we provide further details in implementing the entropy satisfying discontinuous Galerkin (ESDG) method.

4.1. **Reconstruction.** For a high order polynomial approximation, numerical solutions can have negative values. We enforce the solution positivity through some accuracy-preserving reconstruction. Motivated by the definite result on the bound preserving property of cell averages for special cases in Theorem 3.4, we consider the case with positive cell averages.

Let $w_h \in P^k(I_j)$ be an approximation to a smooth function $w(x) \geq 0$, with cell averages $\bar{w}_j > \delta$ for $\delta$ being some small positive parameter or zero. We then reconstruct another polynomial in $P^k(I_j)$ so that

$$(38) \qquad \tilde{w}_h^{\delta}(x) = \bar{w}_j + \frac{\bar{w}_j - \delta}{\bar{w}_j - \min_{I_j} w_h(x)}(w_h(x) - \bar{w}_j), \quad \text{if } \min_{I_j} w_h(x) < \delta.$$

This reconstruction maintains same cell averages and satisfies

$$\min_{I_j} w^{\delta}(x) \geq \delta.$$

It is known that enforcing a maximum principle numerically might damp oscillations in numerical solutions, see, e.g. [16, 30]. Numerical example in Fig.1 confirms such a damping effect near zero from using the positivity preserving limiter (38).

**Lemma 4.1.** *If $\bar{w}_j > \delta$, then the reconstruction satisfies the estimate*

$$|w^\delta(x) - w_h(x)| \le C(k) \left( ||w_h(x) - w(x)||_\infty + \delta \right), \quad \forall x \in I_j,$$

*where $C(k)$ is a constant depending on $k$. This says that the reconstructed $w^\delta(x,t)$ in (38) does not destroy the accuracy when $\delta < h^{k+1}$.*

*Proof.* We have

$$|w^\delta(x) - w_h(x)| = \left| \frac{\delta - \min_{I_j} w_h(x)}{\bar{w}_j - \min_{I_j} w_h(x)} (\bar{w}_j - w_h(x)) \right|$$

$$\le \frac{\max_{I_j} |\bar{w}_j - w_h(x)|}{\max_{I_j} (\bar{w}_j - w_h(x))} \left( ||w_h(x) - w(x)||_\infty + \delta \right).$$

It follows from [23, 30] that

$$\frac{\max_{I_j} |\bar{w}_j - w_h(x)|}{\max_{I_j} (\bar{w}_j - w_h(x))} \le C(k),$$

where $k$ is the degree of the polynomial $w_h(x)$. $\square$

4.2. **Time discretization.** For the time discretization of (9), we use the explicit high order Runge-Kutta method. The explicit time discretization is simple to implement, with entropy dissipation law still preserved under some restriction on the time step.

Let $\{t^n\}, n = 0, 1, \ldots$ be a uniform partition of time interval. Denote $u_h^n \sim u(t_n, x)$, $q_h^n \sim q(t_n, x)$, where $t_n = n\Delta t$ and $\Delta t$ is the uniform temporal step size. The algorithm can be summarized in following steps.

1. Project $u_0(x)$ onto $V_h$ to obtain $u_h(0)$ and solve (9b) to obtain $q_h(0)$.
2. Solve (9a) to obtain $u_h^{n+1}$ with a Runge-Kutta (RK) ODE solver. Perform reconstruction (38) if needed.
3. Solve (9b) to obtain $q_h^{n+1}$ from the obtained $u_h^{n+1}$.
4. Repeat steps 2 and 3 until final time $T$.

In our numerical simulation we choose $\Delta t = C(k)h^2$, where $C(k)$ is smaller for larger $k$. For the case with zero potential and $k = 2$, $C(k)$ is given in Theorem 3.4. The choice of the time step $\Delta t \sim h^2$ suggests that we adopt an $m^{th}$ order Runge-Kutta solver with $m \ge (k+1)/2$, so that in the accuracy test the temporal error is smaller than the spatial error. For polynomials of degree $k = 1, 2, 3$, we use the second order explicit Runge-Kutta method (also called Heun's method) to solve the ODE system $\dot{a} = \mathfrak{L}(\mathbf{a})$:

$$\mathbf{a}^{(1)} = \mathbf{a}^n + \Delta t \mathfrak{L}(\mathbf{a}^n),$$

$$\mathbf{a}^{n+1} = \frac{1}{2}\mathbf{a}^n + \frac{1}{2}\mathbf{a}^{(1)} + \frac{1}{2}\Delta t \mathfrak{L}(\mathbf{a}^{(1)}).$$

The bound preserving property for cell averages in Theorem 3.3, depending on a convex combination of polynomial values in previous time step, works well with the above Runge-Kutta solver since it is simply a convex combination of the forward Euler.

4.3. **Spatial discretization.** In this section, we present some further details on the spatial discretization. The $k$th order basis functions in a 1-D standard reference element $\xi \in [-1, 1]$ are taken as the Legendre polynomials $\{L_i(\xi)\}_{i=0}^k$, then the numerical solutions in each cell $x \in I_j$ can be expressed as

$$u_h(x,t) = \sum_{i=0}^k u_j^i(t) L_i(\xi) =: L^\top(\xi) u_j(t), \quad q_h(x,t) = \sum_{i=0}^k q_j^i(t) L_i(\xi) =: L^\top(\xi) q_j(t),$$

using a uniform mesh size $h$ and the map $x = x_j + \frac{h}{2}\xi$, with notation $L^\top = (L_0, L_1, \cdots, L_k)$ and $u_j = (u_j^0, \cdots, u_j^k)^\top$.

For given $\Phi(x)$, a simple calculation of (9a) with $v = L(\xi)$ gives

$$(39) \qquad M\dot{u}_j = \frac{2}{h}R_1 + \frac{1}{2h}(R_2 + R_3), \quad 2 \le j \le N - 1,$$

where

$$M = \frac{h}{2} \int_{-1}^{1} L(\xi)L^\top(\xi)d\xi,$$

$$R_1 = -\sum_{i=1}^{Q} \omega_i f\left(L^\top(s_i)u_j(t)\right) L_\xi^\top(s_i)q_j L_\xi(s_i),$$

$$R_2 = \left(f\left(L^\top(1)u_j\right) + f\left(L^\top(-1)u_{j+1}\right)\right)(-D^\top q_j + E^\top q_{j+1})L(1)$$
$$- \left(f\left(L^\top(1)u_{j-1}\right) + f\left(L^\top(-1)u_j\right)\right)(-D^\top q_{j-1} + E^\top q_j)L(-1) = R_2^+ - R_2^-,$$

$$R_3 = \left(f\left(L^\top(1)u_j\right) + f\left(L^\top(-1)u_{j+1}\right)\right)(L^\top(1)q_j - L^\top(-1)q_{j+1})L_\xi(1)$$
$$+ \left(f\left(L^\top(1)u_{j-1}\right) + f\left(L^\top(-1)u_j\right)\right)(L^\top(1)q_{j-1} - L^\top(-1)q_j)L_\xi(-1) =: R_3^+ + R_3^-.$$

Here

$$D = \beta_0 L(1) - L_\xi(1) + 4\beta_1 L_{\xi\xi}(1), \quad E = \beta_0 L(-1) + L_\xi(-1) + 4\beta_1 L_{\xi\xi}(-1).$$

In the evaluation of $R_1$, we choose $Q$ Gaussian quadrature points $s_i \in [-1, 1]$ with $1 \le i \le Q$. Here and in what follows, we choose $Q$ quadrature points with $Q \ge \frac{k+2}{2}$ so that the quadrature rule with accuracy of order $\mathcal{O}(h^{2Q})$ does not destroy the scheme accuracy. At two end cells, if the zero flux conditions are specified, we use $R_2 = R_2^+, R_3 = R_3^+$ for $j = 1$ and $R_2 = -R_2^-, R_3 = R_3^-$ for $j = N$.

If Dirichlet boundary conditions, $u(a)$ and $u(b)$, are specified, we modify $R_2$ and $R_3$ according to (11). That is, for $j = 1$,

$$R_2 = R_2^+ - (f(u(a)) + f\left(L^\top(-1)u_1\right))[\beta_0(L^\top(-1)q_1 - \Phi(a) - H'(u(a))) + 2L_\xi^\top(-1)q_1]L(-1),$$

$$R_3 = R_3^+ + (f(u(a)) + f\left(L^\top(-1)u_1\right))[\Phi(a) + H'(u(a)) - L^\top(-1)q_1]L_\xi(-1),$$

and for $j = N$,

$$R_2 = (f\left(L^\top(1)u_N\right) + f(u(b)))[-\beta_0(L^\top(1)q_N - \Phi(b) - H'(u(b))) + 2L_\xi^\top(1)q_N]L(1) - R_2^-,$$

$$R_3 = f\left(L^\top(1)u_N\right) + f(u(b)))[L^\top(1)q_N - \Phi(b) - H'(u(b))]L_\xi(1) + R_3^-.$$

To solve (9b) is, using the $Q$-point Gauss quadrature rule on the interval $(-1, 1)$, to solve

$$(40) \qquad Mq_j = \frac{h}{2}\sum_{i=1}^{Q} \omega_i(\Phi(x(s_i)) + H'(L^\top(s_i)u_j))L(s_i).$$

The collection of (39) and (40) with $1 \le j \le N$ forms a nonlinear ODE system, for which we use a Runge-Kutta method.

## 5. Numerical Tests

In this section, we present a selected set of numerical examples in order to numerically validate our ESDG scheme. Via several physical models from different applications, we examine the order of accuracy by numerical convergence tests, while we quantify $l_1$ errors defined by

$$\|u_h - u_{ref}\|_{l_1} = \sum_{j=1}^{N} \int_{I_j} |u_h(x) - u_{ref}(x)| dx,$$

with the integral on $I_j$ evaluated by a 4-point Gaussian quadrature method and $u_{ref}$ being a reference solution obtained by using a refined mesh size. It is also demonstrated that the scheme captures well the long-time behavior of underlying solutions, as well as the mass concentration phenomenon in certain applications.

5.1. **Porous medium equation.** We consider the porous medium equation of the form

$$(41) \qquad\qquad \partial_t u = \partial_x^2(u^m), \quad m > 1.$$

With this model we will illustrate 1) the scheme's capability in capturing the solution singularity; 2) the positivity preservation proved in Theorem 3.4.

**Example 1. Capturing singularity**
Barenblatt and Pattle independently found an explicit solution of (41) when the Dirac delta function is used as initial condition [3, 25]. A special explicit solution which we will use is

$$(42) \qquad B_m(x,t) = \max\left\{0, t^{-\alpha}\left(0.2 - \frac{\alpha(m-1)}{2m}\frac{|x|^2}{t^{2\alpha}}\right)^{\frac{1}{m-1}}\right\}, \quad \alpha = \frac{1}{m+1}.$$

We compute the solution of (41) with initial data $u_0(x) = B_2(x, 0.1)$, with zero flux boundary conditions $\partial_x u(\pm 2, t) = 0$.

Fig.1 shows the exact solution and $P^2$ numerical solutions without and with reconstruction (38) with $\delta$ set to be 0. This reconstruction is not applied to the cells where the $u_h$ are entirely zero. The scheme with reconstruction gives sharp resolution of expanding fronts, keeping the solution strictly within the initial bounds. The scheme without reconstruction brings visible undershoots near the foot of the numerical solution.

Fig.2 shows a numerical comparison for polynomials with different degrees, $k = 1, 2, 3$. Cell averages are shown in Fig.2 (left) and cell polynomials in Fig.2(right) (zoomed near singularity), we can clearly see that a higher order method gives a more accurate approximation.
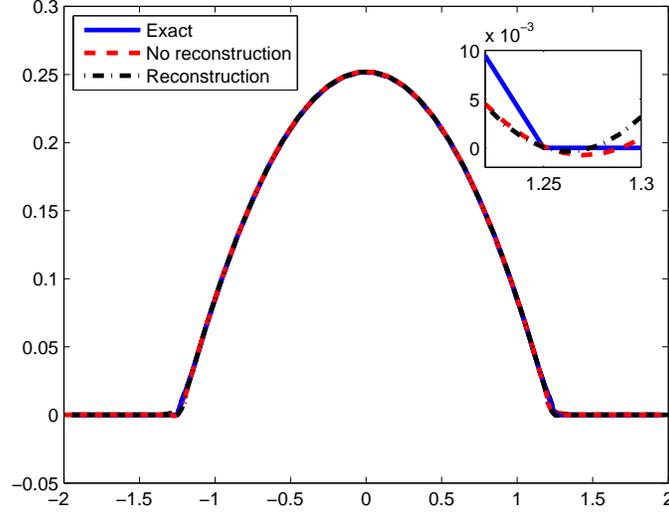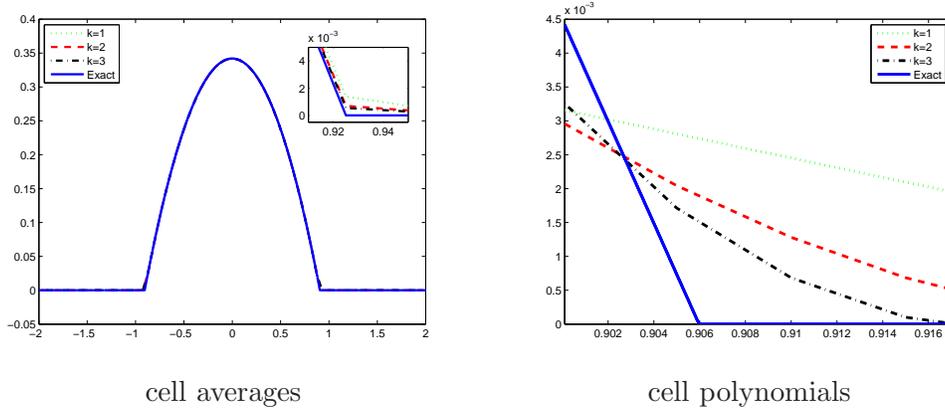
**Example 2. Positivity preservation**
In this example we test the effect of using different parameter $\beta_1$ in terms of the positivity preservation. Equation (41) with $m = 2$, when written in the form

$$\partial_t u = \partial_x(f(u)\partial_x q), \quad f(u) = 2u, \quad q = u,$$

satisfies the requirements in Theorem 3.4. We consider positive initial data with small amplitude,

$$u_0(x) = \epsilon(1 + 30e^{-25x^2}), \quad x \in [-1, 1],$$

and zero flux boundary conditions $\partial_x u(\pm 1, t) = 0$. With $\epsilon = 10^{-5}$, $\delta = 10^{-10}$, $h = 0.2$, $k = 2$ and $\Delta t = 0.25h^2$ in the simulation, our results indicate that cell average $\bar{u}$ remains above $\delta$ at $t = 1000$ when using $(\beta_0, \beta_1) = (2, 1/6)$; while $\bar{u}$ already becomes negative at $t = 41.388$ when taking $(\beta_0, \beta_1) = (2, 0)$. This is consistent with the conclusion in Theorem 3.4 that $\beta_1 \in (1/8, 1/4)$ is sufficient for positivity preservation of cell averages, and for any other $\beta_1$'s such a property is not guaranteed. We note here that the range of $\beta_1$ in Theorem 3.4 is only

FIGURE 1. Capturing singularity in the exact solution at $t = 0.5$



FIGURE 2. Comparison of solutions for $k = 1, 2, 3$



cell averages

cell polynomials

sufficient. Our simulation also indicates that cell average $\bar{u}$ still remains above $\delta$ at $t = 1000$ when using $(\beta_0, \beta_1) = (2, 1/2)$, which does not satisfy the requirement in Theorem 3.4.

We further test the special effect of parameter $\beta_1$ on the positivity preservation for the case with nontrivial potential, $\Phi = 30\epsilon x^2/2$, i.e., we have

$$\partial_t u = \partial_x(f(u)\partial_x q), \quad f(u) = 2u, \quad q = u + 30\epsilon x^2/2.$$

Though Theorem 3.4 is no longer applicable due to the nonzero potential, we still see similar effects of $\beta_1$ through numerical experiments. With the same initial condition and parameters as above, our simulation results in Table 1 show that there is a range for $\beta_1$ in which $\bar{u}$ remains above $\delta$ at $t = 1000$; while $\bar{u}$ becomes negative at $t < 1000$ when $\beta_1 \leq 1/6$ or $\beta_1 \geq 2$. This observation indicates that 1) $\beta_1$ plays a special role for the positivity preservation; 2) the admissibility of $\beta_1$ depends on the underlying problem.

TABLE 1. Time when $\bar{u}$ becomes negative

| $(\beta_0, \beta_1)$ | negative $\bar{u}$ time |
|---|---|
| (2,0) | 35.41 |
| (2, 1/12) | 388.91 |
| (2,1/6) | 845.69 |
| (2,1/3) | >1000 |
| (2,1/2) | >1000 |
| (2,2/3) | >1000 |
| (2,1) | >1000 |
| (2,2) | 917.42 |
| (2,3) | 740.92 |

5.2. **Porous medium equation with linear convection.** We consider the following porous medium equation with linear convection

$$\partial_t u = \partial_x^2(u^m) + \partial_x u, \quad m > 1.$$

This equation corresponds to (1a) with $f(u) = u$, $\Phi = x$ and $H = \frac{u^m}{m-1}$, and has a wide range of applications. With this model equation we shall test the numerical convergence and the scheme accuracy. We note that the case $m = 2$ was tested in [5] with a second order finite volume scheme.

**Example 3 (m=2).** We consider

$$\partial_t u = \partial_x^2(u^2) + \partial_x u,$$

with initial data

$$u_0(x) = 0.5 + 0.5\sin(\pi x), \quad x \in [-1, 1],$$

subject to zero-flux boundary condition, that is $\partial_x u(\pm 1, t) = -\frac{1}{2}$. In Table 2 we observe that the orders of convergence are of $\mathcal{O}(h^{k+1})$ for polynomials of degree $k$ ($k = 1, 2, 3$).

TABLE 2. Error table for porous media equation with $m = 2$ at $t = 1$

| $(k, \beta_0, \beta_1)$ | h | $l_1$ error | order |
|---|---|---|---|
| $(1, 1, -)$ | 0.4 | 0.0056949 | – |
| | 0.2 | 0.0013756 | 2.15 |
| | 0.1 | 0.00034588 | 2.20 |
| | 0.05 | 6.5394e-005 | 2.40 |
| $(2, 4, 1/12)$ | 0.4 | 0.00026132 | – |
| | 0.2 | 3.9026e-005 | 2.86 |
| | 0.1 | 5.3072e-006 | 2.91 |
| | 0.05 | 6.8756e-007 | 2.95 |
| $(3, 9, 1/4)$ | 0.4 | 4.4584e-005 | – |
| | 0.2 | 4.4365e-006 | 3.71 |
| | 0.1 | 3.2099e-007 | 3.91 |
| | 0.05 | 1.9724e-008 | 4.02 |

**Example 4 (m=3).** We further test the case $m = 3$, i.e.,

$$\partial_t u = \partial_x^2(u^3) + \partial_x u,$$

with initial data
$$u_0(x) = 1 + 0.5\sin(\pi x), \quad x \in [-1, 1],$$
subject to zero-flux boundary conditions $(uu_x)(\pm 1, t) = -1/3$. The numerical convergence test is performed with the same flux parameters for each $k$ as in the previous example, both errors and orders of convergence are given in Table 3. These results further confirm the $(k+1)$-th order of accuracy when using $P^k(k = 1, 2, 3)$ elements.

TABLE 3. Error table for porous medium equation with $m = 3$ at $t = 1$

| $(k, \beta_0, \beta_1)$ | h | $l_1$ error | order |
|---|---|---|---|
| | 0.4 | 0.0014749 | – |
| $(1, 1, -)$ | 0.2 | 0.00037363 | 1.99 |
| | 0.1 | 9.5215e-005 | 1.99 |
| | 0.05 | 2.3636e-005 | 2.01 |
| | 0.4 | 7.3404e-005 | – |
| $(2, 4, 1/12)$ | 0.2 | 9.5432e-006 | 2.97 |
| | 0.1 | 1.2268e-006 | 2.98 |
| | 0.05 | 1.5257e-007 | 3.00 |
| | 0.4 | 5.1001e-006 | – |
| $(3, 9, 1/4)$ | 0.2 | 3.4917e-007 | 3.96 |
| | 0.1 | 2.1473e-008 | 4.00 |
| | 0.05 | 1.3609e-009 | 3.98 |

Numerical tests in Example 3 and 4 also indicate that cell averages can be made positive in time when choosing proper parameters $(\beta_0, \beta_1)$, together with reconstruction (38) performed at each time step.

5.3. **Nonlinear diffusion with a double-well potential.** Consider a nonlinear diffusion equation with an external double-well potential of the form
$$\partial_t u = \partial_x(u\partial_x(\nu u^{m-1} + \Phi)), \quad \Phi = \frac{x^4}{4} - \frac{x^2}{2}.$$
This model equation is taken from [6], and it corresponds to system (1) with $H'(u) = \nu u^{m-1}$. With this model we shall test both numerical accuracy and the asymptotic behavior of numerical solutions.

**Example 5. Free energy decay**
In this example, we take $\nu = 1$, $m = 2$ and initial data
$$u_0(x) = \frac{0.1}{\sqrt{0.4\pi}}e^{-\frac{x^2}{0.4}}, \quad x \in [-2, 2],$$
subject to zero-flux boundary conditions $\partial_x u(\pm 2, t) = \mp 6$. Both errors and orders of convergence are given in Table 4, which again demonstrates $\mathcal{O}(h^{k+1})$ order of accuracy for $P^k$ polynomials.
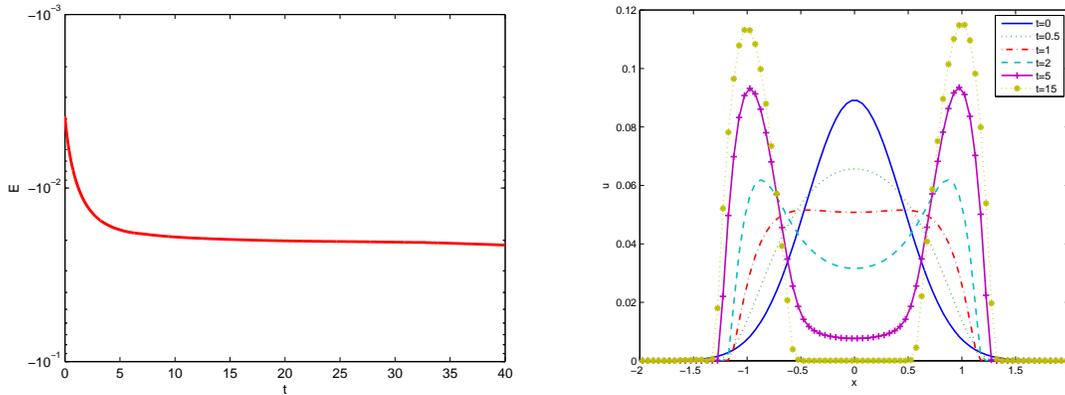
We also examine the decay of the entropy
$$E = \int_{-2}^{2} (\Phi(x)u + H(u))\,dx = \int_{-2}^{2}\left[\left(\frac{x^4}{4} - \frac{x^2}{2}\right)u + \frac{u^2}{2}\right]dx.$$

Figure 3 (left) shows the semilog plot of the free energy decay until final time $T = 40$, and Figure 3 (right) displays the snapshots of $u$ at different times, showing the time-asymptotic convergence of the numerical solutions towards the steady states.

TABLE 4. Error table for nonlinear diffusion with a double-well potential at $t = 1$

| $(k, \beta_0, \beta_1)$ | h | $l_1$ error | order |
|---|---|---|---|
| | 0.4 | 0.082882 | – |
| $(1, 1, -)$ | 0.2 | 0.0051793 | 2.70 |
| | 0.1 | 0.0012178 | 2.06 |
| | 0.05 | 0.00029961 | 2.02 |
| | 0.4 | 0.16726 | – |
| $(2, 4, 1/12)$ | 0.2 | 0.020986 | 3.08 |
| | 0.1 | 0.0023122 | 3.18 |
| | 0.05 | 0.00027875 | 3.05 |
| | 0.8 | 0.09677 | – |
| $(3, 12, 1/24)$ | 0.4 | 0.010059 | 3.82 |
| | 0.2 | 0.00051784 | 4.10 |
| | 0.1 | 3.4058e-005 | 3.93 |

FIGURE 3. Entropy decay of nonlinear diffusion with double well potential



5.4. **The nonlinear Fokker-Planck equation.** We consider the following model for boson gases,

$$(43) \qquad \partial_t u = \partial_x(xu(1 + u^3) + \partial_x u), \quad t > 0,$$

which is a nonlinear Fokker-Planck equation corresponding to (1a) with

$$\Phi = \frac{x^2}{2}, \quad f(u) = u(1 + u^3), \quad H'(u) = \log \frac{u}{\sqrt[3]{1 + u^3}}.$$

This model equation exhibits the critical mass phenomenon (see [1]), that solutions with initial data of large mass blow-up in finite time, whereas solutions with initial data of small mass do not. The authors in [5] numerically verified such critical mass phenomenon using a second order finite volume scheme. With our high order DG scheme, we test the critical mass phenomenon for (43) with initial data
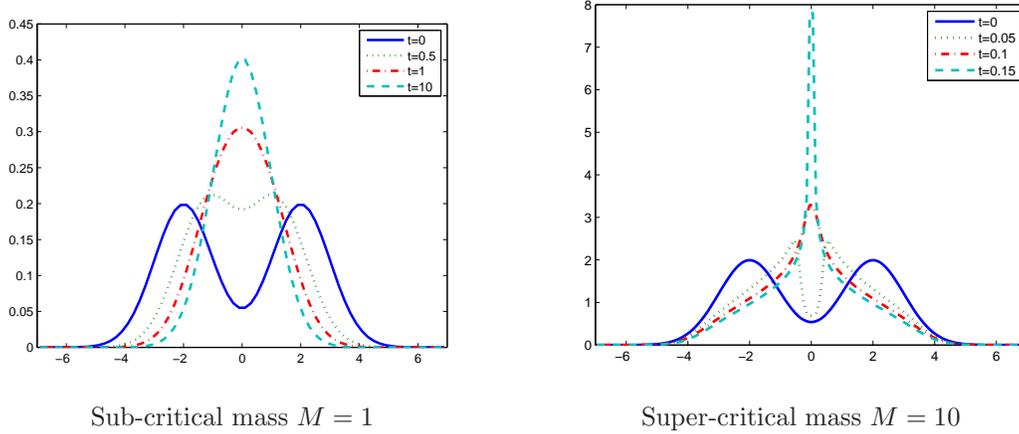
$$u_0(x) = \frac{M}{2\sqrt{2\pi}} \left( \exp\left( -\frac{(x-2)^2}{2} \right) + \exp\left( -\frac{(x+2)^2}{2} \right) \right),$$

which has total mass $M$. This is to illustrate the good performance of the ESDG scheme in capturing complex physical phenomena.

**Example 6. Sub-critical mass $M = 1$ and super-critical mass $M = 10$**
We test the sub-critical mass $M = 1$ with results in Figure 4 (left) and super-critical mass $M = 10$ with results in Figure 4 (right) by $P^2$ polynomial approximations. These results are consistent with the theoretical conclusion made in [1] and the numerical observation in [5], yet our scheme can produce numerical solutions with higher order of accuracy. Note that the reconstruction (38) has to be implemented due to the involvement of log-function in $H'(u)$.

FIGURE 4. Dynamics of the general Fokker-Planck equation



Sub-critical mass $M = 1$                        Super-critical mass $M = 10$

## 6. CONCLUDING REMARKS

In this article, we have developed an entropy satisfying DG method for solving nonlinear Fokker-Planck equations with a gradient flow structure. The idea is to rewrite the equation in the form of a convection equation with flux being $-f(u)\partial_x q$, and $q$ is obtained by a piecewise $L^2$ projection of $\Phi(x) + H'(u)$. Then we apply the numerical flux of the DDG method introduced in [20] to $\partial_x q$. The present scheme is shown to satisfy a discrete version of the entropy dissipation law, therefore preserving steady-states and providing numerical solutions with satisfying long-time behavior. The positivity of numerical solutions is enforced through a reconstruction algorithm, based on positive cell averages. Cell averages can be made positive at each time step by carefully tuning the numerical flux parameter $(\beta_0, \beta_1)$. For the model with trivial potential, a parameter range sufficient for positivity preservation is rigorously established. Numerical examples include the porous medium equation, the nonlinear diffusion equation with a double-well potential, and the general Fokker-Planck equation. Numerical results have demonstrated high-order accuracy of the scheme. Moreover, the long-time solution behavior is also examined to show the robustness of the proposed scheme.

## ACKNOWLEDGMENTS

## REFERENCES

[1] N. Ben Abdallah, I. M. Gamba, and G. Toscani. On the minimization problem of sub-linear convex functionals. *Kinet. Relat. Models*, 4(4):857–871, 2011.
[2] A. Arnold and A. Unterreiter. Entropy decay of discretized Fokker-Planck equations I—Temporal semidiscretization. *Comput. Math. Appl.*, 46(10-11):1683–1690, 2003.

[3] G. I. Barenblatt. On some unsteady fluid and gas motions in a porous medium. *Prikladnaya Matematika i Mekhanika (Applied Mathematics and Mechanics (PMM))*, 16, No. 1, pp. 67-78 , 1952 (in Russian).

[4] M. Burger, J. A. Carrillo, and M.-T. Wolfram. A mixed finite element method for nonlinear diffusion equations. *Kinet. Relat. Models*, 3:59–83, 2010.

[5] M. Bessemoulin-Chatard and F. Filbet. A finite volume scheme for nonlinear degenerate parabolic equations. *SIAM J. Sci. Comput.*, 34(5):B559–B583, 2012.

[6] J. Carrillo, A. Chertock, and Y. H. Huang. A finite-volume method for nonlinear nonlocal equations with a gradient flow structure. *Commun. Comput. Phys.*, 17:233–258, 2015.

[7] J. A. Carrillo, A. Jüngel, P. A. Markowich, G. Toscani, and A. Unterreiter. Entropy dissipation methods for degenerate parabolic problems and generalized Sobolev inequalities. *Monatsh. Math.*, 133(1):1–82, 2001.

[8] J. A. Carrillo, P. Laurençot, and J. Rosado. Fermi-Dirac-Fokker-Planck equation: well-posedness & long-time asymptotics. *J. Differential Equations*, 247(8):2209–2234, 2009.

[9] J. A. Carrillo, R. J. McCann, and C. Villani. Kinetic equilibration rates for granular media and related equations: entropy dissipation and mass transportation estimates. *Rev. Mat. Iberoam.*, 19:971–1018, 2003.

[10] J. A. Carrillo, J. Rosado, and F. Salvarani. 1D nonlinear Fokker-Planck equations for fermions and bosons. *Appl. Math. Lett.*, 21(2):148–154, 2008.

[11] J. A. Carrillo and G. Toscani. Asymptotic $L^1$-decay of solutions of the porous medium equation to self-similarity. *Indiana Univ. Math. J.*, 49(1):113–142, 2000.

[12] J. S. Hesthaven and T. Warburton. Nodal Discontinuous Galerkin Methods: Algorithms, Analysis, and Applications. Springer, New York, 2007.

[13] R. Jordan, D. Kinderlehrer, and F. Otto. The variational formulation of the Fokker–Planck equation. *SIAM J. Math. Anal.*, 29(1): 1–17, 1998.

[14] B. Q. Li. Discontinuous Finite Elements in Fluid Dynamics and Heat Transfer. Computational Fluid and Solid Mechanics, Springer, London, 2006.

[15] H. Liu. Optimal error estimates of the direct discontinuous Galerkin method for convection–diffusion equations. *Math. Comp.*, 84: 2263–2295, 2015.

[16] X. Liu and S. Osher. Nonoscillatory high order accurate self-Similar maximum principle satisfying shock capturing schemes I. *SIAM J. Number. Anal.*, 33(2):760–779, 1996.

[17] H. Liu and M. Pollack. Alternating evolution discontinuous Galerkin methods for convection-diffusion equations. *J. Comput. Phys.*, in press, 2016.

[18] H. Liu and Z. Wang. A free energy satisfying finite difference method for Poisson-Nernst-Planck equations. *J. Comput. Phys.*, 268:363–376, 2014.

[19] H. Liu and J. Yan. The direct discontinuous Galerkin (DDG) methods for diffusion problems. *SIAM J. Numer. Anal.*, 47: 675–698, 2009.

[20] H. Liu and J. Yan. The direct discontinuous Galerkin (DDG) method for diffusion with interface corrections. *Commun. Comput. Phys.*, 8(3):541–564, 2010.

[21] H. Liu and H. Yu. An entropy satisfying conservative method for the Fokker–Planck equation of the finitely extensible nonlinear elastic dumbbell model. *SIAM J. Numer. Anal.*, 50:1207–1239, 2012.

[22] H. Liu and H. Yu. The entropy satisfying dicontinuous Galerkin method for Fokker-Planck equations. *J. Sci. Comput.* 62: 803–830, 2015.

[23] H. Liu and H. Yu. Maximum-Principle-Satisfying third order discontinuous Galerkin schemes for Fokker–Planck equations. *SIAM J. Sci. Comput.*, 36(5):A2296–A2325, 2014.

[24] F. Otto. The geometry of dissipative evolution equations: the porous medium equation. *Comm. Partial Differential Equations*, 26(1-2):101–174, 2001.

[25] R. E. Pattle. Diffusion from an instantaneous point source with a concentration-dependent coefficient. *Quart. J. Mech. Appl. Math.*, 12:407-409, 1959.

[26] B. Rivière. Discontinuous Galerkin Methods for Solving Elliptic and Parabolic Equations: Theory and Implementation, SIAM, Philadelphia, 2008.

[27] C.-W. Shu. Discontinuous Galerkin methods: General approach and stability, in Numerical Solutions of Partial Differential Equations, S. Bertoluzza, S. Falletta, G. Russo, and C.-W. Shu, eds. Advanced Courses in Mathematics, CRM Barcelona, Birkhaüiser, Basel, 2009, pp. 149201.

[28] G. Toscani. Finite time blow up in Kaniadakis-Quarati model of Bose-Einstein particles. *Comm. Partial Differential Equations*, 37(1):77–87, 2012.

[29] T. Warburton and J. S. Hesthaven. On the constants in hp-finite element trace inequalities. *Comput. Methods Appl. Mech. Engin.* 192:2765–2773, 2003.

[30] X.-X. Zhang and C.-W. Shu. On maximum-principle-satisfying high order schemes for scalar conservation laws. *J. Comput. Phys.*, 229(9):3091–3120, 2010.

†Iowa State University, Mathematics Department, Ames, IA 50011
*E-mail address*: hliu@iastate.edu

‡ Florida International University, Department of Mathematics and Statistics, Miami, FL 33199
*E-mail address*: zwang6@fiu.edu