

Michigan Technological University Digital Commons @ Michigan Tech

Dissertations, Master's Theses and Master's Reports

2019

# Discontinuous Galerkin methods for convection-diffusion equations and applications in petroleum engineering

Nattaporn Chuenjarern Michigan Technological University, nchuenja@mtu.edu

Copyright 2019 Nattaporn Chuenjarern

#### **Recommended Citation**

Chuenjarern, Nattaporn, "Discontinuous Galerkin methods for convection-diffusion equations and applications in petroleum engineering", Open Access Dissertation, Michigan Technological University, 2019.

https://doi.org/10.37099/mtu.dc.etdr/786

Follow this and additional works at: https://digitalcommons.mtu.edu/etdr

Part of the Computational Engineering Commons, Numerical Analysis and Computation Commons, and the Partial Differential Equations Commons

## DISCONTINUOUS GALERKIN METHODS FOR CONVECTION-DIFFUSION EQUATIONS AND APPLICATIONS IN PETROLEUM ENGINEERING

By

Nattaporn Chuenjarern

#### A DISSERTATION

Submitted in partial fulfillment of the requirements for the degree of

#### DOCTOR OF PHILOSOPHY

In Mathematical Sciences

#### MICHIGAN TECHNOLOGICAL UNIVERSITY

2019

© 2019 Nattaporn Chuenjarern

### This dissertation has been approved in partial fulfillment of the requirements for the Degree of DOCTOR OF PHILOSOPHY in Mathematical Sciences.

Department of Mathematical Sciences

Dissertation Advisor:	Dr. Yang Yang
Committee Member:	Dr. Zhengfu Xu
Committee Member:	Dr. Cécile M. Piret
Committee Member::	Dr. Zhen Liu
Department Chair:	Dr. Mark S. Gockenbach

# Contents

st of	figures	ii				
st of	tables	ii				
efac	ei	x				
knov	wledgments	i				
ostra	$\mathbf{ct}$	ii				
Intr	$\operatorname{roduction}$	1				
1.1	Introduction to discontinuous Galerkin and local discontinuous					
	Galerkin methods	1				
1.2	Motivation	2				
1.3	Dissertation Outline	6				
2 Conservative local discontinuous Galerkin method for com-						
$\mathbf{pres}$	ssible miscible displacements in porous media	7				
2.1	Introduction	8				
2.2	Compressible miscible displacement problem	2				
	st of st of reface eknow ostra 1.1 1.2 1.3 Con pres 2.1 2.2	st of figures       vii         st of tables       vii         seface       ii         eknowledgments       x         ostract       xi         Introduction       xi         1.1       Introduction to discontinuous Galerkin and local discontinuous Galerkin methods         1.2       Motivation         1.3       Dissertation Outline         Conservative local discontinuous Galerkin method for compressible miscible displacements in porous media         2.1       Introduction         1.2       Compressible miscible displacement problem         1.3       Dissertation Outline         1.4       Introduction         1.5       Introduction         1.6       Introduction         1.7       Introduction         1.8       Displacement problem         1.9       Introduction         1.1       Introduction         1.2       Introduction         1.3       Introduction         1.4       Introduction         1.5       Introduction         1.6       Introduction         1.7       Introduction         1.8       Introduction         1       Introduction				

	2.3	Prelim	inaries	14
		2.3.1	Basic notations	14
		2.3.2	Norms	15
		2.3.3	LDG scheme and the main theorem	17
	2.4	The p	roof of the main theorem	20
		2.4.1	Projections and interpolation properties	20
		2.4.2	A priori error estimates	23
		2.4.3	Error equations	24
		2.4.4	The first energy inequality	24
		2.4.5	The second energy inequality	34
		2.4.6	The third energy inequality	35
		2.4.7	The fourth energy inequality	38
		2.4.8	Proof of Theorem 2.3.2	41
	2.5	Numer	rical example	42
	2.6	Conclu	uding remarks	47
	2.6	Appen	ndix: Proof of Lemma 2.4.5	48
9	U:~	h anda	r hound processing discentinuous Colorkin methods	
3	пigi	n-orae	r bound-preserving discontinuous Galerkin methods	
	for	compr	essible miscible displacements in porous media on	
	tria	ngular	meshes	53
	3.1	Introd	uction	55
	3.2	The D	G scheme	60
	3.3	Second	d-order bound-preserving scheme	64
	3.4	Bound	l-preserving technique for high-order scheme	74

		3.4.1 Flux limiter
		3.4.2 Slope limiter
		3.4.3 High-order time discretization
	3.5	Numerical experiments
	3.6	Concluding remarks
4	Fou	rier analysis of local discontinuous Galerkin methods for
	line	ar parabolic equations on overlapping meshes
	4.1	Introduction
	4.2	LDG method on overlapping meshes
		4.2.1 Overlapping meshes
		4.2.2 LDG scheme
	4.3	Error analysis
		4.3.1 The $P^1$ case $\ldots \ldots \ldots$
		4.3.2 The $P^2$ case $\ldots \ldots \ldots$
	4.4	Superconvergence
	4.5	Numerical experiments
	4.6	Conclusion
5	Cor	nclusion $\ldots$ $\ldots$ $\ldots$ $\ldots$ $\ldots$ $\ldots$ $\ldots$ $\ldots$ $.119$
R	efere	$nces \ldots 134$
$\mathbf{A}$	Cop	$ m by right\ documentations$
	A.1	Copyright documentation of Chapter 2

A.2	Copyright documentation of Chapter 3		•	•		•	•	•	•	•	•	•	•	137
A.3	Copyright documentation of Chapter 4			•		•				•	•	•		139

# List of Figures

2.1	Numerical approximations of c at $t = 0.1$ with $Nx = Ny = 40$ in	
	Example 2.5.3	45
2.2	Numerical approximations of $c$ at $t = 0.1$ with $Nx = Ny = 40$ in	
	Example 2.5.4	46
3.1	Two intersection points for the numerical flux in diffusion part on	
	the triangular mesh	70
3.2	Triangular mesh $(M = 10)$	80
3.3	Example 3.5.2: Numerical approximations of $c_1$ and $c_2$	84
3.4	Example 3.5.3: Concentrations of $c_1, c_2$ and $c_1 + c_2, \ldots, \ldots$	86
3.5	Example 3.5.4: Concentrations of $c_1, c_2$ and $c_3$ with limiters	88
3.6	Example 3.5.4: Concentrations of $c_1, c_2$ and $c_3$ without limiters	89
41	Overlapping meshes	97
1.1	o to hopping mobility i i i i i i i i i i i i i i i i i i	01

# List of Tables

2.1	The numerical results for $c$ with $\alpha = 1$	43
2.2	The numerical results for $c$ with $\alpha = 0.01$	44
2.3	The numerical results for $c$	47
3.1	Example 3.5.1: Accuracy test for $c_1$ and $c_2$ with and without	
	bound-preserving technique.	82
4.1	Example 4.5.1: $\alpha = 0, \ \alpha = 0.05, \ \alpha = 0.25, \ \alpha = 0.5.$	117
4.2	Example 4.5.1: Superconvergence with $\alpha = 0.25$ and $\alpha = 0.5$	118

## Preface

This dissertation contains published and submitted works completed by the author of this dissertation. The contributions of the author are detailed in the following paragraphs.

In the second chapter, the author has collaborated with Fan Yu<sup>1</sup>, Hui Guo<sup>2</sup>, and Yang Yang<sup>3</sup>. The main author's work is considering the problem in the one-dimensional case. The work has been published in the Journal of Scientific Computing. This work was supported by the National Natural Science Foundation of China Grants 11571367 and 11601536, and the Fundamental Research Funds for the Central Universities and Michigan Technological University Research Excellence Fund Scholarship and Creativity Grant 1605052.

In the third chapter, author has collaborated with Ziyao Xu<sup>4</sup> and Yang Yang. The main author's work is the hypercritical part and construction of the secondorder bound-preserving scheme. The work has been published in the Journal of Computational Physics. This work was supported by National Science Foundation DMS-1818467.

In the fourth chapter, the author has collaborated with Yang Yang. Most works in this chapter were performed by the author. First, the author studied

<sup>&</sup>lt;sup>1</sup>College of Science, China University of Petroleum, Qingdao 266580, China.

<sup>&</sup>lt;sup>2</sup>College of Science, China University of Petroleum, Qingdao 266580, China.

<sup>&</sup>lt;sup>3</sup>Department of Mathematical Sciences, Michigan Technological University, Houghton, MI 49931, USA.

<sup>&</sup>lt;sup>4</sup>Department of Mathematical Sciences, Michigan Technological University, Houghton, MI 49931.

C.W. Shu's work, repeated the results, and extended the idea to linear parabolic equations on overlapping meshes. The work has been completed as an article to submit to the Journal of Scientific Computing. This work was supported by National Science Foundation DMS-1818467.

## Acknowledgments

I would like to acknowledge my gratitude to my advisor, Dr. Yang Yang for his valuable advice and support. Dr.Yang provided many opportunities, inspirations, encouragements, patience, and guidance throughout the research and my life as a graduate student.

I would like to thank the committees, Dr. Zhengfu Xu, Dr. Cecile Piret and Dr. Zhen Liu for their advice that makes this work more complete.

Also, I would like to thank Dr. Mark Gockenbach, the Department chair, who gave me the opportunity to teach. Thank Ann Humes and Elizabeth Reed who are the coordinator and mentor for my teaching.

Moreover, I would like to thank Jeanne Meyers, Margaret Perander and Kimberly Puuri who also made my life move smoothly as a graduate student in the department.

I really appreciate the help from my collaborator, Ziyao Xu, for answering all my questions and having helpful discussions.

Moreover, I would like to thank my family who leads me to have the ultimate goals and friends who have encouraged me and joined with me in both academic and non-academic activities.

In addition, I would like to thank the Development and Promotion of Science and Technology talents project (DPST), Institute for the Promotion of Teaching Science and Technology (IPST), Ministry of Education, Thailand for supporting my scholarship.

### Abstract

This dissertation contains research in discontinuous Galerkin (DG) methods applying to convection-diffusion equations. It contains both theoretical analysis and applications. Initially, we develop a conservative local discontinuous Galerkin (LDG) method for the coupled system of compressible miscible displacement problem in two space dimensions. The main difficulty is how to deal with the discontinuity of approximations of velocity,  $\mathbf{u}$ , in the convection term across the cell interfaces. To overcome the problems, we apply the idea of LDG with IMEX time marching using the diffusion term to control the convection term. Optimal error estimates in  $L^{\infty}(0,T;L^2)$  norm for the solution and the auxiliary variables will be derived. Then, a high-order bound-preserving (BP) discontinuous Galerkin (DG) methods for the coupled system of compressible miscible displacements on triangular meshes will be developed. There are three main difficulties to make the concentration of each component between 0 and 1. Firstly, the concentration of each component did not satisfy a maximumprinciple. Secondly, the first-order numerical flux was difficult to construct. Thirdly, the classical slope limiter could not be applied to the concentration of each component. To conquer these three obstacles, we first construct special techniques to preserve two bounds without using the maximum-principlepreserving technique. The time derivative of the pressure was treated as a source of the concentration equation. Next, we apply the flux limiter to obtain highorder accuracy using the second-order flux as the lower order one instead of using the first-order flux. Finally, L<sup>2</sup>-projection of the porosity and constructed special limiters that are suitable for multi-component fluid mixtures were used. Lastly, a new LDG method for convection-diffusion equations on overlapping mesh introduced in [28] showed that the convergence rates cannot be improved if the dual mesh is constructed by using the midpoint of the primitive mesh. They provided several ways to gain optimal convergence rates but the reason for accuracy degeneration is still unclear. We will use Fourier analysis to analyze the scheme for linear parabolic equations with periodic boundary conditions in one space dimension. To investigate the reason for the accuracy degeneration, we explicitly write out the error between the numerical and exact solutions. Moreover, some superconvergence points that may depend on the perturbation constant in the construction of the dual mesh were also found out.

# Chapter 1

# Introduction

# 1.1 Introduction to discontinuous Galerkin and local discontinuous Galerkin methods

The discontinuous Galerkin (DG) methods are a class of finite element methods with completely discontinuous piecewise polynomials as the numerical approximations. The DG method was first introduced in the framework of neutron linear transportation by Reed and Hill [51] in 1973. Subsequently, the Runge-Kutta discontinuous Galerkin (RKDG) methods were proposed for hyperbolic conservation laws in a series of papers [16, 17, 18, 19]. The DG method gained even greater popularity recently for good stability, high order accuracy, and flexibility on h-p adaptivity and on complex geometry. But, it is difficult to apply the DG method directly to the equations with higher order derivatives for example, a convection-diffusion equation. One possible way to form a stable and convergent DG method is to rewrite the equations with higher order derivatives into a first order system, then apply the DG method to the system called local discontinuous Galerkin (LDG) methods . As an extension of DG schemes for hyperbolic conservation laws, the LDG methods share the advantages of the DG methods. Besides, a key advantage of this scheme is the local solvability, i.e. the auxiliary variables approximating the gradient of the solution can be locally eliminated. The first LDG was introduced by Cockburn and Shu in [20] for solving the convection-diffusion equations. Their idea was motivated by Bassi and Rebay [2], where the compressible Navier-Stokes equations were successfully solved. For simplicity, we consider the following linear parabolic equations in one space dimension:

$$u_t - u_{xx} = 0, \quad x \in [0, 2\pi], \quad t > 0,$$
  
 $u(x, 0) = u_0(x), \quad x \in [0, 2\pi],$  (1.1.1)

In [20], the authors introduced an auxiliary variable p to represent the derivative of the primary variable u and thus rewrite (1.1.1) into the following system of first order equations

$$u_t - p_x = 0,$$
  
 $p - u_x = 0.$ 
(1.1.2)

Then one can solve u and p on the same mesh [20].

### 1.2 Motivation

Recently, DG methods have been popular to solve compressible miscible displacements in porous media [21, 22, 71, 72, 37, 73, 77]. Also, there were significant works discussing the DG methods for incompressible miscible displacements, see e.g. [1, 38, 44, 52, 55, 56, 63] and for general porous media flow, see e.g. [3, 30, 29, 57] and the references therein. However, no previous works above focused on the bound-preserving techniques. In many numerical simulations, the approximations of concentration can be placed out of the interval [0, 1]. Especially for problems with large gradients will lead to ill-posedness of the problem, and the numerical approximations will blow up. Therefore, we extend the ideas of [36] to develop high-order bound-preserving (BP) discontinuous Galerkin (DG) methods for the coupled system of multi-component compressible miscible displacements on triangular meshes. The goal was to make the concentration of each component between 0 and 1. There were three main difficulties. Firstly, the concentration of each component did not satisfy a maximum-principle. Secondly, the first-order numerical flux was difficult to construct. Thirdly, the classical slope limiter could not be applied to the concentration of each component. To overcome these three obstacles, special techniques were first constructed to preserve two bounds without using the maximum-principle-preserving technique. The time derivative of the pressure was treated as a source of the concentration equation. Next, the flux limiter was applied to obtain high-order accuracy using the second-order flux as the lower order one instead of using the first-order flux. Finally,  $L^2$ -projection of the porosity and constructed special limiters that are suitable for multi-component fluid mixtures were used.

For the LDG method, it was applied to the one-dimensional coupled system of compressible miscible displacement problem in [37]. But the method in [38] is not conservative. Later in [35], LDG was applied to solve incompressible miscible displacements in porous media. Therefore, we continue to develop a conservative local discontinuous Galerkin (LDG) method for the two-dimensional coupled system of the compressible miscible displacement problem. The main difficulty was the discontinuity of approximations of velocity,  $\mathbf{u}$ , in the convection term across the cell interfaces. Also, if the convection and diffusion terms were considered separately, it would be difficult to obtain error estimates. Due to this difficulty, the traditional error analysis could not be applied directly. To overcome the problems, the idea of LDG with IMEX time marching using the diffusion term to control the convection term was applied. Then, the energy inequalities were rewritten into four parts to obtain optimal error estimates for concentration c,  $-\nabla c$  and velocity  $\mathbf{u}$ .

The LDG method is one of the most important numerical methods for convection diffusion equations. However, for some special convection-diffusion systems, such as chemotaxis model [43, 49] and miscible displacements in porous media [24, 25], the LDG methods are not easy to construct and analyze. In each of the two models, the convection term is the product of one of the primary variables and the derivative of the other primary variable. Because of this obstacle, the upwind fluxes cannot be applied directly. Within the DG framework, there are three main different ways to bridge this gap.

 Combine the convection terms and diffusion terms together and obtain the optimal error estimates. This approach was proposed in [77, 35, 46] However, to make the numerical solutions to be physically relevant, we have to add a very large penalty which depends on the numerical approximations of the derivatives of the primary variables [46, 36, 13].

- 2. Apply the flux-free numerical methods such as the Central DG (CDG) methods [47]. However, for CDG methods, we have to solve each equation in (1.1.2) on both the primary and dual meshes, which may double the computational cost.
- 3. Apply the Staggered DG (SDG) methods [14]. However, the method requires some continuity of the numerical approximations, and hence it is not easy to apply limiters to the numerical solutions.

Recently, a new LDG method was introduced in [28]. The main idea of this method is to compute the primary variable u and auxiliary variable  $p = u_x$  on different meshes. However, the accuracy may not be optimal if odd-order polynomials were applied with the dual mesh constructed by using the midpoint of the primitive mesh. To investigate the reason for accuracy degeneration, Fourier analysis was applied to linear parabolic equations in one space dimension subject to periodic boundary conditions. Then the LDG scheme can be rewritten into an equivalent finite difference scheme, and the numerical solution obtained by finding the eigenvalues and eigenvectors of the amplification matrix. The reason for the accuracy degeneration was discovered by explicitly expressing the error between the numerical and exact solutions. This analysis showed that when the midpoint was used to construct the dual mesh, the nonphysical eigenvalue of the amplification matrix did not decay during mesh refinement. Thus, the scheme generated a spurious wave that caused the accuracy of the scheme to degenerate. Moreover, with the quantitative error estimate, some superconvergence points that may depend on the perturbation constant in the construction of the dual mesh were also found.

### **1.3** Dissertation Outline

The accomplished work will be in three main chapters (Chapter 2 to Chapter 4). First, Chapter 2 describes the work on conservative local discontinuous Galerkin method for compressible miscible displacements in porous media. Second, Chapter 3 presents the research on high-order bound-preserving discontinuous Galerkin methods for compressible miscible displacements in porous media on triangular meshes. Last, the study on Fourier analysis of local discontinuous Galerkin methods for linear parabolic equations on overlapping meshes will be demonstrated in Chapter 4. We will end in Chapter 5 with conclusion.

# Chapter 2

# Conservative local discontinuous Galerkin method for compressible miscible displacements in porous media<sup>1</sup>

#### Abstract

In [H. Guo, Q. Zhang, J. Wang, Applied Mathematics and Computation, 259 (2015), 88-105], a nonconservative local discontinuous Galerkin (LDG) method for both flow and transport equations was introduced for the one-dimensional coupled system of compressible miscible displacement problem. In this paper, we

<sup>&</sup>lt;sup>1</sup>This chapter has been published as an article in Journal of Scientific Computing. Citation: F. Yu, H. Guo, N. Chuenjarern, Y. Yang, J Sci Comput (2017) 73: 1249. https://doi.org/10.1007/s10915-017-0571-z

will continue our effort and develop a conservative LDG method for the problem in two space dimensions. Optimal error estimates in  $L^{\infty}(0,T;L^2)$  norm for not only the solution itself but also the auxiliary variables will be derived. The main difficulty is how to treat the inter-element discontinuities of two independent solution variables (one from the flow equation and the other from the transport equation) at cell interfaces. Numerical experiments will be given to confirm the accuracy and efficiency of the scheme.

**Keywords:** local discontinuous Galerkin method, error estimate, compressible miscible displacement

### 2.1 Introduction

Numerical modeling of miscible displacements in porous media is important and interesting in oil recovery and environmental pollution problem. The miscible displacement problem is described by a coupled system of nonlinear partial differential equations. The need for accurate solutions to the coupled equations challenges numerical analysts to design new methods.

The compressible miscible displacements have been studied intensively in the literature. In [24, 25], Douglas and Roberts presented the mixed finite element method for miscible displacement problem. A variety of numerical techniques have been introduced to obtain better approximations, such as the modified method of characteristic finite element method (MMOC) [26, 31, 79], characteristic finite element method [78], high-order Godunov scheme [4], streamline dif-

fusion method [42], and Mass-conservative characteristic finite element method [45]. Recently, discontinuous Galerkin (DG) for miscible displacement has been investigated by numerical experiments and was reported to exhibit good numerical performance [1, 52]. In [55, 56, 22], primal semi-discrete discontinuous Galerkin methods with interior penalty are proposed to solve the coupled system of flow and reactive transport in porous media.

The DG method gained even greater popularity recently for good stability, high order accuracy, and flexibility on h-p adaptivity and on complex geometry. But, it is difficult to apply the DG method directly to the equations with higher order derivatives. The idea of the local discontinuous Galerkin (LDG) method is to rewrite the equations with higher order derivatives into a first order system, then apply the DG method to the system. As an extension of DG schemes for hyperbolic conservation laws, the LDG methods share the advantages of the DG methods. Besides, a key advantage of this scheme is the local solvability, i.e. the auxiliary variables approximating the gradient of the solution can be locally eliminated. The first LDG method was introduced by Cockburn and Shu in [20] for solving nonlinear convection diffusion equations containing second order spatial derivatives. Their work was motivated by the successful numerical experiments of Bassi and Rebay [2] for the compressible Navier-Stokes equations. The methods were further developed in [66, 67, 69] for solving many nonlinear wave equations with higher order derivatives.

In our previous work [37], we have used the LDG method to the one-dimensional coupled system of compressible miscible displacement problem. But the method in [38] is not conservative. Recently, we [35] applied the LDG methods to solve incompressible miscible displacements in porous media. In this paper we continue our works to develop a conservative LDG method for compressible miscible displacements in two space dimensions. The main difficulty is how to treat the inter-element discontinuities of two independent solution variables (one from the flow equation and the other from the transport equation) at cell interfaces. More precisely, in this problem, the approximations of  $\mathbf{u}$  in the convection term in (2.2.1) is discontinuous across the cell interfaces and it is difficult to obtain error estimates if we analyze the convection and diffusion terms separately. To explain this point, let us consider the following hyperbolic equation

$$u_t + (a(x)u)_x = 0,$$

where a(x) is discontinuous at  $x = x_0$ . In [32, 40], the authors studied such a problem and defined

$$Q = \frac{a(x_0+b) - a(x_0)}{b}.$$

If Q is bounded from below for all b, then the solution exists, but may not be unique. If Q is bounded from above for all b, we can guarantee the uniqueness, but the solution may not exist. Recently, Wang et al. [60, 61] obtained optimal error estimates of the LDG methods with IMEX time marching for linear and nonlinear convection-diffusion problems. The key idea is to explore an important relationship between the gradient and interface jump of the numerical solution polynomial with the numerical approximation of auxiliary variable for the gradient in the LDG methods, which is stated in Lemma 2.4.4. Moreover, the systems are coupled together. Therefore, we will derive four energy inequalities to obtain optimal error estimates in  $L^{\infty}(0,T;L^2)$  for concentration c, in  $L^2(0,T;L^2)$  for  $\mathbf{s} = -\nabla c$  and  $L^{\infty}(0,T;L^2)$  for velocity  $\mathbf{u}$ . Here we should mention the difference between our LDG method and the DG method in [22], where the interior penalty discontinuous Galerkin (IPDG) method was introduced and optimal error estimates in  $L^2(0,T;H^1)$  norm for concentration c were given. In our proof, induction hypothesis is used as a tool, instead of the cut-off operator proposed in [56]. Therefore, it is not necessary to choose the sufficiently large positive constant M, and the possibility of infinite times of loops for extreme cases can be avoided.

The paper is organized as follows. In Section 2.2, we demonstrate the governing equations of the compressible miscible displacements in porous media. In Section 2.3, we present some preliminaries, including the basic notations and norms to be used throughout the paper, the LDG spatial discretization and the error equations. Section 2.4 is the main body of the paper where we present the projections and some essential properties of the finite element spaces, error equations and the details of the optimal error estimates for compressible miscible displacement problem. Then numerical results are given to demonstrate the accuracy and capability of the method in Section 2.5. We will end in Section 2.6 with some concluding remarks.

# 2.2 Compressible miscible displacement problem

In this section, we demonstrate the governing equations of the compressible miscible displacements in porous. Detailed discussion on physical theories can be found in [23]. Let  $\Omega$  be a rectangular domain. The classical equations governing the compressible miscible displacement in porous media in two space dimensions are as follows:

$$d(c)\frac{\partial p}{\partial t} + \nabla \cdot \mathbf{u} = q, \qquad (x, y) \in \Omega, 0 < t \le T,$$
  

$$\mathbf{u} = \frac{-\kappa(x, y)}{\mu(c)} \nabla p, \qquad (x, y) \in \Omega, 0 < t \le T,$$
  

$$\phi\frac{\partial c}{\partial t} + b(c)\frac{\partial p}{\partial t} + \mathbf{u} \cdot \nabla c = \nabla \cdot (\mathbf{D}\nabla c) + (\tilde{c} - c)q, \quad (x, y) \in \Omega, 0 < t \le T,$$
  
(2.2.1)

where the dependent variables p,  $\mathbf{u}$  and c are the pressure in the fluid mixture, the Darcy velocity of the mixture (volume flowing across a unit across-section per unit time), and the concentration of interested species measured in amount of species per unit volume of the fluid mixture, respectively.  $\phi$  and  $\kappa$  are the porosity and the permeability of the rock, respectively.  $\mu$  is the concentrationdependent viscosity. q is the external volumetric flow rate, and  $\tilde{c}$  is the concentration of the fluid in the external flow.  $\tilde{c}$  must be specified at points at which injection (q > 0) takes place, and is assumed to be equal to c at production points (q < 0). We shall also consider only molecular diffusion, so that  $\mathbf{D} = \phi(x, y)d_mI$ with I being the identity matrix. In this paper the tensor matrix  $\mathbf{D}$  is assumed to be positive definite. Moreover, the pressure is uniquely determined up to a constant, thus we assume  $\int_{\Omega} p dx dy = 0$  at t = 0. For simplicity, we confine ourselves to a two component displacement problem. The numerical method can be applied to the multi-component model. The coefficients can be stated as follows:

$$c = c_1 = 1 - c_2,$$
  

$$a(c) = a(x, y, c) = \frac{\kappa(x, y)}{\mu(c)},$$
  

$$b(c) = b(x, y, c) = \phi(x, y)c_1\{m_1 - \sum_{j=1}^2 m_j c_j\},$$
  

$$d(c) = d(x, y, c) = \phi(x, y)\sum_{j=1}^2 m_j c_j,$$

with  $c_i$  being the concentration of i th component of the fluid mixture, and  $m_i$  being the "constant compressibility" factor. In this problem, the initial concentration are pressure are given as

$$c(x, y, 0) = c_0(x, y), \quad p(x, y, 0) = p_0(x, y), \quad (x, y) \in \Omega.$$

Finally, we make the following hypotheses (H) for (2.2.1).

- 1.  $0 < \kappa_* \le \kappa(x, y) \le \kappa^*, 0 < \mu_* \le \mu(c) \le \mu^*, 0 < \phi_* \le \phi(x, y) \le \phi^*,$  $0 < d_* \le d(c) \le d^*, |q| \le C, |b(c)| \le C, |\mu'(c)| \le C \text{ and } |d'(c)| \le C.$
- d(c), μ'(c) and d'(c) are uniformly Lipschtiz continuous with respect to c, respectively.
- 3. **D** is uniformly Lipschtiz continuous, and for any  $\mathbf{v}, \mathbf{w} \in R^2$  there exist two positive constants  $D_*, D^*$  such that  $\mathbf{v}^T \mathbf{D} \mathbf{v} \ge D_* \mathbf{v}^T \mathbf{v} = D_* \|\mathbf{v}\|^2$  and  $\mathbf{v}^T \mathbf{D} \mathbf{w} \le D^* \|\mathbf{v}\| \|\mathbf{w}\|$ .

4.  $\mathbf{u}, \mathbf{u}_t, c, \nabla c, c_t, p_t$  and  $p_{tt}$  are uniformly bounded in  $\mathbb{R}^2$ .

### 2.3 Preliminaries

In this section, we will demonstrate some preliminary results that will be used through out the paper.

#### 2.3.1 Basic notations

In this section, we present the notations. Let  $0 = x_{\frac{1}{2}} < \cdots < x_{N_x + \frac{1}{2}} = 1$ and  $0 = y_{\frac{1}{2}} < \cdots < y_{N_y + \frac{1}{2}} = 1$  be the grid points in the x and y directions, respectively. Define  $I_i = (x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}})$  and  $J_j = (y_{j-\frac{1}{2}}, y_{j+\frac{1}{2}})$ . Let  $K = I_i \times J_j$ ,  $i = 1, \cdots, N_x, \ j = 1, \cdots, N_y$ , be a partition of  $\Omega$  and denote  $\Omega_h = \{K\}$ . The mesh sizes in the x and y directions are given as  $\Delta x_i = x_{i+\frac{1}{2}} - x_{i-\frac{1}{2}}$  and  $\Delta y_j = y_{j+\frac{1}{2}} - y_{j-\frac{1}{2}}$ , respectively and  $h = \max(\max_i \Delta x_i, \max_j \Delta y_j)$ . Moreover, we assume the partition is quasi-uniform. The finite element space is chosen as

$$W_h^k = \{ z : z |_K \in Q^k(K), \forall K \in \Omega_h \},\$$

where  $Q^k(K)$  denotes the space of tensor product polynomials of degrees at most k in K. Note that functions in  $W_h^k$  are discontinuous across element interfaces. This is one of the main differences between the DG method and traditional finite element methods. We choose  $\boldsymbol{\beta} = (1, 1)^T$  to be a fixes vector that is not parallel to any normals of the element interfaces. We denote  $\Gamma_h$  be the set of all element interfaces and  $\Gamma_0 = \Gamma_h \setminus \partial \Omega$ . Let  $e \in \Gamma_0$  be an interior edge shared by elements  $K_\ell$  and  $K_r$ , where  $\boldsymbol{\beta} \cdot \mathbf{n}_\ell > 0$ , and  $\boldsymbol{\beta} \cdot \mathbf{n}_r < 0$ , respectively, with  $\mathbf{n}_\ell$  and  $\mathbf{n}_r$  being the outward normal of  $K_{\ell}$  and  $K_r$ , respectively. For any  $z \in W_h^k$ , we define  $z^- = z|_{\partial K_{\ell}}$  and  $z^+ = z|_{\partial K_r}$ , respectively. The jump is given as  $[z] = z^+ - z^-$ . Moreover, for  $\mathbf{s} \in \mathbf{W}_h^k = W_h^k \times W_h^k$ , we define  $\mathbf{s}^+$  and  $\mathbf{s}^-$  and  $[\mathbf{s}]$  analogously. We also define  $\partial \Omega_- = \{e \in \partial \Omega | \boldsymbol{\beta} \cdot \mathbf{n} < 0\}$ , where  $\mathbf{n}$  is the outer normal of e, and  $\partial \Omega_+ = \partial \Omega \setminus \partial \Omega_-$ . For any  $e \in \partial \Omega_-$ , there exists  $K \in \Omega_h$  such that  $e \in \partial K$ , we define  $z^+|_e = z|_{\partial K}$ , and define  $z^-$  on  $\partial \Omega_+$  analogously. For simplicity, given  $e = \{x_{\frac{1}{2}}\} \times J_j \in \partial \Omega_-$  and  $\tilde{e} = \{x_{N_x + \frac{1}{2}}\} \times J_j \in \partial \Omega_+$ , by periodic boundary condition, we define

$$z^{-}|_{e} = z^{-}|_{\tilde{e}}, \text{ and } z^{+}|_{\tilde{e}} = z^{+}|_{e}$$

Similarly, given  $e = I_i \times \{y_{\frac{1}{2}}\} \in \partial \Omega_-$  and  $\tilde{e} = I_i \times \{y_{N_y + \frac{1}{2}}\} \in \partial \Omega_+$ , we define

$$z^{-}|_{e} = z^{-}|_{\tilde{e}}, \text{ and } z^{+}|_{\tilde{e}} = z^{+}|_{e}.$$

Throughout this paper, the symbol C is used as a generic constant which may appear differently at different occurrences.

#### 2.3.2 Norms

In this subsection, we define several norms that will be used throughout the paper.

Denote  $||u||_{0,K}$  to be the standard  $L^2$  norm of u in cell K. For any natural number  $\ell$ , we consider the norm of the Sobolev space  $H^{\ell}(K)$ , defined by

$$\|u\|_{\ell,K} = \left\{ \sum_{0 \le \alpha + \beta \le \ell} \left\| \frac{\partial^{\alpha + \beta} u}{\partial x^{\alpha} \partial y^{\beta}} \right\|_{0,K}^2 \right\}^{\frac{1}{2}}.$$

Moreover, we define the norms on the whole computational domain as

$$||u||_{\ell} = \left(\sum_{K \in \Omega_h} ||u||^2_{\ell,K}\right)^{\frac{1}{2}}.$$

For convenience, if we consider the standard  $L^2$  norm, then the corresponding subscript will be omitted.

Let  $\Gamma_K$  be the edges of K, and we define

$$\|u\|_{\Gamma_K}^2 = \int_{\partial K} u^2 ds.$$

We also define

$$||u||_{\Gamma_h}^2 = \sum_{K \in \Omega_h} ||u||_{\Gamma_K}^2.$$

Moreover, we define the standard  $L^{\infty}$  norm of u in K as  $||u||_{\infty,K}$ , and define the  $L^{\infty}$  norm on the whole computational domain as

$$||u||_{\infty} = \max_{K \in \Omega_h} ||u||_{\infty,K}.$$

Finally, we define similar norms for vector  $\mathbf{u} = (u_1, u_2)^T$  as

$$\|\mathbf{u}\|_{\ell,K}^{2} = \|u_{1}\|_{\ell,K}^{2} + \|u_{2}\|_{\ell,K}^{2},$$
$$\|\mathbf{u}\|_{\Gamma_{K}}^{2} = \|u_{1}\|_{\Gamma_{K}}^{2} + \|u_{2}\|_{\Gamma_{K}}^{2},$$
$$\|\mathbf{u}\|_{\infty,K} = \max\{\|u_{1}\|_{\infty,K}, \|u_{2}\|_{\infty,K}\}.$$

Similarly, the norms on the whole computational domain are given as

$$\|\mathbf{u}\|_{\ell}^{2} = \sum_{K \in \Omega_{h}} \|\mathbf{u}\|_{\ell}^{2}, \quad \|\mathbf{u}\|_{\Gamma_{h}}^{2} = \sum_{K \in \Omega_{h}} \|\mathbf{u}\|_{\Gamma_{K}}^{2}, \quad \|\mathbf{u}\|_{\infty} = \max_{K \in \Omega_{h}} \|\mathbf{u}\|_{\infty,K}.$$

#### 2.3.3 LDG scheme and the main theorem

To construct the LDG scheme, we introduce some auxiliary variables to approximate the derivatives of the solution which further yields a first order system:

$$\phi \frac{\partial c}{\partial t} + B(c) \frac{\partial p}{\partial t} + \nabla \cdot (\mathbf{u}c) + \nabla \cdot \mathbf{z} = \tilde{c}q, \qquad (2.3.2)$$

$$\mathbf{s} = -\nabla c, \tag{2.3.3}$$

$$\mathbf{z} = \mathbf{D}\mathbf{s},\tag{2.3.4}$$

$$A(c)\mathbf{u} + \nabla p = 0, \tag{2.3.5}$$

$$d(c)\frac{\partial p}{\partial t} + \nabla \cdot \mathbf{u} = q, \qquad (2.3.6)$$

where  $A(c) = \mu(c)\kappa(x,y)^{-1}$ ,  $B(c) = cd(c) + b(c) = c\phi(x,y)m_1$ . We multiply (2.3.2)-(2.3.6) by test functions  $v, \zeta \in W_h^k, \theta, \mathbf{w}, \psi \in \mathbf{W}_h^k$ , respectively. Formally integrate by parts in K to get

$$\begin{aligned} (\phi c_t, v)_K + (B(c)p_t, v)_K &= (\mathbf{u}c + \mathbf{z}, \nabla v)_K - \langle (\mathbf{u}c + \mathbf{z}) \cdot \boldsymbol{\nu}_K, v \rangle_{\partial K} + (\tilde{c}q, v)_K, \\ (\mathbf{s}, \mathbf{w})_K &= (c, \nabla \cdot \mathbf{w})_K - \langle c, \mathbf{w} \cdot \boldsymbol{\nu}_K \rangle_{\partial K}, \\ (\mathbf{z}, \boldsymbol{\psi})_K &= (\mathbf{D}\mathbf{s}, \boldsymbol{\psi})_K, \\ (A(c)\mathbf{u}, \boldsymbol{\theta})_K &= (p, \nabla \cdot \boldsymbol{\theta})_K - \langle p, \boldsymbol{\theta} \cdot \boldsymbol{\nu}_K \rangle_{\partial K}, \\ (d(c)p_t, \zeta)_K &= (\mathbf{u}, \nabla \zeta)_K - \langle \mathbf{u} \cdot \boldsymbol{\nu}_K, \zeta \rangle_{\partial K} + (q, \zeta)_K, \end{aligned}$$

where  $(u, v)_K = \int_K uv dx dy$ ,  $(\mathbf{u}, \mathbf{v})_K = \int_K \mathbf{u} \cdot \mathbf{v} dx dy$ ,  $\langle u, v \rangle_{\partial K} = \int_{\partial K} uv ds$  and  $\boldsymbol{\nu}_K$ is the outer normal of K. Replacing the exact solutions  $c, p, \mathbf{s}, \mathbf{z}, \mathbf{u}$  in the above equations by their numerical approximations  $c_h, p_h \in W_h^k$  and  $\mathbf{s}_h, \mathbf{z}_h, \mathbf{u}_h \in \mathbf{W}_h^k$ , respectively and using numerical fluxes at the cell interfaces to obtain the LDG scheme:

$$(\phi c_{ht}, v)_K + (B(c_h)p_{ht}, v)_K = \mathcal{L}_K^c(\mathbf{u}_h, c_h, v) + \mathcal{L}_K^d(\mathbf{z}_h, v) + (\tilde{c}_h q, v)_K, (2.3.7)$$

$$(\mathbf{s}_h, \mathbf{w})_K = \mathcal{D}_K(c_h, \mathbf{w}), \qquad (2.3.8)$$

$$(\mathbf{z}_h, \boldsymbol{\psi})_K = (\mathbf{D}\mathbf{s}_h, \boldsymbol{\psi})_K, \qquad (2.3.9)$$

$$(A(c_h)\mathbf{u}_h,\boldsymbol{\theta})_K = \mathcal{D}_K(p_h,\boldsymbol{\theta}), \qquad (2.3.10)$$

$$(d(c_h)p_{ht},\zeta)_K = \mathcal{L}_K^d(\mathbf{u}_h,\zeta) + (q,\zeta)_K, \qquad (2.3.11)$$

where

$$\mathcal{L}_{K}^{c}(\mathbf{s}, c, v) = (\mathbf{s}c, \nabla v)_{K} - \langle \widehat{\mathbf{s}c} \cdot \boldsymbol{\nu}_{K}, v \rangle_{\partial K}, \qquad (2.3.12)$$

$$\mathcal{L}_{K}^{d}(\mathbf{s}, v) = (\mathbf{s}, \nabla v)_{K} - \langle \widehat{\mathbf{s}_{h}} \cdot \boldsymbol{\nu}_{K}, v \rangle_{\partial K}, \qquad (2.3.13)$$

$$\mathcal{D}_K(c, \mathbf{w}) = (c, \nabla \cdot \mathbf{w})_K - \langle \hat{c}, \mathbf{w} \cdot \boldsymbol{\nu}_K \rangle_{\partial K}.$$
 (2.3.14)

We use alternating fluxes for the diffusion term and take

$$\widehat{\mathbf{z}}_h = \mathbf{z}_h^-, \quad \widehat{c}_h = c_h^+, \quad \widehat{\mathbf{u}}_h = \mathbf{u}_h^-, \quad \widehat{p}_h = p_h^+.$$

For the convection term, we consider Lax-Friedrichs flux

$$\widehat{\mathbf{u}_h c_h} = \frac{1}{2} (\mathbf{u}_h^+ c_h^+ + \mathbf{u}_h^- c_h^- - \alpha \boldsymbol{\nu}_e (c_h^+ - c_h^-)),$$

where  $\alpha > 0$  can be chosen as any constant and  $\nu_e$  is the unit normal of the  $e \in \Gamma_0$  such that  $\beta \cdot \nu_e > 0$ . Moreover, we define

$$(u,v) = \sum_{K \in \Omega_h} (u,v)_K, \qquad (\mathbf{u},\mathbf{v}) = \sum_{K \in \Omega_h} (\mathbf{u},\mathbf{v})_K,$$

and

$$\mathcal{L}^{c}(\mathbf{s}, c, v) = \sum_{K \in \Omega_{h}} \mathcal{L}^{c}_{K}(\mathbf{s}, c, v),$$
$$\mathcal{L}^{d}(\mathbf{s}, v) = \sum_{K \in \Omega_{h}} \mathcal{L}^{d}_{K}(\mathbf{s}, v),,$$
$$\mathcal{D}(c, \mathbf{w}) = \sum_{K \in \Omega_{h}} \mathcal{D}_{K}(c, \mathbf{w}).$$

It is easy to check the following identities by integration by parts on each cell

Lemma 2.3.1. For any functions v and w,

$$\mathcal{L}^{d}(\mathbf{w}, v) + \mathcal{D}(v, \mathbf{w}) = 0.$$
(2.3.15)

Now we state the main theorem.

**Theorem 2.3.2.** Let  $c \in H^{k+3}$ ,  $s \in (H^{k+2})^2$ ,  $u \in (H^{k+1})^2$  be the exact solutions of the problem (2.3.2)-(2.3.6), and let  $\mathbf{u}_h, p_h, c_h, \mathbf{s}_h, \mathbf{z}_h$  be the numerical solutions of the semi-discrete LDG scheme (2.3.7)-(2.3.11) with initial discretization given as (2.4.4). If the finite element space is the piecewise tensor product polynomials of degree  $k \ge 1$  and h is sufficiently small, then we have the error estimate

$$\begin{aligned} \|c - c_h\|_{L^{\infty}(0,T;L^2)} + \|\mathbf{s} - \mathbf{s}_h\|_{L^{\infty}(0,T;L^2)} \\ + \|\mathbf{u} - \mathbf{u}_h\|_{L^{\infty}(0,T;L^2)} + \|p - p_h\|_{L^{\infty}(0,T;L^2)} + \|(p - p_h)_t\|_{L^{\infty}(0,T;L^2)} \\ + \|(c - c_h)_t\|_{L^2(0,T;L^2)} + \|(\mathbf{u} - \mathbf{u}_h)_t\|_{L^2(0,T;L^2)} \le Ch^{k+1}, \end{aligned}$$
(2.3.16)

where the constant C is independent of h.
# 2.4 The proof of the main theorem

In this section, we proceed to the proof of Theorem 2.3.2. We first introduce several projections and present some auxiliary results. Subsequently, we make an a priori error estimate which provides the boundedness of the numerical approximations. Then we construct the error equations which further yield five main energy inequalities and complete the proof of (2.3.16). Finally, we verify the a priori error estimate at the end of this section.

#### 2.4.1 **Projections and interpolation properties**

In this section, we will demonstrate the projections and several useful lemmas. Let us start with the classical inverse properties [15].

**Lemma 2.4.1.** Assuming  $u \in W_h^k$ , there exists a positive constant C independent of h and u such that

$$h\|u\|_{\infty,K} + h^{1/2}\|u\|_{\Gamma_K} \le C\|u\|_{K}$$

We will use several special projections in this paper. Firstly, we define  $P^+$ into  $W_h^k$  which is, for each cell K

$$(P^+u - u, v)_K = 0, \forall v \in Q^{k-1}(K),$$
$$\int_{J_j} (P^+u - u)(x_{i-\frac{1}{2}}, y)v(y)dy = 0, \forall v \in P^{k-1}(J_j),$$
$$(P^+u - u)(x_{i-\frac{1}{2}}, y_{j-\frac{1}{2}}) = 0$$
$$\int_{I_i} (P^+u - u)(x, y_{j-\frac{1}{2}})v(x)dx = 0, \forall v \in P^{k-1}(I_i),$$

where  $P^k$  denotes the polynomials of degree k. Moreover, we also define  $\Pi_x^-$  and  $\Pi_y^-$  into  $W_h^k$  which are, for each cell K,

$$\begin{split} (\Pi_x^- u - u, v_x)_K &= 0, \forall v \in Q^k(K), \\ \int_{J_j} (\Pi_x^- u - u)(x_{i+\frac{1}{2}}, y)v(y)dy &= 0, \forall v \in P^k(J_j), \\ (\Pi_y^- u - u, v_y)_K &= 0, \forall v \in Q^k(K), \\ \int_{I_i} (\Pi_y^- u - u)(x, y_{j+\frac{1}{2}})v(x)dx &= 0, \forall v \in P^k(I_i), \end{split}$$

as well as a two-dimensional projection  $\Pi^- = \Pi^-_x \otimes \Pi^-_y$ . Finally, we also use the  $L^2$ -projection  $P_k$  into  $W^k_h$  which is, for each cell K

$$(P_k u - u, v)_K = 0, \forall v \in Q^k(K),$$
 (2.4.1)

and its two dimensional version  $\mathbf{P}_k = P_k \otimes P_k$ . For the special projections mentioned above, we give the following lemma by the standard approximation theory [15].

**Lemma 2.4.2.** Suppose  $w \in H^{k+1}(\Omega)$ , then for any project  $P_h$ , which is either  $P^+$ ,  $\Pi_x^-$ ,  $\Pi_y^-$  or  $P_k$ , we have

$$||w - P_h w|| + h^{1/2} ||w - P_h w||_{\Gamma_h} \le C h^{k+1}.$$

Moreover, the projection  $P^+$  on the Cartesian meshes has the following superconvergence property [6].

**Lemma 2.4.3.** Suppose  $w \in H^{k+2}(\Omega)$ , then for any  $\rho \in \mathbf{W}_h$  we have

$$|\mathcal{D}(w - P^+ w, \boldsymbol{\rho})| \le Ch^{k+1} ||w||_{k+2} ||\boldsymbol{\rho}||.$$
(2.4.2)

In this paper, we use e to denote the error between the exact and numerical solutions, i.e.  $e_c = c - c_h$ ,  $e_p = p - p_h$ ,  $\mathbf{e}_u = \mathbf{u} - \mathbf{u}_h$ ,  $\mathbf{e}_s = \mathbf{s} - \mathbf{s}_h$ ,  $\mathbf{e}_z = \mathbf{z} - \mathbf{z}_h$ . As the general treatment of the finite element methods, we split the errors into two terms as

$$e_{c} = \eta_{c} - \xi_{c}, \quad \eta_{c} = c - P^{+}c, \quad \xi_{c} = c_{h} - P^{+}c,$$

$$e_{p} = \eta_{p} - \xi_{p}, \quad \eta_{p} = p - P^{+}p, \quad \xi_{p} = p_{h} - P^{+}p,$$

$$\mathbf{e}_{u} = \boldsymbol{\eta}_{u} - \boldsymbol{\xi}_{u}, \quad \boldsymbol{\eta}_{u} = \mathbf{u} - \boldsymbol{\Pi}^{-}\mathbf{u}, \quad \boldsymbol{\xi}_{u} = \mathbf{u}_{h} - \boldsymbol{\Pi}^{-}\mathbf{u},$$

$$\mathbf{e}_{s} = \boldsymbol{\eta}_{s} - \boldsymbol{\xi}_{s}, \quad \boldsymbol{\eta}_{s} = \mathbf{s} - \mathbf{P}_{k}\mathbf{s}, \quad \boldsymbol{\xi}_{s} = \mathbf{s}_{h} - \mathbf{P}_{k}\mathbf{s},$$

$$\mathbf{e}_{z} = \boldsymbol{\eta}_{z} - \boldsymbol{\xi}_{z}, \quad \boldsymbol{\eta}_{z} = \mathbf{z} - \boldsymbol{\Pi}^{-}\mathbf{z}, \quad \boldsymbol{\xi}_{z} = \mathbf{z}_{h} - \boldsymbol{\Pi}^{-}\mathbf{z}.$$

Based on the above, it is easy to see that

$$\mathcal{L}^{d}(\boldsymbol{\eta}_{u}, v) = \mathcal{L}^{d}(\boldsymbol{\eta}_{z}, v) = 0.$$
(2.4.3)

Following [60, 61, 62, 76] with some minor changes, we have the following lemma

**Lemma 2.4.4.** Suppose  $\xi_c$  and  $\xi_s$  are defined above, we have

$$\|\nabla \xi_c\| \le C(\|\boldsymbol{\xi}_s\| + h^{k+1}), \qquad h^{-\frac{1}{2}} \|[\xi_c]\|_{\Gamma_h} \le C(\|\boldsymbol{\xi}_s\| + h^{k+1}).$$

The proof of the main error estimate requires the following initial discretization, whose detailed construction will be given in the appendix.

Lemma 2.4.5. We choose the initial solution

$$c_h^0 = P^+ c_0, \quad \mathbf{u}_h^0 = \mathbf{\Pi}^- \mathbf{u}_0,$$
 (2.4.4)

where  $\mathbf{u}_0 = -a(c_0)\nabla p_0$ , Then we have

$$||c(x,0) - c_h(x,0)|| \le Ch^{k+1}, \qquad (2.4.5)$$

$$\|\mathbf{u}(x,0) - \mathbf{u}_h(x,0)\| \le Ch^{k+1}, \qquad (2.4.6)$$

$$\|\mathbf{s}(x,0) - \mathbf{s}_h(x,0)\| \le Ch^{k+1},\tag{2.4.7}$$

$$||p_t(x,0) - p_{h_t}(x,0)|| \le Ch^{k+1}, \qquad (2.4.8)$$

$$||p(x,0) - p_h(x,0)|| \le Ch^{k+1}.$$
(2.4.9)

The proof of this lemma will also be given in the appendix.

## 2.4.2 A priori error estimates

In this subsection, we would like to make an a priori error estimate assumption that

$$||c - c_h|| + ||\mathbf{u} - \mathbf{u}_h|| + ||p_t - p_{h_t}|| \le h, \qquad (2.4.10)$$

which further implies

$$||c_h||_{\infty} + ||\mathbf{u}_h||_{\infty} + ||p_{h_t}||_{\infty} \le C$$
(2.4.11)

by hypothesis 4.

#### 2.4.3 Error equations

In this section, we proceed to construct the error equations. From (2.3.7)-(2.3.11), we have the following error equations

$$(B(c)p_t - B(c_h)p_{ht} + \phi e_{ct}, v) = \mathcal{L}^c(\mathbf{u}, c, v) - \mathcal{L}^c(\mathbf{u}_h, c_h, v) \qquad (2.4.12)$$
$$+ \mathcal{L}^d(\mathbf{e}_z, v) + (\tilde{e}_c q, v),$$

$$(\mathbf{e}_s, \mathbf{w}) = \mathcal{D}(e_c, \mathbf{w}), \qquad (2.4.13)$$

$$(\mathbf{e}_z, \boldsymbol{\psi}) = (\mathbf{D}(\mathbf{s} - \mathbf{s}_h), \boldsymbol{\psi}), \qquad (2.4.14)$$

$$((A(c)\mathbf{u} - A(c_h)\mathbf{u}_h), \boldsymbol{\theta}) = \mathcal{D}(e_p, \boldsymbol{\theta}), \qquad (2.4.15)$$

$$(d(c)p_t - d(c_h)p_{ht}, \zeta) = \mathcal{L}^d(\mathbf{e}_u, \zeta), \qquad (2.4.16)$$

 $\forall v, \zeta \in W_h^k, \mathbf{w}, \boldsymbol{\psi}, \boldsymbol{\theta} \in \mathbf{W}_h^k$ , where

$$\tilde{e_c} = \begin{cases} 0, & q > 0, \\ e_c, & q < 0. \end{cases}$$

## 2.4.4 The first energy inequality

Taking the test functions  $v = \xi_c$ ,  $\mathbf{w} = \boldsymbol{\xi}_z$ , and  $\boldsymbol{\psi} = -\boldsymbol{\xi}_s$  in (2.4.12), (2.4.13) and (2.4.14), respectively, and use Lemma 2.3.1 and (2.4.3) to obtain

$$(\phi \frac{\partial \xi_c}{\partial t}, \xi_c) + (\mathbf{D}\boldsymbol{\xi}_s, \boldsymbol{\xi}_s) = R_1 + R_2 - R_3 - R_4 + R_5, \qquad (2.4.17)$$

where

$$\begin{aligned} R_1 &= \left(\phi \frac{\partial \eta_c}{\partial t}, \xi_c\right) + \left(\mathbf{D}\boldsymbol{\eta}_s, \boldsymbol{\xi}_s\right), \\ R_2 &= \left(B(c)p_t - B(c_h)p_{ht}, \xi_c\right), \\ R_3 &= \left(\mathbf{u}c - \mathbf{u}_h c_h, \nabla \xi_c\right) + \sum_{e \in \Gamma_e} \left\langle \left(\mathbf{u}c - \widehat{\mathbf{u}_h c_h}\right) \cdot \boldsymbol{\nu}_e, [\xi_c] \right\rangle_e, \\ R_4 &= \mathcal{D}(\eta_c, \boldsymbol{\xi}_z), \\ R_5 &= \left(\boldsymbol{\eta}_s, \boldsymbol{\xi}_z\right) - \left(\boldsymbol{\eta}_z, \boldsymbol{\xi}_s\right) - \left(\tilde{e_c}q, \xi_c\right), \end{aligned}$$

with  $\Gamma_e = \Gamma_0 \cup \partial \Omega_-$  and  $\langle u, v \rangle_e = \int_e uv ds$ . Now, we estimate  $R'_i s$  term by term. Using hypotheses 1 and 3, Lemma 2.4.2 and the Schwarz inequality, we can get

$$R_{1} \leq C \|\eta_{ct}\| \|\xi_{c}\| + C \|\boldsymbol{\eta}_{s}\| \|\boldsymbol{\xi}_{s}\| \leq Ch^{k+1} \left( \|\xi_{c}\| + \|\boldsymbol{\xi}_{s}\| \right), \qquad (2.4.18)$$

For  $R_2$ , we have

$$R_{2} = \left[ \left( B(c)(p - p_{h})_{t}, \xi_{c} \right) + \left( (B(c) - B(c_{h}))p_{ht}, \xi_{c} \right) \right]$$
  

$$\leq C \| (p - p_{h})_{t} \| \| \xi_{c} \| + C \| c - c_{h} \| \| \xi_{c} \|$$
  

$$\leq C \| \xi_{c} \| \| \xi_{p_{t}} \| + \| \xi_{c} \| + h^{k+1} ), \qquad (2.4.19)$$

where in the second step we use Schwarz inequality and hypothesis 1 and (2.4.11), and the last step requires Lemma 2.4.2. We estimate  $R_3$  by dividing it into three parts

$$R_3 = R_{31} + R_{32} - R_{33}, (2.4.20)$$

where

$$R_{31} = (\mathbf{u}c - \mathbf{u}c_h, \nabla\xi_c) + (\mathbf{u}c_h - \mathbf{u}_h c_h, \nabla\xi_c), \qquad (2.4.21)$$

$$R_{32} = \frac{1}{2} \sum_{e \in \Gamma_e} \langle (2\mathbf{u}c - \mathbf{u}_h^+ c_h^+ - \mathbf{u}_h^- c_h^-) \cdot \boldsymbol{\nu}_e, [\xi_c] \rangle_e, \qquad (2.4.22)$$

$$R_{33} = \frac{1}{2} \sum_{e \in \Gamma_e} \langle \alpha [\eta_c - \xi_c], [\xi_c] \rangle_e.$$
 (2.4.23)

Using hypothesis 4 and (2.4.11), we have

$$R_{31} \leq C \left( \|c - c_h\| + \|\mathbf{u} - \mathbf{u}_h\| \right) \|\nabla \xi_c\|$$
  
$$\leq C \left( h^{k+1} + \|\boldsymbol{\xi}_u\| + \|\boldsymbol{\xi}_c\| \right) \left( \|\boldsymbol{\xi}_s\| + h^{k+1} \right), \qquad (2.4.24)$$

where in the first step, we use Schwarz inequality while the second step follows from Lemmas 2.4.2 and 2.4.4. C depends on  $||u||_{\infty}$  and  $||c_h||_{\infty}$ . The estimate of  $R_{32}$  also requires hypothesis 4 and (2.4.11),

$$R_{32} = \frac{1}{2} \sum_{e \in \Gamma_e} \langle \left( \mathbf{u}(c - c_h^+) + (\mathbf{u} - \mathbf{u}_h^+)c_h^+ + \mathbf{u}(c - c_h^-) + (\mathbf{u} - \mathbf{u}_h^-)c_h^- \right) \cdot \boldsymbol{\nu}_e, [\xi_c] \rangle_e \\ \leq C \left( \|c - c_h\|_{\Gamma_h} + \|\mathbf{u} - \mathbf{u}_h\|_{\Gamma_h} \right) \|[\xi_c]\|_{\Gamma_h} \\ \leq C h^{\frac{1}{2}} (\|\eta_c\|_{\Gamma_h} + \|\xi_c\|_{\Gamma_h} + \|\boldsymbol{\eta}_u\|_{\Gamma_h} + \|\boldsymbol{\xi}_u\|_{\Gamma_h}) (\|\boldsymbol{\xi}_s\| + h^{k+1}) \\ \leq C \left( h^{k+1} + \|\boldsymbol{\xi}_u\| + \|\xi_c\| \right) \left( \|\boldsymbol{\xi}_s\| + h^{k+1} \right), \qquad (2.4.25)$$

where in the second step we use Schwarz inequality, the third step follows from Lemma 2.4.4, the last one requires Lemmas 2.4.1 and 2.4.2. C depends on  $\|\mathbf{u}\|_{\infty}$ and  $\|c_h\|_{\infty}$ . Now we proceed to the estimate of  $R_{33}$ ,

$$R_{33} \leq C(\|\eta_c\|_{\Gamma_h} + \|\xi_c\|_{\Gamma_h})\|[\xi_c]\|_{\Gamma_h}$$
  
$$\leq Ch^{\frac{1}{2}}(\|\eta_c\|_{\Gamma_h} + \|\xi_c\|_{\Gamma_h})(\|\boldsymbol{\xi}_s\| + h^{k+1})$$
  
$$\leq C(h^{k+1} + \|\xi_c\|)(\|\boldsymbol{\xi}_s\| + h^{k+1}), \qquad (2.4.26)$$

where the first step follows from Schwarz inequality, the second step is based on Lemma 2.4.4, the third one requires Lemma 2.4.2. Plug (2.4.24), (2.4.25) and (2.4.26) into (2.4.20) to obtain

$$R_3 \le C \left( h^{k+1} + \|\boldsymbol{\xi}_u\| + \|\boldsymbol{\xi}_c\| \right) \left( \|\boldsymbol{\xi}_s\| + h^{k+1} \right).$$
(2.4.27)

The estimate of  $R_4$  follows from Lemma 2.4.3

$$R_4 \le Ch^{k+1} \|c\|_{k+2} \|\boldsymbol{\xi}_z\|. \tag{2.4.28}$$

Now we begin to deal with  $R_5$ . Using Lemma 2.4.2 and the Schwartz inequality, we easily obtain

$$R_{5} \leq \|\boldsymbol{\eta}_{s}\|\|\boldsymbol{\xi}_{z}\| + \|\boldsymbol{\eta}_{z}\|\|\boldsymbol{\xi}_{s}\| + C\|\tilde{e}_{c}\|\|\boldsymbol{\xi}_{c}\| \leq Ch^{k+1}(\|\boldsymbol{\xi}_{z}\| + \|\boldsymbol{\xi}_{s}\|) + Ch^{k+1}\|\boldsymbol{\xi}_{c}\| + C\|\boldsymbol{\xi}_{c}\|^{2}.$$
(2.4.29)

Substituting the estimation (2.4.18), (2.4.19), (2.4.27), (2.4.28), (2.4.29) into (2.4.17) and use hypothesis 3, we obtain

$$\frac{d}{dt} \|\phi^{\frac{1}{2}} \xi_{c}\|^{2} + \|\boldsymbol{D}^{\frac{1}{2}} \boldsymbol{\xi}_{s}\|^{2} \leq C \left[ \left( h^{k+1} + \|\boldsymbol{\xi}_{u}\| + \|\boldsymbol{\xi}_{c}\| \right) \left( \|\boldsymbol{\xi}_{s}\| + h^{k+1} \right) + h^{k+1} \|\boldsymbol{\xi}_{z}\| + h^{2(k+1)} + \|\boldsymbol{\xi}_{c}\|^{2} + \|\boldsymbol{\xi}_{p_{t}}\|^{2} \right]. (2.4.30)$$

Integrating with the equation with respect to time between 0 and t, we obtain

$$\|\xi_{c}\|^{2} + \int_{0}^{t} \|\boldsymbol{\xi}_{s}\|^{2} dt$$

$$\leq C \int_{0}^{t} (\|\xi_{c}\|^{2} + \|\boldsymbol{\xi}_{u}\|^{2} + \|\boldsymbol{\xi}_{p_{t}}\|^{2} + \|\boldsymbol{\xi}_{z}\|^{2} + \|\boldsymbol{\xi}_{s}\|^{2}) dt + Ch^{2(k+1)}. \quad (2.4.31)$$

We take the time derivative in equation (2.4.13), we have

$$(\mathbf{e}_{st}, \mathbf{w}) = \mathcal{D}(e_{ct}, \mathbf{w}), \qquad (2.4.32)$$

Taking the test functions  $v = \xi_{ct}$ ,  $\mathbf{w} = \boldsymbol{\xi}_z$ , and  $\boldsymbol{\psi} = -\xi_{st}$  in (2.4.12), (2.4.32) and (2.4.14), respectively, and use (2.3.15) and (2.4.3) to obtain

$$(\phi\xi_{ct},\xi_{ct}) + \frac{1}{2}\frac{d}{dt}(\mathbf{D}\boldsymbol{\xi}_s,\boldsymbol{\xi}_s) = \tilde{R}_1 + \tilde{R}_2 + \tilde{R}_3 + \tilde{R}_4 + \tilde{R}_5 + \tilde{R}_6, \qquad (2.4.33)$$

where

$$\begin{split} R_1 &= (\phi \eta_{ct}, \xi_{ct}), \\ \tilde{R}_2 &= (\mathbf{D} \boldsymbol{\eta}_s, \boldsymbol{\xi}_{st}), \\ \tilde{R}_3 &= (B(c)p_t - B(c_h)p_{ht}, \xi_{ct}), \\ \tilde{R}_4 &= -(\mathbf{u}c - \mathbf{u}_h c_h, \nabla \xi_{ct}) - \sum_{e \in \Gamma_e} \langle (\mathbf{u}c - \widehat{\mathbf{u}_h c_h}) \cdot \boldsymbol{\nu}_e, [\xi_{ct}] \rangle_e, \\ \tilde{R}_5 &= -\mathcal{D}(\eta_{ct}, \boldsymbol{\xi}_z), \\ \tilde{R}_6 &= (\boldsymbol{\eta}_{st}, \boldsymbol{\xi}_z) - (\boldsymbol{\eta}_z, \boldsymbol{\xi}_{st}) - (\tilde{e_c}q, \xi_{ct}), \end{split}$$

Now, we estimate  $\tilde{R}'_i s$  term by term. Using the projection and the Schwartz inequality, we can get

$$\tilde{R}_1 \le C \|\eta_{ct}\|^2 + C \|\xi_{ct}\|^2 \le C h^{2(k+1)} + \epsilon \|\xi_{ct}\|^2, \qquad (2.4.34)$$

$$\tilde{R}_{2} = \frac{d}{dt} (\mathbf{D}\boldsymbol{\eta}_{s}, \boldsymbol{\xi}_{s}) - (\mathbf{D}\boldsymbol{\eta}_{st}, \boldsymbol{\xi}_{s})$$

$$\leq \frac{d}{dt} (\mathbf{D}\boldsymbol{\eta}_{s}, \boldsymbol{\xi}_{s}) + C \|\boldsymbol{\xi}_{s}\|^{2} + Ch^{2(k+1)}, \qquad (2.4.35)$$

$$\tilde{R}_{3} = \left[ \left( B(c)(p - p_{h})_{t}, \xi_{ct} \right) + \left( (B(c) - B(c_{h}))p_{ht}, \xi_{ct} \right) \right] \\
\leq C \| (p - p_{h})_{t} \| \| \xi_{ct} \| + C \| c - c_{h} \| \| \xi_{ct} \| \\
\leq C \| \xi_{p_{t}} \|^{2} + C \| \xi_{c} \|^{2} + \epsilon \| \xi_{ct} \|^{2} + C h^{2(k+1)},$$
(2.4.36)

where in the second step we use Schwarz inequality and hypothesis 1, and the last step requires Lemma 2.4.2. We estimate  $R_4$  by dividing it into three parts

$$\tilde{R}_4 = \tilde{R}_{41} + \tilde{R}_{42} + \tilde{R}_{43}, \qquad (2.4.37)$$

where

$$\begin{split} \tilde{R}_{41} &= -(\mathbf{u}c - \mathbf{u}_h c_h, \nabla \xi_{ct}), \\ R_{42} &= -\frac{1}{2} \sum_{e \in \Gamma_e} \langle (2\mathbf{u}c - \mathbf{u}_h^+ c_h^+ - \mathbf{u}_h^- c_h^-) \cdot \boldsymbol{\nu}_e, [\xi_{ct}] \rangle_e, \\ R_{43} &= \frac{1}{2} \sum_{e \in \Gamma_e} \langle \alpha [\eta_c - \xi_c], [\xi_{ct}] \rangle_e. \end{split}$$

Using hypothesis 4 and (2.4.11), we have

$$\tilde{R}_{41} = \frac{d}{dt} \left( \mathbf{u}_h c_h - \mathbf{u}_c, \nabla \xi_c \right) + \left( (\mathbf{u}_c - \mathbf{u}_h c_h)_t, \nabla \xi_c \right) \\
= \frac{d}{dt} \left( \mathbf{u}_h c_h - \mathbf{u}_c, \nabla \xi_c \right) + \left( \mathbf{u}_t c - \mathbf{u}_{ht} c_h, \nabla \xi_c \right) + \left( \mathbf{u}_{ct} - \mathbf{u}_h c_{ht}, \nabla \xi_c \right) \\
= \frac{d}{dt} \left( \mathbf{u}_h c_h - \mathbf{u}_c, \nabla \xi_c \right) + \left( \mathbf{u}_t (c - c_h), \nabla \xi_c \right) + \left( (\mathbf{u} - \mathbf{u}_h)_t c_h, \nabla \xi_c \right) \\
+ (c_t (\mathbf{u} - \mathbf{u}_h), \nabla \xi_c) + ((c - c_h)_t \mathbf{u}_h, \nabla \xi_c) \\
\leq \frac{d}{dt} \left( \mathbf{u}_h c_h - \mathbf{u}_c, \nabla \xi_c \right) + C \|c - c_h\|^2 + \epsilon \| (\mathbf{u} - \mathbf{u}_h)_t \|^2 \\
+ C \| \mathbf{u} - \mathbf{u}_h \|^2 + \epsilon \| (c - c_h)_t \|^2 + C \| \nabla \xi_c \|^2 \\
\leq \frac{d}{dt} \left( \mathbf{u}_h c_h - \mathbf{u}_c, \nabla \xi_c \right) + C h^{2(k+1)} + C \| \xi_c \|^2 + \epsilon \| \xi_{ut} \|^2 \\
+ C \| \xi_u \|^2 + \epsilon \| \xi_{ct} \|^2 + C \| \xi_s \|^2,$$
(2.4.38)

where in the forth step, we use Schwarz inequality while the last step follows from Lemmas 2.4.2 and 2.4.4. The estimate of  $\tilde{R}_{42}$  also requires hypothesis 4 and (2.4.11),

$$\begin{split} \tilde{R}_{42} &= -\frac{1}{2} \sum_{e \in \Gamma_e} \frac{d}{dt} \langle (2\mathbf{u}c - \mathbf{u}_h^+ c_h^+ - \mathbf{u}_h^- c_h^-) \cdot \boldsymbol{\nu}_e, [\boldsymbol{\xi}_c] \rangle_e \\ &+ \sum_{e \in \Gamma_e} \langle (\frac{\mathbf{u}^+ c^+ + \mathbf{u}^- c^-}{2} - \frac{\mathbf{u}_h^+ c_h^+ + \mathbf{u}_h^- c_h^-}{2})_t \cdot \boldsymbol{\nu}_e, [\boldsymbol{\xi}_c] \rangle_e \\ &\leq \frac{1}{2} \sum_{e \in \Gamma_e} \frac{d}{dt} \langle (\mathbf{u}_h^+ c_h^+ + \mathbf{u}_h^- c_h^- - 2\mathbf{u}c) \cdot \boldsymbol{\nu}_e, [\boldsymbol{\xi}_c] \rangle_e + C \| (\mathbf{u}c - \mathbf{u}_h c_h)_t \|_{\Gamma_h} \| [\boldsymbol{\xi}_c] \|_{\Gamma_h} \\ &\leq \frac{1}{2} \sum_{e \in \Gamma_e} \frac{d}{dt} \langle (\mathbf{u}_h^+ c_h^+ + \mathbf{u}_h^- c_h^- - 2\mathbf{u}c) \cdot \boldsymbol{\nu}_e, [\boldsymbol{\xi}_c] \rangle_e \\ &+ Ch^{\frac{1}{2}} (\| c_t (\mathbf{u} - \mathbf{u}_h) \|_{\Gamma_h} + \| (c - c_h)_t \mathbf{u}_h \|_{\Gamma_h}) (\| \boldsymbol{\xi}_{\mathbf{s}} \| + h^{k+1}) \\ &+ Ch^{\frac{1}{2}} (\| \mathbf{u}_t (c - c_h) \|_{\Gamma_h} + \| (\mathbf{u} - \mathbf{u}_h)_t c_h \|_{\Gamma_h}) (\| \boldsymbol{\xi}_{\mathbf{s}} \| + h^{k+1}) \\ &\leq \frac{1}{2} \sum_{e \in \Gamma_e} \frac{d}{dt} \langle (\mathbf{u}_h^+ c_h^+ + \mathbf{u}_h^- c_h^- - 2\mathbf{u}c) \cdot \boldsymbol{\nu}_e, [\boldsymbol{\xi}_c] \rangle_e \\ &+ Ch^{2(k+1)} + C \| \boldsymbol{\xi}_u \|^2 + \epsilon \| \boldsymbol{\xi}_{ct} \|^2 + C \| \boldsymbol{\xi}_c \|^2 + \epsilon \| \boldsymbol{\xi}_{ut} \|^2 + C \| \boldsymbol{\xi}_s \|^2, \quad (2.4.39) \end{split}$$

where in the second step we use Schwarz inequality, the third step follows from and Lemma 2.4.4, the last one requires Lemmas 2.4.1 and 2.4.2. Now we proceed to the estimate of  $\tilde{R}_{43}$ ,

$$\tilde{R}_{43} = \frac{1}{2} \sum_{e \in \Gamma_e} \frac{d}{dt} \langle \alpha [\eta_c - \xi_c], [\xi_c] \rangle_e - \frac{1}{2} \sum_{e \in \Gamma_e} \langle \alpha [\eta_{ct} - \xi_{ct}], [\xi_c] \rangle_e 
\leq \frac{1}{2} \sum_{e \in \Gamma_e} \frac{d}{dt} \langle \alpha [\eta_c - \xi_c], [\xi_c] \rangle_e + C(\|\eta_{ct}\|_{\Gamma_h} + \|\xi_{ct}\|_{\Gamma_h}) \|[\xi_c]\|_{\Gamma_h} 
\leq \frac{1}{2} \sum_{e \in \Gamma_e} \frac{d}{dt} \langle \alpha [\eta_c - \xi_c], [\xi_c] \rangle_e + Ch^{\frac{1}{2}} (\|\eta_{ct}\|_{\Gamma_h} + \|\xi_{ct}\|_{\Gamma_h}) (\|\boldsymbol{\xi}_s\| + h^{k+1}) 
\leq \frac{1}{2} \sum_{e \in \Gamma_e} \frac{d}{dt} \langle \alpha [\eta_c - \xi_c], [\xi_c] \rangle_e + Ch^{2(k+1)} + \epsilon \|\xi_{ct}\|^2 + C\|\boldsymbol{\xi}_s\|^2, \quad (2.4.40)$$

where the second step follows from Schwarz inequality, the third one is based on Lemma 2.4.4, the last one requires Lemmas 2.4.1 and 2.4.2. Plug (2.4.38), (2.4.39) and (2.4.40) into (2.4.37) to obtain

$$\tilde{R}_{4} \leq \frac{d}{dt} \Big( \mathbf{u}_{h} c_{h} - \mathbf{u} c, \nabla \xi_{c} \Big) + \frac{1}{2} \sum_{e \in \Gamma_{e}} \frac{d}{dt} \langle (\mathbf{u}_{h}^{+} c_{h}^{+} + \mathbf{u}_{h}^{-} c_{h}^{-} - 2\mathbf{u} c) \cdot \boldsymbol{\nu}_{e}, [\xi_{c}] \rangle_{e} 
+ \frac{1}{2} \sum_{e \in \Gamma_{e}} \frac{d}{dt} \langle \alpha [\eta_{c} - \xi_{c}], [\xi_{c}] \rangle_{e} + C (h^{2(k+1)} + \|\boldsymbol{\xi}_{u}\|^{2} + \|\boldsymbol{\xi}_{c}\|^{2} + \|\boldsymbol{\xi}_{s}\|^{2}) 
+ \epsilon (\|\xi_{ct}\|^{2} + \|\boldsymbol{\xi}_{ut}\|^{2}).$$
(2.4.41)

The estimate of  $\tilde{R}_5$  follows from Lemma 2.4.3

$$\tilde{R}_5 \le Ch^{k+1} \|c\|_{k+2} \|\boldsymbol{\xi}_z\|.$$
(2.4.42)

Now we begin to deal with  $\tilde{R}_6$ . Using Lemma 2.4.2 and the Schwartz inequality, we easily obtain

$$\tilde{R}_{6} = (\boldsymbol{\eta}_{st}, \boldsymbol{\xi}_{z}) - \frac{d}{dt}(\boldsymbol{\eta}_{z}, \boldsymbol{\xi}_{s}) + (\boldsymbol{\eta}_{zt}, \boldsymbol{\xi}_{s}) - (\tilde{e}_{c}q, \boldsymbol{\xi}_{ct}) 
\leq \|\boldsymbol{\eta}_{st}\| \|\boldsymbol{\xi}_{z}\| - \frac{d}{dt}(\boldsymbol{\eta}_{z}, \boldsymbol{\xi}_{s}) + \|\boldsymbol{\eta}_{zt}\| \|\boldsymbol{\xi}_{s}\| + C \|\tilde{e}_{c}\| \|\boldsymbol{\xi}_{ct}\| 
\leq -\frac{d}{dt}(\boldsymbol{\eta}_{z}, \boldsymbol{\xi}_{s}) + C \left(h^{2(k+1)} + \|\boldsymbol{\xi}_{z}\|^{2} + \|\boldsymbol{\xi}_{s}\|^{2} + \|\boldsymbol{\xi}_{c}\|^{2}\right) + \epsilon \|\boldsymbol{\xi}_{ct}\|^{2} (2.4.43)$$

Substituting the estimation (2.4.34)-(2.4.36) and (2.4.41)-(2.4.43)into (2.4.33) and use hypothesis 3, we obtain

$$\begin{split} \|\phi^{\frac{1}{2}}\xi_{ct}\|^{2} + \frac{1}{2}\frac{d}{dt}\|D^{\frac{1}{2}}\boldsymbol{\xi}_{s}\|^{2} \\ &\leq \frac{d}{dt}(\mathbf{D}\boldsymbol{\eta}_{s},\boldsymbol{\xi}_{s}) - \frac{d}{dt}(\boldsymbol{\eta}_{z},\boldsymbol{\xi}_{s}) + \frac{d}{dt}\left(\mathbf{u}_{h}c_{h} - \mathbf{u}c,\nabla\boldsymbol{\xi}_{c}\right) \\ &+ \frac{1}{2}\sum_{e\in\Gamma_{e}}\frac{d}{dt}\langle(\mathbf{u}_{h}^{+}c_{h}^{+} + \mathbf{u}_{h}^{-}c_{h}^{-} - 2\mathbf{u}c)\cdot\boldsymbol{\nu}_{e},[\boldsymbol{\xi}_{c}]\rangle_{e} + \frac{1}{2}\sum_{e\in\Gamma_{e}}\frac{d}{dt}\langle\alpha[\boldsymbol{\eta}_{c} - \boldsymbol{\xi}_{c}],[\boldsymbol{\xi}_{c}]\rangle_{e} \\ &+ C(h^{2(k+1)} + \|\boldsymbol{\xi}_{u}\|^{2} + \|\boldsymbol{\xi}_{s}\|^{2} + \|\boldsymbol{\xi}_{c}\|^{2} + \|\boldsymbol{\xi}_{p_{t}}\|^{2} + \|\boldsymbol{\xi}_{z}\|^{2}) \\ &+ \epsilon(\|\boldsymbol{\xi}_{ct}\|^{2} + \|\boldsymbol{\xi}_{ut}\|^{2}). \end{split}$$
(2.4.44)

Noticing that

$$(\mathbf{D}\boldsymbol{\eta}_{s},\boldsymbol{\xi}_{s}) - (\boldsymbol{\eta}_{z},\boldsymbol{\xi}_{s}) \leq C \|\boldsymbol{\eta}_{s}\|^{2} + \|\boldsymbol{\eta}_{z}\|^{2} + \epsilon \|\boldsymbol{\xi}_{s}\|^{2} \leq Ch^{2(k+1)} + \epsilon \|\boldsymbol{\xi}_{s}\|^{2}. \quad (2.4.45)$$

and

$$\left( \mathbf{u}_{h}c_{h} - \mathbf{u}c, \nabla\xi_{c} \right) = (c(\mathbf{u}_{h} - \mathbf{u}), \nabla\xi_{c}) + (\mathbf{u}_{h}(c_{h} - c), \nabla\xi_{c})$$

$$\leq C \|\mathbf{u} - \mathbf{u}_{h}\|^{2} + C \|c - c_{h}\|^{2} + C \|\nabla\xi_{c}\|^{2}$$

$$\leq Ch^{2(k+1)} + C \|\boldsymbol{\xi}_{u}\|^{2} + C \|\boldsymbol{\xi}_{c}\|^{2} + \epsilon \|\boldsymbol{\xi}_{s}\|^{2}, \quad (2.4.46)$$

where the last one requires Lemmas 2.4.2 and 2.4.4.

$$\frac{1}{2} \sum_{e \in \Gamma_{e}} \langle (\mathbf{u}_{h}^{+} c_{h}^{+} + \mathbf{u}_{h}^{-} c_{h}^{-} - 2\mathbf{u}c) \cdot \boldsymbol{\nu}_{e}, [\xi_{c}] \rangle_{e} + \frac{1}{2} \sum_{e \in \Gamma_{e}} \langle \alpha [\eta_{c} - \xi_{c}], [\xi_{c}] \rangle_{e} \\
\leq C(\|\mathbf{u}c - \mathbf{u}_{h}c_{h}\|_{\Gamma_{h}} + \|\eta_{c}\|_{\Gamma_{h}} + \|\xi_{c}\|_{\Gamma_{h}})\|[\xi_{c}]\|_{\Gamma_{h}} \\
\leq Ch^{\frac{1}{2}}(\|\mathbf{u}c - \mathbf{u}_{h}c\|_{\Gamma_{h}} + \|\mathbf{u}_{h}c - \mathbf{u}_{h}c_{h}\|_{\Gamma_{h}} + \|\eta_{c}\|_{\Gamma_{h}} + \|\xi_{c}\|_{\Gamma_{h}})(\|\boldsymbol{\xi}_{s}\| + h^{k+1}) \\
\leq Ch^{2(k+1)} + C\|\boldsymbol{\xi}_{u}\|^{2} + C\|\xi_{c}\|^{2} + \epsilon\|\boldsymbol{\xi}_{s}\|^{2}, \qquad (2.4.47)$$

where the second step follows from Schwarz inequality, the third one is based on Lemma 2.4.4, the last one requires Lemmas 2.4.1 and 2.4.2. Integrating (2.4.44)with respect to time between 0 and t, then applying (2.4.45)-(2.4.47), we obtain

$$\int_{0}^{t} \|\xi_{ct}\|^{2} dt + \|\xi_{s}\|^{2} \leq C \int_{0}^{t} (\|\xi_{u}\|^{2} + \|\xi_{s}\|^{2} + \|\xi_{c}\|^{2} + \|\xi_{p_{t}}\|^{2} + \|\xi_{z}\|^{2}) dt 
+ \epsilon \int_{0}^{t} (\|\xi_{ct}\|^{2} + \|\xi_{ut}\|^{2}) dt + Ch^{2(k+1)} 
+ C \|\xi_{u}\|^{2} + C \|\xi_{c}\|^{2} + \epsilon \|\xi_{s}\|^{2}.$$
(2.4.48)

Combining (2.4.48) and (2.4.31), we obtain

$$\int_{0}^{t} \|\xi_{ct}\|^{2} dt + \|\xi_{s}\|^{2} + \|\xi_{c}\|^{2} \leq C \int_{0}^{t} (\|\xi_{u}\|^{2} + \|\xi_{s}\|^{2} + \|\xi_{c}\|^{2} + \|\xi_{p_{t}}\|^{2} + \|\xi_{z}\|^{2}) dt + \epsilon \int_{0}^{t} (\|\xi_{ct}\|^{2} + \|\xi_{ut}\|^{2}) dt + Ch^{2(k+1)} + C\|\xi_{u}\|^{2} + \epsilon \|\xi_{s}\|^{2}.$$

which further yields

$$\int_{0}^{t} \|\xi_{ct}\|^{2} dt + \|\xi_{s}\|^{2} + \|\xi_{c}\|^{2} \leq C \int_{0}^{t} (\|\xi_{u}\|^{2} + \|\xi_{s}\|^{2} + \|\xi_{c}\|^{2} + \|\xi_{p_{t}}\|^{2} + \|\xi_{z}\|^{2}) dt + \epsilon \int_{0}^{t} \|\xi_{ut}\|^{2} dt + Ch^{2(k+1)} + C\|\xi_{u}\|^{2}.$$
(2.4.49)

Now, we proceed to eliminate  $\|\boldsymbol{\xi}_z\|$  on the right-hand side to the above equation. Setting  $\boldsymbol{\psi} = \boldsymbol{\xi}_z$  in (2.4.14) to obtain

$$(\boldsymbol{\xi}_z, \boldsymbol{\xi}_z) = (\boldsymbol{\eta}_z, \boldsymbol{\xi}_z) - (\mathbf{D}(\mathbf{s} - \mathbf{s}_h), \boldsymbol{\xi}_z).$$

Then we have

$$\|\boldsymbol{\xi}_{z}\|^{2} \leq \|\boldsymbol{\eta}_{z}\|\|\boldsymbol{\xi}_{z}\| + C\left(\|\boldsymbol{\eta}_{s}\| + \|\boldsymbol{\xi}_{s}\|\right)\|\boldsymbol{\xi}_{z}\| \leq C(\|\boldsymbol{\xi}_{s}\|^{2} + h^{2(k+1)}) + \epsilon\|\boldsymbol{\xi}_{z}\|^{2},$$

where in the first step we use Schwarz inequality and hypothesis 3, the second step follows from Lemma 2.4.2. We can cancel  $\|\boldsymbol{\xi}_z\|$  in the above equation to obtain

$$\|\boldsymbol{\xi}_{z}\|^{2} \leq C(\|\boldsymbol{\xi}_{s}\|^{2} + h^{2(k+1)}).$$
(2.4.50)

Combining (2.4.49) and (2.4.50), we obtain the first energy Inequality

$$\int_{0}^{t} \|\xi_{ct}\|^{2} dt + \|\xi_{s}\|^{2} + \|\xi_{c}\|^{2} \leq C \int_{0}^{t} (\|\xi_{u}\|^{2} + \|\xi_{s}\|^{2} + \|\xi_{c}\|^{2} + \|\xi_{p_{t}}\|^{2}) dt + \epsilon \int_{0}^{t} \|\xi_{ut}\|^{2} dt + Ch^{2(k+1)} + C \|\xi_{u}\|^{2} (2.4.51)$$

#### 2.4.5 The second energy inequality

We start from an easier case. Take  $\boldsymbol{\theta} = \boldsymbol{\xi}_u$  and  $\zeta = \boldsymbol{\xi}_p$  in (2.4.15) and (2.4.16), respectively and use Lemma 2.3.1 and (2.4.3) to obtain

$$(A(c)\boldsymbol{\xi}_{u},\boldsymbol{\xi}_{u}) + \frac{1}{2}\frac{d}{dt}(d(c)\boldsymbol{\xi}_{p},\boldsymbol{\xi}_{p}) = T_{1} + T_{2} + T_{3} + T_{4} + T_{5} - T_{6}, \qquad (2.4.52)$$

where

$$T_{1} = (A(c)\boldsymbol{\eta}_{u}, \boldsymbol{\xi}_{u}),$$

$$T_{2} = ((A(c) - A(c_{h}))\mathbf{u}_{h}, \boldsymbol{\xi}_{u}),$$

$$T_{3} = \frac{1}{2}(d(c)_{t}\xi_{p}, \xi_{p}),$$

$$T_{4} = (d(c)\eta_{p_{t}}, \xi_{p}),$$

$$T_{5} = ((d(c) - d(c_{h}))p_{h_{t}}, \xi_{p}),$$

$$T_{6} = \mathcal{D}(\eta_{p}, \boldsymbol{\xi}_{u}).$$

Now, we estimate  $T'_is$  term by term. Using Lemma 2.4.2 and Schwarz inequality, we can get

$$T_1 \le C \|\boldsymbol{\eta}_u\|^2 + \epsilon \|\boldsymbol{\xi}_u\|^2 \le Ch^{2(k+1)} + \epsilon \|\boldsymbol{\xi}_u\|^2, \qquad (2.4.53)$$

where we use hypothesis 1 to obtain  $|A(c)| = |\frac{\mu(c)}{\kappa(x,y)}| \le \frac{\mu^*}{\kappa_*}$ . Using 2.4.11, we have

$$T_{2} \leq C \|A(c) - A(c_{h})\|^{2} + \epsilon \|\boldsymbol{\xi}_{u}\|^{2} \leq C \|A_{c}'(c - c_{h})\|^{2} + \epsilon \|\boldsymbol{\xi}_{u}\|^{2}$$
$$\leq C h^{2(k+1)} + C \|\boldsymbol{\xi}_{c}\|^{2} + \epsilon \|\boldsymbol{\xi}_{u}\|^{2}, \qquad (2.4.54)$$

where in the first step we use Schwarz inequality, the second step follows from hypothesis 1, and the last step requires Lemma 2.4.2. Moreover,  $A'_c$  is the mean

value given by  $A'_{c} = A'(\lambda_{c}c + (1 - \lambda_{c})c_{h})$  with  $0 \le \lambda_{c} \le 1$ .

$$T_3 = \frac{1}{2} (d'(c)c_t \xi_p, \xi_p) \le C \|\xi_p\|^2, \qquad (2.4.55)$$

where we use hypothesis 1.

$$T_4 \le C \|\eta_{p_t}\|^2 + C \|\xi_p\|^2 \le Ch^{2(k+1)} + C \|\xi_p\|^2, \qquad (2.4.56)$$

$$T_{5} \leq C \|d(c) - d(c_{h})\|^{2} + C \|\xi_{p}\|^{2} \leq C \|d'_{c}(c - c_{h})\|^{2} + C \|\xi_{p}\|^{2}$$
$$\leq C h^{2(k+1)} + C \|\xi_{c}\|^{2} + C \|\xi_{p}\|^{2}, \qquad (2.4.57)$$

where in the first step we use (2.4.11), the second step follows from hypothesis 1 with  $d'_c$  being the mean value given by  $d'_c = d'(\lambda_c c + (1 - \lambda_c)c_h)$  with  $0 \le \lambda_c \le 1$ . For  $T_6$ , we use Lemma 2.4.3 and Schwarz inequality to obtain

$$T_6 \le Ch^{2(k+1)} + \epsilon \|\boldsymbol{\xi}_u\|^2. \tag{2.4.58}$$

Substituting (2.4.53)-(2.4.58) into (2.4.52), we have the estimate

$$\|A^{\frac{1}{2}}(c)\boldsymbol{\xi}_{u}\|^{2} + \frac{1}{2}\frac{d}{dt}\|d^{\frac{1}{2}}(c)\boldsymbol{\xi}_{p}\|^{2} \le Ch^{2(k+1)} + C\|\boldsymbol{\xi}_{p}\|^{2} + C\|\boldsymbol{\xi}_{c}\|^{2} + \epsilon\|\boldsymbol{\xi}_{u}\|^{2}.$$
 (2.4.59)

Integrating (2.4.59) with respect to time between 0 and t and using the hypothesis 1, we obtain the second energy Inequality

$$\int_0^t \|\boldsymbol{\xi}_u\|^2 dt + \|\boldsymbol{\xi}_p\|^2 \le Ch^{2(k+1)} + C \int_0^t (\|\boldsymbol{\xi}_p\|^2 + \|\boldsymbol{\xi}_c\|^2) dt.$$
(2.4.60)

#### 2.4.6 The third energy inequality

We take the time derivative in equation (2.4.15), we have

$$((A(c)\mathbf{u} - A(c_h)\mathbf{u}_h)_t, \boldsymbol{\theta}) = \mathcal{D}(e_{p_t}, \boldsymbol{\theta}), \qquad (2.4.61)$$

Take  $\boldsymbol{\theta} = \boldsymbol{\xi}_u$  and  $\zeta = \xi_{p_t}$  in (2.4.61) and (2.4.16), respectively and use (2.3.15) and (2.4.3) to obtain

$$\frac{1}{2}\frac{d}{dt}(A(c)\boldsymbol{\xi}_{u},\boldsymbol{\xi}_{u}) + (d(c)\boldsymbol{\xi}_{p_{t}},\boldsymbol{\xi}_{p_{t}}) = \tilde{T}_{1} + \tilde{T}_{2} + \tilde{T}_{3} + \tilde{T}_{4} + \tilde{T}_{5} - \tilde{T}_{6}, \qquad (2.4.62)$$

where

$$\begin{split} \tilde{T}_{1} &= -\frac{1}{2} ((A(c))_{t} \boldsymbol{\xi}_{u}, \boldsymbol{\xi}_{u}), \\ \tilde{T}_{2} &= ((A(c) \boldsymbol{\eta}_{u})_{t}, \boldsymbol{\xi}_{u}), \\ \tilde{T}_{3} &= (((A(c) - A(c_{h})) \mathbf{u}_{h})_{t}, \boldsymbol{\xi}_{u}), \\ \tilde{T}_{4} &= (d(c) \eta_{p_{t}}, \boldsymbol{\xi}_{p_{t}}), \\ \tilde{T}_{5} &= ((d(c) - d(c_{h})) p_{h_{t}}, \boldsymbol{\xi}_{p_{t}}), \\ \tilde{T}_{6} &= \mathcal{D}(\eta_{p_{t}}, \boldsymbol{\xi}_{u}). \end{split}$$

Now, we estimate  $\tilde{T}'_is$  term by term. Using hypothesis 1 and Schwarz inequality, we can get

$$\tilde{T}_{1} = -\frac{1}{2} (A'(c)c_{t}\boldsymbol{\xi}_{u}, \boldsymbol{\xi}_{u}) \leq C \|\boldsymbol{\xi}_{u}\|^{2}, \qquad (2.4.63)$$

and

$$\tilde{T}_{2} = (A'(c)c_{t}\boldsymbol{\eta}_{u}, \boldsymbol{\xi}_{u}) + (A(c)\boldsymbol{\eta}_{ut}, \boldsymbol{\xi}_{u}) 
\leq C \|\boldsymbol{\xi}_{u}\|^{2} + C \|\boldsymbol{\eta}_{u}\|^{2} + C \|\boldsymbol{\eta}_{ut}\|^{2} 
\leq C \|\boldsymbol{\xi}_{u}\|^{2} + Ch^{2(k+1)}.$$
(2.4.64)

The estimate of  $\tilde{T}_3$  is slightly complicated,

$$\tilde{T}_{3} = ((A(c) - A(c_{h}))_{t}\boldsymbol{u}_{h}, \boldsymbol{\xi}_{u}) - ((A(c) - A(c_{h}))(\mathbf{u} - \mathbf{u}_{h})_{t}, \boldsymbol{\xi}_{u}) 
+ ((A(c) - A(c_{h}))\boldsymbol{u}_{t}, \boldsymbol{\xi}_{u}) 
= ((A'(c) - A'(c_{h}))c_{t}\mathbf{u}_{h}, \boldsymbol{\xi}_{u}) + (A'(c_{h})(c - c_{h})_{t}\mathbf{u}_{h}, \boldsymbol{\xi}_{u}) 
- (A'_{c}(c - c_{h})(\mathbf{u} - \mathbf{u}_{h})_{t}, \boldsymbol{\xi}_{u}) + (A'_{c}(c - c_{h})\mathbf{u}_{t}, \boldsymbol{\xi}_{u}) 
\leq C \|c - c_{h}\| \|\boldsymbol{\xi}_{u}\| + C \|(c - c_{h})_{t}\| \|\boldsymbol{\xi}_{u}\| 
+ C \|\boldsymbol{\xi}_{u}\|_{\infty} \|c - c_{h}\| \|(\mathbf{u} - \mathbf{u}_{h})_{t}\| + C \|c - c_{h}\| \|\boldsymbol{\xi}_{u}\| 
\leq C \|c - c_{h}\|^{2} + C \|\boldsymbol{\xi}_{u}\|^{2} + \epsilon \|(c - c_{h})_{t}\|^{2} + \epsilon \|(\mathbf{u} - \mathbf{u}_{h})_{t}\|^{2} 
\leq C \|\xi_{c}\|^{2} + C \|\boldsymbol{\xi}_{u}\|^{2} + \epsilon \|\boldsymbol{\xi}_{ct}\|^{2} + \epsilon \|\boldsymbol{\xi}_{ut}\|^{2} + C h^{2(k+1)}, \qquad (2.4.65)$$

where in the third step we use Schwarz inequality and hypotheses 1 and 2, and the last step requires Lemma 2.4.2. Applying the Schwarz inequality, we have

$$\tilde{T}_4 \le C \|\eta_{p_t}\|^2 + \epsilon \|\xi_{p_t}\|^2 \le Ch^{2(k+1)} + \epsilon \|\xi_{p_t}\|^2, \qquad (2.4.66)$$

$$\tilde{T}_{5} \leq C \|d(c) - d(c_{h})\|^{2} + \epsilon \|\xi_{p_{t}}\|^{2} \leq C \|d_{c}'(c - c_{h})\|^{2} + \epsilon \|\xi_{p_{t}}\|^{2}$$

$$\leq Ch^{2(k+1)} + C \|\xi_{c}\|^{2} + \epsilon \|\xi_{p_{t}}\|^{2}, \qquad (2.4.67)$$

For  $\tilde{T}_6$ , we use Lemma 2.4.3 to obtain

$$\tilde{T}_6 \le Ch^{k+1} \|p\|_{k+2} \|\boldsymbol{\xi}_u\|. \tag{2.4.68}$$

Substituting (2.4.63)-(2.4.68) into (2.4.62), we have the estimate

$$\frac{1}{2} \frac{d}{dt} \|A^{\frac{1}{2}}(c)\boldsymbol{\xi}_{u}\|^{2} + \|d^{\frac{1}{2}}(c)\boldsymbol{\xi}_{p_{t}}\|^{2} \\
\leq Ch^{2(k+1)} + C\|\boldsymbol{\xi}_{u}\|^{2} + C\|\boldsymbol{\xi}_{c}\|^{2} + \epsilon\|\boldsymbol{\xi}_{p_{t}}\|^{2} + \epsilon\|\boldsymbol{\xi}_{ct}\|^{2} + \epsilon\|\boldsymbol{\xi}_{ut}\|^{2}. \quad (2.4.69)$$

Integrating (2.4.69) with respect to time between 0 and t and using the hypothesis 1, we obtain the third energy Inequality

$$\|\boldsymbol{\xi}_{u}\|^{2} + \int_{0}^{t} \|\xi_{p_{t}}\|^{2} dt$$
  

$$\leq Ch^{2(k+1)} + C \int_{0}^{t} (\|\boldsymbol{\xi}_{u}\|^{2} + \|\xi_{c}\|^{2}) dt + \epsilon \int_{0}^{t} (\|\xi_{c_{t}}\|^{2} + \|\boldsymbol{\xi}_{u_{t}}\|^{2}) dt. (2.4.70)$$

#### 2.4.7 The fourth energy inequality

We take the time derivative in equation (2.4.16), we have

$$((d(c)p_t - d(c_h)p_{ht})_t, \zeta) = \mathcal{L}^d(\mathbf{e}_{ut}, \zeta), \qquad (2.4.71)$$

Take  $\boldsymbol{\theta} = \boldsymbol{\xi}_{ut}$  and  $\zeta = \xi_{pt}$  in (2.4.61) and (2.4.71), respectively and use (2.3.15) and (2.4.3) to obtain

$$(A(c)\boldsymbol{\xi}_{ut},\boldsymbol{\xi}_{ut}) + \frac{1}{2}\frac{d}{dt}(d(c_h)\xi_{p_t},\xi_{p_t}) = \tilde{\tilde{T}}_1 + \tilde{\tilde{T}}_2 + \tilde{\tilde{T}}_3 - \tilde{\tilde{T}}_4 + \tilde{\tilde{T}}_5 + \tilde{\tilde{T}}_6 - \tilde{\tilde{T}}_7, \quad (2.4.72)$$

where

$$\begin{split} \tilde{\tilde{T}}_{1} &= -((A(c))_{t} \boldsymbol{\xi}_{u}, \boldsymbol{\xi}_{ut}), \\ \tilde{\tilde{T}}_{2} &= ((A(c) \boldsymbol{\eta}_{u})_{t}, \boldsymbol{\xi}_{ut}), \\ \tilde{\tilde{T}}_{3} &= (((A(c) - A(c_{h}))\mathbf{u}_{h})_{t}, \boldsymbol{\xi}_{ut}), \\ \tilde{\tilde{T}}_{4} &= -\frac{1}{2}((d(c_{h}))_{t} \boldsymbol{\xi}_{p_{t}}, \boldsymbol{\xi}_{p_{t}}), \\ \tilde{\tilde{T}}_{5} &= ((d(c_{h}) \eta_{p_{t}})_{t}, \boldsymbol{\xi}_{p_{t}}), \\ \tilde{\tilde{T}}_{6} &= (((d(c) - d(c_{h}))p_{t})_{t}, \boldsymbol{\xi}_{p_{t}}), \\ \tilde{\tilde{T}}_{7} &= \mathcal{D}(\eta_{p_{t}}, \boldsymbol{\xi}_{ut}). \end{split}$$

Now, we estimate  $\tilde{\tilde{T}}_i's$  term by term. Using hypothesis 1 and Schwarz inequality, we can get

$$\tilde{\tilde{T}}_{1} = -\frac{1}{2} (A'(c)c_{t}\boldsymbol{\xi}_{u}, \boldsymbol{\xi}_{ut}) \leq C \|\boldsymbol{\xi}_{u}\|^{2} + \epsilon \|\boldsymbol{\xi}_{ut}\|^{2}, \qquad (2.4.73)$$

and

$$\tilde{\tilde{T}}_{2} = (A'(c)c_{t}\boldsymbol{\eta}_{u}, \boldsymbol{\xi}_{ut}) + (A(c)\boldsymbol{\eta}_{ut}, \boldsymbol{\xi}_{ut})$$

$$\leq \epsilon \|\boldsymbol{\xi}_{ut}\|^{2} + C\|\boldsymbol{\eta}_{u}\|^{2} + C\|\boldsymbol{\eta}_{ut}\|^{2}$$

$$\leq \epsilon \|\boldsymbol{\xi}_{ut}\|^{2} + Ch^{2(k+1)}.$$
(2.4.74)

Now, we estimate  $\tilde{\tilde{T}}_3$ ,

$$\begin{split} \tilde{T}_{3} &= ((A(c) - A(c_{h}))_{t} \mathbf{u}_{h}, \boldsymbol{\xi}_{ut}) - ((A(c) - A(c_{h}))(\mathbf{u} - \mathbf{u}_{h})_{t}, \boldsymbol{\xi}_{ut}) \\ &+ ((A(c) - A(c_{h}))u_{t}, \boldsymbol{\xi}_{ut}) \\ &= ((A'(c) - A'(c_{h}))c_{t} \mathbf{u}_{h}, \boldsymbol{\xi}_{ut}) + (A'(c_{h})(c - c_{h})_{t} \mathbf{u}_{h}, \boldsymbol{\xi}_{ut}) \\ &+ ((A(c) - A(c_{h}))\boldsymbol{\xi}_{ut}, \boldsymbol{\xi}_{ut}) - ((A(c) - A(c_{h}))\boldsymbol{\eta}_{ut}, \boldsymbol{\xi}_{ut}) + (A'_{c}(c - c_{h})\mathbf{u}_{t}, \boldsymbol{\xi}_{ut}) \\ &\leq C \|c - c_{h}\|\|\boldsymbol{\xi}_{ut}\| + C \|(c - c_{h})_{t}\|\|\boldsymbol{\xi}_{ut}\| \\ &+ \|A^{\frac{1}{2}}(c)\boldsymbol{\xi}_{ut}\|^{2} - \|A^{\frac{1}{2}}(c_{h})\boldsymbol{\xi}_{ut}\|^{2} + C \|\boldsymbol{\eta}_{ut}\|\|\boldsymbol{\xi}_{ut}\| \\ &\leq \|A^{\frac{1}{2}}(c)\boldsymbol{\xi}_{ut}\|^{2} - \|A^{\frac{1}{2}}(c_{h})\boldsymbol{\xi}_{ut}\|^{2} + C \|\boldsymbol{\xi}_{c}\|^{2} \\ &+ C \|\boldsymbol{\xi}_{ct}\|^{2} + \epsilon \|\boldsymbol{\xi}_{ut}\|^{2} + C h^{2(k+1)}, \end{split}$$
(2.4.75)

where in the third step we use Schwarz inequality and hypotheses 1,2, and the last step requires Lemma 2.4.2.

$$\tilde{\tilde{T}}_{4} = \frac{1}{2} \Big( d'(c_{h})(c - c_{h})_{t} \xi_{p_{t}}, \xi_{p_{t}} \Big) - \frac{1}{2} \Big( d'(c_{h})c_{t} \xi_{p_{t}}, \xi_{p_{t}} \Big) \\ \leq C \|\xi_{p_{t}}\|_{\infty} \|(c - c_{h})_{t}\| \|\xi_{p_{t}}\| + C \|\xi_{p_{t}}\|^{2} \\ \leq C \|\xi_{c_{t}}\|^{2} + C \|\xi_{p_{t}}\| \|^{2} + C h^{2(k+1)}, \qquad (2.4.76)$$

where in the second step we use Schwarz inequality and hypothesis 1, and the last step requires Lemma 3.2. C depends on  $||c_t||_{\infty}$ . Similarly, we can estimate  $\tilde{T}_5$  and  $\tilde{T}_6$ 

$$\tilde{\tilde{T}}_{5} = -(d'(c_{h})(c-c_{h})_{t}\eta_{p_{t}},\xi_{p_{t}}) + (d'(c_{h})c_{t}\eta_{p_{t}},\xi_{p_{t}}) + (d(c_{h})\eta_{p_{tt}},\xi_{p_{t}})$$

$$\leq C \|\xi_{p_{t}}\|_{\infty} \|(c-c_{h})_{t}\| \|\eta_{p_{t}}\| + C \|\eta_{p_{t}}\| \|\xi_{p_{t}}\| + C \|\eta_{p_{tt}}\| \|\xi_{p_{t}}\|$$

$$\leq C \|\xi_{c_{t}}\|^{2} + C \|\xi_{p_{t}}\|^{2} + Ch^{2(k+1)}, \qquad (2.4.77)$$

$$\tilde{\tilde{T}}_{6} = \left( (d'(c) - d'(c_{h}))c_{t}p_{t}, \xi_{p_{t}} \right) + \left( d'(c_{h})(c - c_{h})_{t}p_{t}, \xi_{p_{t}} \right) + \left( (d(c) - d(c_{h}))p_{tt}, \xi_{p_{t}} \right) \\ \leq C \|c - c_{h}\|^{2} + C \|(c - c_{h})_{t}\|^{2} + C \|\xi_{p_{t}}\|^{2} \\ \leq C \|\xi_{c}\|^{2} + C \|\xi_{ct}\|^{2} + C \|\xi_{p_{t}}\|^{2} + C h^{2(k+1)}.$$

$$(2.4.78)$$

For  $\tilde{\tilde{T}}_7$ , we use Lemma 2.4.3 to obtain

$$\tilde{\tilde{T}}_{7} \le Ch^{k+1} \|p\|_{k+2} \|\boldsymbol{\xi}_{u_{t}}\|.$$
(2.4.79)

Substituting (2.4.73)-(2.4.79) into (2.4.72), we have the estimate

$$\|A^{\frac{1}{2}}(c_{h})\boldsymbol{\xi}_{ut}\|^{2} + \frac{1}{2}\frac{d}{dt}\|d^{\frac{1}{2}}(c_{h})\boldsymbol{\xi}_{p_{t}}\|^{2}$$
  

$$\leq Ch^{2(k+1)} + C(\|\boldsymbol{\xi}_{u}\|^{2} + C\|\boldsymbol{\xi}_{c}\|^{2} + \|\boldsymbol{\xi}_{p_{t}}\|^{2} + \|\boldsymbol{\xi}_{ct}\|^{2}) + \epsilon\|\boldsymbol{\xi}_{ut}\|^{2}. \quad (2.4.80)$$

Integrating (2.4.80) with respect to time between 0 and t and using the hypothesis 1, we obtain the fourth energy Inequality

$$\int_{0}^{t} \|\boldsymbol{\xi}_{ut}\|^{2} dt + \|\boldsymbol{\xi}_{p_{t}}\|^{2}$$

$$\leq Ch^{2(k+1)} + C \int_{0}^{t} (\|\boldsymbol{\xi}_{u}\|^{2} + \|\boldsymbol{\xi}_{c}\|^{2} + \|\boldsymbol{\xi}_{p_{t}}\|^{2} + \|\boldsymbol{\xi}_{ct}\|^{2}) dt. \qquad (2.4.81)$$

#### 2.4.8 Proof of Theorem 2.3.2

Now we are ready to combine the four energy inequalities and finish the proof of Theorem 2.3.2. Firstly, combing (2.4.51) with (2.4.70), we obtain

$$\int_{0}^{t} \|\xi_{ct}\|^{2} dt + \|\boldsymbol{\xi}_{s}\|^{2} + \|\xi_{c}\|^{2}$$

$$\leq C \int_{0}^{t} (\|\boldsymbol{\xi}_{u}\|^{2} + \|\boldsymbol{\xi}_{s}\|^{2} + \|\xi_{c}\|^{2} + \|\xi_{p_{t}}\|^{2}) dt + \epsilon \int_{0}^{t} \|\boldsymbol{\xi}_{ut}\|^{2} dt + Ch^{2(k+2)}.4.82)$$

Secondly, combing (2.4.81) with (2.4.82), we obtain

$$\int_{0}^{t} \|\boldsymbol{\xi}_{ut}\|^{2} dt + \|\boldsymbol{\xi}_{p_{t}}\|^{2}$$

$$\leq C \int_{0}^{t} (\|\boldsymbol{\xi}_{u}\|^{2} + \|\boldsymbol{\xi}_{s}\|^{2} + \|\boldsymbol{\xi}_{c}\|^{2} + \|\boldsymbol{\xi}_{p_{t}}\|^{2}) dt + Ch^{2(k+1)}. \quad (2.4.83)$$

Then, adding (2.4.60), (2.4.70), (2.4.82) and (2.4.83), we obtain

$$\begin{aligned} \|\boldsymbol{\xi}_{u}\|^{2} + \|\boldsymbol{\xi}_{p}\|^{2} + \|\boldsymbol{\xi}_{p_{t}}\|^{2} + \|\boldsymbol{\xi}_{c}\|^{2} + \|\boldsymbol{\xi}_{s}\|^{2} + \int_{0}^{t} (\|\boldsymbol{\xi}_{u_{t}}\|^{2} + \|\boldsymbol{\xi}_{c_{t}}\|^{2}) dt \\ &\leq Ch^{2(k+1)} + C \int_{0}^{t} (\|\boldsymbol{\xi}_{u}\|^{2} + \|\boldsymbol{\xi}_{p}\|^{2} + \|\boldsymbol{\xi}_{p_{t}}\|^{2} + \|\boldsymbol{\xi}_{c}\|^{2} + \|\boldsymbol{\xi}_{s}\|^{2}) dt \\ &+ \epsilon \int_{0}^{t} (\|\boldsymbol{\xi}_{u_{t}}\|^{2} + \|\boldsymbol{\xi}_{c_{t}}\|^{2}) dt. \end{aligned}$$

$$(2.4.84)$$

Employing Gronwall's lemma, we obtain

$$\|\boldsymbol{\xi}_{u}\|^{2} + \|\boldsymbol{\xi}_{p}\|^{2} + \|\boldsymbol{\xi}_{p_{t}}\|^{2} + \|\boldsymbol{\xi}_{c}\|^{2} + \|\boldsymbol{\xi}_{s}\|^{2} + \int_{0}^{t} (\|\boldsymbol{\xi}_{ut}\|^{2} + \|\boldsymbol{\xi}_{ct}\|^{2}) dt \le Ch^{2(k+1)}. \quad (2.4.85)$$

Finally, by using the standard approximation result, we obtain (2.3.16). To complete the proof, let us verify the a priori assumption (2.4.10). For  $k \ge 1$ , we can consider h small enough so that  $Ch^{k+1} < \frac{1}{2}h$ , where C is the constant determined by the final time T. Then if  $t^* = \inf\{t : ||c - c_h|| + ||\mathbf{u} -$   $\mathbf{u}_h \| + \| p_t - p_{ht} \| \ge h \}$ , we should have  $\| c - c_h \| + \| \mathbf{u} - \mathbf{u}_h \| + \| p_t - p_{ht} \| = h$ by continuity in time at  $t = t^*$ . However, if  $t^* < T$ , theorem 2.3.2 implies that  $\| c - c_h \| + \| \mathbf{u} - \mathbf{u}_h \| + \| p_t - p_{ht} \| \le Ch^{k+1}$  for  $t \le t^*$ , in particular  $h = \| (c - c_h)(t^*) \| + \| (\mathbf{u} - \mathbf{u}_h)(t^*) \| + \| (p_t - p_{ht})(t^*) \| \le Ch^{k+1} < \frac{1}{2}h$ , which is a contradiction. Therefore, there always holds  $t^* \ge T$ , and thus the a priori assumption (2.4.10) is justified.

#### 2.5 Numerical example

In this section we provide numerical examples to illustrate the accuracy and capability of the method. Time discretization is given as the third order strong-stability-preserving Runge-Kutta method [54]. We take the time step to be sufficiently small such that the error in time is negligible compared to spatial error. In the scheme, the numerical flux in the convection term is taken as  $\widehat{\mathbf{u}_h c_h} = \frac{1}{2} (\mathbf{u}_h^+ c_h^+ + \mathbf{u}_h^- c_h^-)$ . Moreover, other parameters are taken as follows

- The solution domain  $\Omega = [0, 1] \times [0, 1]$ , T = 0.01,  $\Delta t = r * h^2$ , here r denotes the grid ratio and r depends on the polynomial degree.
- We take  $\phi(x, y) = 1$ ,  $\kappa(x, y) = 1$ ,  $\mu(c) = 1$ , for simplicity.

**Example 2.5.1.** We first consider the problem with the constant matrix  $\mathbf{D}(\mathbf{u}) = \alpha \mathbf{I}$ , where  $\alpha$  is a constant, in addition, we take the initial and boundary condition  $c_0 = \sin(2\pi(x+y)), p_0 = -2\pi(x^2+y^2), c(0,t) = c(2\pi,t), and the parameters <math>b(c) = 0, d(c) = 1$  and the source term

$$f = 2\pi \cos(2\pi(x+y+t))(4\pi(x+y+t)+1) + 8\alpha\pi^2 \sin(2\pi(x+y+t)) - 2\pi,$$

the exact solution is

$$c = \sin(2\pi(x+y+t)), \mathbf{u} = (4\pi x + 2\pi t, 4\pi y + 2\pi t),$$

The  $L^2$  error and the numerical orders of accuracy at time t = 0.01 with uniform meshes are contained in Tables 2.1 and 2.2. We can see that the method with  $Q^k$  elements gives (k + 1)-th order of accuracy in  $L^2$  norm.

 $Q^3/r = 0.001$  $Q^1/r = 0.01$  $Q^2/r = 0.01$ N $L^2$  error  $L^2$  error  $L^2$  error order order order 102.3021e-028.0016e-04 2.0744e-04\_ \_ 205.8006e-031.999.9746e-053.001.3097e-053.99401.4512e-032.001.2417e-053.018.1846e-074.0080 3.6279e-042.001.5521e-063.005.1097e-084.001609.0695e-052.001.9400e-073.003.1875e-094.00

Table 2.1: The numerical results for c with  $\alpha = 1$ 

**Example 2.5.2.** Next we consider the problem with matrix  $\mathbf{D}(\mathbf{u}) = \mathbf{u} \otimes \mathbf{u} + \mathbf{I}$ , in addition, we take the initial and boundary condition  $c_0 = \sin(2\pi(x+y))$ ,  $p_0 = -2\pi(x^2 + y^2)$ ,  $c(0,t) = c(2\pi,t)$ , and the parameters b(c) = 0, d(c) = 1 and the source term

$$f(x, y, t) = 2\pi \cos(2\pi(x+y+t))(4\pi(x+y+t))(1-12\pi^2) - 2\pi + 4\pi^2(16\pi^2(x+y+t)^2+2)\sin(2\pi(x+y+t)),$$

	$Q^1/r = 0.01$		$Q^2/r = 0.01$		$Q^3/r = 0.001$	
IN .	$L^2$ error	order	$L^2$ error	order	$L^2$ error	order
10	2.3021e-02	_	7.9917e-04	_	2.0744e-04	_
20	5.8006e-03	1.99	9.9612e-05	3.00	1.3097e-05	3.99
40	1.4501e-03	2.00	1.2450e-05	3.00	8.1796e-07	4.00
80	3.6247e-04	2.00	1.5524e-06	3.00	5.1100e-08	4.00
160	9.0603e-05	2.00	1.9355e-07	3.00	3.1875e-09	4.00

Table 2.2: The numerical results for c with  $\alpha = 0.01$ 

the exact solution is

$$c = \sin(2\pi(x+y+t)), \mathbf{u} = (4\pi x + 2\pi t, 4\pi y + 2\pi t),$$

The  $L^2$  error and the numerical orders of accuracy at time t = 0.01 with uniform meshes is contained in Tables 2.3. We can see that the method with  $Q^k$ elements gives (k + 1)-th order of accuracy in  $L^2$  norm.

Example 2.5.3. We choose the initial condition as

$$c_0 = \frac{1}{2}(1 + \cos(2\pi x)\cos(2\pi y)), \qquad p_0 = \cos(2\pi x)\cos(2\pi y) - 1.$$

Other parameters are taken as

$$q(x, y, 0) = 0, m_1 = 0.35, m_2 = 1, \phi(x) = 1, \mathbf{D}(\mathbf{u}) = \begin{pmatrix} |u| & 0 \\ 0 & |u| \end{pmatrix}$$



Figure 2.1: Numerical approximations of c at t = 0.1 with Nx = Ny = 40 in Example 2.5.3.



Figure 2.2: Numerical approximations of c at t = 0.1 with Nx = Ny = 40 in Example 2.5.4.

Table 2.3: The numerical results for c

	λī	Q1/r = 0.01		Q2/r = 0.01		Q3/r = 0.001	
	IN	$L^2$ error	order	$L^2$ error	order	$L^2$ error	order
	10	2.3022e-02	_	7.9948e-04	_	2.0756e-04	_
	20	5.8006e-03	1.99	9.9643e-05	3.00	1.3104e-05	3.99
	40	1.4492e-03	2.00	1.2393e-05	3.01	8.2105e-07	4.00
	80	3.6223e-04	2.00	1.5477e-06	3.00	5.1348e-08	4.00
	160	9.0551e-05	2.00	1.9308e-07	3.00	3.2097e-09	4.00

We choose  $\Delta t = 0.01 \min{\{\Delta x^2, \Delta y^2\}}$  with final time T = 0.1, and the numerical approximation of c is given in Figure 2.1.

Example 2.5.4. We change the initial condition in Example 2.5.3 to

$$c_0 = \begin{cases} 0.001, & (x - 0.5)^2 + (y - 0.5)^2 < 0.09, \\ 0, & otherwise, \end{cases} \qquad p_0 = \sin(\pi x)\sin(\pi y).$$

Other parameters are taken as

$$q(x, y, 0) = 0, m_1 = 1, m_2 = 1, \phi(x) = 1, \mathbf{D}(\mathbf{u}) = \mathbf{I}$$

and the numerical approximation of c is given in Figure 2.2.

# 2.6 Concluding remarks

In this paper, the conservative LDG method for both flow and transport equations is introduced for the coupled system of compressible miscible displacement problem. The optimal order of error estimates hold not only for the solution itself but also for the auxiliary variables. Special projections and a priori assumption help to eliminate the jump terms at the cell interfaces which arise from the discontinuity nature of the numerical method, the nonlinearity and coupling of the model.

#### Acknowledgments

This work is supported by National Natural Science Foundation of China Grants 11571367 and 11601536, and the Fundamental Research Funds for the Central Universities and Michigan Technological University Research Excellence Fund Scholarship and Creativity Grant 1605052.

# 2.6 Appendix: Proof of Lemma 2.4.5

Recall that we have chosen the initial condition  $c_h^0 = P^+c_0$ ,  $\mathbf{u}_h^0 = \mathbf{\Pi}^-\mathbf{u}_0$ , where  $\mathbf{u}_0 = -a(c_0)\nabla p_0$ , and  $\hat{p}_h = p_h^+$ ,  $\widehat{\mathbf{u}}_h = \mathbf{u}_h^-$ ,  $\hat{z}_h = \mathbf{z}_h^-$ ,  $\hat{c}_h = c_h^+$ . For simplicity, we will drop the 0 in the superscripts and subscripts in this section. It is clear that (2.4.5) and (2.4.6) hold. Taking the test function  $\zeta = \xi_{p_t}$  and summing over K in (2.4.16), we have

$$\left(d(c)\xi_{p_t},\xi_{p_t}\right) = \left(d(c)\eta_{p_t},\xi_{p_t}\right) + \left(p_{h_t}(d(c) - d(c_h)),\xi_{p_t})\right), \quad (2.6.1)$$

where we have used  $\mathbf{u}_h = \mathbf{\Pi}^- \mathbf{u}, \widehat{\mathbf{u}_h} = \mathbf{u}_h^-$  and the property of the projection (2.4.3). Using the Schwartz inequality, we can get

$$\|d^{\frac{1}{2}}(c)\xi_{p_t}\|^2 \le C \|\eta_{p_t}\| \|\xi_{p_t}\| + C \|c - c_h\| \|\xi_{p_t}\|, \qquad (2.6.2)$$

By Lemma 2.4.2 and (2.4.5), we easily prove

$$\|\xi_{p_t}\| \le Ch^{k+1}.$$
 (2.6.3)

Similarly, taking the test function  $\mathbf{w} = \boldsymbol{\xi}_s$  and summing over K in (2.4.13), we have

$$(\boldsymbol{\xi}_s, \boldsymbol{\xi}_s) = (\boldsymbol{\eta}_s, \boldsymbol{\xi}_s) - \mathcal{D}(\eta_c, \boldsymbol{\xi}_s), \qquad (2.6.4)$$

where we have used  $c_h = P^+c$ . Using the Schwartz inequality and the Lemma 2.4.3, we can get

$$\|\boldsymbol{\xi}_{s}\|^{2} \leq \|\boldsymbol{\xi}_{s}\|\|\boldsymbol{\eta}_{s}\| + Ch^{k+1}\|c\|_{k+2}\|\boldsymbol{\xi}_{s}\|.$$
(2.6.5)

By Lemma 2.4.2, we easily prove

$$\|\boldsymbol{\xi}_s\| \le Ch^{k+1},\tag{2.6.6}$$

By the standard approximation results, (2.4.7) and (2.4.8) hold. At last we estimate  $p - p_h$ , following the technique in [41]. By (2.3.10) the initial data  $p_h$  is the solution of the following equations

$$(A(c_h)\mathbf{u}_h,\boldsymbol{\theta})_K - (p_h,\nabla\cdot\boldsymbol{\theta})_K + \langle \widehat{p}_h,\boldsymbol{\theta}\cdot\boldsymbol{\nu}_K \rangle_{\partial_K} = 0, \qquad (2.6.7)$$

and also satisfies

$$(p - p_h, 1) = 0. (2.6.8)$$

From (2.4.15), we have

$$(A(c)\mathbf{u} - A(c_h)\mathbf{u}_h, \boldsymbol{\theta})_K - (p - p_h, \nabla \cdot \boldsymbol{\theta})_K + \langle p - \widehat{p}_h, \boldsymbol{\theta} \cdot \boldsymbol{\nu}_K \rangle_{\partial_K} = (2.6.9)$$

We use  $\mathbf{u}_h$  to find a well-defined  $p_h$ , and we only need to prove the uniqueness. If there are two solutions  $p_1$  and  $p_2$  satisfying (2.6.7) and (2.6.8), then we can easily get

$$(p_1 - p_2, \nabla \cdot \boldsymbol{\theta})_K - \langle \hat{p}_1 - \hat{p}_2, \boldsymbol{\theta} \cdot \boldsymbol{\nu}_K \rangle_{\partial_K} = 0, \qquad (2.6.10)$$

$$(p_1 - p_2, 1) = 0. (2.6.11)$$

We consider the elliptic linear problem

$$-\boldsymbol{\zeta}^* = \nabla \boldsymbol{\xi}^*, \text{in} \quad \Omega, \tag{2.6.12}$$

$$\eta^* = \nabla \cdot \boldsymbol{\zeta}^*, \text{in} \quad \Omega, \tag{2.6.13}$$

subject to periodic boundary conditions. To make the problem well-defined, we assume that the average of  $\xi^*$  on  $\Omega$  is a given constant and that of  $\eta^*$  is zero. We have the elliptic regularity result

$$\|\boldsymbol{\zeta}^*\|_{H^1(\Omega)} + \|\boldsymbol{\xi}^*\|_{H^2(\Omega)} \le C \|\boldsymbol{\eta}^*\|.$$
(2.6.14)

Taking  $\eta^* = p_1 - p_2$  and  $\widehat{p_i} = p_i^+, i = 1, 2$ , we get

$$(p_1 - p_2, p_1 - p_2)_K$$

$$= (p_1 - p_2, \nabla \cdot \boldsymbol{\zeta}^*)_K$$

$$= (p_1 - p_2, \nabla \cdot (\boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*))_K + (p_1 - p_2, \nabla \cdot \Pi \boldsymbol{\zeta}^*)_K$$

$$= (p_1 - p_2, \nabla \cdot (\boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*))_K - \langle \hat{p}_1 - \hat{p}_2, (\boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*) \cdot \boldsymbol{\nu}_K \rangle_{\partial_K}$$

$$+ \langle \hat{p}_1 - \hat{p}_2, \boldsymbol{\zeta}^* \cdot \boldsymbol{\nu}_K \rangle_{\partial_K}$$

$$= -(\nabla (p_1 - p_2), \boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*)_K + \langle p_1 - p_2, (\boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*) \cdot \boldsymbol{\nu}_K \rangle_{\partial_K}$$

$$- \langle \hat{p}_1 - \hat{p}_2, (\boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*) \cdot \boldsymbol{\nu}_K \rangle_{\partial_K} + \langle \hat{p}_1 - \hat{p}_2, \boldsymbol{\zeta}^* \cdot \boldsymbol{\nu}_K \rangle_{\partial_K}$$
(2.6.15)

where the third step follows from (2.6.10) and the last equality is based on integration by parts. We take  $\Pi \boldsymbol{\zeta}^* = \Pi^- \boldsymbol{\zeta}^*$  and sum over K. By the continuity of  $\boldsymbol{\zeta}^*$  and the definition of the projection  $\Pi^-$ , we obtain

$$(p_1 - p_2, p_1 - p_2) = 0 (2.6.16)$$

Then we get  $p_1 = p_2$ . We have proved that  $p_h$  is well-defined. In the following, we estimate  $||p - p_h||$ . We use the same technique above and take  $\eta^* = p - p_h$  to obtain

$$(p - p_h, p - p_h)_K$$

$$= (p - p_h, \nabla \cdot \boldsymbol{\zeta}^*)_K$$

$$= (p - p_h, \nabla \cdot (\boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*))_K + (p - p_h, \nabla \cdot \Pi \boldsymbol{\zeta}^*)_K$$

$$= (p - p_h, \nabla \cdot (\boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*))_K - (A(c)\mathbf{u} - A(c_h)\mathbf{u}_h, \boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*)_K$$

$$- \langle p - \hat{p}_h, (\boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*) \cdot \boldsymbol{\nu}_K \rangle_{\partial K} + (A(c)\mathbf{u} - A(c_h)\mathbf{u}_h, \boldsymbol{\zeta}^*)_K$$

$$+ \langle p - \hat{p}_h, \boldsymbol{\zeta}^* \cdot \boldsymbol{\nu}_K \rangle_{\partial K}$$

$$= -(\nabla(p - p_h), \boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*)_K + \langle p - p_h, (\boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*) \cdot \boldsymbol{\nu}_K \rangle_{\partial K}$$

$$-(A(c)\mathbf{u} - A(c_h)\mathbf{u}_h, \boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*)_K - \langle p - \hat{p}_h, (\boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*) \cdot \boldsymbol{\nu}_K \rangle_{\partial K}$$

$$+ (A(c)\mathbf{u} - A(c_h)\mathbf{u}_h, \boldsymbol{\zeta}^*)_K + \langle p - \hat{p}_h, \boldsymbol{\zeta}^* \cdot \boldsymbol{\nu}_K \rangle_{\partial K}$$

$$= -(\nabla(p - p_h), \boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*)_K + \langle \hat{p}_h - p_h, (\boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*) \cdot \boldsymbol{\nu}_K \rangle_{\partial K}$$

$$-(A(c)\mathbf{u} - A(c_h)\mathbf{u}_h, \boldsymbol{\zeta}^* - \Pi \boldsymbol{\zeta}^*)_K + (A(c)\mathbf{u} - A(c_h)\mathbf{u}_h, \boldsymbol{\zeta}^*)_K$$

$$+ \langle p - \hat{p}_h, \boldsymbol{\zeta}^* \cdot \boldsymbol{\nu}_K \rangle_{\partial K}$$

$$(2.6.17)$$

where the third one follows from (2.6.9) and the fourth equality is based on the integrate by parts. Recalling that  $\hat{p}_h = p_h^+$ , we take  $\Pi \boldsymbol{\zeta}^* = \Pi^- \boldsymbol{\zeta}^*$  and sum over

K. By the continuity of  $\pmb{\zeta}^*$  and the definition of the projection  $\pmb{\Pi}^-,$  we obtain

$$\|p - p_{h}\|^{2} = -(\nabla \eta_{p}, \boldsymbol{\zeta}^{*} - \Pi \boldsymbol{\zeta}^{*}) - (A(c)\mathbf{u} - A(c_{h})\mathbf{u}_{h}, \boldsymbol{\zeta}^{*} - \Pi \boldsymbol{\zeta}^{*}) + (A(c)\mathbf{u} - A(c_{h})\mathbf{u}_{h}, \boldsymbol{\zeta}^{*}) = -(\nabla \eta_{p}, \boldsymbol{\zeta}^{*} - \Pi \boldsymbol{\zeta}^{*}) - (A(c)(\mathbf{u} - \mathbf{u}_{h}), \boldsymbol{\zeta}^{*} - \Pi \boldsymbol{\zeta}^{*}) - ((A(c) - A(c_{h}))\mathbf{u}_{h}, \boldsymbol{\zeta}^{*} - \Pi \boldsymbol{\zeta}^{*}) + (A(c_{0})(\mathbf{u} - \mathbf{u}_{h}), \boldsymbol{\zeta}^{*}) + ((A(c) - A(c_{h}))\mathbf{u}_{h}, \boldsymbol{\zeta}^{*}) \leq Ch^{k+1} \|\boldsymbol{\zeta}^{*}\|_{H^{1}(\Omega)} + Ch^{k+2} \|\boldsymbol{\zeta}^{*}\|_{H^{1}(\Omega)} + Ch^{k+1} \|\boldsymbol{\zeta}^{*}\| \leq Ch^{k+1} \|\boldsymbol{\zeta}^{*}\|_{H^{1}(\Omega)}$$
(2.6.18)

which further implies

$$\|p - p_h\| \le Ch^{k+1}.$$
(2.6.19)

# Chapter 3

# High-order bound-preserving discontinuous Galerkin methods for compressible miscible displacements in porous media on triangular meshes<sup>1</sup>

#### Abstract

In this paper, we develop high-order bound-preserving (BP) discontinuous Galerkin (DG) methods for the coupled system of compressible miscible displacements on

<sup>&</sup>lt;sup>1</sup>This chapter has been published as an article in Journal of Computational Physics. Citation: N. Chuenjarern, Z. Xu, Y. Yang, Journal of Computational Physics 378 (2019),110-128. https://doi.org/10.1016/j.jcp.2018.11.003

triangular meshes. We consider the problem with multi-component fluid mixture and the (volumetric) concentration of the *j*th component,  $c_j$ , should be between 0 and 1. There are three main difficulties. Firstly,  $c_j$  does not satisfy a maximum-principle. Therefore, the numerical techniques introduced in (X. Zhang and C.-W. Shu, Journal of Computational Physics, 229 (2010), 3091-3120) cannot be applied directly. The main idea is to apply the positivity-preserving techniques to all  $c'_j s$  and enforce  $\sum_j c_j = 1$  simultaneously to obtain physically relevant approximations. By doing so, we have to treat the time derivative of the pressure dp/dt as a source in the concentration equation and choose suitable fluxes in the pressure and concentration equations. Secondly, it is not easy to construct first-order numerical fluxes for interior penalty DG methods on triangular meshes. One of the key points in the high-order BP technique applied in this paper is the combination of high-order and lower-order numerical fluxes. We will construct second-order BP schemes and use the second-order numerical fluxes as the lower-order one. Finally, the classical slope limiter cannot be applied to  $c_i$ . To construct the BP technique, we will not approximate  $c_i$  directly. Therefore, a new limiter will be introduced. Numerical experiments will be given to demonstrate the high-order accuracy and good performance of the numerical technique.

**Key Words:** compressible miscible displacements, bound-preserving, high-order, discontinuous Galerkin method, triangular meshes, multi-component fluid, flux limiter

# 3.1 Introduction

In this paper, we are interested in constructing high-order bound-preserving discontinuous Galerkin (DG) schemes for compressible miscible displacements in porous media on triangular meshes. We consider the fluid mixture with N components and the governing equations over the computational domain  $\Omega = [0, 1] \times [0, 1]$  read

$$d(\mathbf{c})\frac{\partial p}{\partial t} + \nabla \cdot \mathbf{u} = d(\mathbf{c})\frac{\partial p}{\partial t} - \nabla \cdot \left(\frac{\kappa(x,y)}{\mu(\mathbf{c})}\nabla p\right) = q, \ (x,y) \in \Omega, \ 0 < t \le T, \ (3.1.1)$$
$$\phi\frac{\partial c_j}{\partial t} + \nabla(\mathbf{u} \cdot c_j) - \nabla \cdot (\mathbf{D}\nabla c_j) = \tilde{c}_j q - \phi c_j z_j p_t, \ (x,y) \in \Omega, \ 0 < t \le T, \ j = 1, \cdots, N-1$$
$$(3.1.2)$$

where the dependent variables are the pressure in fluid mixture denoted by p, the Darcy velocity of the mixture (volume flowing across a unit across-section per unit time) denoted by  $\mathbf{u}$  and the concentration of interested species measured in amount of species per unit volume denoted by  $\mathbf{c} = (c_1, \dots, c_N)^T$ , with  $c_j$ being the concentration of the *j*th component.  $\phi$  and  $\kappa$  are the porosity and permeability of the rock, respectively.  $\mu$  refers to the concentration-dependent viscosity. q is the external volumetric flow rate, and  $\tilde{c}_j$  is the concentration of the fluid in the external flow.  $\tilde{c}_j$  must be specified at points where injection (q > 0) takes place, and is assumed to be equal to  $c_j$  at production points (q < 0). The diffusion coefficient  $\mathbf{D}$  is symmetric and arises from two aspects: molecular diffusion, which is rather small for field-scale problems, and dispersion, which is velocity-dependent, in the petroleum engineering literature. Its form is

$$\mathbf{D} = \phi(x, y)(d_{\text{mol}}\mathbf{I} + d_{\text{long}}|\mathbf{u}|\mathbf{E} + d_{\text{tran}}|\mathbf{u}|\mathbf{E}^{\perp}), \qquad (3.1.3)$$
where **E**, a  $2 \times 2$  matrix, represents the orthogonal projection along the velocity vector given as

$$\mathbf{E} = (e_{ij}(\mathbf{u})) = \left(\frac{u_i u_j}{|\mathbf{u}^2|}\right), \quad \mathbf{u} = (u_1, u_2),$$

and  $\mathbf{E}^{\perp} = \mathbf{I} - \mathbf{E}$  is the orthogonal complement. The diffusion coefficient  $d_{\text{long}}$ measures the dispersion in the direction of the flow and  $d_{\text{tran}}$  shows that transverse to the flow. To ensure the stability of the scheme,  $\mathbf{D}$  is assumed to be strictly positive definite in almost all of the previous works. In this paper, we assume  $\mathbf{D}$  to be positive semidefinite. Moreover, the pressure is uniquely determined up to a constant, thus we assume  $\int_{\Omega} p \, dx dy = 0$  at t = 0. However, this assumption is not essential. Other coefficients can be stated as follows:

$$c_N = 1 - \sum_{j=1}^{N-1} c_j, \quad d(\mathbf{c}) = \phi \sum_{j=1}^N z_j c_j,$$

where  $z_j$  is the compressibility factor of the *j*th component of the fluid mixture. In this paper, we consider homogeneous Neumann boundary conditions

$$\mathbf{u} \cdot \mathbf{n} = 0, \quad (\mathbf{D}\nabla c - c\mathbf{u}) \cdot \mathbf{n} = 0,$$

where **n** is the unit outer normal of the boundary  $\partial \Omega$ . Moreover, the initial solutions are given as

$$c_j(x, y, 0) = c_{j_0}(x, y), \quad p(x, y, 0) = p_0(x, y), \quad (x, y) \in \Omega.$$

The miscible displacements in porous media were first presented in [24, 25], where mixed finite element methods were applied. Later, the compressible problem was studied in [23] and the optimal order estimates in  $L^2$ -norm and almost optimal order estimates in  $L^{\infty}$ -norm were given in [11]. Subsequently, many new numerical methods were introduced, such as the finite difference method [81, 82, 83], characteristic finite element method [48], splitting positive definite mixed element method [70] and H1-Galerkin mixed method [7]. Besides the above, in [59], an accurate and efficient simulator was developed for problems with wells. Later, the authors introduced an Eulerian-Lagrangian localized adjoint method to solve the transport partial differential equation for concentration, while a mixed finite element method to solve the pressure equation [58]. Recently, DG methods have been popular to solve compressible miscible displacements in porous media [21, 22, 71, 72, 37, 73, 77]. Some special numerical techniques were introduced to control the jumps of numerical approximations as well as the nonlinearality of the convection term. Besides the above, there were also significant works discussing the DG methods for incompressible miscible displacements, see e.g. [1, 38, 44, 52, 55, 56, 63] and for general porous media flow, see e.g. [3, 30, 29, 57] and the references therein. However, no previous works above focused on the bound-preserving techniques. In many numerical simulations, the approximations of  $c_i$  can be placed out of the interval [0, 1]. Especially for problems with large gradients, the value of  $d(\mathbf{c})$  might be negative, leading to ill-posedness of the problem, and the numerical approximations will blow up. We will use numerical experiments to demonstrate this point in Section 3.5. In [36], we have introduced second-order bound-preserving DG methods on rectangular meshes for two-component miscible displacements in porous media. In this paper, we will extend the idea to multi-component miscible displacements and construct high-order bound-preserving techniques on triangular meshes. Moreover, the idea can be extended to incompressible flows with some minor changes.

The DG method gained even greater popularity for good stability, high-order accuracy, and flexibility on h-p adaptivity and on complex geometry. In 2010, the genuinely maximum-principle-satisfying high-order DG and finite volume schemes were constructed in [85] by Zhang and Shu, the extension to unstructured meshes was given in [88]. After that, the idea was applied to many problems such as compressible Euler equations [86, 87], hyperbolic equations involving  $\delta$ -singularities [74, 75, 90], relativistic hydrodynamics [50] and shallow water equations [64], etc. The basic idea is to take the test function to be 1 in each cell to obtain an equation of the numerical cell average of the target variable, say r, and prove the cell average,  $\bar{r}$ , is within the desired bounds. Then we can apply a slope limiter to the numerical approximation and construct a new one

$$\tilde{r} = \bar{r} + \theta(r - \bar{r}), \quad \theta \in [0, 1].$$
(3.1.4)

If the problem has only one lower bound zero, the technique is also called positivity-preserving technique. Thanks to the limiter, the whole algorithm were proved to be  $L^1$ -stable [75, 50] for some complicated systems. Moreover, the technique does not rely on the trouble cell detector and the limiter keeps the high-order accuracy in regions with smooth solutions for scalar equations [85]. In case of convection-diffusion equations, the same idea was applied to construct genuinely second-order maximum-principle-satisfying DG method on unstructured meshes [89]. Recently, the flux limiter [39, 65, 68] and third-order maximum-principle-preserving direct DG method [8] were also introduced. However, it is not easy to apply the flux limiter to unstructured meshes since the lower order fluxes are not easy to construct, and the only work available is [12] in which the technique for hyperbolic equations was analyzed, and no previous works aimed to discuss convection-diffusion equations. In this paper, we will extend the ideas in [65, 85] and construct high-order bound-preserving DG methods for multi-component compressible miscible displacements. However, there are significant differences from previous techniques. First of all, most of the problems in [65, 85] satisfy maximum-principles while the concentration  $c_i$ in (3.1.2) does not. To solve this problem, we would like to apply the positivitypreserving technique to each  $c_j$  and enforce  $\sum_j c_j = 1$ . Secondly, the high-order positivity-preserving technique in this paper is based on the flux limiter [39, 65]. The basic idea is to combine higher order and lower order fluxes to construct a new one which yield positive numerical cell averages. However, for triangular meshes, first-order fluxes are not easy to construct. Therefore, we will consider the second-order flux as the lower order one. Finally, to obtain the equation satisfied by the cell averages, we need to numerically approximate  $r_j = \phi c_j$  instead of  $c_j$ . By doing so, the upper bound of  $r_j$  is not a constant and the limiter (3.1.4) may fail to work, since such a  $\theta$  may not exist (see the counterexample in [36]). Moreover, the limiter applied in [36] is not straightforward extendable to multi-component problems, since we cannot simply set the upper bound of  $c_j$  to be 1 if the fluid mixture contains more than two components. Therefore, a new bound-preserving limiter will be introduced. In summary, the whole algorithm can be separated into three parts. We first treat  $p_t$  as another source in (3.1.2) to

obtain the positivity of  $c_j$  by the flux limiter [39, 65]. Then we choose consistent fluxes (see Definition 3.2.1) with suitable parameter in the flux limiter in the concentration and pressure equations to obtain the positivity of  $1 - \sum_{j=1}^{N-1} c_j$ . More precisely, in our analysis, instead of solving p and  $c_j$ ,  $j = 1, \dots, N-1$ , we rewrite (3.1.1) and (3.1.2) into a system of  $c_j$ ,  $j = 1, \dots, N$  and enforce  $\sum_{i=j}^{N} c_j = 1$  by choosing consistent fluxes. Finally, we will introduce a new limiter to obtain physically relevant numerical approximations.

The paper is organized as follows: we first discuss the DG scheme in two dimension on triangular mesh in Section 3.2. In Section 3.3, we demonstrate the bound-preserving technique for second-order scheme. The high-order boundpreserving technique with flux limiter will be given in Section 3.4. In Section 3.5, some numerical experiments and results will be shown. We will end in Section 3.6 with concluding remarks.

#### 3.2 The DG scheme

In this section, we will construct the DG scheme for compressible miscible displacements in porous media. We first demonstrate the notations to be used throughout the paper. We consider triangular meshes and denote  $\Omega_h$  to be the set of cells. For any  $K \in \Omega_h$ , we denote the three edges of K to be  $e_K^i$ (i = 1, 2, 3), with corresponding lengths  $\ell_K^i$  (i = 1, 2, 3) and unit outer normal vectors  $\boldsymbol{\nu}_i$  (i = 1, 2, 3). We also denote the neighboring triangle along  $e_K^i$  as  $K_i$ . We use  $\Gamma$  for all the cell interfaces, and  $\Gamma_0 = \Gamma \setminus \partial \Omega$  for all the interior ones. For any  $e \in \Gamma$ , denote |e| to be the length of e. Let  $u^{\pm}$  denote the numerical solution on the edges, evaluated from K or  $K_i$ . The ' $\pm$ ' for each edge  $e_K^i$  is determined by the inner product of  $\boldsymbol{\nu}_i$  and a predetermined constant vector  $\boldsymbol{\nu}_0$  which is not parallel to any edge in the mesh: for each edge  $e_K^i$  in the cell K,

$$\boldsymbol{u}^- = \boldsymbol{u}_K, \quad \boldsymbol{u}^+ = \boldsymbol{u}_{K_i}, \quad \text{if } \boldsymbol{\nu}_0 \cdot \boldsymbol{\nu}_i > 0,$$
  
 $\boldsymbol{u}^+ = \boldsymbol{u}_K, \quad \boldsymbol{u}^- = \boldsymbol{u}_{K_i}, \quad \text{if } \boldsymbol{\nu}_0 \cdot \boldsymbol{\nu}_i < 0.$ 

Moreover, we define  $\mathbf{n}_e$  as the unit outer normal of each edge  $e \in \Gamma_0$  such that  $\mathbf{n}_e \cdot \mathbf{\nu}_0 > 0$  and define the jump and average of any function v at the cell interface e as

$$[v]_e = v_e^+ - v_e^-, \quad \{v\}_e = \frac{1}{2}(v_e^+ + v_e^-).$$

We also denote  $\partial \Omega_+ = \{e \in \partial \Omega : \mathbf{n} \cdot \boldsymbol{\nu}_0 > 0\}$ , where **n** is the unit outer normal of  $\partial \Omega$  and  $\partial \Omega_- = \partial \Omega \setminus \partial \Omega_+$ . The finite element space is chosen as

$$W_h = \{ z : z |_K \in P^k(K), \ \forall K \in \Omega_h \},\$$

where  $P^k(K)$  denotes polynomials of degree at most  $k \ge 1$  in K.

To construct the DG method, we first rewrite the system (3.1.1)-(3.1.2) into the following form

$$d(\mathbf{c})p_t + \nabla \cdot \mathbf{u} = q, \qquad (3.2.5)$$

$$a(\mathbf{c})\mathbf{u} = -\nabla p, \qquad (3.2.6)$$

$$(\phi c_j)_t + \nabla \cdot (\mathbf{u} c_j) - \nabla \cdot (\mathbf{D}(\mathbf{u}) \nabla c_j) = \tilde{c}_j q - \phi c_j z_j p_t, \quad j = 1, 2, \cdots, N-1,$$
(3.2.7)

where  $a(\mathbf{c}) = \frac{\mu(\mathbf{c})}{\kappa}$ .

Next, we would like to demonstrate the key points in this paper that are quite different from most of the previous works.

- 1. Approximate  $r_j = \phi c_j$  instead of  $c_j$ . We cannot simply take the test function to be 1 to obtain the cell average of  $c_j$ .
- 2. Treat  $p_t$  in (3.2.7) as a source to apply the positivity-preserving techniques.
- Apply flux limiters to the high-order scheme by combining the second- and high-order fluxes.
- 4. Suitably choose the parameters in the flux limiter to obtain consistent fluxes for (3.2.5) and (3.2.7) to make  $\bar{r_j} < \bar{\phi}$ , where  $\bar{r_j}$  and  $\bar{\phi}$  are the cell averages of  $r_j$  and  $\phi$ , respectively.
- 5. Take the  $L^2$ -projection of  $\phi$  into  $W_h$ , denoted as  $\Phi$ , and use which as the new approximation of the porosity.
- 6. Construct a new limiter to maintain the cell average  $\bar{r}_j$  and modify the numerical approximations of  $r_j$  such that  $0 < r_j < \Phi$ , which further yields  $c_j = P_k \left\{ \frac{r_j}{\Phi} \right\} \in [0, 1]$ , where  $P_k$  is the  $L^2$ -projection projected into  $W_h$  when  $k \ge 2$  while  $P_1 u|_K$  is the interpolation of u at the three vertices of cell K.

For simplicity, if not otherwise stated, we use  $p, \mathbf{u}, c_j, r_j, j = 1, 2, \dots, N$  as the numerical approximations from now on. Then the DG scheme for (3.2.5) -(3.2.7) is to find  $p, r_j \in W_h$  and  $\mathbf{u} \in \mathbf{W}_h = W_h \times W_h$  such that for any  $\xi, \zeta \in W_h$  and  $\boldsymbol{\eta} \in \mathbf{W}_h$ ,

$$(\tilde{d}(\mathbf{r})p_t,\xi) = (\mathbf{u},\nabla\xi) + \sum_{e\in\Gamma_0} \int_e \hat{\mathbf{u}} \cdot \boldsymbol{n}_e[\xi] ds + (q,\xi), \qquad (3.2.8)$$

$$(a(\mathbf{c})\mathbf{u},\boldsymbol{\eta}) = (p, \nabla \cdot \boldsymbol{\eta}) + \sum_{e \in \Gamma} \int_{e} \hat{p}[\boldsymbol{\eta} \cdot \boldsymbol{n}_{e}] ds, \qquad (3.2.9)$$

$$(r_{j_t}, \zeta) = (\boldsymbol{u}c_j - \boldsymbol{D}(\boldsymbol{u})\nabla c_i, \nabla\zeta) + (\check{c}_j q - r_j z_j p_t, \zeta) + \sum_{e \in \Gamma_0} \int_e \widehat{\boldsymbol{u}c_j} \cdot \boldsymbol{n}_e[\zeta] ds$$
$$- \sum_{e \in \Gamma_0} \int_e \left( \{\boldsymbol{D}(\boldsymbol{u})\nabla c_j \cdot \boldsymbol{n}_e\}[\zeta] + \{\boldsymbol{D}(\boldsymbol{u})\nabla\zeta \cdot \boldsymbol{n}_e\}[c_j] + \frac{\tilde{\alpha}}{|e|}[c_j][\zeta] \right) ds,$$
(3.2.10)

where

$$c_j = P_k \left\{ \frac{r_j}{\Phi} \right\}, \quad \tilde{d}(\mathbf{r}) = \sum_{j=1}^N z_j r_j, \quad (u, v) = \int_K u v dx, \quad \check{c}_j = \left\{ \begin{array}{ll} \tilde{c}_j, & q > 0, \\ \frac{r_j}{\Phi}, & q < 0. \end{array} \right.$$

In (3.2.8)-(3.2.10),  $\hat{p}, \hat{u}$  and  $\widehat{uc_j}$  are the numerical fluxes. We use alternating fluxes for the diffusion term and for any  $e \in \Gamma_0$ 

$$\hat{\boldsymbol{u}}|_e = \boldsymbol{u}^+|_e, \quad \hat{p}|_e = p^-|_e,$$
(3.2.11)

and on  $\partial \Omega$  we take

$$\hat{p}|_e = p^-|_e, \ \forall e \in \partial \Omega^+, \ \hat{p}|_e = p^+|_e, \ \forall e \in \partial \Omega^-.$$

For the convection term, for any  $e \in \Gamma_0$  we take

$$\widehat{\boldsymbol{u}c_j} = \boldsymbol{u}^+ c_j^+ - \alpha[c_j] \mathbf{n}_e. \tag{3.2.12}$$

In (3.2.10) and (3.2.12),  $\alpha$  and  $\tilde{\alpha}$  are two positive constants to be chosen by the bound-preserving technique. Before we complete this subsection, we would like

to introduce the following definition that will be used in the bound-preserving technique.

**Definition 3.2.1.** We say the flux  $\widehat{\mathbf{uc}}_j$  is consistent with  $\hat{\mathbf{u}}$  if  $\widehat{\mathbf{uc}}_j = \hat{\mathbf{u}}$  by taking  $c_j = 1$  in  $\Omega$ .

The numerical flux  $\widehat{\mathbf{u}c_j}$  in (3.2.12) is consistent with the flux  $\hat{\mathbf{u}}$  in (3.2.11), and this is required by the bound-preserving technique.

**Remark 3.2.1.** There are plenty of fluxes can be used following the procedures introduced in the next section. The proofs are basically the same with some minor changes, so we only list some of them below without more details.

- $\hat{\mathbf{u}} = \mathbf{u}^-, \ \hat{p} = p^+, \ \widehat{\mathbf{u}c_j} = \mathbf{u}^-c_j^- \alpha[c_j]\mathbf{n}_e.$
- $\hat{\mathbf{u}} = \frac{1}{2}(\mathbf{u}^+ + \mathbf{u}^-), \ \hat{p} = \frac{1}{2}(p^+ + p^-), \ \widehat{\mathbf{uc}_j} = \frac{1}{2}(\mathbf{u}^+ c_j^+ + \mathbf{u}^- c_j^-) \alpha[c_j]\mathbf{n}_e.$

#### **3.3** Second-order bound-preserving scheme

In this section, we will construct second-order bound-preserving DG scheme with Euler forward time discretization on triangular meshes. For simplicity, we only discuss the technique for cells away from  $\partial\Omega$ , while the boundary cells can be analyzed following the same lines with some minor changes. A similar analysis for the boundary cells can be found in [36]. We use  $o_K$  for the numerical approximation of o in K with cell average  $\bar{o}_K$ . Moreover, we use  $o^n$  as the solution o at time level n. Now, we will demonstrate the bound-preserving technique in detail. For simplicity, we will drop the subindex j in (3.2.10) and use  $r, c, \check{c}, z$ for  $r_j, c_j, \check{c}_j, z_j$ , respectively.

In (3.2.10), we take  $\zeta = 1$  in K to obtain the equation satisfied by the cell average of r

$$\bar{r}_{K}^{n+1} = H_{K}^{c}(r, \boldsymbol{u}, c) + H_{K}^{d}(r, \boldsymbol{u}, c) + H_{K}^{s}(r, \check{c}, q, z, p)$$
(3.3.13)

where

$$H_{K}^{c}(r,\boldsymbol{u},c) = \frac{1}{3}\bar{r}_{K}^{n} - \lambda \sum_{i=1}^{3} \int_{e_{K}^{i}} \widehat{\boldsymbol{u}} \widehat{c} \cdot \boldsymbol{\nu}_{i} ds, \qquad (3.3.14)$$

$$H_{K}^{d}(r,\boldsymbol{u},c) = \frac{1}{3}\bar{r}_{K}^{n} + \lambda \sum_{i=1}^{3} \int_{e_{K}^{i}} \left( \{\boldsymbol{D}(\boldsymbol{u}) \nabla c \cdot \boldsymbol{\nu}_{i}\} + \frac{\tilde{\alpha}}{\ell_{K}^{i}} [c] \mathbf{n}_{e} \cdot \boldsymbol{\nu}_{i} \right) ds, \qquad (3.3.15)$$

$$H_K^s(r,\check{c},q,z,p) = \frac{1}{3}\bar{r}_K^n + \Delta t \overline{\check{c}q - rzp_t}, \qquad (3.3.16)$$

with  $\lambda = \frac{\Delta t}{|K|}$  being the ratio of the time step and the area of triangle K, and  $\overline{cq} - rzp_t$  being the cell average of  $\check{cq} - rzp_t$ . We denote  $V_i$ , i = 1, 2, 3 as the three vertices of cell K. In this section, we will construct the bound-preserving technique in K, hence for any  $w \in W_h$ , we define  $w(V_i)$  to be the limit evaluated in K. We use the (k+1)-point Gaussian quadrature to approximate the integrals along the cell interfaces in (3.3.14)-(3.3.16), and denote  $x_{i,\beta}$ ,  $\beta = 1, 2, \cdots, k+1$  as the quadrature points on  $e_K^i$  with  $w_\beta$  as the corresponding weights on the reference interval  $\left[-\frac{1}{2}, \frac{1}{2}\right]$ . Moreover, we use quadratures discussed in [88] to compute the cell average  $\bar{r}_K^n$ . The quadrature contains  $L = 3(N_G - 2)(k + 1)$  quadrature points, denoted as  $x_\gamma$ , lying in the interior of K with  $2N_G - 3 \geq k$ , and the quadratures points on the cell interfaces are exactly the k + 1 Gaussian

quadratures points. We denote the quadrature weights corresponding to the interior quadrature points as  $\tilde{w}_{\gamma}$  and those on the cell interfaces as  $\hat{w}_{\beta}$ . In [88], it was shown that  $\hat{w}_{\beta} = \frac{2}{3}w_{\beta}\hat{w}$ , where  $\hat{w}$  is the quadrature weight corresponding to the first quadrature point in the  $N_G$ -point Gauss-Lobatto quadrature on the interval  $\left[-\frac{1}{2}, \frac{1}{2}\right]$ . Based on the above notations, we define the values of o ( $o = r, c, p, q, \Phi$ ) at the quadrature points as  $o_K^{i,\beta} = o(x_{i,\beta})$  along the boundary of K and  $o_K^{\gamma} = o(x_{\gamma})$  in cell K. Now, we can demonstrate the bound-preserving techniques. We will consider the source term  $H_K^s$  first, and discuss the high-order bound-preserving technique.

**Lemma 3.3.1.** Suppose  $r^n > 0$   $(c^n > 0)$ , then  $H^s_K(r, \check{c}, q, z, p) > 0$  under the conditions

$$\Delta t \le \frac{1}{6zp_M}, \quad \Delta t \le \frac{\Phi_m}{6q_M}, \tag{3.3.17}$$

where

$$p_{M} = \max_{i,\beta,\gamma} ((p_{t})_{K}^{i,\beta}, (p_{t})_{K}^{\gamma}, 0) \quad \Phi_{m} = \min_{x} \Phi(x), \quad q_{M} = \max_{i,\beta,\gamma} \left\{ -q_{K}^{i,\beta}, -q_{K}^{\gamma}, 0 \right\}.$$
(3.3.18)

*Proof.* We can write  $H_K^s$  as

$$H_K^s(r,\check{c},q,z,p) = \left(\frac{1}{6}\bar{r}_K^n - \triangle t \overline{rzp_t}\right) + \left(\frac{1}{6}\bar{r}_K^n + \triangle t \overline{\check{c}q}\right) := L_1 + L_2.$$

Applying the quadrature in [88], we have

$$\begin{split} L_1 &= \frac{1}{6} \overline{r}_K^n - \bigtriangleup t \overline{rzp_t} \\ &= \frac{1}{6} \left( \sum_{i=1}^3 \sum_{\beta=1}^{k+1} \hat{w}_\beta r_K^{i,\beta} + \sum_{\gamma=1}^L \tilde{w}_\gamma r_K^\gamma \right) \\ &- \bigtriangleup tz \left( \sum_{i=1}^3 \sum_{\beta=1}^{k+1} \hat{w}_\beta r_K^{i,\beta}(p_t)_K^{i,\beta} + \sum_{\gamma=1}^L \tilde{w}_\gamma r_K^\gamma(p_t)_K^\gamma \right) \\ &= \sum_{i=1}^3 \sum_{\beta=1}^{k+1} \hat{w}_\beta \left( \frac{1}{6} - \bigtriangleup tz(p_t)_K^{i,\beta} \right) r_K^{i,\beta} + \sum_{\gamma=1}^L \tilde{w}_\gamma \left( \frac{1}{6} - \bigtriangleup tz(p_t)_K^\gamma \right) r_K^\gamma. \end{split}$$

Then  $L_1 > 0$  under the condition (3.3.17). We apply the same quadrature for  $L_2$  to obtain

$$\begin{split} L_2 = & \frac{1}{6} \left( \sum_{i=1}^3 \sum_{\beta=1}^{k+1} \hat{w}_\beta r_K^{i,\beta} + \sum_{\gamma=1}^L \tilde{w}_\gamma r_K^\gamma \right) + \triangle t \left( \sum_{i=1}^3 \sum_{\beta=1}^{k+1} \hat{w}_\beta \check{c}_K^{i,\beta} q_K^{i,\beta} + \sum_{\gamma=1}^L \tilde{w}_\gamma \check{c}_K^\gamma q_K^\gamma \right) \\ = & \sum_{i=1}^3 \sum_{\beta=1}^{k+1} \hat{w}_\beta \left( \frac{1}{6} r_K^{i,\beta} + \triangle t \check{c}_K^{i,\beta} q_K^{i,\beta} \right) + \sum_{\gamma=1}^L \tilde{w}_\gamma \left( \frac{1}{6} r_K^\gamma + \triangle t \check{c}_K^\gamma q_K^\gamma \right). \end{split}$$

Notice that  $\check{c} = r/\Phi$  if q < 0 while  $\check{c} > 0$  if q > 0. Therefore, under the condition (3.3.17), each term in the summation above is positive.

In the rest part of this section, we will consider second-order scheme only, i.e. k = 1,  $N_G = 2$ , L = 0, then  $\hat{w} = \frac{1}{2}$  and  $w_\beta = 3\hat{w}_\beta$ . Now we can analyze the convection term  $H_K^c$  and the result is given below.

**Lemma 3.3.2.** Suppose  $r^n > 0$  ( $c^n > 0$ ), if  $\alpha$  satisfies

$$\alpha > \max_{i,\beta} \{ |\boldsymbol{u}_{K_i}^{i,\beta}|, 0 \},$$
(3.3.19)

and the time step satisfies

$$\Delta t \le \min_{i,\beta} \left\{ \frac{1}{9\ell_K^i(|\boldsymbol{u}_K^{i,\beta}| + \alpha)} \right\} \Phi_m|K|.$$
(3.3.20)

we have  $H_K^c(r, \boldsymbol{u}, c) > 0$ .

*Proof.* Following the same analysis for the source term, we write

$$H_K^c = \sum_{i=1}^3 \sum_{\beta=1}^2 w_\beta H_{i,\beta}^c, \quad H_{i,\beta}^c = \frac{1}{9} r_K^{i,\beta} - \lambda \ell_K^i \widehat{\boldsymbol{u}} \widehat{\boldsymbol{c}}^{i,\beta} \cdot \boldsymbol{\nu}_i.$$

We only need to show  $H_{i,\beta}^c > 0$ .

Case 1:  $\boldsymbol{\nu}_i = \mathbf{n}_e$ , i.e.  $\boldsymbol{u}^- = \boldsymbol{u}_K$ ,  $\boldsymbol{u}^+ = \boldsymbol{u}_{K_i}$ ,  $c^- = c_K$  and  $c^+ = c_{K_i}$ . Then

$$H_{i,\beta}^{c} = \frac{1}{9} r_{K}^{i,\beta} - \lambda \ell_{K}^{i} (\boldsymbol{u}_{K_{i}}^{i,\beta} c_{K_{i}}^{i,\beta} \cdot \boldsymbol{\nu}_{i} - \alpha c_{K_{i}}^{i,\beta} + \alpha c_{K}^{i,\beta}).$$

Since r and c are both linear functions, we can write the function values of r and c as the interpolation of the values at vertices  $\{V_1, V_2, V_3\}$  of K, i.e. for any point  $x_{\rho}$  in K,

$$r_{K}^{\rho} = \mu_{1}^{\rho} r_{K}(V_{1}) + \mu_{2}^{\rho} r_{K}(V_{2}) + \mu_{3}^{\rho} r_{K}(V_{3}), \quad c_{K}^{\rho} = \mu_{1}^{\rho} c_{K}(V_{1}) + \mu_{2}^{\rho} c_{K}(V_{2}) + \mu_{3}^{\rho} c_{K}(V_{3}),$$
(3.3.21)

with  $\mu_m^{\rho} \ge 0, \, m = 1, 2, 3$ , and  $\sum_{m=1}^{3} \mu_m^{\rho} = 1$ . Then

$$H_{i,\beta}^{c} = \sum_{m=1}^{3} \mu_{m}^{i,\beta} \left( \frac{1}{9} r_{K}(V_{m}) - \lambda \ell_{K}^{i} \alpha c_{K}(V_{m}) \right) + \lambda \ell_{K}^{i} (\alpha - \boldsymbol{u}_{K_{i}}^{i,\beta} \cdot \boldsymbol{\nu}_{i}) c_{K_{i}}^{i,\beta}$$
$$= \sum_{m=1}^{3} \mu_{m}^{i,\beta} \left( \frac{1}{9} \Phi_{K}(V_{m}) - \lambda \ell_{K}^{i} \alpha \right) c_{K}(V_{m}) + \lambda \ell_{K}^{i} (\alpha - \boldsymbol{u}_{K_{i}}^{i,\beta} \cdot \boldsymbol{\nu}_{i}) c_{K_{i}}^{i,\beta}.$$

Then we have  $H_{i,\beta}^c > 0$ , if  $\alpha$  and  $\Delta t$  satisfy (3.3.19) and (3.3.20), respectively. Case 2:  $\boldsymbol{\nu}_i = -\mathbf{n}_e$ , i.e.  $\boldsymbol{u}^+ = \boldsymbol{u}_K$ ,  $\boldsymbol{u}^- = \boldsymbol{u}_{K_i}$ ,  $c^+ = c_K$  and  $c^- = c_{K_i}$ . Then

$$H_{i,\beta}^{c} = \frac{1}{9} r_{K}^{i,\beta} - \lambda \ell_{K}^{i} (\boldsymbol{u}_{K}^{i,\beta} c_{K}^{i,\beta} \cdot \boldsymbol{\nu}_{i} - \alpha c_{K_{i}}^{i,\beta} + \alpha c_{K}^{i,\beta}).$$

Applying (3.3.21) again, we have

$$H_{i,\beta}^{c} = \sum_{m=1}^{3} \mu_{m}^{i,\beta} \left( \frac{1}{9} \Phi_{K}(V_{m}) - \lambda \ell_{K}^{i} \mathbf{u}_{K}^{i,\beta} \cdot \boldsymbol{\nu}_{i} - \lambda \ell_{K}^{i} \alpha \right) c_{K}(V_{m}) + \lambda \ell_{K}^{i} \alpha c_{K_{i}}^{i,\beta}.$$

Then we have  $H_{i,\beta}^c > 0$  under the condition (3.3.20).

Finally, we discuss the diffusion part. We also take k = 1, G = 2, L = 0 and the result is given in the following lemma.

**Lemma 3.3.3.** Assume the minimum angle of each triangle K is uniformly bounded away from zero. Suppose  $r^n > 0$  ( $c^n > 0$ ), then  $H^d_K(r, \boldsymbol{u}, c) > 0$  under the conditions

$$\tilde{\alpha} \ge \frac{(3+\sqrt{3})\Lambda}{2\min_{K,i,j}\left(\sin\left(\theta_K^{i,j}\right)\right)},\tag{3.3.22}$$

and

$$\Delta t \le \frac{\Phi_m |K|}{18\tilde{\alpha}}, \quad \frac{\Delta t}{|K|} \frac{(3+\sqrt{3})\Lambda}{\min_{K,i,j} \left(\sin\left(\theta_K^{i,j}\right)\right)} \le \frac{1}{54} \Phi_m, \tag{3.3.23}$$

where  $\theta_K^{i,j}$ ,  $i, j = 1, 2, 3, i \neq j$  denotes the angle between the edge  $e_K^i$  and  $e_K^j$ , and  $\Lambda$  is the largest absolute value of the eigenvalue of **D**.

*Proof.* First, we will consider the term

$$\int_{e_K^i} \left( \{ \boldsymbol{D}(\boldsymbol{u}) \nabla c \cdot \boldsymbol{\nu}_i \} + \frac{\tilde{\alpha}}{\ell_K^i} [c] \mathbf{n}_e \cdot \boldsymbol{\nu}_i \right) ds.$$

Following [89], we write

$$oldsymbol{D}(oldsymbol{u})
abla c\cdotoldsymbol{
u}_i=
abla c\cdotoldsymbol{D}(oldsymbol{u})oldsymbol{
u}_i=rac{\partial c}{\partialoldsymbol{\eta}_i}\| ilde{oldsymbol{\eta}}_i\|$$

where

$$ilde{\eta_i} = oldsymbol{D}(oldsymbol{u}) oldsymbol{
u}_i, \ oldsymbol{\eta}_i = rac{ ilde{\eta}_i}{\| ilde{\eta}_i\|}.$$



Figure 3.1: Two intersection points for the numerical flux in diffusion part on the triangular mesh.

Define  $\eta_K = \eta_i|_K$  and  $\eta_{K_i} = \eta_i|_{K_i}$ . Likewise for  $\tilde{\eta}_K$  and  $\tilde{\eta}_{K_i}$ . For each quadrature point  $x_{i,\beta}$  on the edge  $e_K^i$ , we can draw a straight line from  $x_{i,\beta}$  with direction  $\eta_{K_i}$  intersects  $\partial K_i$  at  $\tilde{x}_{K_i}^{i,\beta}$ . Similarly, we can draw another straight line from  $x_{i,\beta}$ with direction  $-\eta_K$  intersects  $\partial K$  at  $\tilde{x}_K^{i,\beta}$ . See Figure 3.1 for an illustration. It is easy to verify that at  $x = x_{i,\beta}$ 

$$\begin{aligned} \{\boldsymbol{D}(\boldsymbol{u})\nabla \boldsymbol{c}\cdot\boldsymbol{\nu}_{i}\} &+ \frac{\tilde{\alpha}}{\ell_{K}^{i}}[\boldsymbol{c}]\mathbf{n}_{e}\cdot\boldsymbol{\nu}_{i} \\ &= \frac{1}{2}\boldsymbol{D}(\boldsymbol{u}_{K})\nabla \boldsymbol{c}_{K}\cdot\boldsymbol{\nu}_{i} + \frac{1}{2}\boldsymbol{D}(\boldsymbol{u}_{K_{i}})\nabla \boldsymbol{c}_{K_{i}}\cdot\boldsymbol{\nu}_{i} + \tilde{\alpha}\frac{(\boldsymbol{c}_{K_{i}}-\boldsymbol{c}_{K})}{\ell_{K}^{i}} \\ &= \frac{1}{2}\frac{c_{K}^{i,\beta}-c(\tilde{x}_{K}^{i,\beta})}{\|\boldsymbol{x}_{K}^{i,\beta}-\tilde{x}_{K}^{i,\beta}\|}\|\tilde{\boldsymbol{\eta}}_{K}\| + \frac{1}{2}\frac{c(\tilde{x}_{K_{i}}^{i,\beta})-c_{K_{i}}^{i,\beta}}{\|\tilde{x}_{K_{i}}^{i,\beta}-x_{K}^{i,\beta}\|}\|\tilde{\boldsymbol{\eta}}_{K_{i}}\| + \frac{\tilde{\alpha}}{\ell_{K}^{i}}(\boldsymbol{c}_{K_{i}}^{i,\beta}-\boldsymbol{c}_{K}^{i,\beta}) \\ &= \left(\frac{\|\tilde{\boldsymbol{\eta}}_{K}\|}{2\|\boldsymbol{x}_{K}^{i,\beta}-\tilde{x}_{K}^{i,\beta}\|}-\frac{\tilde{\alpha}}{\ell_{K}^{i}}\right)\boldsymbol{c}_{K}^{i,\beta} + \left(\frac{\tilde{\alpha}}{\ell_{K}^{i}}-\frac{\|\tilde{\boldsymbol{\eta}}_{K_{i}}\|}{2\|\tilde{x}_{K_{i}}^{i,\beta}-x_{K}^{i,\beta}\|}\right)\boldsymbol{c}_{K_{i}}^{i,\beta} \\ &- \frac{\|\tilde{\boldsymbol{\eta}}_{K}\|}{2\|\boldsymbol{x}_{K}^{i,\beta}-\tilde{x}_{K}^{i,\beta}\|}\boldsymbol{c}(\tilde{x}_{K}^{i,\beta}) + \frac{\|\tilde{\boldsymbol{\eta}}_{K_{i}}\|}{2\|\tilde{x}_{K_{i}}^{i,\beta}-x_{K}^{i,\beta}\|}\boldsymbol{c}(\tilde{x}_{K_{i}}^{i,\beta}). \end{aligned}$$

We write the cell average  $\bar{r}_K^n$  as

$$\bar{r}_{K}^{n} = \sum_{i=1}^{3} \sum_{\beta=1}^{2} \hat{w}_{\beta} r_{K}^{i,\beta} = \sum_{i=1}^{3} \sum_{\beta=1}^{2} \sum_{m=1}^{3} \hat{w}_{\beta} \mu_{m}^{i,\beta} \Phi_{K}(V_{m}) c_{K}(V_{m}).$$

we can rewrite  $H^d_K(r, \boldsymbol{u}, c)$  as

$$\begin{aligned} H_{K}^{d} &= \frac{1}{3} \sum_{i=1}^{3} \sum_{\beta=1}^{2} \sum_{m=1}^{3} \hat{w}_{\beta} \mu_{m}^{i,\beta} \Phi_{K}(V_{m}) c_{K}(V_{m}) \\ &+ \lambda \sum_{i=1}^{3} \ell_{K}^{i} \sum_{\beta=1}^{2} w_{\beta} \left[ \{ \boldsymbol{D}(\boldsymbol{u}) \nabla c \cdot \boldsymbol{\nu}_{i} \} + \frac{\tilde{\alpha}}{\ell_{K}^{i}} [c] \mathbf{n}_{e} \cdot \boldsymbol{\nu}_{i} \right]_{x=x_{i,\beta}} \\ &= \sum_{i=1}^{3} \sum_{\beta=1}^{2} w_{\beta} \left( \frac{1}{9} \sum_{m=1}^{3} \mu_{m}^{i,\beta} \Phi_{K}(V_{m}) c_{K}(V_{m}) \right) \\ &+ \sum_{i=1}^{3} \sum_{\beta=1}^{2} w_{\beta} \lambda \ell_{K}^{i} \left[ \{ \boldsymbol{D}(\boldsymbol{u}) \nabla c \cdot \boldsymbol{\nu}_{i} \} + \frac{\tilde{\alpha}}{\ell_{K}^{i}} [c] \mathbf{n}_{e} \cdot \boldsymbol{\nu}_{i} \right]_{x=x_{i,\beta}} \\ &:= \sum_{i=1}^{3} \sum_{\beta=1}^{2} w_{\beta} L_{i,\beta} + L, \end{aligned}$$

where

$$\begin{split} L_{i,\beta} = &\frac{1}{18} \sum_{m=1}^{3} \mu_{m}^{i,\beta} \Phi_{K}(V_{m}) c_{K}(V_{m}) \\ &+ \lambda \ell_{K}^{i} \left[ \left( \frac{\|\tilde{\eta}_{K}\|}{2\|x_{K}^{i,\beta} - \tilde{x}_{K}^{i,\beta}\|} - \frac{\tilde{\alpha}}{\ell_{K}^{i}} \right) c_{K}^{i,\beta} + \left( \frac{\tilde{\alpha}}{\ell_{K}^{i}} - \frac{\|\tilde{\eta}_{K_{i}}\|}{2\|\tilde{x}_{K_{i}}^{i,\beta} - x_{K}^{i,\beta}\|} \right) c_{K_{i}}^{i,\beta} \\ &+ \frac{\|\tilde{\eta}_{K_{i}}\|}{2\|\tilde{x}_{K_{i}}^{i,\beta} - x_{K_{i}}^{i,\beta}\|} c(\tilde{x}_{K_{i}}^{i,\beta}) \right], \\ L = &\frac{1}{6} \tilde{r}_{K}^{n} - \lambda \sum_{i=1}^{3} \sum_{\beta=1}^{2} \frac{\ell_{K}^{i} \|\tilde{\eta}_{K}\|}{2\|x_{K}^{i,\beta} - \tilde{x}_{K}^{i,\beta}\|} c(\tilde{x}_{K}^{i,\beta}). \end{split}$$

We need to make  $L_{i,\beta} > 0$ . In fact

$$\begin{split} L_{i,\beta} &= \frac{1}{18} \sum_{m=1}^{3} \mu_{m}^{i,\beta} \Phi_{K}(V_{m}) c_{K}(V_{m}) + \lambda \ell_{K}^{i} \left( \frac{\|\tilde{\boldsymbol{\eta}}_{K}\|}{2\|x_{K}^{i,\beta} - \tilde{x}_{K}^{i,\beta}\|} - \frac{\tilde{\alpha}}{\ell_{K}^{i}} \right) c_{K}^{i,\beta} \\ &+ \lambda \ell_{K}^{i} \left( \frac{\tilde{\alpha}}{\ell_{K}^{i}} - \frac{\|\tilde{\boldsymbol{\eta}}_{K_{i}}\|}{2\|\tilde{x}_{K_{i}}^{i,\beta} - x_{K}^{i,\beta}\|} \right) c_{K_{i}}^{i,\beta} + \lambda \ell_{K}^{i} \frac{\|\tilde{\boldsymbol{\eta}}_{K_{i}}\|}{2\|\tilde{x}_{K_{i}}^{i,\beta} - x_{K}^{i,\beta}\|} c(\tilde{x}_{K_{i}}^{i,\beta}) \\ &= \sum_{m=1}^{3} \mu_{m}^{i,\beta} \left( \frac{1}{18} \Phi_{K}(V_{m}) + \lambda \ell_{K}^{i} \left( \frac{\|\tilde{\boldsymbol{\eta}}_{K}\|}{2\|x_{K}^{i,\beta} - \tilde{x}_{K}^{i,\beta}\|} - \frac{\tilde{\alpha}}{\ell_{K}^{i}} \right) \right) c_{K}(V_{m}) \\ &+ \lambda \ell_{K}^{i} \left( \frac{\tilde{\alpha}}{\ell_{K}^{i}} - \frac{\|\tilde{\boldsymbol{\eta}}_{K_{i}}\|}{2\|\tilde{x}_{K_{i}}^{i,\beta} - x_{K}^{i,\beta}\|} \right) c_{K_{i}}^{i,\beta} + \lambda \ell_{K}^{i} \frac{\|\tilde{\boldsymbol{\eta}}_{K_{i}}\|}{2\|\tilde{x}_{K_{i}}^{i,\beta} - x_{K}^{i,\beta}\|} c(\tilde{x}_{K_{i}}^{i,\beta}). \end{split}$$

Notice that  $\|\tilde{\boldsymbol{\eta}}\| \leq \Lambda$ . To make  $L_{i,\beta} > 0$ , we need

$$\tilde{\alpha} \geq \frac{\ell_K^i \Lambda}{2 \|\tilde{x}_{K_i}^{i,\beta} - x_{K_i}^{i,\beta}\|}, \quad \lambda \ell_K^i \left(\frac{\tilde{\alpha}}{\ell_K^i} - \frac{\|\tilde{\boldsymbol{\eta}}_K\|}{2 \|x_K^{i,\beta} - \tilde{x}_K^{i,\beta}\|}\right) \leq \frac{1}{18} \Phi_K(V_m).$$

It is easy to compute that

$$\frac{\ell_K^i}{\|\tilde{x}_K^{i,\beta} - x_K^{i,\beta}\|} \le \frac{6}{(3-\sqrt{3})\min_j \sin\left(\theta_K^{i,j}\right)}.$$

and we conclude  $L_{i,\beta} > 0$  under the conditions (3.3.22) and (3.3.23). Finally, we can apply the same idea above to estimate L. Similar to (3.3.21), we write

$$c(\tilde{x}_K^{i,\beta}) = \sum_{m=1}^3 \tilde{\mu}_m^{i,\beta} c_K(V_m),$$

with  $0 \leq \tilde{\mu}_m^{i,\beta} \leq 1$  and  $\sum_{m=1}^3 \tilde{\mu}_m^{i,\beta} = 1$ . Then

$$L = \frac{1}{6} \bar{r}_{K}^{n} - \lambda \ell_{K}^{i} \sum_{i=1}^{3} \sum_{\beta=1}^{2} \frac{\|\tilde{\eta}_{K}\|}{2\|x_{K}^{i,\beta} - \tilde{x}_{K}^{i,\beta}\|} c(\tilde{x}_{K}^{i,\beta})$$
  
$$= \sum_{m=1}^{3} \left( \frac{1}{18} \Phi_{K}(V_{m}) - \lambda \ell_{K}^{i} \sum_{i=1}^{3} \sum_{\beta=1}^{2} \frac{\|\tilde{\eta}_{K}\| \tilde{\mu}_{m}^{i,\beta}}{2\|x_{K}^{i,\beta} - \tilde{x}_{K}^{i,\beta}\|} \right) c_{K}(V_{m})$$
  
$$\geq \sum_{m=1}^{3} \left( \frac{1}{18} \Phi_{K}(V_{m}) - \lambda \sum_{i=1}^{3} \sum_{\beta=1}^{2} \frac{(3 + \sqrt{3})\Lambda}{2\min_{j} \sin\left(\theta_{K}^{i,j}\right)} \right) c_{K}(V_{m})$$

Therefore, we have L > 0 under the condition (3.3.23).

Base on the above three lemmas, we can state the following theorem.

**Theorem 3.3.4.** Suppose  $r^n > 0$  ( $c^n > 0$ ), and the parameters  $\alpha$  and  $\tilde{\alpha}$  satisfy (3.3.19) and (3.3.22), respectively. Then  $\bar{r}^{n+1} > 0$  under the conditions (3.3.17), (3.3.20) and (3.3.23).

Now, we have proved  $\bar{r}_j > 0$  for  $j = 1, 2, \dots, N-1$ . To obtain  $\bar{r}_N > 0$ , we need to subtract (3.2.10) from (3.2.8) to obtain

$$(r_{N_t}, \zeta) = (\boldsymbol{u}c_N - \boldsymbol{D}(\boldsymbol{u})\nabla c_N, \nabla\zeta) + (\check{c}_N q - r_N z_N p_t, \zeta) + \sum_{e \in \Gamma_0} \int_e \widehat{\boldsymbol{u}c_N} \cdot \boldsymbol{n}_e[\zeta] ds$$
$$- \sum_{e \in \Gamma_0} \int_e \left( \{\boldsymbol{D}(\boldsymbol{u})\nabla c_N \cdot \boldsymbol{n}_e\}[\zeta] + \{\boldsymbol{D}(\boldsymbol{u})\nabla\zeta \cdot \boldsymbol{n}_e\}[c_N] + \frac{\tilde{\alpha}}{|e|}[c_N][\zeta] \right) ds.$$
(3.3.24)

Here, we have used the fact that the flux for (3.2.10) is consistent with that in (3.2.8). We can observe that the above equation is similar to (3.2.10). Therefore, following the same analysis above with minor changes we have the following theorem.

**Theorem 3.3.5.** Suppose  $0 \le r^n \le \Phi$ , and the conditions in Theorem 3.3.4 are satisfied. Moreover, if the fluxes  $\widehat{uc_j}$  and  $\hat{u}$  are consistent, then  $\bar{r}^{n+1} \le \bar{\Phi}$ , under the condition

$$\Delta t \le \frac{1}{6z_M p_M},\tag{3.3.25}$$

where  $p_M$  is given in (3.3.18) and  $z_M = \max_{1 \le j \le N} z_j$ .

# 3.4 Bound-preserving technique for high-order scheme

In this section, we will apply the flux limiter to construct high-order boundpreserving technique.

#### 3.4.1 Flux limiter

We use  $P^k$  (k > 2) polynomials and write (3.3.13) as

$$\bar{r}_K^{n+1} = \bar{r}_K^n + \lambda \sum_{i=1}^3 \hat{F}_{e^i} + \Delta t \bar{s}_i$$

where

$$\hat{F}_{e^{i}} = -\int_{e^{i}} \widehat{uc} \cdot \nu_{i} ds + \int_{e^{i}} \left( \{ \boldsymbol{D}(\boldsymbol{u}) \nabla c \cdot \nu_{i} \} + \frac{\tilde{\alpha}}{\ell_{K}^{i}} [c] \right) ds, \quad \bar{s} = \overline{\tilde{c}q - rz_{1}p_{t}}$$

$$(3.4.26)$$

are high-order flux and source, respectively. In Section 3.3, we have demonstrated how to treat the source terms. Therefore, we only discuss the modification of the high-order fluxes only. We will apply the flux limiter [39, 65] and combine the high-order flux  $\hat{F}_{e^i}$  and the second-order fluxes, which was analyzed in Section 3.3, denoted as  $\hat{f}_{e^i}$ . We define the new flux as

$$\tilde{F}_{e^i} = \hat{f}_{e^i} + \theta_{e^i} (\hat{F}_{e^i} - \hat{f}_{e^i}),$$

where  $\theta_{e^i}$  is a parameter that to be chosen. Then the cell average can be written as

$$\bar{r}_{K}^{n+1} = \bar{r}_{K}^{n} + \lambda \sum_{i=1}^{3} \hat{f}_{e^{i}} + \lambda \sum_{i=1}^{3} \theta_{e^{i}} (\hat{F}_{e^{i}} - \hat{f}_{e^{i}}) + \Delta t\bar{s} = \bar{r}_{L}^{n+1} + \lambda \sum_{i=1}^{3} \theta_{e^{i}} (\hat{F}_{e^{i}} - \hat{f}_{e^{i}}),$$

where

$$\bar{r}_L^{n+1} = \bar{r}_K^n + \lambda \sum_{i=1}^3 \hat{f}_{e^i} + \Delta t\bar{s}$$

is the second order cell average which was proved to be positive if  $\Delta t$  is sufficiently small. Notice that, we need the fluxes in (3.2.10) and (3.2.8) to be consistent. Therefore, we have to discuss the fluxes for all components together. We define  $\hat{f}_{e^i}^j$  and  $\hat{F}_{e^i}^j$  as the second- and high-order fluxes for component  $j, j = 1, 2, \dots, N$ , respectively, and the cell average  $\bar{r}$  for the jth component to be  $\bar{r}_j$ . To compute  $\hat{f}_{e^i}^j$ , we only replace the  $c_j$  in  $\hat{F}_{e^i}^j$  in (3.4.26) by a second-order approximation. We cannot change  $\mathbf{u}$ , since we want  $\sum_{j=1}^N \hat{F}_{e^i}^j = \sum_{j=1}^N \hat{f}_{e^i}^j = \hat{\mathbf{u}}_{e^i}$ , which due to the flux consistency requirement. To construct the second-order  $c_j$ , we can simply apply the second-order  $L^2$  projection to the high-order  $c_j$ , and then apply the limiter discussed in 3.4.2 with k = 1 and  $\Phi$  as the second-order  $L^2$  projection of  $\phi$ . We can choose the parameter  $\theta_{e^i}$  as follows:

1. For any  $K \in \Omega_h$ , set  $\beta_K = 0$ .

2. Define 
$$\hat{F}_{e^i}^N = \hat{\mathbf{u}}_{e^i} - \sum_{j=1}^{N-1} \hat{F}_{e^i}^j$$
,  $\hat{f}_{e^i}^N = \hat{\mathbf{u}}_{e^i} - \sum_{j=1}^{N-1} f_{e^i}^j$  and  $\bar{r}_n = \bar{\Phi} - \sum_{j=1}^{N-1} \bar{r}_j$ .

3. For any  $j = 1, 2, \dots, N$ , if  $\hat{F}_{e^i}^j - \hat{f}_{e^i}^j \ge 0$ , take  $\theta_{K,e^i}^j = 1$ , otherwise set  $\beta_K = \beta_K + \hat{F}_{e^i}^j - \hat{f}_{e^i}^j$ .

4. For those edges 
$$e^i$$
 with  $\hat{F}^j_{e^i} - \hat{f}^j_{e^i} < 0$ , we set  $\theta^j_{K,e^i} = \min\left\{-\frac{\bar{r}^{n+1}_{j,L}}{\lambda\beta_K^m}, 1\right\}$ .

- 5. Take  $\theta_{K,e^i} = \min_{1 \le j \le N} \theta_{K,e^i}^j$ .
- 6. For any  $e \in \Gamma_0$ , we can find  $K_1, K_2 \in \Omega_h$  such that  $K_1 \cap K_2 = e$ . We take  $\theta_e = \min\{\theta_{K_1,e}, \theta_{K_2,e}\}.$

Following the same analyses in [12], we have  $\bar{r}_j^{n+1} \ge 0, j = 1, 2, \cdots, N$ . Thus,  $0 \le \bar{r}_j^{n+1} \le \bar{\Phi}$ , since we have the relationship  $\bar{r}_1^{n+1} + \bar{r}_2^{n+1} + \ldots + \bar{r}_N^{n+1} = \bar{\Phi}$ .

**Remark 3.4.1.** In (3.2.8)-(3.2.10), we do not compute  $r_N(c_N)$  directly. Step 2 in the above algorithm is used to compute the fluxes in (3.3.24). Actually, we can simply take  $F_{e^i}^N = -\sum_{j=1}^{N-1} F_{e^i}^j$ ,  $\hat{f}_{e^i}^N = -\sum_{j=1}^{N-1} f_{e^i}^j$ , since we only need the difference of the higher order and lower order fluxes. Moreover, step 5 is used to construct consistent fluxes (See definition 3.2.1).

#### 3.4.2 Slope limiter

In this section, we discuss the limiters to be applied. As discussed in [36], the traditional slope limiter (3.1.4) cannot be applied. In this paper, we will construct a new one. We consider problem with 2 components first and then extend it to N-component ones. The algorithm is given as follows.

1. Define  $\hat{S} = \{x \in K : r(x) \leq 0\}$ . Take

$$\hat{r}_1 = r_1 + \theta \left( \frac{\bar{r}_1}{\bar{\Phi}} \Phi - r_1 \right), \qquad \theta = \max_{y \in \hat{S}} \left\{ \frac{-r_1(y)\bar{\Phi}}{\bar{r}_1 \Phi(y) - r_1(y)\bar{\Phi}}, 0 \right\}.$$
 (3.4.27)

- 2. Set  $r_2 = \Phi \hat{r}_1$ , and repeat the above step for  $r_2$ .
- 3. Take  $\tilde{r}_1 = \Phi \hat{r}_2$  as the new approximation.

**Remark 3.4.2.** In step 1, it is easy to see that  $\hat{r}_1 \ge 0$  which further implies  $r_2 \le \Phi$ . In step 2, we have

$$\hat{r}_2 = r_2 + \theta \left( \frac{\bar{r}_2}{\bar{\Phi}} \Phi - r_2 \right) = (1 - \theta)r_2 + \theta \frac{\bar{r}_2}{\bar{\Phi}} \Phi \le (1 - \theta)\Phi + \theta \Phi = \Phi, \forall \theta \in [0, 1],$$

which means the property  $\hat{r}_2 \leq \Phi$  is inherited naturally from  $r_2 \leq \Phi$ , no matter which parameter  $\theta$  is chosen. This fact gives us enough space to modify  $\hat{r}_2$  such that  $\hat{r}_2 \geq 0$ , as we did in step one. Therefore, after step 3, we have  $0 \leq \tilde{r}_1 \leq$  $\Phi$ . Besides the above, it is easy to check that the limiter does not change the numerical cell averages, i.e.,  $\int_K \tilde{r}(x) dx = \int_K r(x) dx$ .

Moreover, we can also prove that the limiter does not affect the accuracy.

**Theorem 3.4.1.** Let  $R(x) \in C^{k+1}(K)$  and  $r(x), \Phi(x) \in P^k(K)$  with  $0 \le \overline{r} \le \overline{\Phi}$ and  $||r(x) - R(x)||_{\infty} \le Ch^{k+1}$ . Assume there exist two positive constants  $\Phi_m$ and  $\Phi_M$  such that  $0 < \Phi_m \le \Phi(x) \le \Phi_M$ , then  $||\tilde{r}(x) - R(x)||_{\infty} \le Ch^{k+1}$ .

Proof. WLOG, we assume  $\theta > 0$  in (3.4.27) and need to show the modification in step 1 keeps the accurate  $\|\hat{r}(x) - r(x)\|_{\infty} \leq Ch^{k+1}$ . Denote  $r_m = \min_{x \in K} r(x), r_M = \max_{x \in K} r(x)$ . Let  $y \in K$  be the point at which the maximum in (3.4.27) is achieved and define  $r_y = r(y) < 0, \Phi_y = \Phi(y)$ . Then

$$\theta = \frac{-r_y}{\frac{\bar{r}}{\bar{\Phi}}\Phi_y - r_y} \le \frac{-r_y}{\bar{r}\frac{\Phi_m}{\Phi_M} - r_y} \le \frac{-r_y}{\bar{r}\frac{\Phi_m}{\Phi_M} - r_y\frac{\Phi_m}{\Phi_M}} = \frac{-r_y}{\bar{r} - r_y}\frac{\Phi_M}{\Phi_m} \le \frac{-r_m}{\bar{r} - r_m}\frac{\Phi_M}{\Phi_m}$$

which further yields

$$|\hat{r} - r| = \theta |\frac{\bar{r}}{\bar{\Phi}}\Phi - r| \le \frac{\Phi_M}{\Phi_m} \frac{-r_m}{\bar{r} - r_m} |\frac{\bar{r}}{\bar{\Phi}}\Phi - r| = \frac{\Phi_M}{\Phi_m} (-r_m) \frac{|\bar{r}\frac{\Phi}{\bar{\Phi}} - r|}{\bar{r} - r_m}.$$

Since  $\frac{\Phi_M}{\Phi_m}$  is a constant and  $|-r_m| \leq Ch^{k+1}$ , we only need to prove that  $\frac{|\bar{r}\frac{\phi}{\phi}-r|}{\bar{r}-r_m} \leq C$  for some positive constant C independent of x and h. Notice that

$$\bar{r}\frac{\Phi_m}{\Phi_M} - r_M \le \bar{r}\frac{\Phi}{\bar{\Phi}} - r \le \bar{r}\frac{\Phi_M}{\Phi_m} - r_m,$$

we have

$$\left|\bar{r}\frac{\Phi}{\bar{\Phi}}-r\right| \leq \max\left\{\left|\bar{r}\frac{\Phi_M}{\Phi_m}-r_m\right|, \left|\bar{r}\frac{\Phi_m}{\Phi_M}-r_M\right|\right\},\right.$$

which further yields

$$\frac{\left|\bar{r}\frac{\Phi}{\bar{\Phi}}-r\right|}{\bar{r}-r_m} \le \max\left\{\frac{\left|\bar{r}\frac{\Phi_M}{\Phi_m}-r_m\right|}{\bar{r}-r_m}, \frac{\left|\bar{r}\frac{\Phi_m}{\Phi_M}-r_M\right|}{\bar{r}-r_m}\right\}.$$

Next, we will prove the boundedness of  $\frac{|\bar{r}\frac{\Phi_M}{\Phi_m} - r_m|}{\bar{r} - r_m}$ , and  $\frac{|\bar{r}\frac{\Phi_m}{\Phi_M} - r_M|}{\bar{r} - r_m}$ , respectively. For the first term, we have

$$\frac{|\bar{r}\frac{\Phi_M}{\Phi_m} - r_m|}{\bar{r} - r_m} = \frac{\bar{r}\frac{\Phi_M}{\Phi_m} - r_m}{\bar{r} - r_m} \le \frac{\bar{r}\frac{\Phi_M}{\Phi_m} - r_m\frac{\Phi_M}{\Phi_m}}{\bar{r} - r_m} = \frac{\Phi_M}{\Phi_m}$$

while for the second term

$$\begin{aligned} \frac{|\bar{r}\frac{\Phi_m}{\Phi_M} - r_M|}{\bar{r} - r_m} &= -\frac{\bar{r} - r_M + \bar{r}(\frac{\Phi_m}{\Phi_M} - 1)}{\bar{r} - r_m} \\ &\leq -\frac{\bar{r} - r_M}{\bar{r} - r_m} - \frac{\bar{r}(\frac{\Phi_m}{\Phi_M} - 1)}{\bar{r}} \\ &\leq \frac{r_M - \bar{r}}{\bar{r} - r_m} + 1 - \frac{\Phi_m}{\Phi_M}. \end{aligned}$$

In Appendix C of [86], Zhang proved that for any non-constant polynomial of degree k, say p(x), we have

$$\left|\frac{\bar{p} - \max p(x)}{\bar{p} - \min p(x)}\right| \le C_k,$$

where  $C_k$  is a constant only depends on the polynomial degree k. Thus,

$$\frac{\left|\bar{r}\frac{\Phi_m}{\Phi_M} - r_M\right|}{\bar{r} - r_m} \le C_k + 1 - \frac{\Phi_m}{\Phi_M},$$

and we finish the proof.

**Remark 3.4.3.** There are two ways to apply this limiter in an N-component system. One way is to compute the parameter  $\theta_j$  for the *j*th component,  $(j = 1, 2, \dots, N)$  and then take  $\theta = \max_j \theta_j$ . Another way is to modify  $r_1, r_2, \dots, r_{N-1}$  one by one such that  $r_1 \in [0, \Phi], r_2 \in [0, \Phi - r_1], r_3 \in [0, \Phi - r_1 - r_2], \dots, r_{N-1} \in [0, \Phi - r_1 - r_2 \dots - r_{N-2}].$ 

#### 3.4.3 High-order time discretization

In this section, we extend the Euler forward time discretization to high-order ones which are convex combinations of Euler forwards. In this paper, we use third-order strong stability preserving (SSP) high-order time discretization to solve the ODE system  $\mathbf{u_t} = \mathbf{L}(\mathbf{u})$ :

$$\mathbf{u}^{(1)} = \mathbf{u}^{n} + \Delta t \mathbf{L}(\mathbf{u}, t^{n}),$$
  
$$\mathbf{u}^{(2)} = \frac{3}{4}\mathbf{u}^{n} + \frac{1}{4} \left(\mathbf{u}^{(1)} + \Delta t \mathbf{L}(\mathbf{u}^{(1)}, t^{n+1})\right),$$
  
$$\mathbf{u}^{n+1} = \frac{1}{3}\mathbf{u}^{n} + \frac{2}{3} \left(\mathbf{u}^{(2)} + \Delta t \mathbf{L}(\mathbf{u}^{(2)}, t^{n} + \frac{\Delta t}{2})\right).$$

Another choice is third-order SSP multi-step method:

$$\mathbf{u}^{n+1} = \frac{16}{27}(\mathbf{u}^n + 3\Delta t \mathbf{L}(\mathbf{u}^n, t^n)) + \frac{11}{27}(\mathbf{u}^{n-3} + \frac{12}{11}\Delta t \mathbf{L}(\mathbf{u}^{n-3}, t^{n-3})).$$

More details can be found in [33, 34, 53].

#### **3.5** Numerical experiments

In this section, we provide numerical experiments to test the accuracy and stability of the high-order bound-preserving DG scheme. In all the examples, we choose N = 3, and consider fluid mixture with 3 components. Moreover, we use the third-order SSP Runge-Kutta discretization in time and  $P^2$  element in space. The computational domain is set to be  $\Omega = [0, 2\pi] \times [0, 2\pi]$ . To construct  $\Omega_h$ , we first equally divide  $\Omega$  into  $M \times M$  rectangles and the triangles are obtained by equally divide each rectangle into two. See Figure 3.2 for the mesh.



Figure 3.2: Triangular mesh (M = 10)

Example 3.5.1. We set the initial conditions as

$$c_{1,0}(x,y) = \frac{1}{6} \left(1 + \frac{1}{2} (\cos x + \cos y)\right), \quad c_{2,0}(x,y) = \frac{1}{3} \left(1 + \cos x \cos y\right),$$
  
$$c_{3,0}(x,y) = 1 - c_{1,0}(x,y) - c_{2,0}(x,y), \quad p_0(x,y) = \cos x \cos y - 1,$$

and the source variables are taken as

$$\tilde{c}_1(x, y, t) = \frac{1}{6} (1 + \frac{1}{2} e^{-\gamma t} (\cos x + \cos y - \frac{1}{2} \sin x \cos y - \frac{1}{2} \sin y \cos x)),$$
  

$$\tilde{c}_2(x, y, t) = \frac{1}{3} (1 + e^{-2\gamma t} (\cos x \cos y - \frac{1}{2} \sin^2 x \cos^2 y - \frac{1}{2} \cos^2 x \sin^2 y)),$$
  

$$\tilde{c}_3(x, y, t) = 1 - \tilde{c}_1(x, y, t) - \tilde{c}_2(x, y, t), \qquad q(x, y, t) = 2e^{-2t}.$$

Other parameters are chosen as

$$\phi(x, y) = \mu(c_1, c_2) = k(x, y) = a(x, y, c_1, c_2) = z_1 = z_2 = z_3 = 1,$$
  
 $D(u) = diag(\gamma, \gamma).$ 

It is easy to verify that the exact solutions are

$$c_1(x, y, t) = \frac{1}{6} \left(1 + \frac{1}{2} e^{-\gamma t} (\cos x + \cos y)\right), \quad c_2(x, y, t) = \frac{1}{3} \left(1 + e^{-2\gamma t} \cos x \cos y\right),$$
  
$$c_3(x, y, t) = 1 - c_1(x, y, t) - c_2(x, y, t), \quad p(x, y, t) = e^{-2t} (\cos x \cos y - 1).$$

In the numerical simulation, we choose  $\gamma = 0.01$ , final time T = 0.01 and  $\Delta t = 0.001h^2$  to reduce the time error. The computational results are shown in Table 3.1, illustrating the  $L^2$  error and convergence orders for  $c_1$  and  $c_2$  with and without bound-preserving technique. From the table, we observe optimal convergence rates. Therefore, the flux limiter and slope limiter do not degenerate the convergence order.

#### Example 3.5.2. We choose the initial conditions as

$$c_{1,0}(x,y) = \begin{cases} 1, & x \le \frac{\pi}{2}, y \le \frac{\pi}{2}, \\ 0, & otherwise. \end{cases}$$

$$c_{2,0}(x,y) = \begin{cases} 1, & x \ge \frac{3\pi}{2}, y \ge \frac{3\pi}{2}, \\ 0, & otherwise. \end{cases}$$

$$c_{3,0}(x,y) = 1 - c_{1,0}(x,y) - c_{2,0}(x,y) \quad and \quad p_0(x,y) = \cos(\frac{x}{2}) + \cos(\frac{y}{2}).$$

Other parameters are taken as

$$z_1 = z_2 = 1, z_3 = 10, q(x, y, t) = 0, \mathbf{D}(\mathbf{u}) = 0,$$
$$\mu(c_1, c_2) = k(x, y) = a(x, y, c_1, c_2) = \phi(x, y) = 1.$$

	$c_1$				$c_2$			
	no limiter		with limiter		no limiter		with limiter	
M	$L^2$ error	order	$L^2$ error	order	$L^2$ error	order	$L^2$ error	order
5	3.02e-3	_	4.61e-3	_	2.12e-2	_	2.39e-2	_
10	5.00e-4	2.59	5.30e-4	3.12	3.29e-3	2.69	3.47e-3	2.78
20	8.85e-5	2.50	8.86e-5	2.58	5.34e-4	2.63	5.34e-4	2.70
40	1.25e-5	2.82	1.25e-5	2.82	7.25e-5	2.88	7.25e-5	2.88
80	1.71e-6	2.87	1.71e-6	2.87	9.41e-6	2.95	9.41e-6	2.95
160	2.02e-7	3.09	2.02e-7	3.09	1.16e-6	3.02	1.16e-6	3.02

Table 3.1: Example 3.5.1: Accuracy test for  $c_1$  and  $c_2$  with and without boundpreserving technique.

We use this example to demonstrate the stability of the scheme. We choose  $\mathbf{D} = \mathbf{0}$ , then the diffusion term will not provide any dissipation to the scheme. We compute the components  $c_1$  and  $c_2$  at time T = 0.1s and T = 0.6s, respectively, with M = 40 and  $\Delta t = 0.001h^2$  ( $h = \frac{2\pi}{40}$ ). The numerical results are shown as Figure 3.3. From the figure we can see that the concentrations  $c_1$  and  $c_2$  are between 0 and 1. To test the effectiveness of the bound-preserving technique, we simulate the example without the bound-preserving limiters, and the numerical approximations blow up at about 0.003s even though we take time step size as small as  $\Delta t = 0.0001h^2$ . In [36], we demonstrated that the reason for the blow-up of the numerical approximations is the ill-posedness of the system. This example demonstrates the necessity of the bound-preserving technique in solving compressible miscible displacements in porous media.

**Example 3.5.3.** We investigate the displacement of 3-phase porous media flow in the five-spot arrangement of injection and production wells. The computational domain is a square region taken as quarter-of-a-five-spot pattern. The three phases are light oil  $c_1$  (with low viscosity and high compressibility), heavy oil  $c_2$  (with high viscosity and low compressibility) and water  $c_3$  (with medium viscosity and medium compressibility).



Figure 3.3: Example 3.5.2: Numerical approximations of  $c_1 \mbox{ and } c_2$ 

The initial concentrations of oil (water) are

$$c_{1,0}(x,y) = \begin{cases} 1, & x \le \frac{\pi}{2}, y \le \frac{\pi}{2}, \\ 0, & otherwise. \end{cases}$$
$$c_{2,0}(x,y) = \begin{cases} 0, & x \le \frac{\pi}{2}, y \le \frac{\pi}{2}, \\ 1, & otherwise. \end{cases}$$
$$c_{3,0}(x,y) = 0.$$

Therefore, the lower-left part of the region is light oil enrichment area while the other part is heavy oil enrichment area. Moreover, no water exists initially and the initial pressure is taken as 0 in the whole computational domain. To simulate the random perturbation of porosity and permeability around their average value, we choose the porosity and permeability as

$$\phi(x,y) = 0.5 + 0.05\sin(5x)\sin(5y)$$
 and  $k(x,y) = 1.0 + 0.1\cos(5x)\cos(5y)$ ,

respectively. Other parameters are taken as

$$\mu(c_1, c_2, c_3) = 0.4c_1 + 2.0c_2 + 1.0c_3,$$
  
 $z_1 = 1.2, \quad z_2 = 0.8, \quad z_3 = 1.0, \quad \mathbf{D} = diag(|\mathbf{u}|, |\mathbf{u}|)$ 

The injection well is located in lower-left corner and production well is located in upper-right corner, treated as  $\delta$  sources.

This example is used for petroleum production simulations. We compute the components  $c_1$  and  $c_2$  at time T = 0.2, 0.8 with M = 35 and  $\Delta t = 0.001h^2(h = \frac{2\pi}{35})$ . The distributions of  $c_1$ ,  $c_2$  and  $c_1 + c_2$  at different time are shown in figures



Figure 3.4: Example 3.5.3: Concentrations of  $c_1, c_2$  and  $c_1 + c_2$ .

3.4a-3.4f, respectively. From the figure we can see that  $c_1$ ,  $c_2$  and  $c_1 + c_2$  are all between 0 and 1.

**Example 3.5.4.** To show the significance of the bound-preserving technique in real petroleum production simulations, we choose the exact parameters in Example 3.5.3, except  $\mathbf{D} = \mathbf{0}$  in order to avoid any dissipation to the scheme which is resulted from the diffusion term.

This example is used for petroleum production simulations when diffusion effect is negligible. We compute the components  $c_1$  and  $c_2$  at time T = 0.2, 0.8with M = 35 and  $\Delta t = 0.001h^2(h = \frac{2\pi}{35})$ . The distributions of  $c_1$ ,  $c_2$ , and  $c_3$  at different time along diagonal y = x are shown in figures 3.5a-3.5f, respectively. From the figures we can see that the concentrations  $c_1, c_2$ , and  $c_3$  are between 0 and 1.

However, the numerical approximations without bound-preserving limiters blow up at about T = 0.25 if we take the same time step as before. The distribution of components along diagonal at time T = 0.1, 0.2 are shown in figures 3.6a-3.6f, from which we can observe strong oscillations and physically irrelevant values. Further experiments show that, even though we take the time step as small as  $\Delta t = 0.0001h^2$ , the numerical approximations still blow up at about T = 0.26, which implies the necessity of the bound-preserving technique.



Figure 3.5: Example 3.5.4: Concentrations of  $c_1, c_2$  and  $c_3$  with limiters



Figure 3.6: Example 3.5.4: Concentrations of  $c_1, c_2$  and  $c_3$  without limiters

### 3.6 Concluding remarks

In this paper, we constructed high-order bound-preserving DG methods for compressible miscible displacements in porous media on triangular meshes. We have applied the technique to the problem with multi-component fluid mixtures. Numerical simulations shown the accuracy and necessity of the bound-preserving technique

## Chapter 4

# Fourier analysis of local discontinuous Galerkin methods for linear parabolic equations on overlapping meshes<sup>1</sup>

#### Abstract

A new local discontinuous Galerkin (LDG) method for convection-diffusion equations on overlapping mesh was introduced in [28]. In the new method, the primary variable u and auxiliary variable  $p = u_x$  are solved on different meshes. The stability and suboptimal error estimates for problems with periodic boundary conditions were derived. Numerical experiments demonstrated that the con-

<sup>&</sup>lt;sup>1</sup>This chapter has been completed as an article to submit to Journal of Scientific Computing. Citation: N. Chuenjarern, Y. Yang (2019).
vergence rates cannot be improved if the dual mesh is constructed by using the midpoint of the primitive mesh. Several alternatives to gain optimal convergence rates were demonstrated in [28]. However, the reason for accuracy degeneration is still unclear. In this paper, we will use Fourier analysis to analyze the scheme for linear parabolic equations with periodic boundary conditions in one space dimension. We explicitly write out the error between the numerical and exact solutions, and investigate the reason for the accuracy degeneration. Moreover, we also find out some superconvergence points that may depend on the perturbation constant in the construction of the dual mesh. Since the current work is based on Fourier analysis, we only consider uniform meshes. Numerical experiments will be given to verify the theoretical analysis.

**Key Words**: Local Discontinuous Galerkin method, Fourier analysis, Error estimates, Superconvergence, Overlapping meshes

# 4.1 Introduction

In this paper, we apply local discontinuous Galerkin (LDG) method on overlapping meshes [28] for the following linear parabolic equations in one space dimension:

$$u_t - u_{xx} = 0, \quad x \in [0, 2\pi], \quad t > 0,$$
  
 $u(x, 0) = u_0(x), \quad x \in [0, 2\pi],$  (4.1.1)

subject to periodic boundary conditions.

The discontinuous Galerkin (DG) methods are a class of finite element methods with completely discontinuous piecewise polynomials as the numerical approximations. The DG method was first introduced in the framework of neutron linear transportation by Reed and Hill [51] in 1973. Subsequently, the Runge-Kutta discontinuous Galerkin (RKDG) methods were proposed for hyperbolic conservation laws in a series of papers [16, 17, 18, 19]. Later, in [20], Cockburn and Shu introduced the LDG method to solve the convection-diffusion equations. Their idea was motivated by Bassi and Rebay [2], where the compressible Navier-Stokes equations were successfully solved. In [20], the authors introduced an auxiliary variable q to represent the derivative of the primary variable u and thus rewrite (4.1.1) into the following system of first order equations

$$u_t - q_x = 0,$$
  
 $q - u_x = 0.$ 
(4.1.2)

Then one can solve u and p on the same mesh [20].

The LDG method is one of the most important numerical methods for convection diffusion equations. However, for some special convection-diffusion systems, such as chemotaxis model [43, 49] and miscible displacements in porous media [24, 25], the LDG methods are not easy to construct and analyze. In each of the two models, the convection term is the product of one of the primary variables and the derivative of the other primary variable. Most of the well established numerical fluxes for the convection terms, such as the upwind fluxes, cannot be applied, since the coefficients of the convection terms turn out to be discontinuous after the spatial discretization. It is well known that hyperbolic equations

with discontinuous coefficients are in general not well-posed [32, 40]. Therefore, the DG schemes may not be stable when applied to those model equations. Within the DG framework, there are three main different ways to bridge this gap. Firstly, in [77, 35, 46] the authors combined the convection terms and diffusion terms together and obtain the optimal error estimates. The idea was motivated by Wang et. al. [60, 61, 62], where  $u_x$  and the jump of u across the cell interfaces were proved to be bounded by q. Moreover, to make the numerical solutions to be physically relevant, we have to add a very large penalty which depends on the numerical approximations of the derivatives of the primary variables [46, 36, 13]. The second approach is to apply the flux-free numerical methods such as the Central DG (CDG) methods [47]. However, for CDG methods, we have to solve each equation in (4.1.2) on both the primary and dual meshes, which may double the computational cost. The last idea is to apply the Staggered DG (SDG) methods [14]. However, the method requires some continuity of the numerical approximations, and hence it is not easy to apply limiters to the numerical solutions. Recently, one of the authors in this paper introduced a new LDG method in [28], where we solve u and q on the primitive and dual meshes, respectively. To construct the dual mesh, we perturb the midpoint in each cell of the primary mesh, and use them as the cell interfaces of the dual mesh. We denote  $\alpha \in [-1/2, 1/2]$  as the perturbation constance, see [28] for more details. The stability and suboptimal error estimates of the new LDG scheme were also given in [28]. Since q is continuous across the cell interfaces in the primitive mesh, we can apply the upwind fluxes for the convection term for the complicated systems discussed above. Moreover, with the new idea, it is possible to construct thirdorder maximum-principle-preserving LDG methods on the overlapping meshes [27]. However, if the dual mesh is generated by the midpoint in each cell of the primitive mesh and piecewise odd order polynomials are applied, then the new method may not yield optimal convergence rates when applied to the pure linear parabolic equations [28]. This is the main reason why in the SDG method, the numerical approximations are required to be continuous across some of the cell interfaces. Several alternatives to gain the optimal convergence rates were also introduced in [28].

Unfortunately, it is still unclear why the accuracy given in [28] is not optimal. To solve this problem, we would like to apply Fourier analysis to quantitatively analyze the error between the numerical and exact solutions. In [80], the authors applied Fourier analysis to show the conditions of instability of some DG schemes for linear parabolic equations with periodic boundary conditions on uniform meshes. Later, this idea was extended to investigate the superconvergence of the DG scheme for linear hyperbolic equations in [91] and direct DG methods for parabolic equations in [84]. Motivated by the works given above, we take the initial condition as  $u_0(x) = e^{i\omega x}$  and rewrite the LDG scheme on overlapping meshes into an equivalent finite difference scheme. For simplicity, we only consider  $P^1$  and  $P^2$  polynomials, and the extension to high-order polynomials, though quite complicated, can be obtained following the same lines. We will write out the amplification matrix and explore the eigenvalues and eigenvectors. For  $P^1$  case, we anticipate two eigenvalues and only one of them should be physically relevant. We find that if  $\alpha = 0$ , the nonphysical eigenvalue does not decay during mesh refinement, and the scheme will generate a spurious wave that degenerate the accuracy of the scheme. However, if  $\alpha \neq 0$ , the nonphysical eigenvalue will decay exponentially fast during mesh refinement. Hence the nonphysical wave does not contribute much toward the numerical approximations, and keeps the accuracy. For the  $P^2$  case, no matter which  $\alpha$  we choose, both of the two nonphysical eigenvalues decay exponential fast during mesh refinement. Finally, by using Taylor's expansion, we can find out the leading term between the exact and numerical approximations, which gives us the order of accuracy of the scheme.

Moreover, with the quantitative error estimate, we can find some superconvergence points. Superconvergence of DG methods have been studied intensively for parabolic equations, see [9, 10, 76, 5] as an incomplete list. Different from the previous works, we have no idea about the position of the superconvergence points. For simplicity, we take k = 1 as an example. We choose two points in each cell to be determined, denoted as a and b, as the superconvergence points. Then we apply the Fourier analysis and write out the error between the numerical and exact solutions at the two points. The leading terms of the errors should be functions of  $\alpha$ , a and b. By setting the them to be zero, we can find the relationship among  $\alpha$ , a and b. Hence, for fixed  $\alpha$ , we can solve for a and b as the superconvergence points.

The rest of the paper is organized as follows. We first discuss the LDG scheme for one dimensional heat equation on overlapping mesh in Section 4.2. In Section 4.3, we demonstrate the quantitative error estimate using Fourier analysis for piecewise  $P^k$  polynomials with k = 1, 2. The superconvergence of the solution will be given in Section 4.4. In Section 4.5, some numerical experiments will be demonstrated to verify the theoretical results. We will end in Section 4.6 with concluding remarks.

## 4.2 LDG method on overlapping meshes

In this section, we present the formulation of the LDG method on overlapping meshes and study the linear parabolic equation (4.1.2).

#### 4.2.1 Overlapping meshes

Different from the LDG method introduced in [20] where u and q are solved on the same mesh, our new method solves (4.1.2) on two meshes, as shown in Figure 4.1.



Figure 4.1: Overlapping meshes

Let

$$0 = x_{\frac{1}{2}} < x_{\frac{3}{2}} < \dots < x_{N+\frac{1}{2}} = 2\pi$$

be a uniform partition of the domain  $[0, 2\pi]$  with mesh size  $h = \frac{2\pi}{N}$ . We denote

$$I_j = [x_{j-\frac{1}{2}}, x_{j+\frac{1}{2}}]$$
 and  $x_j = \frac{1}{2} \left( x_{j+\frac{1}{2}} + x_{j-\frac{1}{2}} \right), \quad j = 1, ..., N,$ 

as the cells and cell centers of the primitive mesh, respectively.

Based on the primitive mesh, we move each cell center within the corresponding cell to obtain the dual mesh, which is used to solve the auxiliary variable q. Then the cell interfaces of the dual mesh are given as

$$x_j^{\alpha} = x_j + \alpha h, \quad j = 1, ..., N,$$
 (4.2.3)

where  $-\frac{1}{2} \leq \alpha \leq \frac{1}{2}$  is the perturbation constant of the midpoint in the primitive mesh. In this paper, we assume  $\alpha$  to be a constant independent of the cells. Actually, the dual mesh contains all the cell  $J_j = [x_j^{\alpha}, x_{j+1}^{\alpha}]$ , where we define  $x_{N+1}^{\alpha} = x_1^{\alpha} + 2\pi$  due to the periodic boundary condition. For simplicity, we define  $J_0 = J_N = [0, x_1^{\alpha}] \cup [x_N^{\alpha}, 2\pi]$ .

#### 4.2.2 LDG scheme

In this subsection, we proceed to construct the LDG method on the overlapping meshes given above.

The finite element spaces are

$$V_h^k = \{ v : v | _{I_j} \in P^k(I_j), \ j = 1, ..., N \},$$
$$W_h^k = \{ v : v | _{J_j} \in P^k(J_j), \ j = 1, ..., N \},$$

where  $P^k(I_j)$  and  $P^k(J_j)$  denote the set of polynomials of degree up to k on  $I_j$ and  $J_j$ , respectively. It is easy to see that the elements in  $V_h^k$  and  $W_h^k$  are continuous across the cell interfaces on the dual and primitive meshes, respectively. Therefore, it may not be necessary to introduce the numerical fluxes in the LDG scheme. For simplicity, we also use u and q as the numerical approximations. Then the LDG scheme on overlapping meshes is to find  $u \in V_h^k$  and  $q \in W_h^k$  such that for any  $v \in V_h^k$  and  $w \in W_h^k$  we have

$$\int_{I_j} u_t v dx = -\int_{I_j} q v_x dx + q_{j+\frac{1}{2}} v_{j+\frac{1}{2}}^- - q_{j-\frac{1}{2}} v_{j-\frac{1}{2}}^+, \qquad (4.2.4)$$

$$\int_{J_j} qw dx = -\int_{J_j} uw_x dx + u_{j+1}^{\alpha} (w_{j+1}^{\alpha})^- - u_j^{\alpha} (w_j^{\alpha})^+, \qquad (4.2.5)$$

where  $q_{j+\frac{1}{2}} = q(x_{j+\frac{1}{2}}), u_{j+1}^{\alpha} = u(x_{j+1}^{\alpha}), v_{j-\frac{1}{2}}^{-} = v^{-}(x_{j-\frac{1}{2}})$  and  $(w_{j}^{\alpha})^{-} = w^{-}(x_{j}^{\alpha})$ . Likewise for  $v_{j-\frac{1}{2}}^{+}$  and  $(w_{j}^{\alpha})^{+}$ .

To implement the schemes (4.2.4) and (4.2.5), we define  $\phi_j^{\ell}(x)$  and  $\varphi_j^{\ell}(x)$ ,  $\ell = 0, 1, ..., k$ , as the local bases of  $P^k(I_j)$  and  $P^k(J_j)$ , respectively. Then we can represent the numerical solution as

$$u(x) = \sum_{\ell=0}^{k} u_j^{\ell} \phi_j^{\ell}(x), \quad x \in I_j,$$
(4.2.6)

$$q(x) = \sum_{\ell=0}^{k} q_{j}^{\ell} \varphi_{j}^{\ell}(x), \quad x \in J_{j}.$$
(4.2.7)

Substitute (4.2.6) and (4.2.7) into (4.2.4) and (4.2.5) to obtain

$$\frac{d\mathbf{u}_j}{dt} = \frac{1}{h^2} \left( A\mathbf{u}_{j-1} + B\mathbf{u}_j + C\mathbf{u}_{j+1} \right), \qquad (4.2.8)$$

where  $\mathbf{u}_j = (u_j^0, ..., u_j^k)^T$ , and A, B, C are  $(k+1) \times (k+1)$  constant matrices.

Following [91], we define

$$\begin{aligned} x_{j+\frac{2\ell-k}{2(k+1)}} &= x_j + \left(\frac{2\ell-k}{2(k+1)}\right)h, \quad \ell = 0, ..., k, \\ x_{j+\frac{2\ell+1}{2(k+1)}}^{\alpha} &= x_j^{\alpha} + \left(\frac{2\ell+1}{2(k+1)}\right)h, \quad \ell = 0, ..., k, \end{aligned}$$

as the grid points in cell  $I_j$  and  $J_j$ , respectively. Then we can construct Lagrange interpolation polynomials at the grid points as the local bases of  $P^k(I_j)$ , and  $P^k(J_j)$ . With the Lagrange bases,  $\mathbf{u}_j = (u_j^0, ..., u_j^k)^T$  turns out to be the point values of the numerical approximations at the grid points in cell  $I_j$ . Hence, we rewrite the LDG scheme into a finite difference scheme.

**Remark 4.2.1.** To apply Fourier analysis, it is not necessary to choose globally uniformly distributed grid points as we treat the point values at the grid points in each cell as a vector. Therefore, we only need to construct uniform cells. We will choose other grid points to find out the superconvergence points in Section 4.4.

### 4.3 Error analysis

In this section, we proceed to analyze the error between the numerical and exact solutions at the grid points given in Section 4.2. Numerical experiments in [28] demonstrated that, the accuracy may not be optimal only if odd order polynomials were applied. Therefore, we only analyze the LDG scheme with piecewise  $P^1$  and  $P^2$  polynomials in this section to find out the reason of accuracy degeneration.

### 4.3.1 The $P^1$ case

In this subsection, we present the details of error analysis for the piecewise linear case i.e. k = 1. The local basis functions on cell  $I_j$  are  $\phi_{j-\frac{1}{4}}(x)$ ,  $\phi_{j+\frac{1}{4}}(x)$ , which

are Lagrange polynomials based on  $x_{j-\frac{1}{4}}$ ,  $x_{j+\frac{1}{4}}$ . Also, the local basis functions on cell  $J_j$  are  $\varphi_{j+\frac{1}{4}}(x)$ ,  $\varphi_{j+\frac{3}{4}}(x)$ , which are Lagrange polynomials based on  $x_{j+\frac{1}{4}}^{\alpha}$ ,  $x_{j+\frac{3}{4}}^{\alpha}$ . Then the solutions can be written as

$$u(x) = u_{j-\frac{1}{4}}\phi_{j-\frac{1}{4}}(x) + u_{j+\frac{1}{4}}\phi_{j+\frac{1}{4}}(x), \quad x \in I_j,$$
  
$$q(x) = q_{j+\frac{1}{4}}^{\alpha}\varphi_{j+\frac{1}{4}}(x) + q_{j+\frac{3}{4}}^{\alpha}\varphi_{j+\frac{3}{4}}(x), \quad x \in J_j.$$

For  $j = 1, \dots, N$ , the finite difference representation of the LDG scheme (4.2.5) is

$$\begin{pmatrix} q_{j+\frac{1}{4}}^{\alpha} \\ q_{j+\frac{3}{4}}^{\alpha} \end{pmatrix} = \frac{1}{4h} \left[ Q_1 \begin{pmatrix} u_{j-\frac{1}{4}} \\ u_{j+\frac{1}{4}} \end{pmatrix} + Q_2 \begin{pmatrix} u_{j+\frac{3}{4}} \\ u_{j+\frac{5}{4}} \end{pmatrix} \right],$$

where

$$Q_{1} = \begin{pmatrix} -5 + 14\alpha + 12\alpha^{2} & 1 - 26\alpha - 12\alpha^{2} \\ 1 + 2\alpha - 12\alpha^{2} & -5 + 10\alpha + 12\alpha^{2} \end{pmatrix},$$
$$Q_{2} = \begin{pmatrix} 5 + 10\alpha - 12\alpha^{2} & -1 + 2\alpha + 12\alpha^{2} \\ -1 - 26\alpha + 12\alpha^{2} & 5 + 14\alpha - 12\alpha^{2} \end{pmatrix}.$$

Moreover, the finite difference representation of the LDG scheme (4.2.4) can be written as

$$\begin{pmatrix} u'_{j-\frac{1}{4}} \\ u'_{j+\frac{1}{4}} \end{pmatrix} = \frac{1}{4h} \left[ U_1 \begin{pmatrix} q^{\alpha}_{j-\frac{1}{4}} \\ q^{\alpha}_{j-\frac{3}{4}} \end{pmatrix} + U_2 \begin{pmatrix} q^{\alpha}_{j+\frac{1}{4}} \\ q^{\alpha}_{j+\frac{3}{4}} \end{pmatrix} \right],$$

where

$$U_{1} = \begin{pmatrix} -5 - 14\alpha + 12\alpha^{2} & 1 + 26\alpha - 12\alpha^{2} \\ 1 - 2\alpha - 12\alpha^{2} & -5 - 10\alpha + 12\alpha^{2} \end{pmatrix},$$
$$U_{2} = \begin{pmatrix} 5 - 10\alpha - 12\alpha^{2} & -1 - 2\alpha + 12\alpha^{2} \\ -1 + 26\alpha + 12\alpha^{2} & 5 - 14\alpha - 12\alpha^{2} \end{pmatrix}.$$

Here u' denotes the time derivative of u. After some simply algebra, we can obtain

$$\frac{d\mathbf{u}_{j}}{dt} = \frac{1}{h^{2}} \left( A\mathbf{u}_{j-1} + 2B\mathbf{u}_{j} + C\mathbf{u}_{j+1} \right), \qquad (4.3.9)$$

with

$$A = \frac{1}{8} \begin{pmatrix} 13 + 14\alpha - 144\alpha^2 - 168\alpha^3 + 144\alpha^4 & -5 - 2\alpha + 384\alpha^2 + 24\alpha^3 - 144\alpha^4 \\ -5 + 2\alpha + 48\alpha^2 - 24\alpha^3 - 144\alpha^4 & 13 - 14\alpha - 96\alpha^2 + 168\alpha^3 + 144\alpha^4 \end{pmatrix},$$
  

$$B = \frac{1}{8} \begin{pmatrix} -13 - 14\alpha - 168\alpha^2 + 168\alpha^3 - 144\alpha^4 & 5 + 2\alpha + 72\alpha^2 - 24\alpha^3 + 144\alpha^4 \\ 5 - 2\alpha + 72\alpha^2 + 24\alpha^3 + 144\alpha^4 & -13 + 14\alpha - 168\alpha^2 - 168\alpha^3 - 144\alpha^4 \end{pmatrix},$$
  

$$C = \frac{1}{8} \begin{pmatrix} 13 + 14\alpha - 96\alpha^2 - 168\alpha^3 + 144\alpha^4 & -5 - 2\alpha + 48\alpha^2 + 24\alpha^3 - 144\alpha^4 \\ -5 + 2\alpha + 384\alpha^2 - 24\alpha^3 - 144\alpha^4 & 13 - 14\alpha - 144\alpha^2 + 168\alpha^3 + 144\alpha^4 \end{pmatrix}.$$
  
(4.3.10)

Next, we will use the standard Fourier analysis to solve (4.3.9). We consider a general Fourier mode and assume

$$\begin{pmatrix} u_{j-\frac{1}{4}}(t) \\ u_{j+\frac{1}{4}}(t) \end{pmatrix} = \begin{pmatrix} \hat{u}_{-\frac{1}{4}}(t) \\ \hat{u}_{+\frac{1}{4}}(t) \end{pmatrix} e^{i\omega x_j}.$$

Substitute the above into (4.3.9), we get the following ODE system

$$\begin{pmatrix} \hat{u}_{-\frac{1}{4}}'(t) \\ \hat{u}_{+\frac{1}{4}}'(t) \end{pmatrix} = G \begin{pmatrix} \hat{u}_{-\frac{1}{4}}(t) \\ \hat{u}_{+\frac{1}{4}}(t) \end{pmatrix},$$

where the amplification matrix G is

$$G = \frac{1}{h^2} \left( A e^{-i\xi} + 2B + C e^{i\xi} \right), \quad \xi = \omega h, \tag{4.3.11}$$

with the matrices A, B, C given in (4.3.10). For simplicity, we assume  $\omega = 1$ , then  $\xi = h$ . The two eigenvalues of the amplification matrices are

$$\lambda_{1,2} = \frac{1}{8h^2} \left( \gamma \mp \sqrt{\beta} \right), \qquad (4.3.12)$$

where

$$\begin{split} \gamma &= 13 - 26e^{i\xi} + 13e^{2i\xi} + 144\alpha^4 (-1 + e^{i\xi})^2 - 24\alpha^2 (5 + 14e^{i\xi} + 5e^{2i\xi}) \\ \beta &= 25(-1 + e^{i\xi})^4 + 20736\alpha^8 (-1 + e^{i\xi})^4 - 6912\alpha^6 (-1 + e^{i\xi})^2 (5 + 14e^{i\xi} + 5e^{2i\xi}) \\ &- 48\alpha^2 (-1 + e^{i\xi})^2 (41 + 38e^{i\xi} + 41e^{2i\xi}) \\ &+ 288\alpha^4 (55 + 260e^{i\xi} + 522e^{2i\xi} + 260e^{3i\xi} + 55e^{4i\xi}). \end{split}$$

$$(4.3.13)$$

Moreover, the corresponding eigenvectors are

$$V_{1,2} = \begin{pmatrix} \Gamma \pm \sqrt{\beta} \\ \Theta \end{pmatrix}, \qquad (4.3.14)$$

where

$$\Gamma = -14\alpha(-1 + e^{i\xi})^2 + 168\alpha^3(-1 + e^{i\xi})^2 - 24\alpha^2(-1 + e^{2i\xi})$$
$$\Theta = 5(-1 + e^{i\xi})^2 - 2\alpha(-1 + e^{i\xi})^2 + 24\alpha^3(-1 + e^{i\xi})^2 + 144$$
$$\alpha^4(-1 + e^{i\xi})^2 - 48\alpha^2(1 + 3e^{i\xi} + 8e^{2i\xi})$$

with  $\beta$  given in (4.3.13). Then the general solution of the ODE system (4.3.9) is

$$\begin{pmatrix} \hat{u}_{-\frac{1}{4}}(t) \\ \hat{u}_{+\frac{1}{4}}(t) \end{pmatrix} = C_{11}e^{\lambda_1 t}V_1 + C_{12}e^{\lambda_2 t}V_2, \qquad (4.3.15)$$

where the constants  $C_{11}$  and  $C_{12}$  are determined by the initial condition

$$\begin{pmatrix} \hat{u}_{-\frac{1}{4}}(0)\\ \hat{u}_{+\frac{1}{4}}(0) \end{pmatrix} = \begin{pmatrix} e^{-\frac{i\xi}{4}}\\ e^{\frac{i\xi}{4}} \end{pmatrix}.$$

Therefore, we have the explicit solution of the LDG scheme with  $P^1$  polynomials. The quantitative error will arise when we compare the numerical approximations with the exact solutions U(x, t) at the grid points defined by

$$\begin{split} ||e_{-\frac{1}{4}}||_{\infty} &= \max_{1 \le j \le N} |U(x_{j-\frac{1}{4}},t) - u_{j-\frac{1}{4}}(t)|, \\ ||e_{+\frac{1}{4}}||_{\infty} &= \max_{1 \le j \le N} |U(x_{j+\frac{1}{4}},t) - u_{j+\frac{1}{4}}(t)|. \end{split}$$

However, it is not easy to write the analytical form the of errors. Therefore, we would like to apply Taylor's expansion with respect to  $\xi$  at  $\xi = 0$ . Then two eigenvalues of the amplification matrix can be rewritten as

1. For  $\alpha = 0$ ,

$$\lambda_1 = -\frac{9}{4} + \frac{3}{16}\xi^2 - \frac{1}{160}\xi^4 + \frac{1}{8960}\xi^6 + O(\xi^7)$$
  
$$\lambda_2 = -1 + \frac{1}{12}\xi^2 - \frac{1}{360}\xi^4 + \frac{1}{20160}\xi^6 + O(\xi^7).$$

2. For  $\alpha \neq 0$ ,

$$\begin{split} \lambda_1 &= -\frac{9}{4} + 30\alpha^2 - 36\alpha^4 - \frac{144\alpha^2}{\xi^2} + \xi^2 \left(\frac{13}{48} - \frac{5\alpha^2}{2} + 3\alpha^4\right) \\ &- \xi^4 \left(\frac{1}{360} + \frac{5}{6912\alpha^2} - \frac{\alpha^2}{16} + \frac{\alpha^4}{10}\right) \\ &+ \xi^6 \left(\frac{383}{483840} + \frac{25}{3981312\alpha^4} - \frac{1}{13824\alpha^2} - \frac{5\alpha^2}{1008} + \frac{47\alpha^4}{6720}\right) + O(\xi^7), \\ \lambda_2 &= -1 - \xi^4 \left(\frac{1}{160} - \frac{5}{9612\alpha^2} - \frac{\alpha^2}{48}\right) \\ &- \xi^6 \left(\frac{61}{96768} + \frac{25}{3981312\alpha^4} - \frac{1}{13824\alpha^2} - \frac{\alpha^2}{288} + \frac{\alpha^4}{192}\right) + O(\xi^7)). \end{split}$$

It is easy to see that  $\lambda_2$  is the physical eigenvalue, while  $\lambda_1$  is the nonphysical one. For  $\alpha \neq 0$ , the fourth term in  $\lambda_1$  makes the first term in (4.3.15) decay

exponentially fast. In the analysis, we only need to take  $\lambda_2$  into account and omit the contribution of  $\lambda_1$ . However, for  $\alpha = 0$ , the contribution of  $\lambda_1$  is not negligible, leading to a nonphysical wave. With some basic computation, we have the quantitative error:

For 
$$\alpha = 0$$
,

$$\begin{aligned} ||e_{+\frac{1}{4}}||_{\infty} &= \frac{1}{4}e^{-t}(-1+e^{-\frac{5}{4}t})\xi \\ &+ \frac{e^{-3t}\left[(-3+16t^2-6e^{-\frac{5}{4}t}(-1+9t)+3e^{-\frac{5}{2}t})(-1+18t)\right]}{1152(-1+e^{-\frac{5}{4}t})}\xi^3 + O(\xi^4). \end{aligned}$$
(4.3.16)

For  $\alpha \neq 0$ ,

$$\begin{aligned} ||e_{+\frac{1}{4}}||_{\infty} &= \frac{(-1+12\alpha^2)e^{-t}}{9\alpha}\xi^2 \\ &+ \left[75 - 940\alpha^2 - 4080\alpha^4 + 72000\alpha^6 - 103680\alpha^8 - 138240\alpha^7(-1+t) \right. \\ &+ 80\alpha(-1+5t) + 2304\alpha^5(-15+23t) \\ &- 192\alpha^3(-15+43t) \right] \frac{(-1+12\alpha^2)e^{-t}}{552960\alpha(\alpha-12\alpha^3)^2}\xi^4 \\ &+ O(\xi^5). \end{aligned}$$

$$(4.3.17)$$

The error  $||e_{-\frac{1}{4}}||_{\infty}$  is similar, so we omit it here. From the error, we can see that for  $\alpha = 0$  the error is indeed first order accurate, while it is second order accurate for  $\alpha \neq 0$ .

## 4.3.2 The $P^2$ case

In this subsection, we will use the same approach given in Subsection 4.3.1 to demonstrate the error analysis for the  $P^2$  case. Denote the local basis functions for cell  $I_j$  as  $\phi_{j-\frac{1}{3}}(x), \phi_j(x), \phi_{j+\frac{1}{3}}(x)$ , which are Lagrangian polynomials based on the points  $x_{j-\frac{1}{3}}, x_j, x_{j+\frac{1}{3}}$ . The local basis functions for cell  $J_j$  are  $\varphi_{j+\frac{1}{6}}(x)$ ,  $\varphi_{j+\frac{1}{2}}(x), \varphi_{j+\frac{5}{6}}(x)$ , which are Lagrangian polynomials based on the points  $x_{j+\frac{1}{6}}^{\alpha}$ ,  $x_{j+\frac{1}{2}}^{\alpha}, x_{j+\frac{5}{6}}^{\alpha}$ . Then the solutions can be represented as

$$\begin{split} u(x) &= u_{j-\frac{1}{4}}\phi_{j-\frac{1}{4}}(x) + u_{j}\phi_{j}(x) + u_{j+\frac{1}{4}}\phi_{j+\frac{1}{4}}(x), \quad x \in I_{j}, \\ q(x) &= q_{j+\frac{1}{6}}^{\alpha}\varphi_{j+\frac{1}{6}}(x) + q_{j+\frac{1}{2}}^{\alpha}\varphi_{j+\frac{1}{2}}(x) + q_{j+\frac{5}{6}}^{\alpha}\varphi_{j+\frac{5}{6}}(x), \quad x \in J_{j}. \end{split}$$

It is quite complicated to write out the exact forms the eigenvalues and eigenvectors for the  $P^2$  case. Therefore, we will only consider two special cases, namely  $\alpha = 0$  and  $\alpha = \frac{1}{2}$ .

Following the same procedure given in Subsection 4.3.1, the LDG scheme can be written into the matrix form (4.2.8) with

$$\mathbf{u}_{j} = \left(u_{j-\frac{1}{3}}, u_{j}, u_{j+\frac{1}{3}}\right)^{T}, \qquad (4.3.18)$$

and for  $\alpha = 0$ ,

$$A = \frac{1}{512} \begin{pmatrix} -385 & 1674 & 1063 \\ -14 & -318 & 1755 \\ 95 & -310 & 7 \end{pmatrix},$$
  

$$B = \frac{1}{256} \begin{pmatrix} -2211 & 278 & 861 \\ 585 & -2562 & 585 \\ 861 & 278 & -2211 \end{pmatrix},$$
  

$$C = \frac{1}{512} \begin{pmatrix} 7 & -310 & 95 \\ 1755 & -318 & -45 \\ 1755 & -318 & -45 \end{pmatrix},$$
  
(4.3.19)

and for  $\alpha = \frac{1}{2}$ ,

$$A = \frac{1}{16} \begin{pmatrix} 153 & -510 & 765 \\ 9 & -20 & 45 \\ -15 & 50 & -75 \end{pmatrix},$$
  
$$B = \frac{1}{4} \begin{pmatrix} -151 & 42 & 13 \\ 63 & -186 & 171 \\ -13 & 226 & -311 \end{pmatrix},$$
  
$$C = \frac{1}{16} \begin{pmatrix} -29 & 6 & -1 \\ -261 & 54 & -9 \\ 667 & -138 & 23 \end{pmatrix}.$$
  
(4.3.20)

Again, the standard Fourier analysis will be applied and assume

$$\begin{pmatrix} u_{j-\frac{1}{3}}(t) \\ u_{j}(t) \\ u_{j+\frac{1}{3}}(t) \end{pmatrix} = \begin{pmatrix} \hat{u}_{-\frac{1}{3}}(t) \\ \hat{u}_{0}(t) \\ \hat{u}_{+\frac{1}{3}}(t) \end{pmatrix} e^{i\omega x_{j}}.$$
(4.3.21)

For simplicity, we also assume  $\omega = 1$ . Substituting the above into (4.2.8), we can obtain the ODE system

$$\begin{pmatrix} \hat{u}'_{j-\frac{1}{3}}(t) \\ \hat{u}'_{j}(t) \\ \hat{u}'_{j+\frac{1}{3}}(t) \end{pmatrix} = G \begin{pmatrix} \hat{u}_{-\frac{1}{3}}(t) \\ \hat{u}_{0}(t) \\ \hat{u}_{+\frac{1}{3}}(t) \end{pmatrix},$$
(4.3.22)

where the amplification matrix G is given by (4.3.11) with A, B and C defined in (4.3.19) or (4.3.20) for  $\alpha = 0$  and  $\alpha = \frac{1}{2}$ , respectively. Denote  $\lambda_i$  and  $V_i$ , i =1, 2, 3, to be the eigenvalues and corresponding eigenvectors of G, respectively. Then for  $\alpha = 0$ ,

$$\begin{split} \lambda_1 &= -1 - \frac{596651i}{3072} \xi^3 + \frac{4058334841}{3276800} \xi^4 + \frac{3345594197i}{737280} \xi^5 \\ &- \frac{405767495830801}{33030144000} \xi^6 + O(\xi^7) \\ \lambda_{2,3} &= \frac{151}{128} - \frac{15}{\xi^2} \pm \frac{7\sqrt{15}}{8\xi} \mp \frac{2419\sqrt{15}}{20480} \xi \mp \frac{29}{512} \xi^2 + \left(\frac{596651i}{6144} \mp \frac{13228737901}{20971520\sqrt{15}}\right) \xi^3 \\ &\mp \frac{\left(36524902209 + 16115508640\sqrt{15}i\right)}{58982400} \xi^4 \\ &+ \frac{\left(-5481421532364800i \pm 2436959051302733\sqrt{15}\right)}{2415919104000} \xi^5 \\ &+ \frac{\left(405767493603601 \pm 180998522537910\sqrt{15}i\right)}{66060288000} \xi^6 + O(\xi^7). \end{split}$$

and

$$V_{1} = \begin{pmatrix} -720 - 1200i\xi + 1204\xi^{2} + 897i\xi^{3} + O(\xi^{4}) \\ -720 - 1440i\xi + 1644\xi^{2} + 1368i\xi^{3} + O(\xi^{4}) \\ -720 - 1680i\xi + 2164\xi^{2} + 1999i\xi^{3} + O(\xi^{4}) \end{pmatrix}, \quad V_{2,3} = \begin{pmatrix} \Gamma \\ \Theta \\ \Lambda \end{pmatrix};$$

where

$$\begin{split} &\Gamma = -53760(12i \mp \sqrt{15})z + 224(5095 \pm 232i\sqrt{15})\xi^2 + (1198544i \mp 22769\sqrt{15})\xi^3 + O(\xi^4) \\ &\Theta = \mp 161280\sqrt{15}\xi \mp 3360(59 + 96\sqrt{15}i)\xi^2 - (396480i \mp 367571\sqrt{15})\xi^3 + O(\xi^4) \\ &\Lambda = 53760(12i \pm \sqrt{15})\xi - 224(6425i \mp 728\sqrt{15}i)\xi^2 - 3(598128i \pm 81659\sqrt{15})\xi^3 + O(\xi^4) \end{split}$$

and for 
$$\alpha = \frac{1}{2}$$
,  
 $\lambda_1 = -1 - \frac{1144i}{3}\xi^3 + \frac{14300}{9}\xi^4 + \frac{110783530i}{29187}\xi^5 - \frac{42485046399193}{6401682000}\xi^6 + O(\xi^7)$   
 $\lambda_{2,3} = -1 \pm \frac{38}{\sqrt{69}} - \frac{6(13 \mp \sqrt{69})}{\xi^2} + \left(\frac{1}{8} \mp \frac{6821}{1656\sqrt{69}}\right)\xi^2 - \frac{44i}{3}\left(13 \mp 3\sqrt{69}\right)\xi^3$   
 $- \left(\frac{572003}{720} \mp \frac{38588405903}{3427920\sqrt{69}}\right)\xi^4 - 11i\frac{(5035615 \mp 894279\sqrt{69})}{29187}\xi^5$   
 $+ \frac{(502441935138015571557 \mp 74298976612868552411\sqrt{69})}{151416730953936000}\xi^6 + O(\xi^7).$ 

and

$$V_{1} = \begin{pmatrix} 3600 + 6000i\xi - 5246\xi^{2} - 3221i\xi^{3} + O(\xi^{4}) \\ 3600 + 7200i\xi - 7446\xi^{2} - 5313i\xi^{3} + O(\xi^{4}) \\ 3600 + 8400i\xi - 10046\xi^{2} - 8245i\xi^{3} + O(\xi^{4}) \end{pmatrix}, \quad V_{2,3} = \begin{pmatrix} \Gamma \\ \Theta \\ \Lambda \end{pmatrix};$$

where

$$\begin{split} \Gamma &= 1656(141 \mp 7\sqrt{69}) + 138i(1507 \mp 39\sqrt{69})\xi - 10(7222 \pm 749\sqrt{69})\xi^2 + O(\xi^3) \\ \Theta &= 24840(3 \mp \sqrt{69}) + 414i(269 \mp 113\sqrt{69})\xi - 6(6532 \pm 6957\sqrt{69})\xi^2 O(\xi^3) \\ \Lambda &= \frac{1}{3} \left( -4968(171 \mp 17\sqrt{69}) - 414i(3293 \mp 361\sqrt{69})\xi + (937572 \mp 117690\sqrt{69})\xi^2 \right) + O(\xi^3) \end{split}$$

Then the general solution of the ODE system (4.3.22) is

$$\begin{pmatrix} \hat{u}_{-\frac{1}{3}}(t) \\ \hat{u}_{0}(t) \\ \hat{u}_{+\frac{1}{3}}(t) \end{pmatrix} = C_{21}e^{\lambda_{1}t}V_{1} + C_{22}e^{\lambda_{2}t}V_{2} + C_{23}e^{\lambda_{3}t}V_{3}, \qquad (4.3.23)$$

where the constants  $C_{21}$ ,  $C_{22}$  and  $C_{23}$  are determined by the initial condition

$$\begin{pmatrix} \hat{u}_{-\frac{1}{3}}(0) \\ \hat{u}_{0}(0) \\ \hat{u}_{+\frac{1}{3}}(0) \end{pmatrix} = \begin{pmatrix} e^{-\frac{i\xi}{3}} \\ 1 \\ e^{\frac{i\xi}{3}} \end{pmatrix}.$$

We can see that,  $\lambda_1$  is the physical eigenvalue while  $\lambda_{2,3}$  are the nonphysical ones. Moreover, it is easy to observe that the second and third terms in (4.3.23) are decreasing exponentially fast with respect to the mesh size h, hence we can ignore the contribution from them. With some basic computation, we can obtain the quantitative error estimates:

for  $\alpha = 0$ ,

$$\begin{split} ||e_{-\frac{1}{3}}||_{\infty} &:= \max_{1 \leq j \leq N} |U(x_{j-\frac{1}{3}}, t) - u_{j-\frac{1}{3}}(t)| \\ &= \frac{(832 + 80547885t)e^{-t}}{414720} \xi^3 \\ &+ \frac{1}{1019215872000(832 + 80547885t)} [(10979996079226880 \\ &+ 1066737149124583495680t + 48349276106069021512077t^2)e^{-t}] \xi^5 \\ &+ O(\xi^6), \\ ||e_0||_{\infty} &:= \max_{1 \leq j \leq N} |U(x_j, t) - u_j(t)| \\ &= \frac{596651te^{-t}}{3072} \xi^3 \\ &+ \frac{1}{25128767324160000t^2} [(26214400 \\ &+ 976011547208325120t - 14799288676482712431t^2)te^{-t}] \xi^5 + O(\xi^6), \\ ||e_{+\frac{1}{3}}||_{\infty} &:= \max_{1 \leq j \leq N} |U(x_{j+\frac{1}{3}}, t) - u_{j+\frac{1}{3}}(t)| \\ &= \frac{(-832 + 80547885t)e^{-t}}{414720} \xi^3 \\ &+ \frac{1}{1019215872000(832 - 80547885t)} [(-10979996079226880 \\ &+ 1066737149124583495680t + 48349276106069021512077t^2)e^{-t}] \xi^5 \\ &+ O(\xi^6), \end{split}$$

and for 
$$\alpha = \frac{1}{2}$$
,  

$$\begin{aligned} ||e_{-\frac{1}{3}}||_{\infty} &= \frac{(-1+494208t)e^{-t}}{1296}\xi^{3} \\ &+ \frac{1}{3362342400(1-494208t)} \left[ (-85477574647 \\ &+ 42232234477694976t + 806689123688448000t^{2})e^{-t} \right] \xi^{5} \\ &+ O(\xi^{6}), \\ ||e_{0}||_{\infty} &= \frac{(-1+91520t)e^{-t}}{240}\xi^{3} \\ &+ \frac{1}{16811712000(1-91520t)} \left[ (512868994643 \\ &- 47003527618544640t + 746934373785600000t^{2})e^{-t} \right] \xi^{5} \\ &+ O(\xi^{6}), \\ ||e_{+\frac{1}{3}}||_{\infty} &= \frac{(23+2471040t)e^{-t}}{6480}\xi^{3} \\ &+ \frac{1}{16811712000(23+2471040t)} \left[ 13(-151230865483 \\ &+ -16134718463170560t + 1551325237862400000t^{2})e^{-t} \right] \xi^{5} \\ &+ O(\xi^{6}). \end{aligned}$$

We can see that, both cases yield optimal convergence rates.

# 4.4 Superconvergence

In this section, we will consider the one-dimensional linear parabolic equation and investigate the superconvergence of the LDG scheme. We take the perturbation constant  $\alpha \neq 0$ . For simplicity, the finite element spaces are made up of piecewise linear polynomials. The extension to high-order cases, though quite complicated, can be obtain following the same lines. The Fourier analysis technique discussed in Section 4.3 will be used to investigate a relationship between the perturbation constant  $\alpha$  of the dual cells and the superconvergence points. However, the superconvergence property discussed in this section only works for uniform meshes. For general random meshes, the superconvergence points are not easy to derive.

The basis functions in this section are different from those discussed in Section 4.3. We are using  $\phi_{j-\frac{1}{2}}(x)$ ,  $\phi_{j+\frac{1}{2}}(x)$ , which are Lagrange polynomials based on the grid points  $x_{j-\frac{1}{2}}$ ,  $x_{j+\frac{1}{2}}$  as the local basis functions for cell  $I_j$ . Also, the local basis functions for cell  $J_j$  are  $\varphi_j(x)$ ,  $\varphi_{j+1}(x)$ , which are the Lagrange polynomials based on the grid points  $x_j^{\alpha}$ ,  $x_{j+1}^{\alpha}$ . Then the solutions can be represented as

$$u(x) = u_{j-\frac{1}{2}}\phi_{j-\frac{1}{2}}(x) + u_{j+\frac{1}{2}}\phi_{j+\frac{1}{2}}(x), \quad x \in I_j,$$
$$q(x) = q_j^{\alpha}\varphi_j(x) + q_{j+1}^{\alpha}\varphi_{j+1}(x), \quad x \in J_j.$$

Following the same analysis in Section 4.3, the LDG scheme can be written into the matrix form (4.3.9) with

$$A = \frac{1}{8} \begin{pmatrix} 13 + 16\alpha - 24\alpha^2 - 192\alpha^3 + 144\alpha^4 & -5 + 8\alpha + 408\alpha^2 - 96\alpha^3 - 144\alpha^4 \\ -5 - 8\alpha + 24\alpha^2 + 96\alpha^3 - 144\alpha^4 & 13 - 16\alpha - 216\alpha^2 + 192\alpha^3 + 144\alpha^4 \end{pmatrix},$$
  

$$B = \frac{1}{8} \begin{pmatrix} -13 - 16\alpha - 216\alpha^2 + 192\alpha^3 + 144\alpha^4 & 5 - 8\alpha + 72\alpha^2 + 96\alpha^3 + 144\alpha^4 \\ 5 + 8\alpha + 72\alpha^2 - 96\alpha^3 + 144\alpha^4 & -13 + 16\alpha - 168\alpha^2 - 192\alpha^3 - 144\alpha^4 \end{pmatrix},$$
  

$$C = \frac{1}{8} \begin{pmatrix} 13 + 16\alpha - 216\alpha^2 - 192\alpha^3 + 144\alpha^4 & -5 + 8\alpha + 24\alpha^2 - 96\alpha^3 - 144\alpha^4 \\ -5 - 8\alpha + 408\alpha^2 + 96\alpha^3 - 144\alpha^4 & 13 - 16\alpha - 24\alpha^2 + 192\alpha^3 + 144\alpha^4 \end{pmatrix}.$$
  

$$(4.4.24)$$

To observe the superconvergence property, we would like the initial error to be superconvergent at the superconvergence points. Therefore, we can take the initial discretization to be the polynomial interpolation at the superconvergence points. To locate those points, we first map each physical cell into the reference interval  $\left[-\frac{1}{2}, \frac{1}{2}\right]$ , and denote the superconvergence points in the reference interval to be *a* and *b*. Then we map the two points back to the physical cell, and denote them as  $x_j^a$  and  $x_j^b$  in cell  $I_j$ . It is easy to check that

$$x_j^a = x_j + ah, \quad x_j^b = x_j + bh.$$

Then, the initial numerical solution in cell  $I_j$  would be

$$y = \frac{e^{i\omega x_{j}^{b}} - e^{i\omega x_{j}^{a}}}{x_{j}^{b} - x_{j}^{a}}x + \frac{x_{j}^{b}e^{i\omega x_{j}^{a}} - x_{j}^{a}e^{i\omega x_{j}^{b}}}{x_{j}^{b} - x_{j}^{a}}$$

We evaluate the above interpolation at  $x_{j-\frac{1}{2}},\,x_{j+\frac{1}{2}}$  to obtain

$$y(x_{j-\frac{1}{2}}) = \frac{(b+\frac{1}{2})e^{i\xi a} - (a+\frac{1}{2})e^{i\xi b}}{b-a}e^{i\omega x_j}$$
$$y(x_{j+\frac{1}{2}}) = \frac{(b-\frac{1}{2})e^{i\xi a} - (a-\frac{1}{2})e^{i\xi b}}{b-a}e^{i\omega x_j}$$

Then the initial condition of a general Fourier mode

$$\begin{pmatrix} u_{j-\frac{1}{2}}(t) \\ u_{j+\frac{1}{2}}(t) \end{pmatrix} = \begin{pmatrix} \hat{u}_{-\frac{1}{2}}(t) \\ \hat{u}_{+\frac{1}{2}}(t) \end{pmatrix} e^{i\omega x_j},$$
(4.4.25)

can be written as

$$\begin{pmatrix} \hat{u}_{-\frac{1}{2}}(0) \\ \hat{u}_{+\frac{1}{2}}(0) \end{pmatrix} = \begin{pmatrix} \frac{(b+\frac{1}{2})e^{i\xi a} - (a+\frac{1}{2})e^{i\xi b}}{b-a} \\ \frac{(b-\frac{1}{2})e^{i\xi a} - (a-\frac{1}{2})e^{i\xi b}}{b-a} \end{pmatrix}.$$
(4.4.26)

In this problem, the two eigenvalues and the corresponding eigenvectors of the amplification matrix are the same as (4.3.12) and (4.3.14), respectively. Then following the same analysis in Subsection 4.3.1, we can write

$$\begin{pmatrix} \hat{u}_{-\frac{1}{2}}(t) \\ \hat{u}_{\frac{1}{2}}(t) \end{pmatrix} = C_{11}e^{\lambda_1 t}V_1 + C_{12}e^{\lambda_2 t}V_2, \qquad (4.4.27)$$

where the two constants  $C_{11}$  and  $C_{12}$  are determined by the initial condition (4.4.26). After we obtain the numerical approximations at  $x_{j-\frac{1}{2}}$  and  $x_{j+\frac{1}{2}}$  at the final time T, a direct linear function interpolation would yield the numerical solution at  $x_j^a$  and  $x_j^b$ , denoted as  $u_j^a(t)$  and  $u_j^b(t)$ , respectively, which further leads to the quantitative error estimates

$$\begin{split} ||e_a||_{\infty} &\coloneqq \max_{1 \le j \le N} |U(x_j^a, t) - u_j^a(t)| \\ &= \frac{a(1 + 12a\alpha + 12b\alpha - 12\alpha^2)e^{-t}}{24\alpha} \xi^2 \\ &+ \left[ \frac{96a^3\alpha^2 + 384a^2b\alpha^2 + 2\alpha(1 + 12b\alpha - 12\alpha^2)}{576\alpha^2} \right] \\ &+ \frac{a(-5 + 96(1 + b^2)\alpha^2 - 144\alpha^4)e^{-t}}{576\alpha^2} \right] \xi^3 + O(\xi^4), \\ ||e_b||_{\infty} &\coloneqq \max_{1 \le j \le N} |U(x_j^b, t) - u_j^b(t)| \\ &= \frac{b(1 + 12a\alpha + 12b\alpha - 12\alpha^2)e^{-t}}{24\alpha} \xi^2 \\ &+ \left[ \frac{96b^3\alpha^2 + 384ab^2\alpha^2 + 2\alpha(1 + 12a\alpha - 12\alpha^2)}{576\alpha^2} \right] \end{split}$$

+ + 
$$b(-5 + 96(1 + a^2)\alpha^2 - 144\alpha^4)e^{-t}576\alpha^2]\xi^3 + O(\xi^4)$$

To set the coefficients of the leading term to be zero, we have

$$a + b = \frac{12\alpha^2 - 1}{12\alpha} \tag{4.4.28}$$

Then we can state the following theorem.

**Theorem 4.4.1.** Consider the LDG scheme (4.2.4), (4.2.5) on uniform meshes with mesh size h. Suppose the finite element space is made up of piecewise  $P^1$ polynomials and the condition (4.4.28) is satisfied. Assume the initial solution is the interpolation of the exact solution at  $x_j^a = x_j + ah$  and  $x_j^b = x_j + bh$  in cell  $I_j$ , then we have

$$|U(x_j^a) - u_j^a| = O(h^4), \quad |U(x_j^b) - u_j^b| = O(h^4).$$

where U is the exact solution, and  $u_j^a$  and  $u_j^b$  are the numerical solution evaluated at  $x_j^a$  and  $x_j^b$ , respectively.

**Remark 4.4.1.** We choose  $\phi_{j-\frac{1}{2}}(x)$ ,  $\phi_{j+\frac{1}{2}}(x)$  as the local basis only because we would like to demonstrate the general approach to find the superconvergence points. Actually, one may choose any other basis, e.g. those given in Subsection 4.3.1. However, no matter which basis to choose, one has to construct interpolation polynomial at the superconvergence points as the initial discretization and evaluate the error at the same points. Then the superconvergence points can be determined by taking the leading term of the error to be zero.

### 4.5 Numerical experiments

In this section, we will use numerical experiments to demonstrate the accuracy and superconvergence of the LDG method for one dimensional linear heat equation on overlapping meshes. First, we will demonstrate the accuracy using piecewise polynomials of degree k = 1. Next, we will show numerical experiments for superconvergence. Moreover, we use the third-order SSP Runge-Kutta method for time discretization [34] with time step  $\Delta t = 0.01h^2$  to reduce the time error and take the final time T=1.

**Example 4.5.1.** We solve the following heat equation in one space dimension

$$\begin{cases} u_t = u_{xx}, & x \in [0, 2\pi], \\ u(x, 0) = \sin(x). \end{cases}$$
(4.5.29)

Clearly, the exact solution is

$$u(x,t) = e^{-t}\sin(x).$$

We consider uniform meshes and take  $\alpha = 0$  in (4.2.3), i.e, the dual mesh is generated by using the midpoint of the primitive mesh. Moreover, we also take  $\alpha = 0.05$  which is closed to 0,  $\alpha = 0.25$  which is away from 0, and  $\alpha = 0.5$  that the dual mesh agrees with the primitive mesh. We compute the error between the numerical and exact solutions and the results under  $L^2$ -norm are given in Table 4.1. From the table, we can observe suboptimal accuracy when taking  $\alpha = 0$  with piecewise linear polynomials. To obtain optimal accuracy, we can choose  $\alpha \neq 0$ .

Next, we proceed to verify the superconvergence property discussed in Section 4.4. We first take  $\alpha = 0.25$ , then  $a + b = -\frac{1}{12}$ . One example would be  $a = -\frac{1}{6}$  and  $b = \frac{1}{12}$ , and the result is given in Table 4.2. We can observe third-order convergence, which verifies Theorem 4.4.1. Next, we take  $\alpha = 0.5$ , then  $a + b = \frac{1}{3}$ . In this case, the dual mesh agrees with the primitive mesh. In [76] we have demonstrated third-order superconvergence at the right-biased Radau

k	number of cells	$\alpha = 0$		$\alpha = 0.05$	
		$L^2$ norm	order	$L^2$ norm	order
1	10	1.19E-01	-	9.05E-02	-
	20	5.96E-02	0.96	2.62E-02	1.79
	40	2.98E-02	0.99	5.86E-03	2.16
	80	1.49E-02	1.00	1.37E-03	2.10
	160	7.46E-03	1.00	3.35E-04	2.03
k	number of cells	$\alpha = 0.25$		$\alpha = 0.5$	
		$L^2$ norm	order	$L^2$ norm	order
1	10	5.73E-03	-	1.77E-02	-
	20	1.19E-03	2.27	4.39E-03	2.01
	40	2.80E-04	2.09	1.10E-03	2.00
	80	6.88E-05	2.02	2.74E-04	2.00
	160	1.71E-05	2.00	6.84E-05	2.00

Table 4.1: Example 4.5.1:  $\alpha = 0, \alpha = 0.05, \alpha = 0.25, \alpha = 0.5$ .

points  $(a = -\frac{1}{6}, b = \frac{1}{2})$ . We will choose some other superconvergence points, for example,  $a = -\frac{1}{8}$  and  $b = \frac{11}{24}$ , and the results are given in Table 4.2. From the table, we can also observe third-order superconvergence which verifies Theorem 4.4.1.

k	number of cells	$\alpha = 0.25$		$\alpha = 0.5$	
		a=-1/6	b=1/12	a=-1/8	b=11/24
		$L^2$ norm	order	$L^2$ norm	order
1	10	8.855857E-04	-	1.920474E-03	-
	20	1.054519E-04	3.07	2.358132E-04	3.03
	40	1.299687E-05	3.02	2.934797E-05	3.01
	80	1.617833E-06	3.01	3.664505E-06	3.00
	160	2.020093E-07	3.00	4.579387E-07	3.00

Table 4.2: Example 4.5.1: Superconvergence with  $\alpha = 0.25$  and  $\alpha = 0.5$ 

# 4.6 Conclusion

In this paper, we applied Fourier analysis to demonstrate the quantitative error estimates of the LDG methods on overlapping meshes with piecewise  $P^k$ polynomials (k = 1, 2) for linear parabolic equations in one space dimension. We analyzed the reason for the accuracy degeneration. Some superconvergence points were also investigated.

# Chapter 5

# Conclusion

In the first work, the conservative LDG method for both flow and transport equations was introduced for the coupled system of compressible miscible displacement problem that is important and interesting in oil recovery and environmental pollution problem. The optimal order of error estimates hold not only for the solution itself but also for the auxiliary variables. Special projections and a priori assumption help to eliminate the jump terms at the cell interfaces which arise from the discontinuity nature of the numerical method, the non-linearity and coupling of the model.

In the second study, we expanded the idea of the previous work to construct high-order bound-preserving DG methods for compressible miscible displacements in porous media on triangular meshes. The technique have been applied to the problem with multi-component fluid mixtures. Numerical simulations shown the accuracy and necessity of the bound-preserving technique.

In the third research, Fourier analysis was applied to demonstrate the quanti-

tative error estimates of the LDG methods on overlapping meshes with piecewise  $P^k$  polynomials (k = 1, 2) for linear parabolic equations in one space dimension. We analyzed the reason for the accuracy degeneration. Some superconvergence points were also investigated.

# References

- S. Bartels, M. Jensen and R. Müller, Discontinuous Galerkin finite element convergence for incompressible miscible displacement problem of low regularity, SIAM Journal on Numerical Analysis, 47 (2009), 3720-3743.
- [2] F. Bassi, S. Rebay, A high-order accurate discontinuous finite element method for the numerical solution of the compressible Navier-Stokes equations, Journal of Computational Physics, 131 (1997), 267-279.
- [3] P. Bastian, A fully-coupled discontinuous Galerkin method for two-phase flow in porous media with discontinuous capillary pressure, Computational Geosciences, 18 (2014), 779-796.
- [4] J. Bell, C.N. Dawson, G.R. Shubin, An unsplit high-order Godunov scheme for scalar conservation laws in two dimensions, Journal of Computational Physics, 74 (1988), 1-24.
- [5] W. Cao and Z. Zhang, Superconvergence of Local Discontinuous Galerkin method for one-dimensional linear parabolic equations, Mathematics of Computation, 85 (2016), 63-84.

- [6] P. Castillo, B. Cockburn, I. Perugia and D. Schötzau, Superconvergence of the local discontinuous Galerkin method for elliptic problems on cartesian grids, SIAM Journal on Numerical Analysis, 39 (2001), 264-285.
- [7] H.-Z. Chen and H. Wang, An optimal-order error estimate on an H<sup>1</sup>-Galerkin mixed method for a nonlinear parabolic equation in porous medium flow, Numerical Methods for Partial Differential Equations, 26 (2010), 188-205.
- [8] Z. Chen, H. Huang and J. Yan, Third order Maximum-principle-satisfying direct discontinuous Galerkin methods for time dependent convection diffusion equations on unstructured triangular meshes, Journal of Computational Physics, 308 (2016), 198-217.
- Y. Cheng and C.-W. Shu, Superconvergence of local discontinuous Galerkin methods for convection-diffusion equations, Computers and Structures, 87 (2009), pp. 630-641.
- [10] Y. Cheng and C.-W. Shu, Superconvergence of Discontinuous Galerkin and Local Discontinuous Galerkin Schemes for Linear Hyperbolic and Convection-Diffusion Equations in One Space Dimension, SIAM Journal on Numerical Analysis, 47 (2010), 4044-4072.
- [11] S.-H. Chou and Q. Li, Mixed finite element methods for compressible miscible displacement in porous media, Mathematics of Computation, 57 (1991), 507-527.

- [12] A. Christlieb, Y. Liu, Q. Tang and Z. Xu, Parametrized Maximum-principlepreserving and positivity-preserving flux limiter for WENO schemes on unstructured meshes, Journal of Computational Physics, 281 (2015), 334-351.
- [13] N. Chuenjarern, Z. Xu and Y. Yang, High-order bound-preserving discontinuous Galerkin methods for compressible miscible displacements in porous media on triangular meshes, Journal of Computational Physics, accepted.
- [14] E. Chung and C.S. Lee, A staggered discontinuous Galerkin method for convection-diffusion equations, Journal of Numerical Mathematics, 20 (2012), 1-31.
- [15] P. Ciarlet, The finite element method for elliptic problem, North Holland, 1975.
- [16] B. Cockburn, S. Hou, C. W. Shu, The Runge-Kutta local projection discontinuous Galerkin finite element method for conservation laws. IV: The multidimensional case, Math. Comp. 54 (1990) 545-581.
- [17] B. Cockburn, S. Y. Lin, C. W. Shu, TVB Runge-Kutta local projection discontinuous Galerkin finite element method for conservation laws. III: Onedimensional systems, J. Comput. Phys. 84 (1989) 90-113.
- B. Cockburn, C. W. Shu, TVB Runge-Kutta local projection discontinuous Galerkin finite element method for conservation laws. II: General framework, Math. Comp. 52 (1989) 411-435.

- [19] B. Cockburn, C. W. Shu, The Runge-Kutta discontinuous Galerkin method for conservation laws. V: Multidimensional systems, J. Comput. Phys. 141 (1998) 199-224.
- [20] B. Cockburn, C.-W. Shu, The local discontinuous Galerkin method for timedependent convection-diffusion systems, SIAM Journal on Numerical Analysis, 35 (1998), 2440-2463.
- [21] M. Cui, A combined mixed and discontinuous Galerkin method for compressible miscible displacement problem in porous media, Journal of Computational and Applied Mathematics, 198 (2007), 19-34.
- [22] M. Cui, Analysis of a semidiscrete discontinuous Galerkin scheme for compressible miscible displacement problem, Journal of Computational and Applied Mathematics, 214 (2008), 617-636.
- [23] J. Douglas, Jr. and J. Roberts, Numerical methods for a model for compressible miscible displacement in porous media, Mathematics of Computation, 41 (1983), 441-459.
- [24] J. Douglas Jr., R.E. Ewing, M.F. Wheeler, A time-discretization procedure for a mixed finite element approximation of miscible displacement in porous media, RAIRO Analyse Numérique, 17 (1983), 249-256.
- [25] J. Douglas Jr., R.E. Ewing, M.F. Wheeler, The approximation of the pressure by a mixed method in the simulation of miscible displacement, RAIRO Analyse Numérique 17 (1983), 17-33.

- [26] J. Douglas Jr., T.F. Russell, Numerical methods for convection-dominated diffusion problems based on combining the method of characteristic with finite element of finite difference procedures, SIAM Journal on Numerical Analysis, 19 (1982), 871-885.
- [27] J. Du and Y. Yang, Maximum-principle-preserving third-order local discontinuous Galerkin methods on overlapping meshes, Journal of Computational Physics, 377 (2019), 117-141.
- [28] J. Du, Y. Yang and E. Chung, Stability analysis and error estimates of local discontinuous Galerkin method for convection-diffusion equations on overlapping meshes, BIT Numerical Mathematics, accepted.
- [29] A. Ern, I. Mozolevski and L. Schuh, Discontinuous Galerkin approximation of two-phase flows in heterogeneous porous media with discontinuous capillary pressures, Computer Methods in Applied Mechanics and Engineering, 199 (2010), 1491-1501.
- [30] A. Ern, I. Mozolevski and L. Schuh, Accurate velocity reconstruction for Discontinuous Galerkin approximations of two-phase porous media flows, Comptes Rendus Mathematique, 347 (2009), 551-554.
- [31] R.E. Ewing, T.F. Russell, M.F. Wheeler, Convergence analysis of an approximation of miscible displacement in porous media by mixed finite elements and a modified method of characteristics, Computer Methods in Applied Mechanics and Engineering, 47 (1984), 73-92.

- [32] I.M. Gelfand, Some questions of analysis and differential equations, American Mathematical Society Translations, 26 (1963), 201-219.
- [33] S. Gottlieb, D. Ketcheson and C.-W. Shu, *High order strong stability pre*serving time discretizations, Journal of Scientific Computing, 38 (2009), 251-289.
- [34] S. Gottlieb, C.-W. Shu and E. Tadmor, Strong stability-preserving highorder time discretization methods, SIAM review 43.1 (2001), pp. 89-112.
- [35] H. Guo, F. Yu and Y. Yang, Local discontinuous Galerkin method for incompressible miscible displacement problem in porous media, Journal of Scientific Computing, 71 (2017), 615-633.
- [36] H. Guo and Y. Yang, Bound-Preserving Discontinuous Galerkin Method for Compressible Miscible Displacement in Porous Media, SIAM Journal on Scientific Computing, 39 (2017), A1969-A1990.
- [37] H. Guo and Q. Zhang, Error analysis of the semi-discrete local discontinuous Galerkin method for compressible miscible displacement problem in porous media, Applied Mathematics and Computation, 259 (2015), 88-105.
- [38] H. Guo, Q. Zhang and Y. Yang, A combined mixed finite element method and local discontinuous Galerkin method for miscible displacement problem in porous media, Science China Mathematics, 57 (2014), 2301-2320.

- [39] L. Guo and Y. Yang, Positivity preserving high-order local discontinuous Galerkin method for parabolic equations with blow-up solutions, Journal of Computational Physics, 289 (2015), 181-195.
- [40] A.E. Hurd and D.H. Sattinger, Questions of existence and uniqueness for hyperbolic equations with discontinuous coefficients, Transactions of the American Mathematical Society, 132 (1968), 159-174.
- [41] L. Ji, Y. Xu, Optimal error estimates of the local discontinuous Galerkin method for Willmore flow of graphs on Cartesian meshes, International Journal of Numerical Analysis and Modeling, 8 (2011), 252-283.
- [42] C. Johnson, Streamline diffusion methods for problems in fluid mechanics, Finite Element in Fluids VI, Wiley, New York, 1986.
- [43] E. F. Keller and L. A. Segel, Initiation on slime mold aggregation viewed as instability, Journal of Theoretical Biology, 26 (1970), 399-415.
- [44] S. Kumar, A mixed and discontinuous Galerkin finite volume element method for incompressible miscible displacement problems in porous media, Numerical Methods for Partial Differential Equations, 28 (2012), 1354-1381.
- [45] X. Li, H. Rui, W. Xu, A new MCC-MFE method for compressible miscible displacement in porous media, Journal of Computational and Applied Mathematics, 302 (2016), 139-156.
- [46] X. Li, C.-W. Shu and Y. Yang, Local discontinuous Galerkin method for the Keller-Segel chemotaxis model, Journal of Scientific Computing, 73 (2017), 943-967.
- [47] Y. Liu, C.-W. Shu, E. Tadmor and M. Zhang, Central local discontinuous Galerkin method on overlapping cells for diffusion equations, ESAIM: Mathematical Modeling and Numerical Analysis (M2AN), 45 (2011), 1009-1032.
- [48] N. Ma, D. Yang and T. Lu, L<sup>2</sup>-norm error bounds of characteristics collocation method for compressible miscible displacement in porous media, International Journal of Numerical Analysis and Modeling, 2 (2005), 28-42.
- [49] C. Patlak, Random walk with persistence and external bias, The bulletin of mathematical biophysics, 15 (1953), 311338.
- [50] T. Qin, C.-W. Shu and Y. Yang, Bound-preserving discontinuous Galerkin methods for relativistic hydrodynamics, Journal of Computational Physics, 315 (2016), 323-347.
- [51] W.H. Reed and T. R. Hill, Triangular mesh methods for the neutron transport equation, Los Alamos Scientific Laboratory Report LA-UR-73-479, Los Alamos, NM, 1973.
- [52] B. Rivière, Discontinuous Galerkin finite element methods for solving the miscible displacement problem in porous media, Ph.D. Thesis, The University of Texas at Austin, 2000.

- [53] C.-W. Shu, Total-variation-diminishing time discretizations, SIAM Journal on Scientific and Statistical Computing 9 (1988), 1073-1084.
- [54] C.-W. Shu, S. Osher, Efficient implementation of essentially non-oscillatory shock-capturing schemes, Journal of Computational Physics, 77 (1988), 439-471.
- [55] S. Sun, B. Rivière and M.F. Wheeler, A combined mixed finite element and discontinuous Galerkin method for miscible displacement problem in porous media, Recent Progress in Computational and Applied PDEs, Tony Chan et al. (Eds.), Kluwer Academic Publishers, Plenum Press, Dordrecht, NewYork, 2002, 323-351.
- [56] S. Sun and M.F. Wheeler, Discontinuous Galerkin methods for coupled flow and reactive transport problems, Applied Numerical Mathematics, 52 (2005), 273-298.
- [57] S. Sun and M.F. Wheeler, Symmetric and nonsymmetric discontinuous Galerkin methods for reactive transport in porous media, SIAM Journal on Numerical Analysis, 43 (2005), 195-219.
- [58] H. Wang, D. Liang, R.E. Ewing, S.L. Lyons and G. Qin, An approximation to miscible fluid flows in porous media with point sources and sinks by an Eulerian-Lagrangian localized adjoint method and mixed finite element methods, SIAM Journal on Scentific Computing, 22 (2000), 561-581.

- [59] H. Wang, D. Liang, R.E. Ewing, S.L. Lyons and G. Qin, An accurate approximation to compressible flow in porous media with wells, Numerical Treatment of Multiphase Flows in Porous Media, Lecture Notes in Physics, 552 (2000), 324-332.
- [60] H. Wang, C.-W. Shu, and Q. Zhang, Stability and error estimates of local discontinuous Galerkin methods with implicit-explicit time-marching for advection-diffusion problems, SIAM Journal on Numerical Analysis, 53 (2015), 206-227.
- [61] H. Wang, C.-W. Shu and Q. Zhang, Stability analysis and error estimates of local discontinuous Galerkin methods with implicit-explicit time-marching for nonlinear convection-diffusion problems, Applied Mathematics and Computation, 272 (2016), 237-258.
- [62] H. Wang, S. Wang, Q. Zhang and C.-W. Shu, Local discontinuous Galerkin methods with implicit-explicit time marching for multi-dimensional convection diffusion problems, ESAIM: M2AN, 50 (2016), 1083-1105.
- [63] M. F. Wheeler and B. L. Darlow, Interiori penalty Galerkin methods for miscible displacement problems in porous media, Computational Methods in Nonlinear Mechanics, North-Holland, Amsterdam, 1980, 458-506.
- [64] Y. Xing, X. Zhang and C.-W. Shu, Positivity preserving high order well balanced discontinuous Galerkin methods for the shallow water equations, Advances inWater Resources, 33 (2010), 1476-1493.

- [65] T. Xiong, J.-M. Qiu and Z. Xu, High order maximum-principle-preserving discontinuous Galerkin method for convection-diffusion equations, SIAM Journal on Scientific Computing, 37 (2015), A583-A608.
- [66] Y. Xu, C.-W. Shu, Local discontinuous Galerkin methods for nonlinear Schrodinger equations, Journal of Computational Physics, 205 (2005), 72-97.
- [67] Y. Xu, C.-W. Shu, Local discontinuous Galerkin methods for the Kuramoto-Sivashinsky equations and the Ito-type coupled KdV equations, Computer Methods in Applied Mechanics and Engineering, 195 (2006), 3430-3447.
- [68] Z. Xu, Parametrized maximum principle preserving flux limiters for high order schemes solving hyperbolic conservation laws: One-dimensional scalar problem, Mathematics of Computation, 83 (2014), 310-331.
- [69] J. Yan, C.-W. Shu, A local discontinuous Galerkin method for KdV type equations, SIAM Journal on Numerical Analysis, 40 (2002), 769-791.
- [70] D. Yang A splitting positive definite mixed element method for miscible displacement of compressible flow in porous media, Numerical Methods for Partial Differential Equations, 17 (2001), 229-249.
- [71] J. Yang and Y. Chen, A priori error estimates of a combined mixed finite element and discontinuous Galerkin method for compressible miscible displacement with molecular diffusion and dispersion, Journal of Computational Mathematics, 28 (2010), 1005-1022.

- [72] J. Yang, A posteriori error of a discontinuous Galerkin scheme for compressible miscible displacement problems with molecular diffusion and dispersion, International Journal for Numerical Methods in Fluids, 65 (2011), 781-797.
- [73] J. Yang and Y. Chen, A priori error analysis of a discontinuous Galerkin approximation for a kind of compressible miscible displacement problems, Science China Mathematics, 53 (2010), 2679-2696.
- [74] Y. Yang and C.-W. Shu, Discontinuous Galerkin method for hyperbolic equations involving δ-singularities: Negative-order norm error estimates and applications, Numerische Mathematik, 124,(2013),753-781.
- [75] Y. Yang, D. Wei and C.-W. Shu, Discontinuous Galerkin method for Krause's consensus models and pressureless Euler equations, Journal of Computational Physics, 252 (2013), 109-127.
- [76] Y. Yang and C.-W. Shu, Analysis of optimal superconvergence of local discontinuous Galerkin method for one-dimensional linear parabolic equations, Journal of Computational Mathematics, 33 (2015), 323-340.
- [77] F. Yu, H. Guo, N. Chuenjarern and Y. Yang, Conservative local discontinuous Galerkin method for compressible miscible displacements in porous media, Journal of Scientific Computing, 73 (2017), 1249-1275.
- [78] Y. Yuan, Characteristic finite element methods for positive semidefinite problem of two phase miscible flow in three dimensions, Chinese Science Bulletin, 22 (1996), 2027-2032.

- [79] J. Zhang, D. Yang, S. Shen, J. Zhu, A new MMOCAA-MFE method for compressible miscible displacement in porous media, Applied Numerical Mathematics, 80 (2014), 65-80.
- [80] M. Zhang and C. W. Shu, An analysis of three different formulations of the discontinuous Galerkin method for diffusion equations, Mathematical Models and Methods in Applied Sciences, 13 (2003), 395-413.
- [81] Y. Yuan, The characteristic finite difference fractional steps methods for compressible two-phase displacement problem, Science in China Series A-Mathematics, 42 (1999), 48-57.
- [82] Y. Yuan, The upwind finite difference fractional steps methods for two-phase compressible flow in porous media, Numerical Methods Partial Differential Equations, 19 (2003), 67-88.
- [83] Y. Yuan, The modified upwind finite difference fractional steps method for compressible two-phase displacement problem, Acta Mathematicae Applicatae Sinica, 20 (2004), 381-396.
- [84] M. Zhang and J. Yan, Fourier Type Error Analysis of the Direct Discontinuous Galerkin Method and Its Variations for Diffusion Equations, Journal of Scientific Computing, 52 (2012),638-655.
- [85] X. Zhang and C.-W. Shu, On maximum-principle-satisfying high order schemes for scalar conservation laws, Journal of Computational Physics, 229 (2010), 3091-3120.

- [86] X. Zhang and C.-W. Shu, On positivity preserving high order discontinuous Galerkin schemes for compressible Euler equations on rectangular meshes, Journal of Computational Physics, 229 (2010), 8918-8934.
- [87] X. Zhang and C.-W. Shu, Positivity-preserving high order discontinuous Galerkin schemes for compressible Euler equations with source terms, Journal of Computational Physics, 230 (2011), 1238-1248.
- [88] X. Zhang, Y. Xia and C.-W. Shu, Maximum-Principle-Satisfying and Positivity-Preserving High Order Discontinuous Galerkin Schemes for Conservation Laws on Triangular Meshes, Journal of Scientific Computing, (2012), 50: 29-32.
- [89] Y. Zhang, X. Zhang and C.-W. Shu, Maximum-principle-satisfying second order discontinuous Galerkin schemes for convection-diffusion equations on triangular meshes, Journal of Computational Physics, 234 (2013), 295-316.
- [90] X. Zhao, Y. Yang and C. Seyler, A positivity-preserving semi-implicit discontinuous Galerkin scheme for solving extended magnetohydrodynamics equations, Journal of Computational Physics, 278 (2014), 400-415.
- [91] X. Zhong and C.-W. Shu, Numerical resolution of discontinuous Galerkin methods for time dependent wave equations, Computer Methods in Applied Mechanics and Engineering, 200 (2011), 2814-2827.

# Appendix A

# **Copyright documentations**

## A.1 Copyright documentation of Chapter 2

Springer Science and Bus Media B V LICENSE TERMS AND CONDITIONS

Jan 09, 2019

This is a License Agreement between Nattaporn Chuenjarern ("You") and Springer Science and Bus Media B V ("Springer Science and Bus Media B V") provided by Copyright Clearance Center ("CCC"). The license consists of your order details, the terms and conditions provided by Springer Science and Bus Media B V, and the payment terms and conditions.

All payments must be made in full to CCC. For payment instructions, please see information listed at the bottom of this form.

License Number	4504930587962
License date	Jan 09, 2019
Licensed content publisher	Springer Science and Bus Media B V
Licensed content title	Journal of scientific computing
Licensed content date	Jan 1, 1986
Type of Use	Thesis/Dissertation
Requestor type	Academic institution
Format	Print, Electronic

Portion	chapter/article
The requesting person/organization is:	Nattaporn Chuenjarern
Title or numeric reference of the portion(s)	Entire article
Title of the article or chapter the portion is from	Conservative local discontinuous Galerkin method for compressible miscible displacements in porous media
Editor of portion(s)	N/A
Author of portion(s)	Fan Yu, Hui Guo, Nattaporn Chuenjarern, Yang Yang
Volume of serial or monograph.	73
Issue, if republishing an article from a serial	2-3
Page range of the portion	
Publication date of portion	09 October 2017
Rights for	Main product
Duration of use	Current edition and up to 5 years
Creation of copies for the disabled	no
With minor editing privileges	no
For distribution to	Worldwide
In the following language(s)	Original language of publication
With incidental promotional use	no
The lifetime unit quantity of new product	Up to 4,999
Title	Discontinuous Galerkin methods for convection-diffusion equations and applications in petroleum engineering
Institution name	Michigan Technological University
Expected presentation date	Feb 2019
Billing Type	Invoice
Billing Address	Nattaporn Chuenjarern 1906 A WOODMAR DRIVE
	Houghton, MI 49931 United States Attn: Nattaporn Chuenjarern

Total (may include CCC user ~0.00~USD fee)

## A.2 Copyright documentation of Chapter 3

#### Elsevier Science and Technology Journals LICENSE TERMS AND CONDITIONS

Jan 09, 2019

This is a License Agreement between Nattaporn Chuenjarern ("You") and Elsevier Science and Technology Journals ("Elsevier Science and Technology Journals") provided by Copyright Clearance Center ("CCC"). The license consists of your order details, the terms and conditions provided by Elsevier Science and Technology Journals, and the payment terms and conditions.

All payments must be made in full to CCC. For payment instructions, please se	e
information listed at the bottom of this form.	

License Number	4504930589262
License date	Jan 09, 2019
Licensed content publisher	Elsevier Science and Technology Journals
Licensed content title	Journal of computational physics
Licensed content date	Jan 1, 1966
Type of Use	Thesis/Dissertation
Requestor type	Academic institution
Format	Print, Electronic
Portion	chapter/article
Number of pages in chapter/article	19
The requesting person/organization is:	Nattaporn Chuenjarern
Title or numeric reference of the portion(s)	Entire article
Title of the article or chapter the portion is from	High-order bound-preserving discontinuous Galerkin methods for compressible miscible displacements in porous media on triangular meshes
Editor of portion(s)	N/A
Author of portion(s)	Nattaporn Chuenjarern,Ziyao Xu Yang Yang
Volume of serial or monograph.	378
Page range of the portion	110-128
Publication date of portion	7 November 2018
Rights for	Main product

Duration of use	Life of current edition
Creation of copies for the disabled	no
With minor editing privileges	no
For distribution to	Worldwide
In the following language(s)	Original language of publication
With incidental promotional use	no
The lifetime unit quantity of new product	Up to 4,999
Title	Discontinuous Galerkin methods for convection-diffusion equations and applications in petroleum engineering
Institution name	Michigan Technological University
Expected presentation date	Feb 2019
Billing Type	Invoice
Billing Address	Nattaporn Chuenjarern 1906 A WOODMAR DRIVE
	Houghton, MI 49931 United States Attn: Nattaporn Chuenjarern
Total (may include CCC user fee)	0.00 USD

## A.3 Copyright documentation of Chapter 4

Fourier analysis of local discontinuous Galerkin methods for linear parabolic equations on overlapping meshes has been completed as an article to submit to Journal of Scientific Computing. As it has not yet been published (at the time of this dissertation's publication), we (Nattaporn Chuenjarern, and Yang Yang) still retain the copyright and as such, there is no need to obtain copyright documentation.