

# KINETIC MODELS OF DILUTE POLYMERS:

ANALYSIS, APPROXIMATION AND COMPUTATION

*John W. Barrett*

*David J. Knezevic*

*Endre Süli*

2009

**John W. Barrett**

Department of Mathematics,  
Imperial College London,  
South Kensington,  
London SW7 2AZ,  
UK.

e-mail: [j.barrett@imperial.ac.uk](mailto:j.barrett@imperial.ac.uk)

**David J. Knezevic**

Massachusetts Institute of Technology,  
School of Engineering,  
Boston, MA,  
USA.

e-mail: [dknez@MIT.EDU](mailto:dknez@MIT.EDU)

**Endre Süli**

Oxford Centre for Nonlinear PDE,  
Mathematical Institute,  
University of Oxford,  
Gibson Building,  
Radcliffe Observatory Quarter,  
Woodstock Road,  
Oxford, OX2 6HA,  
UK.

e-mail: [endre.suli@maths.ox.ac.uk](mailto:endre.suli@maths.ox.ac.uk)

NEČAS CENTER FOR MATHEMATICAL MODELING:  
11TH SCHOOL ON MATHEMATICAL THEORY IN FLUID MECHANICS  
22-29 MAY 2009  
KACOV, CZECH REPUBLIC

© John W. Barrett, David John Knezevic, Endre Süli  
London/Boston/Oxford  
2009

**Acknowledgement**

ES is grateful to Josef Málek and Mirko Rokyta (Nečas Center, Prague) for the kind invitation, and to the Nečas Center and OxpDE, the Oxford Centre for Nonlinear Partial Differential Equations, for financial support. OxpDE is funded by the UK Engineering and Physical Sciences Research Council under the EPSRC grant number EP/E035027/1.

**ABSTRACT**

We review recent analytical and computational results for macroscopic-microscopic bead-spring models that arise from the kinetic theory of dilute solutions of incompressible polymeric fluids with noninteracting polymer chains, involving the coupling of the unsteady Navier–Stokes system in a bounded domain  $\Omega \subset \mathbb{R}^d$ ,  $d = 2$  or  $3$ , with an elastic extra-stress tensor as right-hand side in the momentum equation, and a (possibly degenerate) Fokker–Planck equation over the  $(2d + 1)$ -dimensional region  $\Omega \times D \times [0, T]$ , where  $D \subset \mathbb{R}^d$  is the configuration domain and  $[0, T]$  is the temporal domain. The Fokker–Planck equation arises from a system of (Itô) stochastic differential equations, which models the evolution of a  $2d$ -component vectorial stochastic process comprised by the  $d$ -component centre-of-mass vector and the  $d$ -component orientation (or configuration) vector of the polymer chain. We show the existence of global-in-time weak solutions to the coupled Navier–Stokes–Fokker–Planck system for a general class of spring potentials including, in particular, the widely used finitely extensible nonlinear elastic (FENE) potential. The numerical approximation of this high-dimensional coupled system is a formidable computational challenge, complicated by the fact that for practically relevant spring potentials, such as the FENE potential, the drift term in the Fokker–Planck equation is unbounded on  $\partial D$ . We present numerical simulations for this coupled high-dimensional micro-macro model and we consider the analysis of the algorithms.



# Contents

<b>1</b>	<b>Introduction</b>	<b>9</b>
1.1	Overview of Newtonian fluid dynamics . . . . .	10
1.2	Modelling polymeric fluids . . . . .	11
1.2.1	The freely rotating chain model . . . . .	11
1.2.2	The bead-rod chain model . . . . .	12
1.2.3	The bead-spring chain model . . . . .	12
1.2.4	The dumbbell model . . . . .	12
1.2.5	Spring force laws . . . . .	12
1.3	The micro-macro model . . . . .	14
1.3.1	Derivation of the Fokker–Planck equation . . . . .	15
1.3.2	Properties of the probability density function . . . . .	19
1.3.3	Polymeric extra-stress . . . . .	20
1.3.4	The coupled Navier–Stokes–Fokker–Planck system . . . . .	21
1.4	Numerics for polymeric fluids: literature review . . . . .	22
1.4.1	Fully macroscopic methods . . . . .	22
1.4.2	Stochastic multiscale methods . . . . .	23
1.4.3	Deterministic multiscale methods . . . . .	23
1.5	Outlook and goals . . . . .	26
<b>2</b>	<b>The Fokker–Planck equation in configuration space</b>	<b>29</b>
2.1	Introduction . . . . .	29
2.2	Properties of Maxwellian-weighted spaces . . . . .	33
2.3	Analysis of the backward Euler semidiscretisation . . . . .	34
2.4	The Chauvière–Lozinski transformed FENE model . . . . .	40
2.5	The fully-discrete method . . . . .	43
2.6	Approximation results . . . . .	45
2.7	Convergence analysis of the numerical method . . . . .	55
2.8	Numerical results . . . . .	57
2.8.1	Numerical methods in the two-dimensional case . . . . .	58
2.8.2	The semi-implicit numerical method . . . . .	70
2.8.3	Three-dimensional implementation of the spectral method . . . . .	71
2.9	Conclusions . . . . .	75
<b>3</b>	<b>ADI methods for the full Fokker–Planck equation</b>	<b>79</b>
3.1	Introduction . . . . .	79
3.2	Weak formulation and spatial discretisation . . . . .	80

3.3	The alternating-direction numerical method . . . . .	82
3.3.1	The hybrid alternating-direction scheme . . . . .	86
3.3.2	Method I: Semi-implicit scheme . . . . .	89
3.3.3	Method II: Fully-implicit scheme . . . . .	93
3.4	Stability of methods I and II . . . . .	94
3.5	Convergence analysis for method I: Part 1 . . . . .	97
3.6	Approximation results on $\Omega \times D$ . . . . .	102
3.7	Convergence analysis for method I: Part 2 . . . . .	105
3.8	Implementation of methods I and II . . . . .	109
3.8.1	The $q$ -direction stage . . . . .	109
3.8.2	The $\tilde{x}$ -direction stage . . . . .	110
3.8.3	The $\tilde{x}$ -direction quadrature rule . . . . .	110
3.8.4	Parallel implementation of the alternating-direction method . . . . .	112
3.9	Numerical results . . . . .	113
3.10	Conclusions . . . . .	119
<b>4</b>	<b>The coupled Navier–Stokes–Fokker–Planck system</b>	<b>121</b>
4.1	Introduction . . . . .	121
4.2	Numerical method for the micro-macro model . . . . .	121
4.3	Numerical results . . . . .	125
4.3.1	4-to-1 planar contraction flow . . . . .	125
4.3.2	Flow around a sphere . . . . .	126
4.4	Conclusions . . . . .	127
<b>5</b>	<b>Global weak solutions to NSFP systems</b>	<b>131</b>
5.1	Introduction . . . . .	131
5.2	The polymer model . . . . .	139
5.3	Existence of global weak solutions . . . . .	140
5.3.1	Existence for $(P_{\varepsilon,L,\delta})$ . . . . .	143
5.3.2	Existence for $(P_{\varepsilon,L})$ . . . . .	155
5.4	Appendix: Compact embeddings . . . . .	158
5.4.1	Step 1: Compact embedding of $H_M^1(D)$ into $L_M^2(D)$ , completeness, separability . . . . .	158
5.4.2	Step 2: Isometric isomorphisms . . . . .	159
5.4.3	Step 3: Compact embedding of $H_M^1(\Omega \times D)$ into $L_M^2(\Omega \times D)$ . . . . .	162
<b>6</b>	<b>Finite element approximation of NSFP systems</b>	<b>163</b>
6.1	Introduction . . . . .	163
6.2	Polymer models . . . . .	166
6.2.1	Microscopic-macroscopic polymer models . . . . .	166
6.2.2	FENE model . . . . .	168
6.2.3	General structural assumptions on the potential . . . . .	168
6.2.4	Formal estimates . . . . .	169
6.3	Function spaces . . . . .	171
6.4	Finite element approximation . . . . .	174
6.5	Appendix: Maxwellian Sobolev norm quasi-interpolation . . . . .	197
6.5.1	The univariate case . . . . .	198

6.5.2	Multiple dimensions . . . . .	201
6.5.3	Two dimensions: flat boundary . . . . .	201
6.5.4	Two dimensions: curved boundary . . . . .	207
6.5.5	Three dimensions . . . . .	212
6.5.6	Stability of the Maxwellian-weighted $L^2$ projector in the Maxwellian-weighted $H^1$ norm . . . . .	212

## Structure of the course

- Lecture 1.** Introduction. Derivation of the model. The Fokker–Planck equation in configuration space. [Chapter 1 and Chapter 2].
- Lecture 2.** Alternating direction method for the full Fokker–Planck equation. The coupled Navier–Stokes–Fokker–Planck system. [Chapter 3 and Chapter 4].
- Lecture 3.** Existence of global weak solutions to coupled Navier–Stokes–Fokker–Planck systems with microscopic cut-off. [Chapter 5].
- Lecture 4.** Finite element approximation of Navier–Stokes–Fokker–Planck systems with microscopic cut-off. [Chapter 6].

*These Lecture Notes have been written as supporting material for a lecture-series delivered by ES at the 11th School on Mathematical Theory in Fluid Mechanics, held between 22-29 May 2009 in Kacov, Czech Republic, organised by the Nečas Center for Mathematical Modeling, Prague. The material in these Lecture Notes is based on the following sources: Chapters 1–4 of the Doctoral Thesis of DJK defended at the University of Oxford in October 2008 (cf. [70]), two publications by DJK and ES (cf. [71, 72]) and a series of articles by JWB and ES (cf. [10, 11, 12, 13, 14]), the first of which was written in collaboration with Christoph Schwab (ETH Zürich).*

# Chapter 1

## Introduction

The study of the dynamics of polymeric fluids has been an area of active research since the 1950's and has undergone significant evolution since that time. In the early work in this field, analytical techniques were developed with the goal of deriving exact solutions for idealised flow problems. With the increasing availability of computational power in subsequent years, it was natural for researchers to apply numerical methods to more complicated flow problems for polymeric fluids (and non-Newtonian fluids in general) than were tractable with analytical methods. This line of research, known as *computational rheology*, took root in the 1970's and it remains an exciting and challenging area of scientific computing today. Simultaneously, there have been significant and exciting advances on the mathematical analysis of nonlinear partial differential equations that arise as mathematical models of polymeric fluids.

In these lectures we investigate a particular class of problems from rheology: bead-spring models for dilute polymeric fluids. We explore the existence of weak solutions, the approximation of bead-spring models using deterministic multiscale algorithms based on the Galerkin method, the rigorous analysis of the numerical algorithms, and we also present computational results, which demonstrate the effectiveness of the methods in practice.

The essence of the subject of modelling dilute polymeric fluids is encapsulated in the coupled Navier–Stokes–Fokker–Planck system (discussed in detail in Section 1.3). This system of equations is often referred to as the “micro-macro” model to emphasise that it is fundamentally multiscale in nature. It is worth highlighting at the outset that there is an extensive literature on numerical methods for simulating polymeric fluids, but most of the previous work uses either a fully macroscopic approach in order to circumvent the multiscale nature of the Navier–Stokes–Fokker–Planck system (see the text [103] for an overview of this field) or a stochastic approach in which the micro-macro system is treated using Monte Carlo type methods (*cf.* [101] and [80]). The direction pursued in this work is rather different; our goal is to solve the micro-macro system using deterministic methods (*e.g.* finite element or spectral methods). This will subsequently be referred to as the *deterministic multiscale* approach. The various advantages and disadvantages of fully macroscopic, stochastic and deterministic multiscale methods will be discussed in detail later, but it should be noted at the outset that the deterministic multiscale method has received far less attention in the literature than the other approaches, probably because this approach can be highly computationally intensive. The central goal of this work, therefore, is to develop multiscale numerical methods for the micro-macro model of dilute polymeric fluids, to address some of the questions related to numerical analysis of such methods, which, up to now, have not been considered in

the literature, and to develop the mathematical theory of the underlying partial differential equations.

In this introductory chapter, we discuss background material on the mathematical modelling of polymer fluids. Newtonian fluids are briefly considered in Section 1.1, and then in Section 1.2 some “coarse-grained” mechanical models for polymer molecules are introduced. Next, in Section 1.3, we derive the Fokker–Planck equation and define the coupled Navier–Stokes–Fokker–Planck system. Section 1.4 contains a literature review of the many and varied numerical methods that have been used for simulating polymeric fluids (these methods fall into the three categories mentioned in the previous paragraph), and the chapter concludes with an overview of the outlook and goals of this work.

## 1.1 Overview of Newtonian fluid dynamics

The success of classical fluid dynamics in accurately describing the properties of a wide range of fluids (typically with low molecular weight, *e.g.* water) using macroscopic continuum models is well established. We begin with a very brief review of some basic principles of classical fluid dynamics (for a full discussion see [15]) as this will be useful for elucidating important ideas in the theory of polymeric fluids.

In the case of Newtonian fluids it has been experimentally established that in a shear flow, *i.e.*  $u_1 = u_1(x_2)$ ,  $u_2 = 0$  where  $u_1$  and  $u_2$  are the components of a two-dimensional velocity field  $\underline{u} = (u_1, u_2)$ , the fluid stress can be related to shear rate by “Newton’s law of viscosity”:

$$\sigma_{21} = \mu \frac{du_1}{dx_2}, \quad (1.1.1)$$

where  $\sigma_{21}$  denotes the force per unit area acting in the  $x_1$ -direction, on a surface normal to the  $x_2$ -direction. That is, stress is proportional to shear rate and the viscosity,  $\mu$ , is the constant of proportionality. This relationship can be generalised to a tensor equation for the stress tensor,  $\underline{\underline{\sigma}}$ , and the strain tensor as follows:

$$\underline{\underline{\sigma}} = -p\underline{\underline{I}} + \mu (\underline{\underline{\nabla}}_x \underline{u} + (\underline{\underline{\nabla}}_x \underline{u})^T). \quad (1.1.2)$$

This equation provides a relationship between the stress and strain of a fluid (in this case, a simple linear equation) and is known as a *constitutive equation*.

Combining the Newtonian constitutive equation (1.1.2) with the equations of conservation of mass:

$$\underline{\underline{\nabla}}_x \cdot \underline{u} = 0, \quad (1.1.3)$$

and momentum:

$$\rho \left( \frac{\partial \underline{u}}{\partial t} + (\underline{u} \cdot \underline{\underline{\nabla}}_x) \underline{u} \right) = \underline{\underline{\nabla}}_x \cdot \underline{\underline{\sigma}}, \quad (1.1.4)$$

where  $\rho$  is the fluid density (assumed to be constant), gives rise to the Navier–Stokes equations for an incompressible, viscous, isothermal fluid:

$$\frac{\partial \underline{u}}{\partial t} + (\underline{u} \cdot \underline{\underline{\nabla}}_x) \underline{u} - \nu \Delta_x \underline{u} + \underline{\underline{\nabla}}_x p = 0, \quad (1.1.5)$$

$$\underline{\underline{\nabla}}_x \cdot \underline{u} = 0, \quad (1.1.6)$$

where the momentum equation has been divided through by  $\rho$ , the pressure in (1.1.5) has implicitly been rescaled by  $\rho$  and  $\nu := \mu/\rho$  is the kinematic viscosity. These equations (which involve only macroscopic quantities) form the cornerstone of classical fluid dynamics.

The situation with polymeric fluids, however, is quite different. In general the contributions to the stress tensor  $\underline{\underline{\sigma}}$  from microscopic polymer molecules cannot be averaged out into purely macroscopic quantities and therefore in order to faithfully simulate a polymeric fluid, the microscopic and macroscopic length-scales must be coupled together. This coupling is achieved by the Navier–Stokes–Fokker–Planck system alluded to above.

In the next section, mechanical models (*i.e.* systems containing masses, rigid rods and/or springs) for microscopic polymer molecules are considered. From the perspective of polymer fluid dynamics, the purpose of these models is to capture the most important characteristics of polymer molecules in systems with many fewer degrees of freedom and in order to yield mathematical models for polymeric fluids that are analytically and computationally tractable.

## 1.2 Modelling polymeric fluids

Polymer molecules consist of long chains of repeated basic structural units, or *monomers*. Polymers of interest typically contain on the order of  $10^3$  to  $10^6$  monomers and the presence of these long chain molecules in a fluid can dramatically affect the fluid’s macroscopic properties. In particular, polymer molecules introduce elastic properties and, as a result, polymeric fluids are often described as *viscoelastic*. Viscoelasticity gives rise to a range of exotic phenomena, such as shear-thinning, rod-climbing, the “tubeless siphon”, and elastic recoil [22].

Most approaches to the mathematical modelling of polymeric fluids are based on kinetic theory, in which the behaviour of the microscopic polymer molecules is characterised in a statistical sense. The starting point in deriving kinetic-theory-based equations is to propose a simple mechanical model that represents an individual polymer molecule. A mechanical model that would faithfully capture the microscopic properties of an actual polymer would be extremely complicated, with a very high number of degrees of freedom, and would be prohibitively difficult to deal with and as a result, a range of simplifications and idealisations have been proposed.

The following “coarse-grained” models for polymer molecules are discussed below: the freely rotating chain model; the bead-rod chain model; the bead-spring chain model; and the dumbbell model (see Chapter 10 of Bird *et. al.* [23] for more details on each of these). This hierarchy of models is depicted in Figure 1.1(a).

### 1.2.1 The freely rotating chain model

It was observed by Flory [49] that bond angles between monomers in a polymer chain are restricted to quite narrow ranges about their average values (up to  $\sim 3\%$  deviation). This motivated the freely rotating chain model, which represents each monomer unit as a bead, where adjacent beads are joined by a rigid, massless rod and where rods are set at a fixed angle (the average bond angle) but are free to rotate. This model has been used in a number of kinetic theory studies by Kirkwood [69]. For the purposes of multiscale computations, though, this model is far too complex. It requires one degree of freedom for each monomer,

so that the number of degrees of freedom in a single chain would be on the order of  $10^3$  to  $10^6$ .

### 1.2.2 The bead-rod chain model

The bead-rod chain model is significantly simpler. It lumps a group of monomers into a single bead and adjacent beads are connected by a massless rod. The restriction on the bond angle (cf. Section 1.2.1) is dropped so that this model is referred to as “freely jointed”. The number of degrees of freedom for this model is typically around 100. The bead-rod chain was first analysed in a seminal paper by Kramers in 1944 [76], and the model is often referred to as a *Kramers chain*. While clearly a considerable simplification from the freely rotating chain, this model still reflects a number of the important characteristics of a polymer molecule – in particular the bead-rod chain has a large number of internal degrees of freedom, it can be oriented and deformed by a flow and it has a constant contour length.

### 1.2.3 The bead-spring chain model

The bead-spring chain is a yet coarser model; a polymer is modelled by a chain of typically around 10 beads joined by springs. The model is completed by specifying a force law for the springs (see below). This model has been the basis of a number of kinetic-theory-based investigations of polymer fluids, *e.g.* the seminal papers of Rouse and Zimm [107, 130].

### 1.2.4 The dumbbell model

The dumbbell model is the simplest in the hierarchy of coarse-grained mechanical models for polymers; it consists of only two masses, which are connected by a spring (or sometimes a rigid rod, although we only consider the spring case in this work). A dumbbell is fully specified by the position of its centre of mass,  $\underline{x}$ , and its configuration (or end-to-end) vector,  $\underline{q}$  (see Figure 1.1(b)). Despite the simplicity of the dumbbell model, it is still very useful for simulating polymeric fluids in many flow regimes because dumbbells can be stretched and oriented by a flow, and these two actions determine the main contributions from polymer molecules to the macroscopic properties of a viscoelastic fluid.

### 1.2.5 Spring force laws

As indicated above, a force law,  $\underline{F}$ , must also be defined for the coarse-grained models that contain one or more springs. In general, the elastic force is assumed to be defined by a (sufficiently smooth) potential  $U : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$  via

$$\underline{F}(\underline{q}) = H U'(\frac{1}{2}|\underline{q}|^2)\underline{q}, \quad (1.2.1)$$

where  $\underline{q}$  is the configuration vector (as illustrated in Figure 1.1(b)) of a given spring and  $H \in \mathbb{R}_{>0}$  is the spring constant. The simplest force law is that of a Hookean spring:

$$U(s) = s \quad \text{and} \quad \underline{F}(\underline{q}) = H\underline{q}. \quad (1.2.2)$$

Many interesting analytical results have been derived for dilute solutions of Hookean dumbbells; indeed the simple linear relationship in (1.2.2) makes this model attractive from the

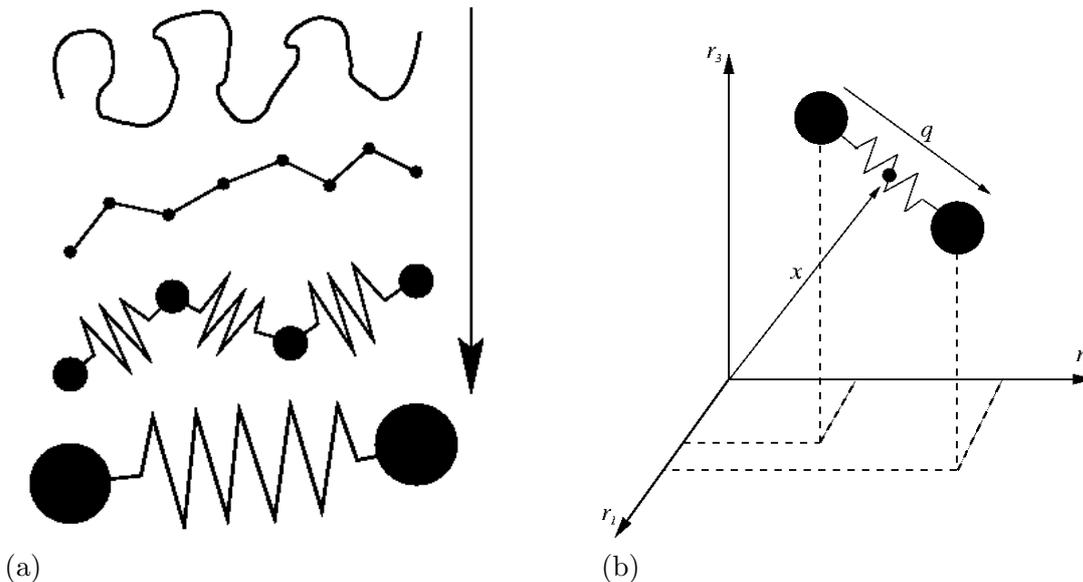


Figure 1.1: (a) Diagram of the hierarchy of mechanical models for polymer molecules, descending from a polymer molecule with on the order of  $10^3$  to  $10^6$  monomers to the dumbbell model, containing only two masses connected by a spring. (b) A more detailed depiction of the dumbbell model. The state of a dumbbell is defined by the position of its centre of mass,  $\underline{x}$ , and its configuration (or end-to-end) vector,  $\underline{q}$ . The dumbbell shown in this schematic can move in  $\mathbb{R}^3$ , and therefore has six degrees-of-freedom.

mathematical point of view. For example, it is well known that the Oldroyd-B macroscopic model for dilute polymeric fluids (originally derived from continuum mechanics considerations [100]) is equivalent to the Hookean dumbbell micro-macro model (*e.g.* see [10]). However, due to the physically unrealistic ability of Hookean springs to be infinitely stretched these models can break down in certain cases, such as strong extensional flows. A remedy is to use the Finitely Extensible Nonlinear Elastic (FENE) force law, suggested by Warner [124], for which we have,

$$U(\frac{1}{2}|\underline{q}|^2) = -\frac{l_{\max}^2}{2} \ln \left( 1 - \frac{|\underline{q}|^2}{l_{\max}^2} \right) \quad \text{and} \quad \underline{F}(\underline{q}) = \frac{H\underline{q}}{1 - |\underline{q}|^2/l_{\max}^2}. \quad (1.2.3)$$

As the name suggests, FENE springs can only be stretched a finite amount because the spring potential is unbounded as  $|\underline{q}| \rightarrow l_{\max}$ . Unlike with Hookean springs, there is no equivalent macroscopic formulation for suspensions of FENE dumbbells; the FENE dumbbell model requires a truly multiscale approach. Note also that for any  $\underline{q}$  fixed in the open ball  $\{\underline{q} : |\underline{q}| < l_{\max}\}$ , the FENE force converges to the Hookean spring force as  $l_{\max} \rightarrow \infty$ .

In this work, the focus is on developing deterministic multiscale methods for simulating the flow of a suspension of FENE-type dumbbells<sup>1</sup> in a Newtonian solvent. This is an imposing challenge in itself because (as discussed in Section 1.3) for a  $d$ -dimensional flow, the Fokker–Planck equation is posed in  $2d$  spatial dimensions, where we consider  $d = 2$  or  $3$ . Solving

<sup>1</sup>In Chapter 2, we consider a more general class of spring potentials that include the FENE potential as a special case.

this high-dimensional equation is a large-scale computational problem, which requires highly specialised numerical methods. Replacing dumbbells with bead-spring chains would clearly make the problem far more challenging still. The development of methods to treat the bead-spring chain case efficiently using deterministic algorithms (as opposed to Monte Carlo approaches) has received attention in the literature recently (see Section 1.4). The extension of the work herein to the bead-spring case is the subject of ongoing research.

### 1.3 The micro-macro model

With the background material developed in the previous two sections it is now possible to derive the Navier–Stokes–Fokker–Planck model for dilute polymeric fluids. As indicated above, we consider a dilute solution of polymer chains suspended in a Newtonian solvent, and we assume that individual polymer chains do not interact with one another, but can be convected, stretched and oriented by the macroscopic velocity field, and are also subject to thermal agitation due to the motion of the solvent molecules.

Suppose the fluid is confined to a physical domain  $\Omega$ , assumed to be a bounded open set in  $\mathbb{R}^d$ ,  $d = 2$  or  $3$ , and that appropriate boundary conditions are imposed on  $\partial\Omega$ . The conservation equations for polymeric fluids are the same as for the Newtonian case, but the presence of polymer molecules contributes a *polymeric extra-stress*, represented by the tensor  $\underline{\underline{\tau}}$ . That is, the total stress tensor  $\underline{\underline{\sigma}}$  is given by

$$\underline{\underline{\sigma}} = -p\underline{\underline{I}} + \mu_s(\underline{\underline{\nabla}}_x \underline{u} + (\underline{\underline{\nabla}}_x \underline{u})^T) + \underline{\underline{\tau}}, \quad (1.3.1)$$

where in this case the viscosity is labelled with a subscript  $s$  to indicate that it comes from the solvent. Combining (1.3.1) with the conservation of mass and momentum equations yields a modified form of the Navier–Stokes equations in which the divergence of  $\underline{\underline{\tau}}$  arises as a source term. Thus, the model problem takes the following form.

Find  $\underline{u} : (\underline{x}, t) \in \Omega \times \mathbb{R} \mapsto \underline{u}(\underline{x}, t) \in \mathbb{R}^d$  and  $p : (\underline{x}, t) \in \Omega \times \mathbb{R} \mapsto p(\underline{x}, t) \in \mathbb{R}$  such that

$$\frac{\partial \underline{u}}{\partial t} + (\underline{u} \cdot \underline{\underline{\nabla}}_x) \underline{u} - \nu_s \Delta_x \underline{u} + \underline{\underline{\nabla}}_x p = \frac{1}{\rho} \underline{\underline{\nabla}}_x \cdot \underline{\underline{\tau}} \quad \text{in } \Omega \times (0, T], \quad (1.3.2)$$

$$\underline{\underline{\nabla}}_x \cdot \underline{u} = 0 \quad \text{in } \Omega \times (0, T], \quad (1.3.3)$$

$$\underline{u}(\underline{x}, 0) = \underline{u}^0(\underline{x}) \quad \forall \underline{x} \in \Omega, \quad (1.3.4)$$

where  $\nu_s$  is the kinematic solvent viscosity,  $\nu_s := \mu_s/\rho$ . The system is completed by specifying appropriate boundary conditions on  $\partial\Omega$ .

The system (1.3.2)–(1.3.4) models the macroscopic flow of a polymeric fluid, and the contributions of microscopic polymer molecules enter through the extra-stress tensor,  $\underline{\underline{\tau}}$ . In the case that the polymer molecules are represented by coarse-grained objects (*e.g.* dumbbells), it turns out that  $\underline{\underline{\tau}}$  can be computed in terms of a statistical average of the probability density function describing the distribution of configurations of polymer molecules within the fluid.<sup>2</sup> The probability density function for dumbbell configurations will henceforth be denoted  $\psi$ , and the idea of the deterministic multiscale method is to compute  $\psi$  directly by solving a partial differential equation (the high-dimensional Fokker–Planck equation alluded to above) so that  $\underline{\underline{\tau}}$  can be computed and fed into the macroscopic system (1.3.2)–(1.3.4).

---

<sup>2</sup>The precise equation for computing  $\underline{\underline{\tau}}$  is known as *Kramers expression*, and it is discussed below in Section 1.3.3.

### 1.3.1 Derivation of the Fokker–Planck equation

In this section the Fokker–Planck equation for polymeric fluids that governs  $\psi$  is derived from first principles. For the purposes of the derivation, it suffices to consider the general spring force law (1.2.1). Similar derivations can be found in Bird *et. al.* [23], the Ph.D. thesis of Lozinski [90] or the paper by Barrett & Süli [11].

First of all, consider an isolated dumbbell immersed in a Newtonian solvent with fluid velocity given at point  $\underline{x} \in \Omega$  and time  $t \in [0, T]$ , where  $T \in \mathbb{R}_{>0}$ , by  $\underline{u}(\underline{x}, t)$ . Denote by  $\underline{r}_1(t), \underline{r}_2(t) \in \Omega \subset \mathbb{R}^d$  the position vectors of the two masses of the dumbbell at time  $t$ , where  $\Omega$  is referred to as *physical space*. For the purpose of this derivation we assume that  $\Omega = \mathbb{R}^d$ ; this allows us to avoid complications associated with the behaviour of dumbbells at the domain boundary. From Section 1.3.2 onwards, we shall assume that  $\Omega$  is a bounded subset of  $\mathbb{R}^d$ .

As in Figure 1.1(b), the centre of mass,  $\underline{x}(t)$ , and configuration vector,  $\underline{q}(t)$ , are defined as:

$$\underline{x}(t) = (\underline{r}_1(t) + \underline{r}_2(t)) / 2 \quad \text{and} \quad \underline{q}(t) = \underline{r}_2(t) - \underline{r}_1(t). \quad (1.3.5)$$

Assuming that  $\Omega$  is convex, we then have that  $\underline{x}(t) \in \Omega$ . Also, let the *configuration space* be the set of all admissible configuration vectors (which we assume to be a time-invariant domain), *i.e.*,

$$D = \{\underline{q} \in \mathbb{R}^d : \underline{q} = \underline{r}_2 - \underline{r}_1, \text{ for all admissible } \underline{r}_1, \underline{r}_2 \in \Omega\}.$$

For example, for Hookean dumbbells, the configuration space is all of  $\mathbb{R}^d$ , whereas for FENE dumbbells we have  $D = B(0, l_{\max})$ , where  $B(0, s) \subset \mathbb{R}^d$  is the ball centered at the origin with radius  $s$ . It is more natural to treat the Fokker–Planck equation in  $(\underline{x}, \underline{q})$ -coordinates than in  $(\underline{r}_1, \underline{r}_2)$ -coordinates because with the FENE model for example, for a given  $\underline{r}_1$ , we have  $\underline{r}_2 \in B(\underline{r}_1, l_{\max})$ , *i.e.* in contrast to the vectors  $(\underline{x}, \underline{q}) \in \Omega \times D$ , the domains of  $\underline{r}_1$  and  $\underline{r}_2$  cannot be decoupled in this case.

Considering an isolated dumbbell, Newton’s Second Law can be applied to the  $i^{\text{th}}$  bead such that  $\underline{F}_i^{\text{total}} = m_i \underline{a}_i$ , where  $\underline{a}_i$  is the acceleration of bead  $i = 1, 2$  and  $\underline{F}_i^{\text{total}}$ , the total force on bead  $i$ , is the sum of the following components:

- $\underline{F}_i^{\text{drag}}$ : Drag force due to bead  $i$  moving through the viscous solvent;
- $\underline{B}_i$ : Brownian force due to random collisions of solvent molecules with bead  $i$ ;
- $\underline{F}_i$ : The spring force on bead  $i$ , *e.g.* (1.2.3).

Hence, we have the following force balance equations for beads 1 and 2:

$$\begin{aligned} m_1 \underline{a}_1(t) &= \underline{F}_1^{\text{drag}}(t) + \underline{B}_1(t) + \underline{F}(\underline{r}_2(t) - \underline{r}_1(t)), \\ m_2 \underline{a}_2(t) &= \underline{F}_2^{\text{drag}}(t) + \underline{B}_2(t) + \underline{F}(\underline{r}_1(t) - \underline{r}_2(t)). \end{aligned}$$

We model the hydrodynamic drag force,  $\underline{F}^{\text{drag}}$ , using Stokes’ law for the viscous drag on a sphere at low Reynolds number [1], *i.e.*

$$\underline{F}_i^{\text{drag}} = -\zeta \left( \frac{d\underline{r}_i}{dt}(t) - \underline{u}(\underline{r}_i(t), t) \right),$$

where the term inside the brackets is the velocity of bead  $i$  relative to the velocity of the solvent, and  $\zeta$  is the friction/drag coefficient.

Following Schieber & Öttinger [109] we consider the zero-mass limit for the dumbbell beads and therefore multiplying through by  $dt$  we obtain the following two equations:

$$\zeta (d\mathfrak{r}_1(t) - \mathfrak{u}(\mathfrak{r}_1(t), t) dt) = \mathfrak{B}_1(t) dt + \mathfrak{F}(\mathfrak{r}_2(t) - \mathfrak{r}_1(t)) dt, \quad (1.3.6)$$

$$\zeta (d\mathfrak{r}_2(t) - \mathfrak{u}(\mathfrak{r}_2(t), t) dt) = \mathfrak{B}_2(t) dt + \mathfrak{F}(\mathfrak{r}_1(t) - \mathfrak{r}_2(t)) dt. \quad (1.3.7)$$

Equations (1.3.6) and (1.3.7) are referred to as Langevin's equations [36] for the dumbbell. The Brownian force is defined as

$$\mathfrak{B}_i(t) dt := \sqrt{2k_B\mathcal{T}\zeta} d\mathfrak{W}_i(t), \quad (1.3.8)$$

where  $\mathfrak{W}_i(t)$  is a  $d$ -component Wiener process [104],  $k_B = 1.38 \times 10^{-23} \text{ m}^2\text{kg s}^{-2}\text{K}^{-1}$  is Boltzmann's constant and  $\mathcal{T}$  is the absolute temperature measured in Kelvin, K. The coefficient  $\sqrt{2k_B\mathcal{T}\zeta}$  in (1.3.8) is due to the Einstein–Smoluchowski relation, which determines the diffusion coefficient in Brownian motion [98]. Therefore, (1.3.6), (1.3.7) can be rewritten as follows:

$$d \begin{bmatrix} \mathfrak{r}_1(t) \\ \mathfrak{r}_2(t) \end{bmatrix} = \begin{bmatrix} \mathfrak{u}(\mathfrak{r}_1(t), t) + \zeta^{-1}\mathfrak{F}(\mathfrak{r}_2(t) - \mathfrak{r}_1(t)) \\ \mathfrak{u}(\mathfrak{r}_2(t), t) + \zeta^{-1}\mathfrak{F}(\mathfrak{r}_1(t) - \mathfrak{r}_2(t)) \end{bmatrix} dt + \sqrt{\frac{2k_B\mathcal{T}}{\zeta}} d \begin{bmatrix} \mathfrak{W}_1(t) \\ \mathfrak{W}_2(t) \end{bmatrix}. \quad (1.3.9)$$

Defining

$$\mathfrak{X}(t) := \begin{bmatrix} \mathfrak{r}_1(t) \\ \mathfrak{r}_2(t) \end{bmatrix}, \quad \mathfrak{W}(t) := \begin{bmatrix} \mathfrak{W}_1(t) \\ \mathfrak{W}_2(t) \end{bmatrix}, \quad \mathfrak{\sigma} := \sqrt{\frac{2k_B\mathcal{T}}{\zeta}} \mathfrak{I},$$

$$\mathfrak{b}(\mathfrak{X}(t)) := \begin{bmatrix} \mathfrak{u}(\mathfrak{r}_1(t), t) + \zeta^{-1}\mathfrak{F}(\mathfrak{r}_2(t) - \mathfrak{r}_1(t)) \\ \mathfrak{u}(\mathfrak{r}_2(t), t) + \zeta^{-1}\mathfrak{F}(\mathfrak{r}_1(t) - \mathfrak{r}_2(t)) \end{bmatrix},$$

and allowing, more generally, the possible dependence of  $\mathfrak{\sigma}$  on  $\mathfrak{X}(t)$ , (1.3.9) can be written as the following stochastic differential equation:

$$d\mathfrak{X}(t) = \mathfrak{b}(\mathfrak{X}(t)) + \mathfrak{\sigma}(\mathfrak{X}(t)) d\mathfrak{W}(t), \quad \mathfrak{X}(0) = \mathfrak{X}. \quad (1.3.10)$$

We can now use the *forward Kolmogorov equation* to obtain a partial differential equation for the evolution of the probability density function of the stochastic process  $t \mapsto \mathfrak{X}(t)$  (see Corollary 5.2.10 in [78]).

**Theorem 1.3.1 (Forward Kolmogorov (Fokker–Planck) equation)** *Let the random variable  $\mathfrak{X}(t)$  have a density function  $(z, t) \mapsto \psi(z, t)$  of class  $C^{2,1}(\mathbb{R}^d \times \mathbb{R}^d, [0, \infty))$  (i.e. twice continuously differentiable with respect to  $z \in \mathbb{R}^d \times \mathbb{R}^d$  and once with respect to  $t$ ), and let  $\mathfrak{X}(0) = \mathfrak{X}$  be a square-integrable random variable with density function  $\psi^0 \in C^2(\mathbb{R}^d \times \mathbb{R}^d)$ . Also, suppose that  $\mathfrak{b}$  and  $\mathfrak{\sigma}$  in (1.3.10) are globally Lipschitz continuous, and  $\mathfrak{a}(z) = \mathfrak{\sigma}(z)\mathfrak{\sigma}(z)^\mathbb{T}$ . Then,*

$$\frac{\partial \psi}{\partial t} + \sum_{j=1}^{2d} \frac{\partial}{\partial z_j} (b_j \psi) = \frac{1}{2} \sum_{i,j=1}^{2d} \frac{\partial^2}{\partial z_i \partial z_j} (a_{ij} \psi), \quad (1.3.11)$$

in  $\mathbb{R}^{2d} \times (0, \infty)$ , with  $\psi(z, 0) = \psi^0(z)$  for  $z \in \mathbb{R}^d$ .

**Remark 1.3.2** The Hookean spring force satisfies the global Lipschitz continuity assumption in Theorem 1.3.1, whereas the FENE force does not. Indeed, the FENE force is only locally Lipschitz on  $D$ , and it is not defined on all of  $\mathbb{R}^d$ . Nevertheless, we shall proceed on the basis of the conjecture that Theorem 1.3.1 applies in the case of the FENE model as well.  $\diamond$

Applying Theorem 1.3.1 to (1.3.10) yields:

$$\begin{aligned} \frac{\partial \psi^{12}}{\partial t} + \nabla_{r_1} \cdot \left[ \underline{u}(r_1, t) \psi^{12} + \frac{1}{\zeta} \underline{F}(r_2 - r_1) \psi^{12} \right] \\ + \nabla_{r_2} \cdot \left[ \underline{u}(r_2, t) \psi^{12} + \frac{1}{\zeta} \underline{F}(r_1 - r_2) \psi^{12} \right] = \frac{k_B \mathcal{T}}{\zeta} \Delta_{r_1} \psi^{12} + \frac{k_B \mathcal{T}}{\zeta} \Delta_{r_2} \psi^{12}, \end{aligned} \quad (1.3.12)$$

where  $\psi^{12}$  denotes the probability density function with respect to  $(r_1, r_2)$ -coordinates. Changing to  $(\underline{x}, \underline{q})$ -coordinates and letting  $\psi(\underline{x}, \underline{q}, t) := \psi^{12}(r_1, r_2, t)$ , we obtain

$$\begin{aligned} \frac{\partial \psi}{\partial t} + \nabla_{\underline{q}} \cdot \left( \left[ \underline{u}(\underline{x} + \underline{q}/2, t) - \underline{u}(\underline{x} - \underline{q}/2, t) \right] \psi - \frac{2}{\zeta} \underline{F}(\underline{q}) \psi \right) \\ + \nabla_{\underline{x}} \cdot \left( \frac{\underline{u}(\underline{x} + \underline{q}/2, t) + \underline{u}(\underline{x} - \underline{q}/2, t)}{2} \psi \right) = \frac{k_B \mathcal{T}}{2\zeta} \Delta_{\underline{x}} \psi + \frac{2k_B \mathcal{T}}{\zeta} \Delta_{\underline{q}} \psi, \end{aligned} \quad (1.3.13)$$

where we have used the fact that  $\underline{F}(\underline{q}) = -\underline{F}(-\underline{q})$  (cf. (1.2.1)).

In order to simplify (1.3.13) further, we adopt the *local homogeneity assumption*, which states that  $\underline{u}$  and  $\psi$  are linear in  $\underline{x}$  on the length-scale of a dumbbell. This is a plausible assumption because the dumbbell length-scale is typically orders of magnitude smaller than the macroscopic length-scale. Using, in (1.3.13), Taylor series expansions of  $\underline{u}(\underline{x} + \underline{q}/2)$  and  $\underline{u}(\underline{x} - \underline{q}/2)$  about the point  $\underline{x}$  up the linear terms and ignoring quadratic and higher-order terms yields:

$$\frac{\partial \psi}{\partial t} + \nabla_{\underline{x}} \cdot (\underline{u} \psi) + \nabla_{\underline{q}} \cdot \left( \underline{\kappa} \underline{q} \psi - \frac{2}{\zeta} \underline{F}(\underline{q}) \psi \right) = \frac{k_B \mathcal{T}}{2\zeta} \Delta_{\underline{x}} \psi + \frac{2k_B \mathcal{T}}{\zeta} \Delta_{\underline{q}} \psi, \quad (1.3.14)$$

where  $\underline{\kappa} := \nabla_{\underline{x}} \underline{u}$  is a standard short-hand notation for  $\nabla_{\underline{x}} \underline{u}$ . Note that by incompressibility of  $\underline{u}$ ,  $\text{tr}(\underline{\kappa}) = 0$ .

The next step is to put (1.3.14) into nondimensional form by scaling as follows:

$$\underline{x} := L_0 \hat{\underline{x}}, \quad \underline{q} := l_0 \hat{\underline{q}}, \quad \underline{u} := U_0 \hat{\underline{u}}, \quad t := (L_0/U_0) \hat{t}, \quad \psi := \hat{\psi}/l_0^d, \quad (1.3.15)$$

where  $l_0 := \sqrt{k_B \mathcal{T}/H}$  is the characteristic length-scale of a dumbbell and  $L_0, U_0$  are the characteristic length and velocity of the macroscopic flow, respectively; we let  $\hat{\underline{F}}(\hat{\underline{q}}) := \hat{U}'(\frac{1}{2}|\hat{\underline{q}}|^2)\hat{\underline{q}}$ , where  $\hat{U}(s) := l_0^{-2}U(l_0^2 s)$ .

Applying (1.3.15) to (1.3.14) and omitting the hat superscripts for notational convenience yields:

$$\frac{U_0}{L_0} \frac{\partial \psi}{\partial t} + \frac{U_0}{L_0} \nabla_{\underline{x}} \cdot (\underline{u} \psi) + \nabla_{\underline{q}} \cdot \left( \frac{U_0}{L_0} \underline{\kappa} \underline{q} \psi - \frac{1}{2\lambda} \underline{F}(\underline{q}) \psi \right) = \frac{1}{2\lambda} \Delta_{\underline{q}} \psi + \frac{1}{8\lambda} \left( \frac{l_0}{L_0} \right)^2 \Delta_{\underline{x}} \psi, \quad (1.3.16)$$

where  $\lambda := \zeta/4H$  is the characteristic relaxation time of a dumbbell.

Note that for the FENE case  $|\hat{q}| \in [0, \sqrt{b})$ , where  $b := Hl_{\max}^2/k_B\mathcal{T}$  and therefore the configuration space in nondimensional form is  $D = B(0, \sqrt{b}) \subset \mathbb{R}^d$ , and (again, omitting the hat superscripts,) (1.2.3) becomes:

$$U(\tfrac{1}{2}|q|^2) := -\frac{b}{2} \ln \left( 1 - \frac{|q|^2}{b} \right), \quad \mathcal{F}(q) = \frac{q}{1 - |q|^2/b}. \quad (1.3.17)$$

The dimensionless parameter  $b$  is typically in the range  $[10, 100]$ . In [62], Jourdain, Lelièvre and Le Bris showed that for the stochastic differential equation modelling a suspension of FENE dumbbells (which corresponds to the deterministic Fokker–Planck-based model considered here), the solution exists and has trajectorial uniqueness if, and only if,  $b > 2$  (cf. also Example 1.2 in [12]). Hence, throughout the rest of this work, we assume that  $b \in (2, \infty)$  for the FENE potential.

Multiplying (1.3.16) through by  $L_0/U_0$  gives:

$$\frac{\partial \psi}{\partial t} + \nabla_x \cdot (u\psi) + \nabla_q \cdot \left( \kappa q \psi - \frac{1}{2\text{Wi}} \mathcal{F}(q)\psi \right) = \frac{1}{2\text{Wi}} \Delta_q \psi + \frac{1}{8\text{Wi}} \left( \frac{l_0}{L_0} \right)^2 \Delta_x \psi, \quad (1.3.18)$$

where  $\text{Wi} := \lambda U_0/L_0$  is the nondimensional *Weissenberg number*, which is the ratio of the microscopic to macroscopic time-scales, and is typically on the order of 1 or 10.

Equation (1.3.18) contains an  $x$ -diffusion term, referred to as the center-of-mass diffusion term. The standard approach in the literature has been to discard this term outright because its coefficient,  $\varepsilon := (1/8\text{Wi})(l_0/L_0)^2$ , is typically on the order of  $10^{-8}$  [19]. However, it has been recognised by Barrett & Süli [11] that, from the point of view of analysis, this simplification is counterproductive because when the centre-of-mass diffusion term is neglected (1.3.18) becomes a degenerate parabolic equation that exhibits hyperbolic behaviour in physical space. We shall return to this point in Chapters 5 and 6. Since in Chapters 2 to 4 at least the emphasis is on other, largely algorithmic, questions, due to its negligible size the centre-of-mass diffusion coefficient is simply set to zero. Hence, in Chapters 2 to 4, we consider the Fokker–Planck equation with no  $x$ -diffusion, *i.e.*,

$$\frac{\partial \psi}{\partial t} + \nabla_x \cdot (u\psi) + \nabla_q \cdot \left( \kappa q \psi - \frac{1}{2\text{Wi}} \mathcal{F}(q)\psi \right) = \frac{1}{2\text{Wi}} \Delta_q \psi. \quad (1.3.19)$$

We shall re-introduce the  $x$ -diffusion term in Chapters 5 and 6. Notice that (at least in the case of FENE or Hookean dumbbells) the Fokker–Planck equation (1.3.19) contains an unbounded advection coefficient  $\mathcal{F}$ . This is inconvenient from the point of view of analysis. Therefore we shall focus on the following Kolmogorov symmetrisation [73] of the Fokker–Planck equation, in which the spring force,  $\mathcal{F}$ , has been absorbed into a weighted diffusion term,

$$\frac{\partial \psi}{\partial t} + \nabla_x \cdot (u\psi) + \nabla_q \cdot (\kappa q \psi) = \frac{1}{2\text{Wi}} \nabla_q \cdot \left( M \nabla_q \left( \frac{\psi}{M} \right) \right), \quad (1.3.20)$$

where  $M$  is the (normalised) *Maxwellian* defined by

$$q \mapsto M(q) := \frac{1}{Z} \exp \left( -U(\tfrac{1}{2}|q|^2) \right) \in L^1(D), \quad Z := \int_D \exp \left( -U(\tfrac{1}{2}|q|^2) \right) dq. \quad (1.3.21)$$

The Maxwellian transformation used in (1.3.20) allows us to circumvent analytical difficulties introduced by the unbounded convection term,  $\underline{F}$ . In Section 2, we shall also consider an alternative transformation of (1.3.19) due to Chauvière & Lozinski [33] that allows us to deal with the unbounded convection term in a different manner, and hence a range of theoretical results can be proved for the Chauvière–Lozinski transformed equation also.

The function  $(\underline{x}, \underline{q}, t) \mapsto \psi(\underline{x}, \underline{q}, t)$  represents the probability, at time  $t$ , of finding a dumbbell with center of mass in the volume element  $\underline{x} + d\underline{x}$  and orientation vector in the volume element  $\underline{q} + d\underline{q}$ . Recall that the above derivation of the Fokker–Planck equation assumed that  $\Omega = \mathbb{R}^d$ , but we shall henceforth assume that  $\Omega$  is a bounded subset of  $\mathbb{R}^d$ . Also, it is crucial to note that (1.3.20) is posed in  $2d$  spatial dimensions, plus time. Since the computational complexity of classical numerical methods grows exponentially with the dimension of the spatial domain, the high-dimensionality of (1.3.20) poses a significant computational challenge. Developing a fully practical computational framework for this high-dimensional equation is one of the central goals of this work.

### 1.3.2 Properties of the probability density function

Since  $\psi$  is a probability density function (pdf) for each  $\underline{x} \in \Omega$ , the initial condition should be nonnegative:

$$\psi(\underline{x}, \underline{q}, 0) = \psi^0(\underline{x}, \underline{q}) \geq 0, \quad \text{for a.e. } (\underline{x}, \underline{q}) \in \Omega \times D, \quad (1.3.22)$$

and should also satisfy the following normalisation property:

$$\int_D \psi^0(\underline{x}, \underline{q}) d\underline{q} = 1, \quad \text{for a.e. } \underline{x} \in \Omega. \quad (1.3.23)$$

Assuming that  $\psi$  is sufficiently smooth to ensure that the manipulations below are meaningful, we now show that (1.3.23) is preserved for  $t \in (0, T]$  for solutions of (1.3.20) provided that a suitable boundary condition is imposed. Specifically, suppose that

$$\left( \frac{1}{2\text{Wi}} M \nabla_{\underline{q}} \left( \frac{\psi}{M} \right) - \underline{\kappa} \underline{q} \psi \right) \cdot \underline{n}_{\partial D} = 0 \quad \text{on } \Omega \times \partial D \times (0, T], \quad (1.3.24)$$

where  $\underline{n}_{\partial D}$  is the outward unit normal on  $\partial D$ . Then, integrating (1.3.20) in configuration space and in time and applying the divergence theorem gives, for  $\underline{x} \in \Omega$  and  $t \in (0, T]$ :

$$\begin{aligned} & \int_D \psi(\underline{x}, \underline{q}, t) d\underline{q} - \int_D \psi(\underline{x}, \underline{q}, 0) d\underline{q} + \int_0^t \nabla_{\underline{x}} \cdot \left( \underline{u}(\underline{x}, \tau) \left( \int_D \psi(\underline{x}, \underline{q}, \tau) d\underline{q} \right) \right) d\tau \\ &= \int_0^t \int_D \nabla_{\underline{q}} \cdot \left( -\underline{\kappa} \underline{q} \psi(\underline{x}, \underline{q}, \tau) + \frac{1}{2\text{Wi}} M \nabla_{\underline{q}} \left( \frac{\psi(\underline{x}, \underline{q}, \tau)}{M} \right) \right) d\underline{q} d\tau \\ &= \int_0^t \int_{\partial D} \left( -\underline{\kappa} \underline{q} \psi(\underline{x}, \underline{q}, \tau) + \frac{1}{2\text{Wi}} M \nabla_{\underline{q}} \left( \frac{\psi(\underline{x}, \underline{q}, \tau)}{M} \right) \right) \cdot \underline{n}_{\partial D} ds d\tau = 0. \end{aligned} \quad (1.3.25)$$

Let  $\varrho(\underline{x}, t)$  be defined as follows:

$$\varrho(\underline{x}, t) := \int_D \psi(\underline{x}, \underline{q}, t) d\underline{q}.$$

Then (1.3.25) can be rewritten, for  $\underline{x} \in \Omega$  and  $t \in (0, T]$ , as

$$\varrho(\underline{x}, t) - \varrho(\underline{x}, 0) + \int_0^t \nabla_x \cdot (\underline{y}(\underline{x}, \tau) \varrho(\underline{x}, \tau)) \, d\tau = 0.$$

Now, assuming that  $\underline{y} \cdot \underline{n}_{\partial\Omega} = 0$  on  $\partial\Omega \times (0, T]$ , where  $\underline{n}_{\partial\Omega}$  is the unit outward normal to  $\partial\Omega$ , it follows on integration over  $\Omega$  and using the divergence theorem that

$$\int_{\Omega} \varrho(\underline{x}, t) \, d\underline{x} = \int_{\Omega} \varrho(\underline{x}, 0), \quad t \in (0, T), \quad (1.3.26)$$

and hence the following result has been established.

**Lemma 1.3.3** *For  $t \in [0, T]$  and  $\varrho^0(\underline{x}) := \int_D \psi^0(\underline{x}, \underline{q}) \, d\underline{q}$ , we have that*

$$\int_{\Omega} \varrho(\underline{x}, t) \, d\underline{x} = \int_{\Omega} \varrho^0(\underline{x}) \, d\underline{x} \quad (1.3.27)$$

for all, sufficiently smooth, solutions of (1.3.20) satisfying the boundary condition (1.3.24).

An important consideration that will be returned to in subsequent sections is whether results analogous to Lemma 1.3.3 can be established for solutions (both continuous and discrete) based on the weak formulation of (1.3.20).

It is also desirable to preserve the property (1.3.22) for  $t \in (0, T]$ . This nonnegativity property is considered for weak solutions of the Fokker–Planck equation (*cf.* Lemma 2.3.4) as well as for approximate solutions obtained via a Galerkin spectral approach (*cf.* Remark 2.8.2) in Section 2.

### 1.3.3 Polymeric extra-stress

As indicated above, in the context of the coupled Navier–Stokes–Fokker–Planck system, the purpose of solving (1.3.20) is so that the polymeric extra-stress tensor,  $\underline{\tau}$ , can be computed and fed into the right-hand side of (1.3.2). The polymeric extra-stress tensor is determined by the following equality, known as the *Kramers expression*:

$$\underline{\tau}(\underline{x}, t) = n_p \left( \int_D \underline{F}(\underline{q}) \otimes \underline{q} \psi \, d\underline{q} - \underline{I} \right), \quad (\underline{x}, t) \in \Omega \times (0, T], \quad (1.3.28)$$

where  $n_p$  is the polymer number density, *i.e.* the number of polymer molecules per unit volume. For a derivation of (1.3.28), see, for example, [82]. Note that it follows from (1.3.28) that  $\underline{\tau}$  is symmetric. Since  $\underline{\tau}$  enters into (1.3.2) only via its divergence, the constant  $n_p \underline{I}$  in (1.3.28) has no effect in the coupled system and therefore we ignore it from now on. Nondimensionalising (1.3.28) according to (1.3.15) (and omitting the hat superscripts as before) gives

$$\underline{\tau}(\underline{x}, t) = n_p k_B \mathcal{T} \int_D \underline{F}(\underline{q}) \otimes \underline{q} \psi(\underline{x}, \underline{q}, t) \, d\underline{q}. \quad (1.3.29)$$

At this point we make the specific assumption that  $\underline{F}$  is the FENE force in order to derive the full Navier–Stokes–Fokker–Planck system, in nondimensional form, for a suspension of FENE dumbbells.

It can be shown that for a dilute solution of FENE dumbbells in shear flow, the  $(1, 2)$ -component of  $\underline{\underline{\tau}}$  is approximated by

$$\tau_{12} \approx \dot{\gamma} \lambda n_p k_B \mathcal{T} \left( \frac{b}{b+d+2} \right), \quad (1.3.30)$$

where  $\dot{\gamma}$  is the shear rate (see [23]). Equation (1.3.30) is an asymptotic expression for  $\tau_{12}$  that is valid for small  $\dot{\gamma}$ . Therefore, by analogy with Newtonian fluids, the *polymeric viscosity*,  $\mu_p$ , for FENE dumbbell suspensions is defined as

$$\mu_p := \lambda n_p k_B \mathcal{T} \left( \frac{b}{b+d+2} \right), \quad (1.3.31)$$

so that (1.3.29) can be rewritten:

$$\frac{1}{\rho} \underline{\underline{\tau}}(\underline{x}, t) = \frac{\nu_p}{\lambda} \frac{b+d+2}{b} \int_D \underline{\underline{F}}(\underline{q}) \otimes \underline{q} \psi(\underline{x}, \underline{q}, t) d\underline{q}, \quad (1.3.32)$$

where the equation has been divided through by the density  $\rho$  as in (1.3.2), and  $\nu_p := \mu_p/\rho$ .

Equation (1.3.32) provides a bridge between the Fokker–Planck equation and the Navier–Stokes equation. The fully coupled form of the micro-macro system is discussed in the next section.

### 1.3.4 The coupled Navier–Stokes–Fokker–Planck system

The Fokker–Planck equation and Kramers expression have been written in terms of nondimensional variables in (1.3.20) and (1.3.32), respectively. Thus, it remains to nondimensionalise the Navier–Stokes equations, (1.3.2), (1.3.3), in the same manner. The mass conservation equation, (1.3.3), contains only one nonzero term and therefore rescaling is trivial. Applying (1.3.15) in the momentum equation, letting  $\nu = \nu_s + \nu_p$ , rescaling the pressure as  $p = U_0^2 \hat{p}$  and using (1.3.32) yields (on omitting the hat superscripts):

$$\frac{\partial \underline{u}}{\partial t} + (\underline{u} \cdot \underline{\nabla}_x) \underline{u} + \underline{\nabla}_x p = \frac{\gamma}{\text{Re}} \Delta_x \underline{u} + \frac{b+d+2}{b} \frac{1-\gamma}{\text{Re Wi}} \underline{\nabla}_x \cdot \underline{\underline{\tau}}, \quad (1.3.33)$$

where  $\text{Re} := L_0 U_0 / \nu$  (*i.e.* the Reynolds number) and  $\gamma := \nu_s / \nu$  are nondimensional parameters.<sup>3</sup> Note that we have absorbed the coefficients on the right-hand side of (1.3.32) into the momentum equation in order to perform nondimensionalisation.

Combining the equations heretofore derived gives the following system:

$$\frac{\partial \underline{u}}{\partial t} + (\underline{u} \cdot \underline{\nabla}_x) \underline{u} + \underline{\nabla}_x p = \frac{\gamma}{\text{Re}} \Delta_x \underline{u} + \frac{b+d+2}{b} \frac{1-\gamma}{\text{Re Wi}} \underline{\nabla}_x \cdot \underline{\underline{\tau}}, \quad (\underline{x}, t) \in \Omega \times (0, T], \quad (1.3.34)$$

$$\underline{\nabla}_x \cdot \underline{u} = 0, \quad (\underline{x}, t) \in \Omega \times (0, T], \quad (1.3.35)$$

$$\frac{\partial \psi}{\partial t} + \underline{\nabla}_x \cdot (\underline{u} \psi) + \underline{\nabla}_q \cdot (\underline{\kappa} \underline{q} \psi) = \frac{1}{2\text{Wi}} \underline{\nabla}_q \cdot \left( M \underline{\nabla}_q \frac{\psi}{M} \right), \quad (\underline{x}, \underline{q}, t) \in \Omega \times D \times (0, T], \quad (1.3.36)$$

$$\underline{\underline{\tau}}(\underline{x}, t) = \int_D \underline{\underline{F}} \otimes \underline{q} \psi(\underline{x}, \underline{q}, t) d\underline{q}, \quad (\underline{x}, t) \in \Omega \times (0, T], \quad (1.3.37)$$

$$\underline{u}(\underline{x}, 0) = \underline{u}^0(\underline{x}), \quad \underline{x} \in \Omega, \quad \psi(\underline{x}, \underline{q}, 0) = \psi^0(\underline{x}, \underline{q}), \quad (\underline{x}, \underline{q}) \in \Omega \times D. \quad (1.3.38)$$

<sup>3</sup>Hat superscripts have again been dropped in (1.3.33) for notational simplicity; the variables are to be understood as nondimensional.

Equations (1.3.34)–(1.3.38), when supplemented with appropriate boundary conditions, are the coupled Navier–Stokes–Fokker–Planck model for dilute polymeric fluids. Note that the nondimensionalisation used above is the same as the one introduced on page 8 of [82]. In Chapter 4, we also consider a Stokes–Fokker–Planck model in which (1.3.34) is replaced by a simpler linear equation (*cf.* (4.2.4)) that is relevant for modelling creeping flows, *i.e.* in the limit  $\text{Re} \rightarrow 0_+$ .

In the discussion above, we have assumed that both  $\Omega$  and  $D$  are domains in  $\mathbb{R}^d$  so that the Fokker–Planck equation is posed on  $\Omega \times D \subset \mathbb{R}^{2d}$ . However, it is not essential that this is the case and, for example, in [32] the authors considered a micro-macro model in which  $\Omega \subset \mathbb{R}^2$  and  $D \subset \mathbb{R}^3$ . No significant complications are introduced from the theoretical or implementational point of view by allowing the dimensionality of  $D$  and  $\Omega$  to be different, but for the rest of this work we shall restrict our attention to the case when these domains have the same dimensionality.

## 1.4 Literature review of numerics for polymeric fluids

As indicated in the opening of this section, the techniques for numerically simulating polymeric fluids can be grouped into three categories: fully macroscopic methods, stochastic multiscale methods and deterministic multiscale methods. A survey of some of the key literature for each method is presented below.

### 1.4.1 Fully macroscopic methods

Continuum numerical simulations of polymeric fluids have been popular since the 1970’s. In some sense, this is the most natural approach to simulating polymeric fluids because by avoiding consideration of the microscopic length-scales, one can save an enormous amount of computational effort. However, except in certain simple cases (*e.g.* a suspension of Hookean dumbbells, see Section 1.2) in order to derive a closed-form macroscopic model for a polymeric fluid, it is necessary to resort to an ad hoc “closure approximation”, and the shortcomings of such approximations are well documented [65, 83, 128]. Nevertheless, in many situations, macroscopic models are sufficiently accurate to capture the relevant characteristics of polymer flows and in such cases these methods are preferable to using multiscale methods.

A macroscopic computation typically employs standard algorithms of computational fluid dynamics, such as finite element, finite volume, finite difference or spectral methods, but specialised considerations are usually necessary in practice in order to ensure convergence. The challenges of developing continuum numerical methods for polymeric fluids are epitomised by the well-known “high Weissenberg number” problem, which refers to the difficulty of developing numerical methods that remain stable as  $\text{Wi}$  is increased. The development of macroscopic numerical methods for polymer fluids is clearly a very important field of research; a vast literature has been developed and yet there remain many unresolved issues in this area that are the focus of ongoing research. However, since the focus of this work is on multiscale methods, we shall not consider fully macroscopic methods any further here (for a detailed discussion, see the book by Owens & Phillips [103]).

### 1.4.2 Stochastic multiscale methods

An alternative approach that has gained popularity since the early 1990's is to treat the micro-macro model directly by solving the stochastic differential equation (1.3.10) using Monte Carlo-type methods and coupling with deterministic numerical methods for solving the Navier–Stokes equations (1.3.2), (1.3.3). The Monte Carlo method involves distributing a large number of model polymer molecules throughout the computational domain and tracking their motion as they are convected along streamlines and stretched and oriented by a flow. The stress field,  $\underline{\tau}$ , can then be determined by computing ensemble averages, so that the Navier–Stokes equations can then be solved (with source term  $\nabla_x \cdot \underline{\tau}$ ) to determine the macroscopic velocity field, typically using finite elements or some other standard CFD method. In 1992 Öttinger & Laso [79] proposed the first scheme of this type, which is referred to by the acronym CONNFESSIT for “Calculation of Non-Newtonian Flow: Finite Elements and Stochastic Simulation Technique”. Many other flavours of stochastic multiscale methods have subsequently been developed, such as the method of Brownian configuration fields [58] and the Lagrangian particle method [54]. Note also that there has been a lot of interest in the mathematical properties of multiscale stochastic methods. For example, in certain simple flow regimes, existence and uniqueness of solutions have been established for stochastic approximations of suspensions of Hookean and FENE dumbbells in papers by Jourdain, Lelièvre & Le Bris [60, 61, 62].

The stochastic multiscale approach is a computationally intensive procedure – it is little wonder, therefore, that there was no work done in this direction prior to the 1990's. Moreover, a drawback of the stochastic approach is that it introduces a slowly decaying stochastic error (typically  $\mathcal{O}(N^{-1/2})$  as  $N \rightarrow \infty$ , where  $N$  is the number of sample points). Variance reduction techniques were developed to ameliorate this error term and reduce the number of polymer molecules one must track in order to compute an ensemble average to within a given error tolerance (see [67] for an overview of variance reduction in this context). However, even with variance reduction techniques, the presence of stochastic error is a significant limitation of the stochastic approaches and circumventing this is an important motivation for moving to deterministic methods. On the other hand, an important advantage of the stochastic approach is that it scales well with the number of degrees of freedom in the polymer model – this ensures that stochastic methods remain effective when applied to bead-spring chain polymer models [67].

### 1.4.3 Deterministic multiscale methods

As indicated earlier, the deterministic multiscale approach involves solving the coupled Navier–Stokes–Fokker–Planck system directly. This approach has received comparatively little attention, most likely because solving the high-dimensional Fokker–Planck equation is an imposing computational challenge. Nevertheless, literature on this method extends back to the 1970's although the early works in which the Fokker–Planck equation was solved directly were not truly multiscale since simplified flow regimes were considered for which  $\psi$  was assumed to be a function of  $q$  and  $t$  only (problems in which  $\psi$  does not depend on  $\underline{x}$  are often referred to as *homogeneous flows*). For example, Stewart & Sørensen in 1972 [116] used spherical harmonics to solve the Fokker–Planck equation for a steady shear flow of a dilute suspension of rigid dumbbells. Warner [124] applied a similar approach to the study of shear flows of FENE dumbbells. The first work in which a deterministic approach was utilised to simulate

a nonhomogeneous velocity flow was by Fan in 1989 [48], who considered a planar channel flow using a rigid dumbbell polymer model, and also made the simplifying assumption that the physical space convection term,  $\underline{u} \cdot \nabla_x \psi$ , vanished. Fan's work was subsequently built upon by Nayak [99] and Grosso *et al.* [52] who eliminated this assumption on  $\underline{u} \cdot \nabla_x \psi$ .

Recently, the deterministic multiscale approach has been further developed by Lozinski, Chauvière and collaborators, who proposed a spectral method for simulating the micro-macro model for dilute solutions of FENE dumbbells [32, 33, 90, 91, 92]. Similarly, Helzel & Otto [55] solved the micro-macro model arising in the simulation of suspensions of rod-like polymers using finite difference and finite volume methods.

In the papers of Lozinski, Chauvière *et al.* and Helzel & Otto, the authors decomposed the Fokker–Planck equation (1.3.19) (*i.e.* in the nonsymmetrised form) according to

$$\frac{\partial \psi}{\partial t} + (L_x + L_q) \psi = 0, \quad (1.4.1)$$

where

$$L_q \psi = \nabla_q \cdot (\underline{\kappa} \underline{q} \psi) - \frac{1}{2\text{Wi}} \left( \nabla_q \cdot \underline{F}(\underline{q}) \psi + \Delta_q \psi \right), \quad (1.4.2)$$

$$L_x \psi = \nabla_x \cdot (\underline{u} \psi), \quad (1.4.3)$$

and then they used an alternating-direction approach based on the operators  $L_q$  and  $L_x$  to compute numerical solutions.

That is, suppose that  $0 = t^0 < t^1 < \dots < t^n < \dots \leq T$  is a uniform partition of spacing  $\Delta t$  of the interval  $[0, T]$ . A (two-stage) alternating-direction scheme involves approximating the solution,  $\psi$ , by  $\psi_2$  in the following manner: given  $\psi_2(t^n)$ ,  $n \geq 0$ , with  $\psi_2(t^0) = \psi^0$ , find  $\psi_1$  and  $\psi_2$  such that,

$$\frac{\partial \psi_1}{\partial t} + L_q \psi_1 = 0, \quad t \in (t^n, t^{n+1}], \quad \psi_1(t^n) = \psi_2(t^n), \quad (1.4.4)$$

$$\frac{\partial \psi_2}{\partial t} + L_x \psi_2 = 0, \quad t \in (t^n, t^{n+1}], \quad \psi_2(t^n) = \psi_1(t^{n+1}). \quad (1.4.5)$$

A practical alternating-direction numerical method is based on spatial and temporal discretisation of (1.4.4) and (1.4.5).

In the case of the Fokker–Planck equation (with the centre-of-mass diffusion term  $\varepsilon \Delta_x \psi$  omitted), (1.4.4) is a convection-diffusion equation posed on  $D$  and (1.4.5) is a first-order hyperbolic equation on  $\Omega$ . After discretising in space and time, the two-stage scheme described above can be implemented by alternating between applying  $L_x$  to  $\Omega$  cross-sections of  $\Omega \times D$  and  $L_q$  to  $D$  cross-sections of  $\Omega \times D$ . This type of scheme is also referred to as a dimension-splitting or operator-splitting approach. We shall use the three terms (*i.e.* alternating direction/dimension-splitting/operator-splitting) interchangeably in this work, but our preference will be for the name ‘alternating-direction method’, since we believe it is more descriptive than the alternatives.

Using this operator-splitting, the ‘curse of dimensionality’ associated with the numerical solution of the Fokker–Planck equation in  $2d$  dimensions is ameliorated, as the splitting leads to a sequence of  $d$ -dimensional solves at each time step rather than a single  $2d$ -dimensional solve. Also, this splitting of  $L$  allows different numerical methods to be used in  $\Omega$  and  $D$ . In Chapter 3 we consider alternating-direction numerical methods for the FENE Fokker–Planck equation on  $\Omega \times D$  and we use a *heterogeneous* alternating-direction method based

on a finite element method in  $\Omega$  and a single-domain Galerkin spectral in  $D$ . These are appropriate choices because a finite element method is flexible enough to deal with the general domain  $\Omega$ , whereas  $D$  is always a ball in  $\mathbb{R}^d$  and therefore the  $L_q$  operator is well suited to a spectral discretisation via a polar or spherical co-ordinate transformation to a cartesian product domain. Note also that we shall primarily focus on the Maxwellian-transformed Fokker–Planck equation and therefore instead of  $L_q$  as defined in (1.4.2), we shall generally consider the following  $q$ -direction operator:

$$L_q \psi = \nabla_q \cdot (\kappa \underline{q} \psi) - \frac{1}{2\text{Wi}} \nabla_q \cdot \left( M \nabla_q \left( \frac{\psi}{M} \right) \right). \quad (1.4.6)$$

The operators (1.4.2) and (1.4.6) are identical. However, as we discuss in Chapter 2, the natural weak formulation of (1.4.6), in which we use test functions  $\varphi/M$ , is not identical to the standard weak formulation of (1.4.2) in which unweighted test functions,  $\varphi$ , are used.

Lozinski & Chauvière [32, 33, 91] demonstrated that compared to a stochastic method for the FENE dumbbell model, their deterministic multiscale scheme was more efficient in terms of computational cost, and was also more accurate due to the absence of stochastic error for the benchmark problem of laminar flow around a cylindrical obstacle in a channel.

A further interesting observation by Lozinski & Chauvière was that the direct discretisation of (1.3.19) did not lead to a stable numerical method, and instead they used a substitution of the form  $\psi/(1-|q|^2/b)^s$ , for some  $s$  that was chosen on computational grounds (for example, the authors recommended  $s = 2$  and  $s = 2.5$  for  $d = 2$  and  $d = 3$ , respectively [32, 33]). We return to this point in Section 2.4 where we show that the bilinear form corresponding to the Chauvière–Lozinski-transformed FENE Fokker–Planck equation is coercive for  $s > 1/2$ ; thus it is not surprising that Lozinski & Chauvière’s method was unstable when no substitution was used.

Based on the results of Lozinski & Chauvière, it is clear that the deterministic multiscale approach can be effective for models with low-dimensional configuration space. However, it is still an open question whether this approach can be extended to bead-spring chain dumbbell models in which the configuration space has dimension greater than three. There has been some recent work in this direction using numerical methods that were developed for high-dimensional (*i.e.*  $d \gg 3$ ) PDEs. For example, Ammar, Mokdad, Chinesta & Keunings developed a reduced-basis approach and used it to solve the Fokker–Planck equation in configuration space of dimension up to 20 [3, 4]. An alternative idea is to use sparse grids, which have been shown to be effective for solving elliptic and parabolic PDEs in high-dimensional domains [111, 122]. This idea was applied to the Fokker–Planck equation by Delaunay, Lozinski & Owens [38]. Attempts to solve the Fokker–Planck equation for configuration spaces for  $d \gg 3$  are still at an early stage, and indeed the numerical results presented in the literature so far have been for homogeneous flows only. Nevertheless, reduced-basis and sparse-grid methods appear to be promising approaches for this problem and may enable the development of efficient deterministic multiscale methods for simulating suspensions of bead-spring chains.

Clearly the well-posedness of the Navier–Stokes–Fokker–Planck system is a prerequisite for the success of the deterministic multiscale approach. The analysis of the question of existence of solutions to these equations has been the subject of active research in recent years. We shall review some of the recent results and ongoing research in this direction in

Chapter 5 and will consider this question further; see also [10, 12, 11, 86, 87]) and the survey article of Li & Zhang [82], which provides an informative overview of this literature.

## 1.5 Outlook and goals

We are now in a position to give more details on the aims of this work. Our focus is on the deterministic multiscale method. As discussed in Section 1.4, several different deterministic multiscale numerical methods have been developed in the literature, but the numerical analysis of these methods has not previously been considered in detail. The central goal of this work, therefore, is to develop rigorous analysis of deterministic multiscale methods in order to ensure that there is a firm theoretical foundation for this approach, and to explore the existence of weak solutions to the underlying Navier–Stokes–Fokker–Planck system.

We begin in Chapter 2, by focusing on the analysis of a Galerkin spectral method for discretising (1.4.4), *i.e.* the  $q$ -direction part of the Fokker–Planck equation (or equivalently, the Fokker–Planck equation for a homogeneous flow problem). The focus in Chapter 2, is on the Maxwellian-transformed Fokker–Planck equation (*cf.* (1.4.6)), but we also consider the transformation proposed by Chauvière & Lozinski for (1.4.2) in some detail. Numerical methods based on either transformation require careful analysis; the Maxwellian weight arising in the principal part of the symmetrised formulation is degenerate in the sense that it vanishes on  $\partial D$ , and the Chauvière–Lozinski-transformed scheme contains the unbounded convection coefficient  $\underline{F}$ . We also pay particular attention to the practical implementation of the spectral method on  $D$ , and we present numerical results for the cases  $d = 2$  and  $d = 3$ .

In Chapter 3, the Galerkin spectral method developed in Chapter 2 is combined with a finite element method in  $\Omega$  to yield the alternating-direction scheme with which we obtain approximate solutions of (1.3.20). We show that some subtle issues arise in the numerical analysis of such alternating-direction schemes and, as a result, we develop a specialised quadrature-based Galerkin alternating-direction method for the Fokker–Planck equation that is amenable to stability and convergence analysis; this analysis builds upon the arguments in Chapter 2. We also present some computational results in order to provide experimental support for our theoretical results, and to demonstrate the effectiveness of our alternating-direction approach in practice.

The focus in Chapter 4 is on obtaining computational results for the Navier–Stokes–Fokker–Planck system. Our approach is to couple a standard finite element scheme for solving the Navier–Stokes equations with an alternating-direction method from Chapter 3 for the Fokker–Planck equation. Solving the Fokker–Planck equation is the bottleneck step in this algorithm, due to the fact that it is posed on  $\Omega \times D$ . The numerical results in Chapter 4, and indeed in Chapters 2 and 3 as well, are for the FENE dumbbell case only. However, it would be straightforward to apply the methods developed in this work to more general dumbbell spring potentials, such as potentials that satisfy Hypotheses A and B stated in Chapter 2.

We emphasise an important innovation developed in this work: the application of parallel computation to alternating-direction numerical methods for the Fokker–Planck equation. Alternating-direction algorithms are well suited to implementation on parallel computers since they involve solving a large number of independent equations in each time-step. We show in Chapter 3 and 4 that our alternating direction approach can be efficiently implemented in parallel, and this enables us to solve large-scale deterministic multiscale problems that may

otherwise have been computationally intractable (*e.g.* an important large-scale case is when  $\Omega \times D \subset \mathbb{R}^6$ ).

In Chapter 5 we study the question of existence of global-in-time weak solutions to a coupled microscopic-macroscopic bead-spring model with microscopic cut-off, which arises from the kinetic theory of dilute solutions of polymeric liquids with noninteracting polymer chains. The model consists of the unsteady incompressible Navier–Stokes equations in a bounded domain of  $\mathbb{R}^d$ ,  $d = 2$  or  $3$ , for the velocity and the pressure of the fluid, with an elastic extra-stress tensor as the right-hand side in the momentum equation. The extra-stress tensor stems from the random movement of the polymer chains and is defined through the associated probability density function that satisfies a Fokker–Planck-type parabolic equation, a crucial feature of which is the presence of a center-of-mass diffusion term and a cut-off function  $\beta^L(\psi) = \min(\psi, L)$  in the drag term, where  $L \gg 1$ . We establish the existence of global-in-time weak solutions to the model for a general class of spring force potentials including, in particular, the widely used finitely extensible nonlinear elastic (FENE) potential. A key ingredient of the argument is a special testing procedure in the weak formulation of the Fokker–Planck equation, based on the convex entropy function

$$s \in \mathbb{R}_{\geq 0} \mapsto F(s) := s(\ln s - 1) + 1 \in \mathbb{R}_{\geq 0}.$$

In the case of a corotational drag term, passage to the limit as  $L \rightarrow \infty$  recovers the Navier–Stokes–Fokker–Planck model with centre-of-mass diffusion, without cut-off.

In Chapter 6 we construct a Galerkin finite method for the numerical approximation of weak solutions to Navier–Stokes–Fokker–Planck systems with microscopic cut-off. We perform a rigorous passage to the limit as the spatial and temporal discretization parameters tend to zero, and show that a (sub)sequence of these finite element approximations converges to a weak solution of this coupled Navier–Stokes–Fokker–Planck system. The passage to the limit is performed under minimal regularity assumptions on the data. The convergence proof rests on several auxiliary technical results, including the stability, in the Maxwellian-weighted  $H^1$  norm, of the orthogonal projector, in the Maxwellian-weighted  $L^2$  inner product, onto finite element spaces consisting of continuous piecewise linear functions. We establish optimal-order quasi-interpolation error bounds in the Maxwellian-weighted  $L^2$  and  $H^1$  norms, and prove a new elliptic regularity result in the Maxwellian-weighted  $H^2$  norm.



## Chapter 2

# The Fokker–Planck equation in configuration space

### 2.1 Introduction

This section is concerned with the numerical approximation of the  $d$ -dimensional Fokker–Planck equation posed in configuration space  $D := B(\mathbf{0}, \sqrt{b})$  with  $b \in (2, \infty)$ :

$$\frac{\partial \psi}{\partial t} + \nabla_{\underline{q}} \cdot (\underline{\kappa} \underline{q} \psi) = \frac{1}{2\mathbf{Wi}} \nabla_{\underline{q}} \cdot \left( M \nabla_{\underline{q}} \frac{\psi}{M} \right), \quad (\underline{q}, t) \in D \times (0, T], \quad (2.1.1)$$

where the  $d \times d$  tensor  $\underline{\kappa}$  is assumed to belong to  $\mathcal{C}[0, T] := (C[0, T])^{d \times d}$  (*i.e.* it is independent of  $\underline{x}$ ) and is such that  $\text{tr}(\underline{\kappa})(t) = 0$  for all  $t \in [0, T]$ , where  $T \in \mathbb{R}_{>0}$ . It will be assumed throughout that (2.1.1) is supplemented with the following initial and boundary conditions:

$$\psi(\underline{q}, 0) = \psi^0(\underline{q}), \quad \text{for all } \underline{q} \in D, \quad (2.1.2)$$

$$\left( \frac{1}{2\mathbf{Wi}} M \nabla_{\underline{q}} \left( \frac{\psi}{M} \right) - \underline{\kappa} \underline{q} \psi \right) \cdot \underline{n}_{\partial D} = 0 \quad \text{on } \Omega \times \partial D \times (0, T]. \quad (2.1.3)$$

The initial datum  $\psi^0$  is such that  $\psi^0 \geq 0$  and  $\int_D \psi^0(\underline{q}) \, d\underline{q} = 1$ , as in (1.3.22) and (1.3.23). We will henceforth use the notation  $\mathfrak{d}(\underline{q}) := \text{dist}(\underline{q}, \partial D) = \sqrt{b} - |\underline{q}|$ .

The motivation for studying this subproblem is that, as indicated in Chapter 1, an efficient approach to the numerical solution of (1.3.36) in  $2d+1$  variables is based on operator-splitting with respect to  $(\underline{q}, t)$  and  $(\underline{x}, t)$  as in (1.4.4), (1.4.5). Thereby, the resulting time-dependent transport equation with respect to  $(\underline{x}, t)$  is completely standard,  $\psi_t + \nabla_{\underline{x}} \cdot (\underline{u}(\underline{x}, t)\psi) = 0$ , while the transport-diffusion equation with respect to  $(\underline{q}, t)$  is (2.1.1).

The focus of this chapter is on the analysis and implementation of spectral methods for computing numerical solutions of (2.1.1). We emphasise rigour in establishing the analytical properties of the weak formulation of (2.1.1) and also in developing spectral convergence estimates for the numerical methods based on this weak formulation. Most of the material in this chapter follows the paper [71].

As indicated in Chapter 1, we are primarily interested in solving the micro-macro equations for FENE dumbbells. However, the analysis in this chapter is valid for a more general class of spring force laws. Therefore, the following structural hypotheses, which generalise the relevant properties of the FENE spring potential, are adopted.

**Hypothesis A.** The spring potential  $U \in C^1([0, \frac{b}{2}))$  is a nonnegative monotonic increasing function, with  $U(0) = 0$ ,  $\lim_{s \rightarrow b/2_-} U(s) = +\infty$ ,  $\lim_{s \rightarrow b/2_-} (\frac{b}{2} - s)U'(s) < \infty$ .  $\diamond$

Hypothesis A is consistent with the physical requirement that, in order to faithfully model *finite* stretching of polymer chains, the spring force  $\underline{F}(\underline{q})$  should have infinite intensity when the maximum admissible elongation  $|\underline{q}| = \sqrt{b}$  is reached; *i.e.*, the function  $\underline{q} \mapsto U'(\frac{1}{2}|\underline{q}|^2)$  should tend to  $+\infty$  as  $\mathfrak{d}(\underline{q}) \rightarrow 0_+$ .

Recall the definition of the Maxwellian  $M$  for a spring potential  $U$ , (1.3.21). Since, by Hypothesis A,  $U(\frac{1}{2}|\underline{q}|^2) \rightarrow +\infty$  as  $\mathfrak{d}(\underline{q}) \rightarrow 0_+$ , it follows that  $M(\underline{q}) \rightarrow 0_+$  as  $\mathfrak{d}(\underline{q}) \rightarrow 0_+$ .

**Hypothesis B.**  $\sqrt{M} \in H_0^1(D)$ , and  $M$  is a *weight function of type 3* on  $D$  in the sense of Triebel [120], p.247, Definition 3.2.1.3c; *i.e.*, there exist positive constants  $c_1$ ,  $c_2$  and  $\lambda$ , and a positive monotonic increasing function  $\tau$  defined on the interval  $(0, \lambda)$ , such that  $c_1 \tau(\mathfrak{d}(\underline{q})) \leq M(\underline{q}) \leq c_2 \tau(\mathfrak{d}(\underline{q}))$  for all  $\underline{q} \in D$  satisfying  $\mathfrak{d}(\underline{q}) < \lambda$ .  $\diamond$

Hypotheses A and B will be assumed throughout this chapter.

**Example 2.1.1** Consider the function  $U$  defined by

$$U(s) := -f(s) \ln \left( 1 - \frac{2s}{b} \right), \quad s \in [0, \frac{b}{2}), \quad \text{with } b > 2,$$

where  $f \in C^\infty[0, \frac{b}{2}]$  is a monotonic nondecreasing function, positive on  $(0, \frac{b}{2}]$ , with  $f(\frac{b}{2}) > 1$ ; then  $U$  and the associated Maxwellian  $M$  satisfy hypotheses A and B, respectively. When  $f(s) \equiv b/2$ , the FENE potential is recovered.

The central difficulty of (2.1.1), (2.1.2), (2.1.7), from both the analytical and the computational point of view, is the presence in (2.1.1) of the degenerate Maxwellian  $M(\underline{q})$ , with  $\lim_{\mathfrak{d}(\underline{q}) \rightarrow 0_+} M(\underline{q}) = 0$ .

Most numerical methods developed for the Fokker–Planck equation have been based on the ‘original’ form of the equation,

$$\frac{\partial \psi}{\partial t} + \nabla_{\underline{q}} \cdot (\underline{\kappa} \underline{q} \psi) = \frac{1}{2\mathbb{W}_i} \nabla_{\underline{q}} \cdot (\nabla_{\underline{q}} \psi + \underline{F}(\underline{q}) \psi), \quad (2.1.4)$$

see, for example, [32,33,91] or [3,4]. From the theoretical viewpoint at least, the advantage of (2.1.1) over (2.1.4), is that on transformation into weak form the diffusion operator becomes symmetric (see (2.1.5)), which facilitates the analysis of the Fokker–Planck equation for a general class of Maxwellians. Notwithstanding this potential theoretical advantage, the computational benefits, or otherwise, of discretising (2.1.1) rather than (2.1.4) remain to be understood.

The aims of the analysis in this chapter are therefore two-fold:

- (a) The principal objective is to develop the mathematical and numerical analysis of equation (2.1.1) for the class of Maxwellians satisfying Hypotheses A and B. The discretisation of the equation is based on a spectral Galerkin method in the spatial variable  $\underline{q}$  coupled with backward Euler time-stepping. One can, of course, consider more accurate time discretisation schemes, such as an  $n$ th-order backward differentiation formula, BDF $n$ ,  $n \in \{2, \dots, 6\}$ , for example. High-order time discretisation of the problem is, however, a secondary consideration to the central theme of this chapter, and it is not discussed here.

- (b) In the special case of the FENE model, it will be shown how the results under (a) can be adapted to the case of an alternative discretisation proposed by Chauvière & Lozinski [32,33,90,91], which applies a transformation, different from the symmetrising transformation considered under (a), to the ‘original’ form (2.1.4) of the Fokker–Planck equation. The transformed equation is then approximated in the same way as in (a), using a spectral Galerkin method in space and a backward Euler discretisation in time.

Since the analytical arguments under (b) are almost identical to those under (a), for the sake of brevity, attention will be focused on (a), but the key adjustments that need to be made in order to obtain the corresponding results under (b) will be systematically indicated.

First of all, we define the function spaces relevant to the weak formulation of (2.1.1). Note that since only the configuration space is considered in this chapter,  $\|\cdot\|$  and  $(\cdot, \cdot)$  will denote the  $L^2(D)$  norm and inner-product, respectively. In subsequent chapters when numerical methods for the Fokker–Planck equation on physical space as well as configuration space are considered, the unsubscripted norm and inner-product will imply the domain  $\Omega \times D$ .

Let

$$\begin{aligned} \mathfrak{H} &:= \left\{ \varphi \in L^2_{\text{loc}}(D) : \int_D \left( \frac{\varphi}{\sqrt{M}} \right)^2 dq < \infty \right\}, \\ \mathfrak{K} &:= \left\{ \varphi \in \mathfrak{H} : \int_D \left( \left( \frac{\varphi}{\sqrt{M}} \right)^2 + \left| \sqrt{M} \nabla_q \left( \frac{\varphi}{M} \right) \right|^2 \right) dq < \infty \right\}, \end{aligned}$$

and define  $\mathfrak{K}_0$  as the closure of  $\sqrt{M}C_0^\infty(D)$  in the norm of  $\mathfrak{K}$ . Taking test functions as  $\varphi/M$  with  $\varphi \in \mathfrak{K}_0$ , we get the following weak formulation of the initial-boundary-value problem (2.1.1).

Given  $\psi^0 \in \mathfrak{H}$ , find  $\psi \in L^\infty(0, T; \mathfrak{H}) \cap L^2(0, T; \mathfrak{K}_0)$  such that

$$\begin{aligned} \frac{d}{dt} \int_D \frac{\psi \varphi}{M} dq - \int_D \frac{\psi}{\sqrt{M}} \cdot \sqrt{M} \nabla_q \left( \frac{\varphi}{M} \right) dq \\ + \frac{1}{2\text{Wi}} \int_D \sqrt{M} \nabla_q \left( \frac{\psi}{M} \right) \cdot \sqrt{M} \nabla_q \left( \frac{\varphi}{M} \right) dq = 0 \quad \forall \varphi \in \mathfrak{K}_0, \end{aligned} \quad (2.1.5)$$

in the sense of distributions on  $(0, T)$ , and  $\psi(\cdot, 0) = \psi^0(\cdot)$ .

Now, by introducing the notation

$$\hat{\varphi} := \frac{\varphi}{\sqrt{M}} \quad \text{and} \quad \nabla_M \hat{\varphi} := \sqrt{M} \nabla_q \left( \frac{\hat{\varphi}}{\sqrt{M}} \right),$$

(2.1.5) can be reformulated on observing that, by the definition of  $\mathfrak{K}$ ,  $\varphi \in \mathfrak{K}_0$  if, and only if,  $\hat{\varphi} \in H_0^1(D; M)$ , where  $H_0^1(D; M)$  is the closure of  $C_0^\infty(D)$  in the norm of  $H^1(D; M)$ , and

$$H^1(D; M) := \left\{ \zeta \in L^2(D) : \|\zeta\|_{H^1(D; M)}^2 := \int_D \left( |\zeta|^2 + |\nabla_M \zeta|^2 \right) dq < \infty \right\}.$$

When applied to an element of  $H_0^1(D; M)$  the norm  $\|\cdot\|_{H^1(D; M)}$  will be written  $\|\cdot\|_{H_0^1(D; M)}$ . As a matter of fact, it will be shown in Section 2.2 that  $C_0^\infty(D)$  is dense in  $H^1(D; M)$  and therefore, perhaps somewhat counter-intuitively,  $H_0^1(D; M) = H^1(D; M)$ , and also  $\mathfrak{K}_0 = \mathfrak{K}$ .

**Remark 2.1.2** We note in passing that the substitution  $\hat{\varphi} = \varphi/\sqrt{M}$  also appears in the recent paper by Du, Liu and Yu [42], though the operator  $\nabla_M$  does not.

In the case of the FENE Maxwellian (*cf.* Example 2.1.1), Chauvière & Lozinski [32, 33, 90, 91] used a spectral method to approximate  $\psi/M^{2s/b}$  instead of  $\psi/\sqrt{M}$ , where  $s$  is a parameter that was chosen on the basis of numerical experiments. Clearly, the two expressions coincide when  $s = b/4$ ; on the other hand, the values  $s = 2$  and  $s = 2.5$  were recommended in [32, 33, 90, 91] on computational grounds for  $d = 2$  and  $d = 3$ , respectively. More will be said in Sections 2.3, 2.5 and 2.7 about the analytical implications of using, in the special case of the FENE model, the substitution  $\hat{\psi} := \psi/M^{2s/b}$  instead of the substitution  $\hat{\psi} := \psi/\sqrt{M}$ . In particular, we shall show that both substitutions result in stable and convergent numerical methods, although in the case of the Chauvière & Lozinski type substitution it will be necessary to assume for this purpose that  $b \geq 4s^2/(2s - 1)$  with  $s > 1/2$ , while the symmetrised formulation based on (2.1.1) will be seen to result in a stable and optimally convergent scheme for all  $b > 2$ . In Section 2.8 we shall perform quantitative comparisons of the two approaches through numerical experiments.  $\diamond$

With these notational conventions, (2.1.5) has the following form.

Given  $\hat{\psi}^0 := \psi^0/\sqrt{M} \in L^2(D)$ , find  $\hat{\psi} \in L^\infty(0, T; L^2(D)) \cap L^2(0, T; H_0^1(D; M))$  such that

$$\frac{d}{dt} \int_D \hat{\psi} \hat{\varphi} \, dq - \int_D \underline{\kappa} q \hat{\psi} \cdot \nabla_M \hat{\varphi} \, dq + \frac{1}{2\text{Wi}} \int_D \nabla_M \hat{\psi} \cdot \nabla_M \hat{\varphi} \, dq = 0 \quad \forall \hat{\varphi} \in H_0^1(D; M), \quad (2.1.6)$$

in the sense of distributions on  $(0, T)$ , and  $\hat{\psi}(\cdot, 0) = \hat{\psi}^0(\cdot)$ .

No boundary condition on the function  $\hat{\psi} := \psi/\sqrt{M}$  will be *directly/explicitly* imposed along  $\partial D$  in the weak formulation. However,  $\hat{\psi}$  will be sought in the weighted Sobolev space  $H^1(D; M) = H_0^1(D; M)$ ; thereby, indirectly,  $\psi/\sqrt{M}$  will be forced to satisfy a homogeneous Dirichlet boundary condition on  $\partial D$ . This is consistent with the recent results of Zhang & Zhang [126], Liu & Liu [88] and Masmoudi [96]; see in particular Theorem 1.1 in [88] and Remark 3.6 in [96]. The implied homogeneous Dirichlet boundary condition on  $\psi/\sqrt{M}$  can be seen as an asymptotic decay condition for  $\psi$  as  $q$  approaches  $\partial D$ ; *viz.*,

$$\psi(q, t) = o\left(\sqrt{M(q)}\right), \quad \text{as } \text{dist}(q, \partial D) \rightarrow 0_+, \text{ for all } t \in (0, T], \quad (2.1.7)$$

i.e.  $\hat{\psi}(q, t) = o(1)$  as  $\text{dist}(q, \partial D) \rightarrow 0_+$ , for all  $t \in (0, T]$ . The role of the subscript 0 in the notation  $H_0^1(D; M)$  is to emphasize this *indirect/implicit* imposition of a boundary condition on elements of the function space  $H^1(D; M)$ .

The function space  $H^1(D; M)$  may appear exotic. However, it will be shown in Section 2.2 that, under Hypotheses A and B,  $H^1(D; M) = H_0^1(D; M)$  and  $H_0^1(D) \subset H_0^1(D; M)$ . The connection between  $H_0^1(D; M)$  and  $H_0^1(D)$  will prove helpful in the development of Galerkin methods for (2.1.6), since the construction of finite-dimensional subspaces of  $H_0^1(D)$  and the analysis of their approximation properties are well understood.

In Section 2.3 the weak formulation (2.1.6) of the initial-boundary-value problem will be revisited. A backward Euler semidiscretisation of the weak formulation will be constructed, and the stability of the temporal semidiscretisation in the  $\ell^\infty(0, T; L^2(D))$  and  $\ell^2(0, T; H_0^1(D; M))$  norms will be established. Also, in the case of the FENE model with  $b \geq 4s^2/(2s - 1)$  and  $s > 1/2$ , it will be demonstrated that these results can be carried across,

independent of the spatial dimension  $d$ , to a weak formulation that results from using the alternative substitution  $\hat{\psi} := \psi/M^{2s/b}$ ; the cases of  $s = 2$  and  $s = 2.5$  correspond to the methods proposed by Chauvière & Lozinski for  $d = 2$  and  $d = 3$ , respectively.

In Section 2.5 the fully-discrete method is developed and, using the stability results from Section 2.3, a bound on the global error in terms of the approximation error in a suitably defined spectral projection operator is derived.

In Section 2.6, the precise definition of the projection operator is given: its nonstandard form stems from a *decomposition lemma*, Lemma 2.6.2, for elements of the Sobolev space  $H^1(D)$  transformed to polar coordinates. For ease of presentation, we confine ourselves to the case of two space dimensions ( $d = 2$ ) in Section 2.6; analogous arguments could be developed in the  $d = 3$  case.

The convergence analysis is completed in Section 2.7 by showing that, under Hypotheses A and B, the method exhibits optimal-order convergence in the Maxwellian-weighted norm  $\|\cdot\|_{\ell^2(0,T;H_0^1(D;M))}$  with respect to the spatial and temporal discretisation parameters.

Section 2.8 is devoted to numerical experiments that illustrate the performance of the method. We focus solely on the FENE potential in this section. First of all, we discuss the implementation of our Galerkin spectral method for the case of  $d = 2$  in Section 2.8.1, and we also present a range of computational results in order to illustrate the behaviour of the method in practice, as well as to provide experimental verification of the convergence analysis from Section 2.7. In Section 2.8.2, we compare the behaviour of the numerical method based on the backward Euler temporal discretisation with a semi-implicit scheme in which the transport term in (2.1.6) is treated explicitly in time. The semi-implicit scheme is used in Chapter 3, and the results of Section 2.8.2 have important implications there. Finally, we consider the implementation of the spectral method in three spatial dimensions in Section 2.8.3 and we demonstrate that, as expected, the behaviour of the Galerkin spectral method is essentially the same as in the case of  $d = 2$ .

## 2.2 Properties of Maxwellian-weighted spaces

In this section, density results are derived for the Maxwellian-weighted function spaces that were defined above. Since the density results below are not specific to the FENE model, they will be stated more generally, for any potential  $U$  and associated Maxwellian  $M$  that satisfy Hypotheses A and B, respectively.

(a) Suppose that the Maxwellian  $M$  satisfies Hypothesis B;  $M$  is then a weight-function of Type 3 in the sense of Triebel. According to [120], Theorem 3.2.2a, the weighted Sobolev space  $H_M^1(D) = \{v \in L_M^2(D) : \nabla_q v \in L_M^2(D) := [L_M^2(D)]^d\}$  is a Hilbert space with respect to the norm  $\|\cdot\|_{H_M^1(D)}$  defined by

$$\|v\|_{H_M^1(D)} := \left( \|v\|_{L_M^2(D)}^2 + \|\nabla_q v\|_{L_M^2(D)}^2 \right)^{\frac{1}{2}},$$

and  $L_M^2(D) = (1/\sqrt{M})L^2(D)$  is a Hilbert space with norm  $\|\cdot\|_{L_M^2(D)}$  defined by  $\|v\|_{L_M^2(D)} := \|\sqrt{M}v\|$ , where  $\|\cdot\|$  denotes the  $L^2(D)$  norm induced by the  $L^2(D)$  inner product  $(\cdot, \cdot)$ . By [120], Theorem 3.2.2c,  $C^\infty(\bar{D})$  is dense in both  $H_M^1(D)$  and  $L_M^2(D)$ ; see also Ch. I, Sec. 7, in Kufner [77], or one of [17, 18]. Thus, since  $v \in H_M^1(D)$  if, and only if,  $\sqrt{M}v \in H^1(D; M)$ , it follows that  $\sqrt{M}C^\infty(\bar{D})$  is dense in the Hilbert spaces  $H^1(D; M)$  and  $L^2(D)$ , whereby  $H^1(D; M)$  is dense in  $L^2(D)$ .

(b) Now suppose that  $U$  satisfies Hypothesis A and the associated Maxwellian  $M$  satisfies Hypothesis B. It follows from Hardy’s inequality (see, for example, [7, 95]) that

$$\int_D \left(1 - \frac{|q|^2}{b}\right)^{-2} |\hat{\psi}(q)|^2 dq \leq 4b \|\nabla_q \hat{\psi}\|^2 \quad \forall \hat{\psi} \in \mathbf{H}_0^1(D). \quad (2.2.1)$$

Since  $\nabla_M \hat{\psi} = \nabla_q \hat{\psi} + \frac{1}{2} q U' \left(\frac{1}{2} |q|^2\right) \hat{\psi}$ , Hypothesis A implies that there exists  $C_1 \in \mathbb{R}_{>0}$  (for the FENE model  $C_1 = 1$ ) such that  $(1 - |q|^2/b)^2 |U'(\frac{1}{2} |q|^2)|^2 \leq C_1^2$  for all  $q \in D$ , whereby

$$\|\nabla_M \hat{\psi}\| \leq (1 + C_1 b) \|\nabla_q \hat{\psi}\| \quad \forall \hat{\psi} \in \mathbf{H}_0^1(D). \quad (2.2.2)$$

Now, (2.2.2) implies that  $\mathbf{H}_0^1(D) \subset \mathbf{H}^1(D; M)$ .

Finally, we show that  $\mathbf{H}^1(D; M) = \mathbf{H}_0^1(D; M)$ . As  $\sqrt{M} C^\infty(\bar{D}) \subset \mathbf{H}_0^1(D) \subset \mathbf{H}^1(D; M)$  and  $\sqrt{M} C^\infty(\bar{D})$  is dense in  $\mathbf{H}^1(D; M)$  (cf. (a) above), we deduce that  $\mathbf{H}_0^1(D)$  is dense in  $\mathbf{H}^1(D; M)$ . Since  $C_0^\infty(D)$  is dense in  $\mathbf{H}_0^1(D)$ , it follows from (2.2.2) that  $C_0^\infty(D)$  is also dense in  $\mathbf{H}^1(D; M)$ . By definition,  $\mathbf{H}_0^1(D; M)$  is the closure of  $C_0^\infty(D)$  in  $\mathbf{H}^1(D; M)$ ; thus  $\mathbf{H}^1(D; M) = \mathbf{H}_0^1(D; M)$ , and therefore also  $\mathfrak{K} = \mathfrak{K}_0$ . As  $\mathbf{H}^1(D; M)$  is continuously and densely embedded into  $L^2(D)$ , it follows that  $\mathbf{H}_0^1(D; M)$  is continuously and densely embedded into  $L^2(D)$ .

**Remark 2.2.1** A third hypothesis (referred to as Hypothesis C) was introduced in [71], which enabled the inequalities:

$$\inf_{c \in \text{Ker}(\nabla_M)} \int_D |\hat{\psi} - c|^2 dq \leq \int_D |\nabla_M \hat{\psi}|^2 dq, \quad (2.2.3)$$

and

$$\inf_{c \in \text{Ker}(\nabla_M)} \int_D \frac{|\hat{\psi} - c|^2}{1 - \frac{|q|^2}{b}} dq \leq \frac{b}{b-2} \int_D |\nabla_M \hat{\psi}|^2 dq \quad (2.2.4)$$

to be established for all  $\hat{\psi} \in \mathbf{H}^1(D; M)$ . These can be seen as counterparts of Poincaré’s inequality in the Maxwellian-weighted Sobolev space  $\mathbf{H}^1(D; M) = \mathbf{H}_0^1(D; M)$ .  $\diamond$

## 2.3 Analysis of the backward Euler semidiscretisation

As noted in the opening of this chapter, by setting  $\hat{\psi}(\cdot, t) := \psi(\cdot, t)/\sqrt{M}$  for  $t \in [0, T]$  and  $\hat{\varphi} := \varphi/\sqrt{M}$  in (2.1.5) and writing  $\hat{\psi}^0 := \psi^0/\sqrt{M}$ , the following weak formulation of the initial-boundary-value problem (2.1.1), (2.1.2), (2.1.7) is obtained:

Given  $\hat{\psi}^0 \in L^2(D)$ , find  $\hat{\psi} \in L^\infty(0, T; L^2(D)) \cap L^2(0, T; \mathbf{H}_0^1(D; M))$  such that (2.1.6) holds in the sense of distributions on  $(0, T)$ , and  $\hat{\psi}(\cdot, 0) = \hat{\psi}^0(\cdot)$ .

The function  $\psi$ , representing a weak solution to the problem (2.1.5), is then recovered from  $\hat{\psi}$  through the substitution  $\psi := \sqrt{M} \hat{\psi}$ . Thus, instead of constructing a Galerkin approximation to  $\psi$ , the aim is to construct a Galerkin approximation to  $\hat{\psi}$  from a finite-dimensional subspace of  $\mathbf{H}_0^1(D; M)$ , from which an approximation to  $\hat{\psi}$  can be obtained straightforwardly.

Let  $N_T \geq 1$  be an integer,  $\Delta t = T/N_T$ , and  $t^n = n\Delta t$ , for  $n = 0, 1, \dots, N_T$ . Discretising (2.1.6) in time using the backward Euler method yields the following semi-discrete numerical scheme.

Given  $\hat{\psi}^0 := \psi^0 / \sqrt{M} \in L^2(D)$ , find  $\hat{\psi}^{n+1} \in H_0^1(D; M)$ ,  $n = 0, \dots, N_T - 1$ , such that

$$\int_D \frac{\hat{\psi}^{n+1} - \hat{\psi}^n}{\Delta t} \hat{\varphi} \, d\mathbf{q} - \int_D (\underline{\kappa}^{n+1} \mathbf{q} \hat{\psi}^{n+1}) \cdot \nabla_M \hat{\varphi} \, d\mathbf{q} + \frac{1}{2\text{Wi}} \int_D \nabla_M \hat{\psi}^{n+1} \cdot \nabla_M \hat{\varphi} \, d\mathbf{q} = 0, \quad (2.3.1)$$

for all  $\hat{\varphi} \in H_0^1(D; M)$ .

Let us first show that for any  $\Delta t$ , sufficiently small, problem (2.3.1) has a unique solution. To this end, we consider the bilinear form  $\mathfrak{B}(\cdot, \cdot)$  defined on  $H_0^1(D; M) \times H_0^1(D; M)$  by

$$\mathfrak{B}(\hat{\psi}, \hat{\varphi}) := \frac{1}{\Delta t} \int_D \hat{\psi} \hat{\varphi} \, d\mathbf{q} - \int_D (\underline{\kappa}^{n+1} \mathbf{q} \hat{\psi}) \cdot \nabla_M \hat{\varphi} \, d\mathbf{q} + \frac{1}{2\text{Wi}} \int_D \nabla_M \hat{\psi} \cdot \nabla_M \hat{\varphi} \, d\mathbf{q},$$

and, for  $\hat{\psi}^n \in L^2(D)$  fixed, we define the linear functional  $\ell(\hat{\psi}^n; \cdot)$  on  $H_0^1(D; M)$  by

$$\ell(\hat{\psi}^n; \hat{\varphi}) := \frac{1}{\Delta t} \int_D \hat{\psi}^n \hat{\varphi} \, d\mathbf{q}.$$

Clearly,

$$\mathfrak{B}(\hat{\psi}, \hat{\psi}) \geq \frac{1}{\Delta t} \left(1 - \Delta t \text{Wi} b \|\underline{\kappa}\|_{L^\infty(0,T)}^2\right) \int_D |\hat{\psi}|^2 \, d\mathbf{q} + \frac{1}{4\text{Wi}} \int_D |\nabla_M \hat{\psi}|^2 \, d\mathbf{q},$$

and hence, on assuming that  $\Delta t \text{Wi} b \|\underline{\kappa}\|_{L^\infty(0,T)}^2 < 1$  and letting

$$c_{\Delta t} := \frac{1}{\Delta t} \left(1 - \Delta t \text{Wi} b \|\underline{\kappa}\|_{L^\infty(0,T)}^2\right),$$

we deduce that

$$\mathfrak{B}(\hat{\psi}, \hat{\psi}) \geq \min\left(c_{\Delta t}, \frac{1}{4\text{Wi}}\right) \|\hat{\psi}\|_{H_0^1(D; M)}^2. \quad (2.3.2)$$

Also, by a simple application of the Cauchy–Schwarz inequality,  $\mathfrak{B}(\cdot, \cdot)$  is a bounded bilinear functional on  $H_0^1(D; M) \times H_0^1(D; M)$  and, for any  $\hat{\psi}^n \in L^2(D)$ ,  $\ell(\hat{\psi}^n; \cdot)$  is a bounded linear functional on  $H_0^1(D; M)$ . Since  $H_0^1(D; M)$  is a Hilbert space with norm  $\|\cdot\|_{H_0^1(D; M)}$ , the Lax–Milgram theorem implies the existence of a unique solution  $\hat{\psi}^{n+1} \in H_0^1(D; M)$  such that

$$\mathfrak{B}(\hat{\psi}^{n+1}, \hat{\varphi}) = \ell(\hat{\psi}^n; \hat{\varphi}) \quad \forall \hat{\varphi} \in H_0^1(D; M), \quad n = 0, 1, \dots, N_T - 1. \quad (2.3.3)$$

As  $\hat{\psi}^0 \in L^2(D)$ , we have thus shown that, for any  $\Delta t = T/N_T$  such that  $\Delta t \text{Wi} b \|\underline{\kappa}\|_{L^\infty(0,T)}^2 < 1$ , the problem (2.3.1) has a unique solution  $\{\hat{\psi}^n \in H_0^1(D; M) : n = 1, \dots, N_T\}$ .

For the purposes of the convergence analysis that will be carried out below, we consider an extended version of the scheme (2.3.1) with a nonzero right-hand side:

$$\begin{aligned} & \int_D \frac{\hat{\psi}^{n+1} - \hat{\psi}^n}{\Delta t} \hat{\varphi} \, d\mathbf{q} - \int_D (\underline{\kappa}^{n+1} \mathbf{q} \hat{\psi}^{n+1}) \cdot \nabla_M \hat{\varphi} \, d\mathbf{q} + \frac{1}{2\text{Wi}} \int_D \nabla_M \hat{\psi}^{n+1} \cdot \nabla_M \hat{\varphi} \, d\mathbf{q} \\ & = \int_D \mu^{n+1} \hat{\varphi} \, d\mathbf{q} + \int_D \nu^{n+1} \cdot \nabla_M \hat{\varphi} \, d\mathbf{q} \quad \forall \hat{\varphi} \in H_0^1(D; M), \end{aligned} \quad (2.3.4)$$

for  $n = 0, \dots, N_T - 1$ , where  $\mu^{n+1} \in L^2(D)$  and  $\nu^{n+1} \in \mathbb{L}^2(D)$  for all  $n \geq 0$ . We have the following stability result for (2.3.4).

**Lemma 2.3.1 (The first stability inequality)** *Let  $\Delta t = T/N_T$ ,  $N_T \geq 1$ ,  $\kappa \in \mathcal{C}[0, T]$ ,  $\hat{\psi}^0 \in L^2(D)$ , and define  $c_0 := 1 + 4\text{Wi}b\|\kappa\|_{L^\infty(0, T)}^2$ . Let, further,  $\Delta t$  be such that  $0 < c_0\Delta t \leq 1/2$ ; then, we have, for all  $m$  such that  $1 \leq m \leq N_T$ , that*

$$\begin{aligned} \|\hat{\psi}^m\|^2 + \sum_{n=0}^{m-1} \Delta t \left\| \frac{\hat{\psi}^{n+1} - \hat{\psi}^n}{\sqrt{\Delta t}} \right\|^2 + \sum_{n=0}^{m-1} \frac{\Delta t}{2\text{Wi}} \|\nabla_M \hat{\psi}^{n+1}\|^2 \\ \leq e^{2c_0 m \Delta t} \left\{ \|\hat{\psi}^0\|^2 + \sum_{n=0}^{m-1} 2\Delta t (\|\mu^{n+1}\|^2 + 4\text{Wi}\|\mathcal{L}^{n+1}\|^2) \right\}. \end{aligned}$$

**Proof.** Let  $0 \leq n \leq N_T - 1$ . Setting  $\hat{\varphi} = \hat{\psi}^{n+1}$ , we write the first term in (2.3.4) as

$$\int_D \frac{\hat{\psi}^{n+1} - \hat{\psi}^n}{\Delta t} \hat{\psi}^{n+1} \, dq = \frac{1}{2\Delta t} (\|\hat{\psi}^{n+1}\|^2 - \|\hat{\psi}^n\|^2) + \frac{1}{2\Delta t} \|\hat{\psi}^{n+1} - \hat{\psi}^n\|^2$$

using the identity  $(\alpha - \beta)\alpha = \frac{1}{2}(\alpha^2 - \beta^2) + \frac{1}{2}(\alpha - \beta)^2$ .

Applying the Cauchy–Schwarz inequality to the transport term in (2.3.4), we have

$$\int_D (\kappa^{n+1} q \hat{\psi}^{n+1}) \cdot \nabla_M \hat{\psi}^{n+1} \, dq \leq \sqrt{b} |\kappa^{n+1}| \|\hat{\psi}^{n+1}\| \|\nabla_M \hat{\psi}^{n+1}\|.$$

Combining these results and applying the Cauchy–Schwarz inequality to the right-hand side terms in (2.3.4) gives

$$\begin{aligned} \|\hat{\psi}^{n+1}\|^2 + \|\hat{\psi}^{n+1} - \hat{\psi}^n\|^2 + \frac{\Delta t}{\text{Wi}} \|\nabla_M \hat{\psi}^{n+1}\|^2 \\ \leq \|\hat{\psi}^n\|^2 + 2\Delta t \sqrt{b} |\kappa^{n+1}| \|\hat{\psi}^{n+1}\| \|\nabla_M \hat{\psi}^{n+1}\| \\ + 2\Delta t \|\mu^{n+1}\| \|\hat{\psi}^{n+1}\| + 2\Delta t \|\mathcal{L}^{n+1}\| \|\nabla_M \hat{\psi}^{n+1}\| \\ =: \|\hat{\psi}^n\|^2 + \text{T}_1 + \text{T}_2 + \text{T}_3. \end{aligned}$$

Using Cauchy’s inequality  $2\alpha\beta \leq \varepsilon\alpha^2 + \varepsilon^{-1}\beta^2$  with  $\varepsilon > 0$  on each of  $\text{T}_1$  and  $\text{T}_3$ , we deduce that

$$\text{T}_1 \leq \varepsilon \|\nabla_M \hat{\psi}^{n+1}\|^2 + \frac{1}{\varepsilon} \Delta t^2 b |\kappa^{n+1}|^2 \|\hat{\psi}^{n+1}\|^2, \quad \text{T}_3 \leq \varepsilon \|\nabla_M \hat{\psi}^{n+1}\|^2 + \frac{1}{\varepsilon} \Delta t^2 \|\mathcal{L}^{n+1}\|^2.$$

Choosing  $\varepsilon = \Delta t/(4\text{Wi})$  then gives

$$\begin{aligned} \|\hat{\psi}^{n+1}\|^2 + \|\hat{\psi}^{n+1} - \hat{\psi}^n\|^2 + \frac{\Delta t}{2\text{Wi}} \|\nabla_M \hat{\psi}^{n+1}\|^2 \\ \leq \|\hat{\psi}^n\|^2 + 4\Delta t \text{Wi} b |\kappa^{n+1}|^2 \|\hat{\psi}^{n+1}\|^2 + 4\Delta t \text{Wi} \|\mathcal{L}^{n+1}\|^2 + \text{T}_2. \end{aligned}$$

Similarly, we have  $\text{T}_2 \leq \Delta t \|\hat{\psi}^{n+1}\|^2 + \Delta t \|\mu^{n+1}\|^2$ , and therefore, on defining  $c_0 := 1 + 4\text{Wi}b\|\kappa\|_{L^\infty(0, T)}^2$ , we get

$$\begin{aligned} (1 - c_0\Delta t) \|\hat{\psi}^{n+1}\|^2 + \|\hat{\psi}^{n+1} - \hat{\psi}^n\|^2 + \frac{\Delta t}{2\text{Wi}} \|\nabla_M \hat{\psi}^{n+1}\|^2 \\ \leq \|\hat{\psi}^n\|^2 + \Delta t \|\mu^{n+1}\|^2 + 4\Delta t \text{Wi} \|\mathcal{L}^{n+1}\|^2. \end{aligned}$$

As  $c_0\Delta t \leq \frac{1}{2}$ , dividing through by  $(1 - c_0\Delta t)$  and using that  $1 \leq \frac{1}{1 - c_0\Delta t} \leq 1 + 2c_0\Delta t \leq 2$ , we have

$$\begin{aligned} & \|\hat{\psi}^{n+1}\|^2 + \|\hat{\psi}^{n+1} - \hat{\psi}^n\|^2 + \frac{\Delta t}{2\mathbb{W}i} \|\nabla_M \hat{\psi}^{n+1}\|^2 \\ & \leq \frac{1}{1 - c_0\Delta t} \left( \|\hat{\psi}^n\|^2 + \Delta t \|\mu^{n+1}\|^2 + 4\Delta t \mathbb{W}i \|\zeta^{n+1}\|^2 \right) \\ & \leq (1 + 2c_0\Delta t) \|\hat{\psi}^n\|^2 + 2\Delta t (\|\mu^{n+1}\|^2 + 4\mathbb{W}i \|\zeta^{n+1}\|^2). \end{aligned} \quad (2.3.5)$$

Summing over  $n = 0, \dots, m-1$  in (2.3.5) we obtain

$$\begin{aligned} & \|\hat{\psi}^m\|^2 + \sum_{n=0}^{m-1} \Delta t \left\| \frac{\hat{\psi}^{n+1} - \hat{\psi}^n}{\sqrt{\Delta t}} \right\|^2 + \sum_{n=0}^{m-1} \frac{\Delta t}{2\mathbb{W}i} \|\nabla_M \hat{\psi}^{n+1}\|^2 \\ & \leq \left\{ \|\hat{\psi}^0\|^2 + \sum_{n=0}^{m-1} 2\Delta t (\|\mu^{n+1}\|^2 + 4\mathbb{W}i \|\zeta^{n+1}\|^2) \right\} + 2c_0 \sum_{n=0}^{m-1} \Delta t \|\hat{\psi}^n\|^2, \end{aligned} \quad (2.3.6)$$

for all  $m \in \{1, \dots, N_T\}$ . By induction (or by a discrete Gronwall lemma) we deduce that

$$\begin{aligned} & \|\hat{\psi}^m\|^2 + \sum_{n=0}^{m-1} \Delta t \left\| \frac{\hat{\psi}^{n+1} - \hat{\psi}^n}{\sqrt{\Delta t}} \right\|^2 + \sum_{n=0}^{m-1} \frac{\Delta t}{2\mathbb{W}i} \|\nabla_M \hat{\psi}^{n+1}\|^2 \\ & \leq e^{2c_0 m \Delta t} \left\{ \|\hat{\psi}^0\|^2 + \sum_{n=0}^{m-1} 2\Delta t (\|\mu^{n+1}\|^2 + 4\mathbb{W}i \|\zeta^{n+1}\|^2) \right\}, \quad 1 \leq m \leq N_T, \end{aligned}$$

and that completes the proof.  $\square$

**Theorem 2.3.2** *Suppose that  $\hat{\psi}^0 \in L^2(D)$  and that  $\underline{\kappa} \in \underline{\mathbb{C}}[0, T]$ . Then, there exists a function  $\hat{\psi}$  in  $L^\infty(0, T; L^2(D)) \cap L^2(0, T; \mathbb{H}_0^1(D; M))$  such that*

$$\begin{aligned} & -(\hat{\psi}^0, \hat{\varphi}(\cdot, 0)) - \int_0^T \int_D \hat{\psi} \frac{\partial \hat{\varphi}}{\partial t} \, d\underline{q} \, dt - \int_0^T \int_D (\underline{\kappa} \underline{q} \hat{\psi}) \cdot \nabla_M \hat{\varphi} \, d\underline{q} \, dt \\ & + \frac{1}{2\mathbb{W}i} \int_0^T \int_D \nabla_M \hat{\psi} \cdot \nabla_M \hat{\varphi} \, d\underline{q} \, dt = 0, \quad \forall \hat{\varphi} \in \mathbb{H}^1(0, T; \mathbb{H}_0^1(D; M)), \quad \hat{\varphi}(\cdot, T) = 0. \end{aligned} \quad (2.3.7)$$

Further,  $\hat{\psi} \in \mathbb{H}^1(0, T; \mathbb{H}_0^1(D; M)') \cap C([0, T]; L^2(D))$  and  $(\hat{\psi}(\cdot, 0) - \hat{\psi}^0, \varphi) = 0$  for all  $\varphi \in L^2(D)$ ; moreover,  $\hat{\psi}$  is the unique such function. The function  $\psi = \sqrt{M} \hat{\psi}$  will be called the weak solution of the initial-boundary-value problem (2.1.1), (2.1.2), (2.1.7).

**Proof.** This theorem is proved in Section 3 of [71]; the interested reader is referred to that paper for details. The argument makes use of the stability result in Lemma 2.3.1 in order to use compactness results for the bounded sequence of solutions to (2.3.1) as  $\Delta t \rightarrow 0_+$ .  $\square$

In the next lemma, a configuration space analogue of Lemma 1.3.3 is established and also it is shown that a weak form of (1.3.22) is preserved on  $D$ . In the remark below, a result is stated that is necessary for the proof of Lemma 2.3.4.

**Remark 2.3.3** Suppose  $\hat{\varphi} \in \mathbf{H}_0^1(D; M)$  and  $L \geq 0$ , and let  $[\hat{\psi}^n]_-$  be the pointwise negative part of  $\hat{\psi}^n$ , i.e.  $[x]_{\pm} := (x \pm |x|)/2$  for  $x \in \mathbb{R}$ . Then, it is shown in Lemma 3.5 of [71] that

$$\nabla_M[\hat{\varphi} - L\sqrt{M}]_+ = \begin{cases} \nabla_M(\hat{\varphi} - L\sqrt{M}) = \nabla_M\hat{\varphi} & \text{if } \hat{\varphi} > L\sqrt{M}, \\ 0 & \text{if } \hat{\varphi} \leq L\sqrt{M}; \end{cases} \quad (2.3.8)$$

and

$$\nabla_M[\hat{\varphi} - L\sqrt{M}]_- = \begin{cases} \nabla_M(\hat{\varphi} - L\sqrt{M}) = \nabla_M\hat{\varphi} & \text{if } \hat{\varphi} < L\sqrt{M}, \\ 0 & \text{if } \hat{\varphi} \geq L\sqrt{M}; \end{cases} \quad (2.3.9)$$

i.e. that the  $[\cdot]_{\pm}$  operators act on functions in  $\mathbf{H}_0^1(D; M)$  as one would expect. Moreover,  $[\hat{\varphi} - L\sqrt{M}]_+$  and  $[\hat{\varphi} - L\sqrt{M}]_-$  belong to  $\mathbf{H}_0^1(D; M)$ . The proofs of these results is omitted here, for the sake of brevity; we refer to [71] for details.  $\diamond$

**Lemma 2.3.4** Let  $\psi^0 \in \mathfrak{H}$  and  $\psi = \sqrt{M}\hat{\psi}$  where  $\hat{\psi} \in C([0, T]; L^2(D)) \cap L^2(0, T; \mathbf{H}_0^1(D; M)) \cap \mathbf{H}^1(0, T; \mathbf{H}_0^1(D; M)')$  is the weak solution to (2.3.7) subject to the initial condition  $\hat{\psi}^0 = \psi^0/\sqrt{M}$  (i.e., the function  $\psi$  is the weak solution of the initial-boundary-value problem (2.1.1), (2.1.2), (2.1.7)). Then,

$$\int_D \psi(\underline{q}, t) \, d\underline{q} = \int_D \psi^0(\underline{q}) \, d\underline{q} \quad \forall t \in [0, T].$$

Furthermore if  $\psi^0 \geq 0$  a.e. on  $D$ , then  $\psi(\cdot, t) \geq 0$  a.e. on  $D$  for all  $t \in [0, T]$ .

**Proof.** Fix any  $t \in (0, T)$ , and let  $\varepsilon \in (0, T - t]$ . Consider the function  $\hat{\varphi}_\varepsilon$  defined by

$$\hat{\varphi}_\varepsilon(\underline{q}, s) := \begin{cases} \sqrt{M} & \text{for } s \in [0, t], \\ \sqrt{M}(t + \varepsilon - s)/\varepsilon & \text{for } s \in [t, t + \varepsilon), \\ 0 & \text{for } s \in [t + \varepsilon, T]. \end{cases}$$

Clearly,  $\hat{\varphi}_\varepsilon \in \mathbf{H}^1(0, T; \mathbf{H}_0^1(D; M))$  and  $\hat{\varphi}_\varepsilon(\cdot, T) = 0$ . Taking  $\hat{\varphi}_\varepsilon$  as test function in (2.3.7) yields

$$-(\hat{\psi}^0, \sqrt{M}) + \frac{1}{\varepsilon} \int_t^{t+\varepsilon} (\hat{\psi}(\cdot, s), \sqrt{M}) \, ds = 0.$$

Passing to the limit  $\varepsilon \rightarrow 0_+$  yields  $-(\hat{\psi}^0, \sqrt{M}) + (\hat{\psi}(\cdot, t), \sqrt{M}) = 0$ , whereby  $(\psi(\cdot, t), 1) = (\psi^0, 1)$ , as required, for all  $t \in (0, T)$ ; for  $t = 0$  the equality holds trivially.

Now, suppose that  $\psi^0 \in \mathfrak{H}$  and  $\psi^0 \geq 0$ ; then,  $\hat{\psi}^0 \in L^2(D)$  and  $\hat{\psi}^0 \geq 0$ . For  $\Delta t$  as in Lemma 2.3.1, consider the sequence of functions  $(\hat{\psi}^n)_{n=0}^{N_T} \subset \mathbf{H}_0^1(D; M)$  defined by (2.3.3). Let  $[\hat{\psi}^n]_-$  be the pointwise negative part of  $\hat{\psi}^n$ , where  $[x]_{\pm} := (x \pm |x|)/2$  for  $x \in \mathbb{R}$ . Then, by Remark 2.3.3,  $([\hat{\psi}^n]_-)_{n=0}^{N_T} \subset \mathbf{H}_0^1(D; M)$ . It follows that

$$\mathfrak{B}([\hat{\psi}^{n+1}]_-, [\hat{\psi}^{n+1}]_-) = \mathfrak{B}(\hat{\psi}^{n+1}, [\hat{\psi}^{n+1}]_-) = \ell(\hat{\psi}^n; [\hat{\psi}^{n+1}]_-),$$

where the first equality is due to the fact that  $[\hat{\psi}^{n+1}]_-$  vanishes when  $\hat{\psi}^{n+1} > 0$ , and the second equality is due to (2.3.3). Suppose, for induction, that  $\hat{\psi}^n \geq 0$ ; this is certainly true for  $n = 0$ , since  $\hat{\psi}^0 \geq 0$ . Hence,

$$\ell(\hat{\psi}^n; [\hat{\psi}^{n+1}]_-) = \frac{1}{\Delta t} \int_D \hat{\psi}^n(\underline{q}) [\hat{\psi}^{n+1}(\underline{q})]_- \, d\underline{q} \leq 0.$$

Therefore,  $\mathfrak{B}([\hat{\psi}^{n+1}]_-, [\hat{\psi}^{n+1}]_-) \leq 0$ ; thus, (2.3.2) implies that  $\|[\hat{\psi}^{n+1}]_-\|_{\mathbb{H}_0^1(D;M)} \leq 0$ , whereby  $[\hat{\psi}^{n+1}]_- = 0$  and hence  $\hat{\psi}^{n+1} \geq 0$ . By induction,  $\hat{\psi}^n \geq 0$  for all  $n = 0, 1, \dots, N_T$ . Then, passing to the limit  $\Delta t \rightarrow 0_+$ , it follows from Theorem 2.3.2 that the weak solution  $\hat{\psi}$  is nonnegative on  $D \times [0, T]$  (see [71]).  $\square$

**Remark 2.3.5** The same argument used above to establish the nonnegativity of  $\hat{\psi}$  can be used to derive a weak maximum principle in the case that  $\underline{q}^T \underline{\kappa}(t) \underline{q} \leq 0$  for a.e.  $t \in [0, T]$  and  $\underline{q} \in D$ .

Let

$$L = \text{ess.sup}_{\underline{q} \in D} \hat{\psi}^0(\underline{q}) / \sqrt{M(\underline{q})},$$

where it is assumed that the essential supremum above is finite. Suppose that  $\hat{\psi}^n \leq L\sqrt{M}$ ; this is certainly true for  $n = 0$ . Then, following the argument above:

$$\begin{aligned} \mathfrak{B}([\hat{\psi}^{n+1} - L\sqrt{M}]_+, [\hat{\psi}^{n+1} - L\sqrt{M}]_+) &= \mathfrak{B}(\hat{\psi}^{n+1} - L\sqrt{M}, [\hat{\psi}^{n+1} - L\sqrt{M}]_+) \\ &= \mathfrak{B}(\hat{\psi}^{n+1}, [\hat{\psi}^{n+1} - L\sqrt{M}]_+) - L\mathfrak{B}(\sqrt{M}, [\hat{\psi}^{n+1} - L\sqrt{M}]_+) \\ &= \ell(\hat{\psi}^n; [\hat{\psi}^{n+1} - L\sqrt{M}]_+) - L\mathfrak{B}(\sqrt{M}, [\hat{\psi}^{n+1} - L\sqrt{M}]_+) \\ &= \frac{1}{\Delta t} \int_D (\hat{\psi}^n(\underline{q}) - L\sqrt{M}) [\hat{\psi}^{n+1} - L\sqrt{M}]_+ \, d\underline{q} \\ &\quad + L \text{Wi} \int_D (\underline{\kappa} \underline{q} \sqrt{M}) \cdot \nabla_M [\hat{\psi}^{n+1} - L\sqrt{M}]_+ \, d\underline{q}, \end{aligned}$$

where the diffusion term in  $\mathfrak{B}(\cdot, \cdot)$  vanishes because  $\sqrt{M} \in \text{Ker}(\nabla_M)$ . The term on the second-last line above is nonpositive by the inductive hypothesis and, after integrating by parts, we deduce that the term on the last line is also nonpositive when  $\underline{q}^T \underline{\kappa} \underline{q} \leq 0$ .<sup>1</sup> Therefore,  $[\hat{\psi}^{n+1} - L\sqrt{M}]_+ = 0$ ; i.e.,  $\hat{\psi}^{n+1} \leq L\sqrt{M}$ . Then, in the same way as in Lemma 2.3.4, on passage to the limit  $\Delta t \rightarrow 0_+$ , this implies that

$$\text{ess.sup}_{(\underline{q}, t) \in D \times [0, T]} \psi(\underline{q}, t) / M(\underline{q}) \leq \text{ess.sup}_{\underline{q} \in D} \psi^0(\underline{q}) / M(\underline{q}),$$

which can be thought of as a maximum principle for the initial-boundary value problem in the case that  $\underline{q}^T \underline{\kappa} \underline{q} \leq 0$ .  $\diamond$

By the next lemma, if  $\underline{\kappa} \in \mathbb{H}^1(0, T)$  and  $\hat{\psi}^0 \in \mathbb{H}_0^1(D; M)$ , then stability can be established in stronger norms than in Lemma 2.3.1.

**Lemma 2.3.6 (The second stability inequality)** *Let  $\Delta t = T/N_T$ ,  $N_T \geq 1$ , suppose that  $\underline{\kappa} \in \mathbb{H}^1(0, T)$ ,  $\hat{\psi}^0 \in \mathbb{H}_0^1(D; M)$ , and define  $c_0 := 1 + 4\text{Wi} \, b \|\underline{\kappa}\|_{L^\infty(0, T)}^2$ . Let us assume, further,*

<sup>1</sup>In fact, if  $\underline{q}^T \underline{\kappa}(t) \underline{q} \leq 0$  for all  $\underline{q} \in \mathbb{R}^d$  and  $t \in [0, T]$ , and  $\text{tr}(\underline{\kappa}(t)) = 0$  for all  $t \in [0, T]$ , then it must be the case that  $\underline{q}^T \underline{\kappa}(t) \underline{q} = 0$  for all  $\underline{q} \in \mathbb{R}^d$  and  $t \in [0, T]$ .

that  $\Delta t$  is such that  $0 < c_0 \Delta t \leq 1/2$ ; then, for all  $m$  such that  $1 \leq m \leq N_T$ ,

$$\begin{aligned} & \Delta t \sum_{n=0}^{m-1} \left\| \frac{\hat{\psi}^{n+1} - \hat{\psi}^n}{\Delta t} \right\|^2 + \frac{1}{4\text{Wi}} \|\nabla_M \hat{\psi}^m\|^2 + \frac{1}{2\text{Wi}} \sum_{n=0}^{m-1} \Delta t \left\| \nabla_M \frac{\hat{\psi}^{n+1} - \hat{\psi}^n}{\sqrt{\Delta t}} \right\|^2 \\ & \leq e^{2c_1 m \Delta t} \left\{ 2\Delta t \sum_{n=0}^{m-1} \|\mu^{n+1}\|^2 + 12\text{Wi} \max_{1 \leq n \leq m} \|\zeta^n\|^2 + \Delta t \sum_{n=1}^{m-1} \left\| \frac{\zeta^{n+1} - \zeta^n}{\Delta t} \right\|^2 \right. \\ & \quad \left. + \frac{1}{\text{Wi}} \|\nabla_M \hat{\psi}^0\|^2 + \left( b \|\kappa_t\|_{L^2(0,T)}^2 + 12\text{Wi} b \|\kappa\|_{L^\infty(0,T)}^2 \right) \mathfrak{S}(\hat{\psi}^0, \mu, \zeta, \text{Wi}, m\Delta t) \right\}, \end{aligned}$$

where  $\mathfrak{S}(\hat{\psi}^0, \mu, \zeta, \text{Wi}, m\Delta t)$  is the right-hand side of the inequality from Lemma 2.3.1 and  $c_1 = 4\text{Wi}(1 + b \|\kappa\|_{L^\infty(0,T)}^2)$ .

**Proof.** The proof is similar to that of Lemma 2.3.1, except one uses the test function  $\hat{\varphi} = (\hat{\psi}^{n+1} - \hat{\psi}^n)/\Delta t$ .  $\square$

It follows from Lemma 2.3.6, by an identical argument as in the proof of Theorem 2.3.2, that the weak solution  $\hat{\psi}$  of (2.3.7) belongs to  $H^1(0, T; L^2(D)) \cap L^\infty(0, T; H_0^1(D; M))$ , provided that  $\kappa \in \mathbb{H}^1(0, T)$  and  $\hat{\psi}^0 \in H_0^1(D; M)$ .

The stability result in Lemma 2.3.1 will be useful in Section 2.5, but for now, note that setting  $\mu = 0$  and  $\zeta = \mathcal{Q}$  in Lemmas 2.3.1 and 2.3.6 demonstrates the stability of the time semidiscretisation in various norms. Also note that, evidently, any fully-discrete method based on the semidiscrete scheme (2.3.1) and conforming Galerkin discretisation in  $q$  using a finite-dimensional subspace  $\mathcal{P}_N(D)$  of  $H_0^1(D; M)$  will be stable in the norms appearing on the left-hand sides of the bounds in Lemmas 2.3.1 and 2.3.6.

## 2.4 The Chauvière–Lozinski transformed FENE model

In this section we show that, in the case of the FENE model, the weak formulation resulting from the substitution  $\hat{\psi} := \psi/M^{2s/b}$  with  $b \geq 4s^2/(2s-1)$  and  $s > 1/2$  also leads to a well-posed problem and a stable semidiscretisation in any number of space dimensions. The minimum value of the function  $s \in (0, \infty) \mapsto 4s^2/(2s-1)$  is attained at  $s = 1$ , yielding the maximum range of  $b$  values,  $b \geq 4$ . This transformation was proposed by Chauvière & Lozinski [90, 33, 32, 91] in the special cases  $s = 2$  and  $s = 2.5$ , where these values were chosen on the basis of numerical experiments in two and three space dimensions, respectively. For the sake of brevity, we shall confine ourselves to establishing an energy estimate analogous to our first stability inequality in Lemma 2.3.1, and the discussion in this section is restricted to the FENE model.

Inserting  $\psi(q) = [M(q)]^{2s/b} \hat{\psi}(q)$  into our model problem (2.1.1), where now  $M$  is the

FENE Maxwellian, yields, on noting that  $\text{tr}(\underline{\kappa})(t) = 0$  for all  $t \in [0, T]$ ,

$$\begin{aligned} \frac{\partial \hat{\psi}}{\partial t} - \frac{1}{2\text{Wi}} \Delta_{\underline{q}} \hat{\psi} &= \frac{1}{2\text{Wi}} \left[ \left(1 - \frac{4s}{b}\right) \left(1 - \frac{|\underline{q}|^2}{b}\right)^{-1} \underline{q} - 2\text{Wi}(\underline{\kappa} \underline{q}) \right] \cdot \nabla_{\underline{q}} \hat{\psi} \\ &+ \frac{1}{2\text{Wi}} \left(1 - \frac{|\underline{q}|^2}{b}\right)^{-2} \left[ d \left(1 - \frac{2s}{b}\right) \left(1 - \frac{|\underline{q}|^2}{b}\right) \right. \\ &\quad \left. + \frac{2(s-1)(2s-b)}{b^2} |\underline{q}|^2 + \frac{4s\text{Wi}}{b} (\underline{q}^T \underline{\kappa} \underline{q}) \left(1 - \frac{|\underline{q}|^2}{b}\right) \right] \hat{\psi}. \end{aligned} \quad (2.4.1)$$

Denoting by  $A(\underline{q}, t)$  the expression in the first square bracket on the right-hand side of (2.4.1) and by  $B(\underline{q}, t)$  the expression in the second square bracket, multiplying (2.4.1) by any  $\hat{\varphi} \in \mathbf{H}_0^1(D)$ , integrating the resulting expression over  $D$ , and integrating by parts in the second term on the left-hand side, yields the following weak formulation.

Let  $\hat{\psi}^0 = \psi^0 / M^{2s/b} \in \mathbf{L}^2(D)$ ; find  $\hat{\psi} \in \mathbf{C}([0, T]; \mathbf{L}^2(D)) \cap \mathbf{L}^2(0, T; \mathbf{H}_0^1(D)) \cap \mathbf{H}^1(0, T; \mathbf{H}_0^1(D)')$  such that

$$\begin{aligned} \frac{d}{dt} \int_D \hat{\psi} \hat{\varphi} \, d\underline{q} + \frac{1}{2\text{Wi}} \int_D \nabla_{\underline{q}} \hat{\psi} \cdot \nabla_{\underline{q}} \hat{\varphi} \, d\underline{q} \\ = \frac{1}{2\text{Wi}} \int_D (A(\underline{q}, t) \cdot \nabla_{\underline{q}} \hat{\psi}) \hat{\varphi} \, d\underline{q} + \frac{1}{2\text{Wi}} \int_D \left(1 - \frac{|\underline{q}|^2}{b}\right)^{-2} B(\underline{q}, t) \hat{\psi} \hat{\varphi} \, d\underline{q}, \end{aligned} \quad (2.4.2)$$

for all  $\hat{\varphi} \in \mathbf{H}_0^1(D)$ , in the sense of distributions on  $(0, T)$ , and with  $\hat{\psi}(\cdot, 0) = \hat{\psi}^0$ .

The backward Euler semidiscretisation of this weak formulation is as follows.

Given  $\hat{\psi}^0 := \psi^0 / M^{2s/b} \in \mathbf{L}^2(D)$ , find  $\hat{\psi}^{n+1} \in \mathbf{H}_0^1(D)$ ,  $n = 0, 1, \dots, N_T - 1$ , such that

$$\begin{aligned} \int_D \frac{\hat{\psi}^{n+1} - \hat{\psi}^n}{\Delta t} \hat{\varphi} \, d\underline{q} + \frac{1}{2\text{Wi}} \int_D \nabla_{\underline{q}} \hat{\psi}^{n+1} \cdot \nabla_{\underline{q}} \hat{\varphi} \, d\underline{q} \\ = \frac{1}{2\text{Wi}} \int_D (A(\underline{q}, t^{n+1}) \cdot \nabla_{\underline{q}} \hat{\psi}^{n+1}) \hat{\varphi} \, d\underline{q} \\ + \frac{1}{2\text{Wi}} \int_D \left(1 - \frac{|\underline{q}|^2}{b}\right)^{-2} B(\underline{q}, t^{n+1}) \hat{\psi}^{n+1} \hat{\varphi} \, d\underline{q}, \end{aligned} \quad (2.4.3)$$

for all  $\hat{\varphi} \in \mathbf{H}_0^1(D)$ .

We begin by showing that, for  $\Delta t$  sufficiently small and all  $b \geq 4s^2/(2s-1)$  and  $s > 1/2$ , this problem has a unique solution. To this end, for  $t \in [0, T]$  fixed, we consider the bilinear form defined on  $\mathbf{H}_0^1(D) \times \mathbf{H}_0^1(D)$  by

$$\begin{aligned} \mathfrak{C}(\hat{\psi}, \hat{\varphi}) &:= \frac{1}{\Delta t} \int_D \hat{\psi} \hat{\varphi} \, d\underline{q} + \frac{1}{2\text{Wi}} \int_D \nabla_{\underline{q}} \hat{\psi} \cdot \nabla_{\underline{q}} \hat{\varphi} \, d\underline{q} \\ &\quad - \frac{1}{2\text{Wi}} \int_D (A(\underline{q}, t) \cdot \nabla_{\underline{q}} \hat{\psi}) \hat{\varphi} \, d\underline{q} - \frac{1}{2\text{Wi}} \int_D \left(1 - \frac{|\underline{q}|^2}{b}\right)^{-2} B(\underline{q}, t) \hat{\psi} \hat{\varphi} \, d\underline{q}. \end{aligned}$$

Now, taking  $\hat{\varphi} = \hat{\psi} \in C_0^\infty(D)$ , integration by parts in the third integral in the definition of  $\mathfrak{C}$ , and then merging the resulting integral with the fourth integral in the definition of  $\mathfrak{C}$ , yields

$$\begin{aligned} \mathfrak{C}(\hat{\psi}, \hat{\psi}) &= \frac{1}{\Delta t} \|\hat{\psi}\|^2 + \frac{1}{2\text{Wi}} \|\nabla_q \hat{\psi}\|^2 + \frac{1}{2\text{Wi}} \left(2s - 1 - \frac{4s^2}{b}\right) \int_D \frac{|q|^2}{b} \left(1 - \frac{|q|^2}{b}\right)^{-2} |\hat{\psi}|^2 dq \\ &\quad - \frac{1}{4\text{Wi}} \int_D \left[ d + \frac{8s\text{Wi}}{b} (q^T \kappa q) \right] \left(1 - \frac{|q|^2}{b}\right)^{-1} |\hat{\psi}|^2 dq. \end{aligned}$$

Assuming that  $b \geq 4s^2/(2s - 1)$  with  $s > 1/2$ , and recalling that  $|q| < \sqrt{b}$  for  $q \in D$ , we then have that

$$\mathfrak{C}(\hat{\psi}, \hat{\psi}) \geq \frac{1}{\Delta t} \|\hat{\psi}\|^2 + \frac{1}{2\text{Wi}} \|\nabla_q \hat{\psi}\|^2 - \frac{1}{4\text{Wi}} (d + 8s\text{Wi} \|\kappa\|_{L^\infty(0,T)}) \int_D \left(1 - \frac{|q|^2}{b}\right)^{-1} |\hat{\psi}|^2 dq.$$

Let us note that for, any  $\beta > 0$ ,

$$\int_D \left(1 - \frac{|q|^2}{b}\right)^{-1} |\hat{\psi}|^2 dq \leq \frac{1}{4\beta} \int_D |\hat{\psi}|^2 dq + \beta \int_D \left(1 - \frac{|q|^2}{b}\right)^{-2} |\hat{\psi}|^2 dq. \quad (2.4.4)$$

Hence, by (2.2.1) and fixing  $\beta$  as the unique solution of the equation  $4b(d + 8s\text{Wi} \|\kappa\|_{L^\infty(0,T)})\beta = 1$ , we have that

$$\mathfrak{C}(\hat{\psi}, \hat{\psi}) \geq \frac{1}{\Delta t} \left(1 - \frac{b\Delta t}{4\text{Wi}} (d + 8s\text{Wi} \|\kappa\|_{L^\infty(0,T)})^2\right) \|\psi\|^2 + \frac{1}{4\text{Wi}} \|\nabla_q \hat{\psi}\|^2 \quad \forall \hat{\psi} \in C_0^\infty(D).$$

Recalling that  $C_0^\infty(D)$  is dense in  $H_0^1(D)$  and, by [17] and [18], also in the  $(1 - |q|^2/b)^{-2}$ -weighted  $L^2$  space,  $L_{M^{-4/b}}^2(D)$ , we deduce that, for any  $\Delta t < 4\text{Wi}/(b(d + 8s\text{Wi} \|\kappa\|_{L^\infty(0,T)})^2)$ , the bilinear form  $\mathfrak{C}$  is coercive on  $H_0^1(D) \times H_0^1(D)$ . The existence of a unique solution  $\{\hat{\psi}^n\}_{n=0}^{N_T}$  to the semidiscretisation (2.4.3) in  $H_0^1(D)$  then follows from the Lax–Milgram theorem, as in the previous section. Using the above coercivity argument, the proof of stability of (2.4.3), stated in Lemma 2.4.1 below, is completely analogous to the proof of Lemma 2.3.1 and is therefore omitted.<sup>2</sup>

**Lemma 2.4.1 (Stability inequality)** *Let  $\Delta t = T/N_T$ ,  $N_T \geq 1$ ,  $\kappa \in \mathbb{C}[0, T]$ ,  $\hat{\psi}^0 \in L^2(D)$ ,  $b \geq 4s^2/(2s - 1)$  with  $s > 1/2$ , and define  $c_0 := b(d + 8s\text{Wi} \|\kappa\|_{L^\infty(0,T)})^2/(2\text{Wi})$ . Suppose that  $\Delta t$  is such that  $0 < c_0\Delta t \leq 1/2$ ; then, we have, for all  $m$  such that  $1 \leq m \leq N_T$ ,*

$$\|\hat{\psi}^m\|^2 + \sum_{n=0}^{m-1} \Delta t \left\| \frac{\hat{\psi}^{n+1} - \hat{\psi}^n}{\sqrt{\Delta t}} \right\|^2 + \sum_{n=0}^{m-1} \frac{\Delta t}{2\text{Wi}} \|\nabla_q \hat{\psi}^{n+1}\|^2 \leq e^{2c_0 m \Delta t} \|\hat{\psi}^0\|^2.$$

Using Lemma 2.4.1, the existence of a unique weak solution to (2.4.2) can be established in the same way as for the symmetrised formulation.

<sup>2</sup>Note that the weak solution here was shown to exist and be unique in  $H_0^1(D)$  (rather than  $H_0^1(D; M)$  as in our earlier, symmetric formulation); clearly,  $H_0^1(D) \subset H^1(D; M) = H_0^1(D; M)$ .

## 2.5 The fully-discrete method

We now return to the semidiscrete method (2.3.1) based on the symmetrised version of the Fokker–Planck equation and describe the construction of a fully-discrete numerical method that stems from this semidiscretisation. At the end of the section we shall comment on the extension of our results to a fully-discrete method based on the semidiscretisation (2.4.3) of the Chauvière–Lozinski-transformed Fokker–Planck equation (2.4.1) for the FENE model.

Let  $\mathcal{P}_N(D)$  be a finite-dimensional subspace of  $H_0^1(D; M)$ , to be chosen below, and let  $\hat{\psi}_N^n \in \mathcal{P}_N(D)$  be the solution at time level  $n$  of our fully-discrete Galerkin method:

$$\int_D \frac{\hat{\psi}_N^{n+1} - \hat{\psi}_N^n}{\Delta t} \hat{\varphi} \, d\mathfrak{q} - \int_D (\mathfrak{k}^{n+1} \mathfrak{q} \hat{\psi}_N^{n+1}) \cdot \nabla_M \hat{\varphi} \, d\mathfrak{q} + \frac{1}{2\text{Wi}} \int_D \nabla_M \hat{\psi}_N^{n+1} \cdot \nabla_M \hat{\varphi} \, d\mathfrak{q} = 0$$

$$\forall \hat{\varphi} \in \mathcal{P}_N(D), \quad n = 0, \dots, N_T - 1, \quad (2.5.1)$$

$$\hat{\psi}_N^0(\cdot) := \text{the } L^2(D) \text{ orthogonal projection of } \hat{\psi}^0(\cdot) = \hat{\psi}(\cdot, 0) \text{ onto } \mathcal{P}_N(D). \quad (2.5.2)$$

**Remark 2.5.1** If the linear space  $\mathcal{P}_N(D)$  is selected so that  $\sqrt{M} \in \mathcal{P}_N(D)$ , then, since  $\sqrt{M} \in \text{Ker}(\nabla_M)$ , it follows on taking  $\hat{\varphi} = \sqrt{M}$  in (2.5.1) that

$$\int_D \sqrt{M(\mathfrak{q})} \hat{\psi}_N^n(\mathfrak{q}) \, d\mathfrak{q} = \int_D \sqrt{M(\mathfrak{q})} \hat{\psi}_N^0(\mathfrak{q}) \, d\mathfrak{q}, \quad n = 1, \dots, N_T,$$

whereby, on letting  $\psi_N^n := \sqrt{M} \hat{\psi}_N^n$ , we have that

$$\int_D \psi_N^n(\mathfrak{q}) \, d\mathfrak{q} = \int_D \psi_N^0(\mathfrak{q}) \, d\mathfrak{q}, \quad n = 1, \dots, N_T.$$

The function  $\psi_N^n$  represents an approximation to the probability density function  $\psi = \sqrt{M} \hat{\psi}$  at  $t = t^n$ . Since, by Lemma 2.3.4,  $\int_D \psi(\mathfrak{q}, t) \, d\mathfrak{q} = \int_D \psi^0(\mathfrak{q}) \, d\mathfrak{q} = 1$  for all  $t \geq 0$ , we deduce, by choosing  $\mathcal{P}_N(D)$  so that  $\sqrt{M} \in \mathcal{P}_N(D)$ , that this integral identity is preserved under discretisation. The integral  $\int_D \psi(\mathfrak{q}, t) \, d\mathfrak{q}$  will sometimes be referred to as the *volume* of  $\psi$ .  $\diamond$

Our objective is to derive a bound on the global error  $e_N^n := \hat{\psi}(\cdot, t^n) - \hat{\psi}_N^n$ . Clearly,

$$e_N^n = (\hat{\psi}(\cdot, t^n) - \hat{\Pi}_N \hat{\psi}(\cdot, t^n)) + (\hat{\Pi}_N \hat{\psi}(\cdot, t^n) - \hat{\psi}_N^n) =: \eta^n + \xi^n,$$

where  $\hat{\Pi}_N \hat{\psi}(\cdot, t^n) \in \mathcal{P}_N(D)$  is a certain projection of  $\hat{\psi}(\cdot, t^n)$  onto  $\mathcal{P}_N(D)$  that will be defined below. For the moment, the specific choices of  $\mathcal{P}_N \subset H_0^1(D; M)$  and  $\hat{\Pi}_N$  are irrelevant. Note also that  $\eta$  is defined for a.e.  $t \in (0, T)$ , *i.e.* not only at the discrete time-levels.

We begin by bounding norms of  $\xi$  in terms of suitable norms of  $\eta$ . Substituting  $\xi$  into (2.5.1), setting  $\hat{\varphi} = \xi^{n+1}$ , and noting that  $\xi^n = \hat{\psi}(\cdot, t^n) - \hat{\psi}_N^n - \eta^n$ , we have

$$\int_D \frac{\xi^{n+1} - \xi^n}{\Delta t} \xi^{n+1} \, d\mathfrak{q} - \int_D (\mathfrak{k}^{n+1} \mathfrak{q} \xi^{n+1}) \cdot \nabla_M \xi^{n+1} \, d\mathfrak{q} + \frac{1}{2\text{Wi}} \int_D \nabla_M \xi^{n+1} \cdot \nabla_M \xi^{n+1} \, d\mathfrak{q}$$

$$= \int_D \mu^{n+1} \xi^{n+1} \, d\mathfrak{q} + \int_D \nu^{n+1} \cdot \nabla_M \xi^{n+1} \, d\mathfrak{q}, \quad (2.5.3)$$

for  $n = 0, \dots, N_T - 1$ , where

$$\mu^{n+1} := \left( \frac{\hat{\psi}(\cdot, t^{n+1}) - \hat{\psi}(\cdot, t^n)}{\Delta t} - \frac{\partial \hat{\psi}}{\partial t}(\cdot, t^{n+1}) \right) - \frac{\eta^{n+1} - \eta^n}{\Delta t}, \quad (2.5.4)$$

$$\varrho^{n+1} := \underline{\kappa}^{n+1} \underline{q} \eta^{n+1} - \frac{1}{2\text{Wi}} \nabla_M \eta^{n+1}. \quad (2.5.5)$$

Since  $\mathcal{P}_N(D) \subset \text{H}_0^1(D; M)$ , (2.5.3) is in the form of (2.3.4); hence, applying Lemma 2.3.1, we obtain

$$\|\xi^m\|^2 + \frac{1}{2\text{Wi}} \sum_{n=0}^{m-1} \Delta t \|\nabla_M \xi^{n+1}\|^2 \leq e^{2c_0 m \Delta t} \left\{ \|\xi^0\|^2 + \sum_{n=0}^{m-1} 2\Delta t (\|\mu^{n+1}\|^2 + 4\text{Wi} \|\varrho^{n+1}\|^2) \right\}, \quad (2.5.6)$$

for  $m = 1, \dots, N_T$ . Let us first consider the term  $\|\xi^0\|$  on the right-hand side of (2.5.6). Since  $\hat{\psi}_N^0$  is the  $L^2(D)$  orthogonal projection of  $\hat{\psi}(\cdot, 0) = \hat{\psi}^0$  onto  $\mathcal{P}_N(D)$ , we have  $(\xi^0, \hat{\varphi}_N) = -(\eta^0, \hat{\varphi}_N)$  for all  $\hat{\varphi}_N \in \mathcal{P}_N(D)$ . Setting  $\hat{\varphi}_N = \xi^0$  here and applying the Cauchy-Schwarz inequality on the right-hand side yields  $\|\xi^0\| \leq \|\eta^0\|$ .

By the triangle inequality we have the following bound on  $\|\varrho^{n+1}\|$ :

$$\|\varrho^{n+1}\| \leq \sqrt{b} |\underline{\kappa}^{n+1}| \|\eta^{n+1}\| + \frac{1}{2\text{Wi}} \|\nabla_M \eta^{n+1}\|, \quad n = 0, \dots, N_T - 1.$$

Hence for the third term on the right-hand-side of (2.5.6), we have

$$\begin{aligned} \sum_{n=0}^{m-1} 8\text{Wi} \Delta t \|\varrho^{n+1}\|^2 &\leq \sum_{n=0}^{m-1} \Delta t \left( 16\text{Wi} b |\underline{\kappa}^{n+1}|^2 \|\eta^{n+1}\|^2 + \frac{4}{\text{Wi}} \|\nabla_M \eta^{n+1}\|^2 \right) \\ &\leq 4c_2 \sum_{n=0}^{m-1} \Delta t \|\eta^{n+1}\|_{\text{H}_0^1(D; M)}^2 = 4c_2 \|\eta\|_{\ell^2(0, t^m; \text{H}_0^1(D; M))}^2, \end{aligned}$$

for  $m = 1, \dots, N_T$ , where  $c_2 := \max(1/\text{Wi}, 4\text{Wi} b |\underline{\kappa}|_{L^\infty(0, T)}^2)$ .

It remains to bound  $\|\mu^{m+1}\|$ . We begin by observing that

$$\|\mu^{m+1}\| \leq \left\| \frac{\hat{\psi}(\cdot, t^{n+1}) - \hat{\psi}(\cdot, t^n)}{\Delta t} - \frac{\partial \hat{\psi}}{\partial t}(\cdot, t^{n+1}) \right\| + \left\| \frac{\eta^{n+1} - \eta^n}{\Delta t} \right\| =: I + II.$$

Bounding both  $I$  and  $II$  by Taylor's theorem with integral remainder yields

$$\begin{aligned} I^2 &\leq \Delta t \int_{t^n}^{t^{n+1}} \left\| \frac{\partial^2 \hat{\psi}}{\partial t^2}(\cdot, t) \right\|^2 dt, \\ II^2 &\leq \int_D \frac{1}{\Delta t} \int_{t^n}^{t^{n+1}} \left| \frac{\partial \eta}{\partial t}(q, t) \right|^2 dt dq = \frac{1}{\Delta t} \int_{t^n}^{t^{n+1}} \left\| \frac{\partial \eta}{\partial t}(\cdot, t) \right\|^2 dt. \end{aligned}$$

Therefore, we now have that

$$\begin{aligned} \sum_{n=0}^{m-1} 2\Delta t \|\mu^{n+1}\|^2 &\leq 4 \sum_{n=0}^{m-1} \Delta t^2 \int_{t^n}^{t^{n+1}} \left\| \frac{\partial^2 \hat{\psi}}{\partial t^2}(\cdot, t) \right\|^2 dt + 4 \sum_{n=0}^{m-1} \int_{t^n}^{t^{n+1}} \left\| \frac{\partial \eta}{\partial t}(\cdot, t) \right\|^2 dt \\ &= 4\Delta t^2 \left\| \frac{\partial^2 \hat{\psi}}{\partial t^2} \right\|_{L^2(0, t^m; L^2(D))}^2 + 4 \left\| \frac{\partial \eta}{\partial t} \right\|_{L^2(0, t^m; L^2(D))}^2. \end{aligned}$$

Combining the bounds on the three terms on the right-hand side of (2.5.6) we deduce that

$$\begin{aligned} & \|\xi^m\|^2 + \frac{1}{2\text{Wi}} \sum_{n=0}^{m-1} \Delta t \|\nabla_M \xi^{n+1}\|^2 \\ & \leq e^{2c_0 m \Delta t} \left( \|\eta^0\|^2 + 4c_2 \|\eta\|_{\ell^2(0,t^m; \mathbf{H}_0^1(D;M))}^2 \right. \\ & \quad \left. + 4 \left\| \frac{\partial \eta}{\partial t} \right\|_{\mathbf{L}^2(0,t^m; \mathbf{L}^2(D))}^2 + 4\Delta t^2 \left\| \frac{\partial^2 \hat{\psi}}{\partial t^2} \right\|_{\mathbf{L}^2(0,t^m; \mathbf{L}^2(D))}^2 \right). \end{aligned} \quad (2.5.7)$$

It remains to bound the first three terms in the bracket on the right-hand side of (2.5.7). To do so we need to make a specific choice of the finite-dimensional space  $\mathcal{P}_N(D)$  from which approximations to  $\hat{\psi} \in \mathbf{H}_0^1(D; M)$  are sought, and we also need to specify the projector  $\hat{\Pi}_N$ . These issues will be discussed in the next section. We shall then return, in Section 2.7, to (2.5.7) and complete the convergence analysis of the numerical method.

**Remark 2.5.2** In the case of the FENE model with  $b \geq 4s^2/(2s-1)$  and  $s > 1/2$  a bound analogous to (2.5.7) can be shown to hold for the fully-discrete version of the semidiscretisation (2.4.3) based on a Chauvière–Lozinski-type transformation, with suitable fixed positive constants  $c_0$  and  $c_2$ , except that  $\mathcal{P}_N(D)$  is then taken to be a finite-dimensional subspace of  $\mathbf{H}_0^1(D)$ ,  $\nabla_M \xi^{n+1}$  on the left-hand side of the bound (2.5.7) is replaced by  $\nabla_q \xi^{n+1}$ , and the norm  $\|\cdot\|_{\ell^2(0,t^m; \mathbf{H}_0^1(D;M))}$  on the right-hand side of (2.5.7) is replaced by  $\|\cdot\|_{\ell^2(0,t^m; \mathbf{H}_0^1(D))}$ . The main steps of the proof are identical to those above: the Cauchy–Schwarz inequality and inequalities (2.2.1) and (2.4.4) are used in the course of bounding the terms on the right-hand side of an error identity analogous to (2.5.3) relating the sequence  $\{\xi^m\}_{m=0}^{N_T}$  to the sequence  $\{\eta^m\}_{m=0}^{N_T}$ , while the terms on the left-hand side of the error identity are bounded below as in the proof the stability inequality stated in Lemma 2.4.1.

We note in particular that the fully-discrete version of the semidiscretisation (2.4.3) based on a Chauvière–Lozinski type transformation  $\hat{\psi} = \psi/M^{2s/b}$  and the finite-dimensional Galerkin subspace  $\mathcal{P}_N(D) \subset \mathbf{H}_0^1(D)$  is stable in the sense that the sequence of numerical solutions  $\{\hat{\psi}_N^n\}_{n=0}^{N_T}$  generated by the fully-discrete scheme satisfies the stability inequality stated in Lemma 2.4.1, with  $\Delta t = T/N_T$ ,  $N_T \geq 1$ ,  $\kappa \in \mathcal{C}[0, T]$ ,  $\hat{\psi}_N^0 \in \mathcal{P}_N(D)$ ,  $b \geq 4s^2/(2s-1)$ ,  $s > 1/2$ ,  $c_0 := b(d + 8s\text{Wi}\|\kappa\|_{\mathbf{L}^\infty(0,T)})^2/(2\text{Wi})$ ,  $0 < c_0 \Delta t \leq 1/2$ , and  $\psi^m$ ,  $\psi^{m-1}$  and  $\psi^0$  replaced by  $\psi_N^m$ ,  $\psi_N^{m-1}$  and  $\psi_N^0$ , respectively, without any conditions relating  $\Delta t$  to  $N$ . The proof of this is identical to that of Lemma 2.4.1, *mutatis mutandis*. We thus deduce that for  $b \gg 1$  a time-step limitation of the form  $\Delta t = \mathcal{O}(b^{-1})$  is needed in order to ensure that  $0 < c_0 \Delta t \leq 1/2$ , and thereby the stability of the method. In this respect the scheme behaves identically to the fully-discrete numerical method (2.5.1), (2.5.2), based on the symmetrised form of the Fokker–Planck equation (*cf.* the conditions of Lemma 2.3.1, for example).  $\diamond$

## 2.6 Approximation results

It was shown in Section 2.2(b) that, under Hypotheses A and B,  $\mathbf{H}_0^1(D) \subset \mathbf{H}^1(D; M) = \mathbf{H}_0^1(D; M)$ . Therefore, any finite-dimensional space  $\mathcal{P}_N(D) \subset \mathbf{H}_0^1(D)$  is, trivially, also contained in  $\mathbf{H}_0^1(D; M)$ . The aim now is to make a specific choice of  $\mathcal{P}_N(D)$  and to explore the approximation properties of the chosen space.

**Remark 2.6.1** As in Remark 2.5.1, if, in addition,  $\sqrt{M} \in \mathcal{P}_N(D)$ , then

$$\int_D \psi_N^n(\underline{q}) \, d\underline{q} = \int_D \psi_N^0(\underline{q}) \, d\underline{q}.$$

In the notation of Lemma 1.3.3, this can be written as  $\varrho_N^n = \varrho_N^0$ . Since, by Hypothesis B,  $\sqrt{M} \in \mathbf{H}_0^1(D)$ , one can ensure that this integral identity holds by including  $\sqrt{M}$  in the finite-dimensional space  $\mathcal{P}_N(D)$ .  $\diamond$

The definition of  $\mathcal{P}_N(D)$  and the choice of the projector  $\hat{\Pi}_N : \mathbf{H}_0^1(D; M) \rightarrow \mathcal{P}_N(D)$  will depend on the number  $d$  of space dimensions. Since the case of  $d = 2$  is sufficiently representative, for the sake of brevity and ease of presentation we shall confine ourselves to two space dimensions in this section, that is, when  $D$  is a disc of radius  $\sqrt{b}$  in  $\mathbb{R}^2$ .

Let  $D_0$  denote the slit disc  $D_0 := D \setminus \{(q_1, 0) : 0 \leq q_1 < \sqrt{b}\}$ . It is natural to transform  $D_0$  into the rectangle  $(r, \theta) \in R := (0, 1) \times (0, 2\pi)$  in a polar co-ordinate system, using the (bijective) change of variables  $\underline{q} = (q_1, q_2) = (\sqrt{b}r \cos \theta, \sqrt{b}r \sin \theta) \in D_0$  where  $(r, \theta) \in R$ . Given  $f \in \mathbf{H}^1(D)$ , define  $\tilde{f}$  on  $R$  by

$$\tilde{f}(r, \theta) := f(q_1, q_2), \quad \underline{q} = (q_1, q_2) \in D_0, \quad (r, \theta) \in R, \quad q_1 = \sqrt{b}r \cos \theta, \quad q_2 = \sqrt{b}r \sin \theta. \quad (2.6.1)$$

Thus,

$$\|f\|_{\mathbf{H}^1(D)}^2 = \|f\|_{\mathbf{H}^1(D_0)}^2 = \int_0^1 r \int_0^{2\pi} \left( b|\tilde{f}|^2 + |\mathbf{D}_r \tilde{f}|^2 + \left| \frac{\mathbf{D}_\theta \tilde{f}}{r} \right|^2 \right) \, d\theta \, dr.$$

where  $\mathbf{D}_r$  denotes differentiation with respect to  $r$ . Motivated by this identity and writing, here and henceforth,  $\tilde{w}(r) := r$  for the weight-function on the interval  $(0, 1)$ , the space  $\tilde{\mathbf{H}}_{\tilde{w}}^1(R)$  is defined as:

$$\tilde{\mathbf{H}}_{\tilde{w}}^1(R) := \{ \tilde{f} \in L_{\text{loc}}^2(0, 1; \mathbf{H}_p^1(0, 2\pi)) : \tilde{f} \in L_{\tilde{w}}^2(R), \quad \mathbf{D}_r \tilde{f} \in L_{\tilde{w}}^2(R) \quad \text{and} \quad \frac{1}{r} \mathbf{D}_\theta \tilde{f} \in L_{\tilde{w}}^2(R) \}, \quad (2.6.2)$$

equipped with the norm  $\|\cdot\|_{\tilde{\mathbf{H}}_{\tilde{w}}^1(R)}$  defined by

$$\|\tilde{f}\|_{\tilde{\mathbf{H}}_{\tilde{w}}^1(R)}^2 := \int_0^1 \tilde{w}(r) \int_0^{2\pi} \left( |\tilde{f}|^2 + |\mathbf{D}_r \tilde{f}|^2 + \left| \frac{\mathbf{D}_\theta \tilde{f}}{r} \right|^2 \right) \, d\theta \, dr, \quad (2.6.3)$$

where  $L_{\tilde{w}}^2(R)$  is the  $\tilde{w}$ -weighted space of square-integrable functions on  $R$ , with norm  $\|\cdot\|_{L_{\tilde{w}}^2(R)}$  defined by

$$\|\tilde{f}\|_{L_{\tilde{w}}^2(R)}^2 := \int_0^1 \tilde{w}(r) \int_0^{2\pi} |\tilde{f}(r, \theta)|^2 \, d\theta \, dr = \int_R |\tilde{f}(r, \theta)|^2 r \, dr \, d\theta,$$

and, for a nonnegative integer  $t$ , the periodic Sobolev space  $\mathbf{H}_p^t(0, 2\pi)$  is given by

$$\mathbf{H}_p^t(0, 2\pi) := \{ \tilde{f} \in \mathbf{H}_{\text{loc}}^t(\mathbb{R}) : \tilde{f}(\theta + 2\pi) = \tilde{f}(\theta) \quad \forall \theta \in \mathbb{R} \}.$$

$\tilde{\mathbf{H}}_{\tilde{w},0}^1(R)$  denotes the subspace of  $\tilde{\mathbf{H}}_{\tilde{w}}^1(R)$  consisting of all functions  $\tilde{f}$  such that the trace  $\tilde{f}(1, \cdot) = 0$ .

We shall also require weighted Sobolev spaces of the form  $H_{\tilde{w}}^{s,t}(R) := H_{\tilde{w}}^s(0, 1; H_p^t(0, 2\pi))$ , equipped (for nonnegative integers  $s$  and  $t$ ) with the norm  $\|\cdot\|_{H_{\tilde{w}}^{s,t}(R)}$  defined by

$$\|\tilde{f}\|_{H_{\tilde{w}}^{s,t}(R)}^2 := \sum_{0 \leq i \leq s, 0 \leq j \leq t} \int_0^1 \tilde{w}(r) \int_0^{2\pi} |D_r^i D_\theta^j \tilde{f}(r, \theta)|^2 d\theta dr.$$

Similarly, for integers  $s \geq 1$  and  $t \geq 0$ , we define  $H_{\tilde{w},0}^{s,t}(R) := H_{\tilde{w},0}^s(0, 1; H_p^t(0, 2\pi))$ , where  $H_{\tilde{w},0}^s(0, 1) := H_{\tilde{w}}^s(0, 1) \cap H_{\tilde{w},0}^1(0, 1)$ , and  $H_{\tilde{w},0}^1(0, 1)$  denotes the set of all  $\tilde{u} \in H_{\tilde{w}}^1(0, 1)$  such that  $\tilde{u}(1) = 0$ .  $H_{\tilde{w},0}^1(0, 1)$  is endowed with the following inner product and norm:

$$(\tilde{u}, \tilde{v})_{H_{\tilde{w},0}^1(0,1)} := \int_0^1 \tilde{w}(r) D_r \tilde{u} D_r \tilde{v} dr \quad \text{and} \quad \|\tilde{u}\|_{H_{\tilde{w},0}^1(0,1)} := \{(\tilde{u}, \tilde{u})_{H_{\tilde{w},0}^1(0,1)}\}^{\frac{1}{2}}.$$

Note that  $\tilde{w}$  is a Jacobi weight function when transformed to  $s \in (-1, 1)$ , since  $\tilde{w}(r(s)) = \frac{1}{2}(1+s)$ .<sup>3</sup> This fact will be important later in this section.

Next, the projection operators are introduced. Due to the cartesian product structure of the set  $R$  it is natural to define distinct projection operators in the  $r$  and  $\theta$  co-ordinate directions. In the  $\theta$ -direction, the orthogonal projection in the  $L^2(0, 2\pi)$  inner product is used (*i.e.*, truncation of the Fourier series). This is denoted by  $P_N^F : L^2(0, 2\pi) \rightarrow \mathbb{S}_N(0, 2\pi)$ , for  $N \geq 1$ , where  $\mathbb{S}_N(0, 2\pi)$  is the space of all trigonometric polynomials in  $\theta \in [0, 2\pi]$  of degree  $N$  or less.<sup>4</sup> Also, let  $\mathbb{S}_{N,0}(0, 2\pi)$  be the orthogonal complement in  $\mathbb{S}_N(0, 2\pi)$ , with respect to the  $L^2(0, 2\pi)$  inner product, of the one-dimensional subspace spanned by constant functions.

The appropriate choice of projector in the  $r$ -direction is less immediate. First of all, for  $N \geq 1$ , let the operator  $P_N^J : H_{\tilde{w},0}^1(0, 1) \rightarrow \mathbb{P}_{N,0}(0, 1)$  be the orthogonal projection in the  $H_{\tilde{w},0}^1(0, 1)$  inner product,<sup>5</sup> where  $\mathbb{P}_{N,0}(0, 1)$  is the space of all algebraic polynomials in  $r \in [0, 1]$ , of degree  $N$  or less, that vanish at  $r = 1$ .

It is tempting to define a two-dimensional projector onto  $\mathbb{S}_N(0, 2\pi) \otimes \mathbb{P}_{N,0}(0, 1)$  as the tensor product of the projectors  $P_N^F$  and  $P_N^J$ . Unfortunately, this choice is inadequate due to the presence of the singular factor  $1/r$  in the weighted Sobolev norm  $\|\cdot\|_{\tilde{H}_{\tilde{w}}^1(R)}$ , and a different definition is required. The lemma below motivates the choice of the two-dimensional projector.

**Lemma 2.6.2 (Decomposition Lemma)** *Let  $\tilde{g} \in \tilde{H}_{\tilde{w}}^1(R)$  and, for  $\varepsilon \in (0, 1)$ , define  $R_\varepsilon := (\varepsilon, 1) \times (0, 2\pi)$ . There exist  $\tilde{g}_1 \in H_{\tilde{w}}^1(0, 1)$  and  $\tilde{g}_2 \in H_{\tilde{w}}^{0,1}(R)$ , with  $\tilde{g}_2 \in H^1(R_\varepsilon)$  for each  $\varepsilon \in (0, 1)$  and  $r\tilde{g}_2 \in \tilde{H}_{\tilde{w}}^1(R)$ , such that*

$$\tilde{g}(r, \theta) = \tilde{g}_1(r) + r\tilde{g}_2(r, \theta) \quad \text{for a.e. } (r, \theta) \in R \quad \text{and} \quad \tilde{g}_1(r) := \frac{1}{2\pi}(\tilde{g}(r, \cdot), 1)_{L^2(0, 2\pi)}.$$

*This is the unique such decomposition of  $\tilde{g}$ . If  $\tilde{g} \in \tilde{H}_{\tilde{w},0}^1(R)$ , then  $\tilde{g}_1 \in H_{\tilde{w},0}^1(0, 1)$  and  $r\tilde{g}_2 \in \tilde{H}_{\tilde{w},0}^1(R)$ , with  $\tilde{g}_2(1, \cdot) = 0$  in the sense of the trace theorem on  $H^1(R_\varepsilon)$ ,  $\varepsilon \in (0, 1)$ .*

<sup>3</sup>Jacobi weight functions are of the form  $(1-s)^\alpha(1+s)^\beta$ ,  $s \in (-1, 1)$  with  $\alpha, \beta > -1$ .

<sup>4</sup>The superscript  $F$  indicates Fourier projection.

<sup>5</sup>The  $J$  superscript indicates projection in a Jacobi-weighted inner-product.

**Proof.** Let  $\tilde{g} \in \tilde{H}_{\tilde{w}}^1(R)$ ; then, by virtue of Fubini's theorem,  $\tilde{g}(r, \cdot) \in H_p^1(0, 2\pi)$  for a.e.  $r \in (0, 1)$ . Let us define, for  $r \in (0, 1)$ , the Fourier coefficients of  $\tilde{g}(r, \cdot)$  by

$$\tilde{\gamma}_n(r) := \frac{1}{\sqrt{2\pi}} \int_0^{2\pi} \tilde{g}(r, \theta) \exp(-in\theta) d\theta, \quad n = 0, 1, \dots$$

According to Parseval's identity,

$$\|\tilde{g}\|_{\tilde{H}_{\tilde{w}}^1(R)}^2 = \sum_{n \in \mathbb{Z}} \int_0^1 \left( |\tilde{\gamma}_n(r)|^2 + |\tilde{\gamma}'_n(r)|^2 + n^2 \left| \frac{\tilde{\gamma}_n(r)}{r} \right|^2 \right) r dr < \infty,$$

whereby, in particular,  $\tilde{\gamma}_0 \in H_{\tilde{w}}^1(0, 1)$  and

$$\tilde{\gamma}_n \in H^1(0, 1; r^{-1}, r) := \left\{ \tilde{f} \in H_{\text{loc}}^1(0, 1) : \int_0^1 \left( r^{-1} |\tilde{f}(r)|^2 + r |\tilde{f}'(r)|^2 \right) dr < \infty \right\},$$

for all  $n \in \mathbb{Z} \setminus \{0\}$ .

For any  $\varepsilon \in (0, 1)$  and  $n \in \mathbb{Z} \setminus \{0\}$ ,  $\tilde{\gamma}_n \in H^1(\varepsilon, 1)$ , and hence by a standard Sobolev embedding,  $\tilde{\gamma}_n \in C(0, 1]$ . Also, for  $0 < r_1 < r_2 < 1$ ,

$$\begin{aligned} \tilde{\gamma}_n(r_2)^2 - \tilde{\gamma}_n(r_1)^2 &= \int_{r_1}^{r_2} \frac{d}{ds} (\tilde{\gamma}_n(s)^2) ds = 2 \int_{r_1}^{r_2} \frac{\tilde{\gamma}_n(s)}{\sqrt{s}} \sqrt{s} \tilde{\gamma}'_n(s) ds \\ &\leq 2 \left( \int_{r_1}^{r_2} s^{-1} |\tilde{\gamma}_n(s)|^2 ds \right)^{\frac{1}{2}} \left( \int_{r_1}^{r_2} s |\tilde{\gamma}'_n(s)|^2 ds \right)^{\frac{1}{2}}, \end{aligned}$$

which is finite by the definition of  $H^1(0, 1; r^{-1}, r)$ , and hence the left-most integral above is finite also. Since the integral is a continuous function of its limits, it follows that  $\tilde{\gamma}_n^2 \in C[0, 1]$ , and hence that  $|\tilde{\gamma}_n| = \sqrt{\tilde{\gamma}_n^2} \in C[0, 1]$ . We now show that  $\tilde{\gamma}_n \in C(0, 1]$  and  $|\tilde{\gamma}_n| \in C[0, 1]$  together imply that  $\tilde{\gamma}_n \in C[0, 1]$ .

There are two cases to consider; (i)  $|\tilde{\gamma}_n(0)| = 0$ , and (ii)  $|\tilde{\gamma}_n(0)| > 0$ . In case (i), we set  $\tilde{\gamma}_n(0) := 0$ . Then  $|\tilde{\gamma}_n(r) - \tilde{\gamma}_n(0)| = |\tilde{\gamma}_n(r)| = ||\tilde{\gamma}_n(r)| - |\tilde{\gamma}_n(0)|| \rightarrow 0_+$  as  $r \rightarrow 0_+$ , by the continuity of  $|\tilde{\gamma}_n|$  on  $[0, 1]$ . In case (ii), there exists  $\delta > 0$  such that  $|\tilde{\gamma}_n(r)| > 0$  for  $r \in [0, \delta]$ . Hence the sign of  $\tilde{\gamma}_n$  does not change on  $(0, \delta]$ , so that  $\tilde{\gamma}_n$  is either  $|\tilde{\gamma}_n|$  or  $-|\tilde{\gamma}_n|$  on the interval  $(0, \delta]$ . Since  $|\tilde{\gamma}_n|, -|\tilde{\gamma}_n| \in C[0, 1]$ , we can define  $\tilde{\gamma}_n(0)$  to be one of  $|\tilde{\gamma}_n(0)|$  or  $-|\tilde{\gamma}_n(0)|$  so that  $\tilde{\gamma}_n \in C[0, 1]$  also.

Now, since  $\tilde{\gamma}_n \in C[0, 1]$ , Parseval's identity above then implies that, necessarily,  $\tilde{\gamma}_n(0) = 0$  for all  $n \in \mathbb{Z} \setminus \{0\}$ .

Let  $\tilde{G}_n(r) := \tilde{\gamma}_n(r)/r$  for  $n \in \mathbb{Z} \setminus \{0\}$ ,  $r \in (0, 1]$  and  $\tilde{E}_n(\theta) := (\exp(in\theta))/\sqrt{2\pi}$ ,  $n \in \mathbb{Z}$ ,  $\theta \in [0, 2\pi]$ . By Parseval's identity, again,  $\sqrt{r^2 + n^2} \tilde{G}_n \in L_{\tilde{w}}^2(0, 1)$ ,  $n \in \mathbb{Z} \setminus \{0\}$ . The following Fourier series expansion of  $\tilde{g}$  can be written as follows:

$$\tilde{g} = \frac{1}{\sqrt{2\pi}} \tilde{\gamma}_0 + r \sum_{n \in \mathbb{Z} \setminus \{0\}} \tilde{G}_n \tilde{E}_n,$$

with equality in the sense of  $\tilde{H}_{\tilde{w}}^1(R)$ . We define  $\tilde{g}_1 := \tilde{\gamma}_0/\sqrt{2\pi}$  and  $\tilde{g}_2 = \sum_{n \in \mathbb{Z} \setminus \{0\}} \tilde{G}_n \tilde{E}_n$  to deduce the stated decomposition  $\tilde{g}(r, \theta) = \tilde{g}_1(r) + r\tilde{g}_2(r, \theta)$ , and we note that  $\tilde{g}_1 = \frac{1}{2\pi}(\tilde{g}, 1)_{L^2(0, 2\pi)} \in H_{\tilde{w}}^1(0, 1)$  and  $\tilde{g}_2 \in H_{\tilde{w}}^{0,1}(R)$ ; moreover, trivially,  $r\tilde{g}_2 = \tilde{g} - \tilde{g}_1 \in \tilde{H}_{\tilde{w}}^1(R)$ .

Also, since  $\tilde{g} \in \tilde{H}_{\tilde{w}}^1(R)$  it follows that  $\tilde{g} \in H^1(R_\varepsilon)$  and  $\tilde{g}_1 \in H^1(\varepsilon, 1)$  for any  $\varepsilon \in (0, 1)$ . Hence,  $\tilde{g}_2 = (\tilde{g} - \tilde{g}_1)/r \in H^1(R_\varepsilon)$  for any  $\varepsilon \in (0, 1)$ .

For  $\tilde{g}_1 = \tilde{\gamma}_0/\sqrt{2\pi}$  fixed, as in the statement of the lemma, the uniqueness of  $\tilde{g}_2$  follows easily by *reductio ad absurdum*: suppose that  $\tilde{h}_2$  is another function, with the same regularity properties as  $\tilde{g}_2$ , and such that  $\tilde{g} = \tilde{g}_1 + r\tilde{h}_2$ . Then,  $r(\tilde{h}_2 - \tilde{g}_2) = 0$  a.e. on  $R$ , and therefore  $\tilde{h}_2 = \tilde{g}_2$  a.e. on  $R$ .

The final statement of the lemma follows directly from the definitions of  $\tilde{\gamma}_n$ ,  $n \in \mathbb{Z}$  and the definitions of  $\tilde{g}_1$  and  $\tilde{g}_2$  via the  $\tilde{\gamma}_n$ ,  $n \in \mathbb{Z}$ .  $\square$

Suppose that  $\tilde{g} \in \tilde{H}_{\tilde{w},0}^1(R)$ . On applying Lemma 2.6.2 we deduce that  $\tilde{g}$  has the (unique) decomposition

$$\tilde{g}(r, \theta) = \tilde{g}_1(r) + r\tilde{g}_2(r, \theta), \quad (2.6.4)$$

where  $\tilde{g}_1 := \frac{1}{2\pi}(\tilde{g}, 1)_{L^2(0,2\pi)} \in H_{\tilde{w},0}^1(0,1)$ ,  $\tilde{g}_2 \in H_{\tilde{w}}^{0,1}(R)$  and  $\tilde{g}_2(1, \cdot) = 0$ . Note also that  $(g_2(r, \cdot), 1)_{L^2(0,2\pi)} = 0$  for a.e.  $r \in (0, 1)$ . We shall assume in addition that  $\tilde{g}_2(\cdot, \theta) \in H_{\tilde{w},0}^1(0,1)$  for a.e.  $\theta \in (0, 2\pi)$ ; by virtue of Fubini's theorem, a convenient sufficient condition for this is that  $\tilde{g}_2 \in H_{\tilde{w},0}^{1,0}(R)$ , for example. We then define

$$\tilde{P}_N^J \tilde{g}(\cdot, \theta) := P_N^J \tilde{g}_1(\cdot) + r P_N^J \tilde{g}_2(\cdot, \theta), \quad \theta \in (0, 2\pi),$$

where  $P_N^J : H_{\tilde{w},0}^1(0,1) \rightarrow \mathbb{P}_{N,0}(0,1)$  is the orthogonal projector defined above.

There are a number of approximation results available in the literature related to projectors in Jacobi-weighted inner products (see for example [16] or [28]). Since the setting here is specific, we shall establish the required approximation properties of the univariate projector  $P_N^J$  from first principles. The approximation properties of  $\tilde{P}_N^J$  and of our two-dimensional projector  $P_N^F \tilde{P}_N^J$  will then follow. The relevant results are stated in the next two lemmas.

**Lemma 2.6.3** *Suppose that  $\tilde{g} \in H_{\tilde{w},0}^k(0,1)$  with  $k \geq 1$ ; then,*

$$\|\tilde{g} - P_N^J \tilde{g}\|_{H_{\tilde{w}}^1(0,1)} \leq cN^{1-k} \|\tilde{g}\|_{H_{\tilde{w}}^k(0,1)} \quad (2.6.5)$$

and

$$\|\tilde{g} - P_N^J \tilde{g}\|_{L_{\tilde{w}}^2(0,1)} \leq cN^{-k} \|\tilde{g}\|_{H_{\tilde{w}}^k(0,1)}. \quad (2.6.6)$$

**Proof.** First consider (2.6.5). Note that by Pythagoras' theorem,

$$\|\tilde{g} - P_N^J \tilde{g}\|_{H_{\tilde{w},0}^1(0,1)} = \left( \|\tilde{g}\|_{H_{\tilde{w},0}^1(0,1)}^2 - \|P_N^J \tilde{g}\|_{H_{\tilde{w},0}^1(0,1)}^2 \right)^{\frac{1}{2}} \leq \|\tilde{g}\|_{H_{\tilde{w},0}^1(0,1)} \leq \|\tilde{g}\|_{H_{\tilde{w}}^k(0,1)}.$$

If  $k = 1$ , the right-most term in this chain is equal to  $1 \cdot N^{1-k} \|\tilde{g}\|_{H_{\tilde{w}}^k(0,1)}$ , while if  $k \geq 2$  and  $1 \leq N < k - 1$ , then it is bounded by  $(k - 1)^{k-1} N^{1-k} \|\tilde{g}\|_{H_{\tilde{w}}^k(0,1)}$ .

Finally, if  $k \geq 2$  and  $N \geq \max(2, k - 1)$ , then recall that, by the definition of  $P_N^J$ ,

$$\|\tilde{g} - P_N^J \tilde{g}\|_{H_{\tilde{w},0}^1(0,1)} \leq \|\tilde{g} - \tilde{v}\|_{H_{\tilde{w},0}^1(0,1)} \quad \forall \tilde{v} \in \mathbb{P}_{N,0}(0,1).$$

Select, in particular,

$$\tilde{v}(r) = - \int_r^1 Q_{N-1}^J D_s \tilde{g}(s) ds, \quad r \in [0, 1],$$

where  $Q_{N-1}^J$  is the orthogonal projector in  $L_{\tilde{w}}^2(0, 1)$  onto  $\mathbb{P}_{N-1}(0, 1)$ , the set of all algebraic polynomials of degree  $N - 1$  or less on the interval  $[0, 1]$ . Thus,

$$\|\tilde{g} - P_N^J \tilde{g}\|_{\mathbf{H}_{\tilde{w},0}^1(0,1)} \leq \|D_r \tilde{g} - D_r \tilde{v}\|_{L_{\tilde{w}}^2(0,1)} = \|D_r \tilde{g} - Q_{N-1}^J(D_r \tilde{g})\|_{L_{\tilde{w}}^2(0,1)} \leq c(N-1)^{1-k} \|\tilde{g}\|_{\mathbf{H}_{\tilde{w}}^k(0,1)},$$

where the last bound (scaled from the standard interval  $(-1, 1)$  to  $(0, 1)$ ) comes from Sec. 5.7.1 of Canuto *et al.* [28], and is valid for  $N \geq \max(2, k - 1)$ ,  $k \geq 2$ . Hence, after bounding  $(N - 1)^{1-k}$  by  $2^{k-1} N^{1-k}$  (recall that  $N \geq 2$  by hypothesis), it follows that

$$\|\tilde{g} - P_N^J \tilde{g}\|_{\mathbf{H}_{\tilde{w},0}^1(0,1)} \leq c 2^{k-1} N^{1-k} \|\tilde{g}\|_{\mathbf{H}_{\tilde{w}}^k(0,1)}.$$

Now choosing  $\hat{c} = \max\{(k - 1)^{k-1}, c 2^{k-1}\}$  for  $k \geq 1$ , with the convention that  $0^0 := 1$ ,

$$\|\tilde{g} - P_N^J \tilde{v}\|_{\mathbf{H}_{\tilde{w},0}^1(0,1)} \leq \hat{c} N^{1-k} \|\tilde{g}\|_{\mathbf{H}_{\tilde{w}}^k(0,1)}$$

for all  $N \geq 1$  (regardless of whether or not  $N \geq k - 1$ ).

For any  $\tilde{v} \in \mathbf{H}_{\tilde{w},0}^1(0, 1)$ , we have:

$$\begin{aligned} \|\tilde{v}\|_{L_{\tilde{w}}^2(0,1)}^2 &= \int_0^1 \tilde{v}^2(r) r \, dr = \int_0^1 \left( \int_r^1 (\sqrt{s} D_s \tilde{v}(s) \frac{1}{\sqrt{s}} \, ds) \right)^2 r \, dr \\ &\leq \int_0^1 r \left( \int_r^1 |D_s \tilde{v}(s)|^2 s \, ds \right) \left( \int_r^1 \frac{1}{s} \, ds \right) \, dr \\ &\leq \left( \int_0^1 r |\log r| \, dr \right) \|\tilde{v}\|_{\mathbf{H}_{\tilde{w},0}^1(0,1)}^2 = \frac{1}{4} \|\tilde{v}\|_{\mathbf{H}_{\tilde{w},0}^1(0,1)}^2, \end{aligned} \quad (2.6.7)$$

where we make the substitution  $r = e^t$  to evaluate  $\int_0^1 r |\log r| \, dr$ . It follows from the Friedrichs inequality (2.6.7) that  $\|\cdot\|_{\mathbf{H}_{\tilde{w},0}^1(0,1)}$  and  $\|\cdot\|_{\mathbf{H}_{\tilde{w}}^1(0,1)}$  are equivalent norms on  $\mathbf{H}_{\tilde{w},0}^1(0, 1)$ , and therefore (2.6.5) holds for any  $N \geq 1$ .

The proof of (2.6.6) is based on a duality argument. Let  $e := \tilde{g} - P_N^J \tilde{g}$  and note that, by the hypotheses of the lemma on  $\tilde{g}$ , we have  $e \in L_{\tilde{w}}^2(0, 1)$ . Consider the mixed Neumann–Dirichlet boundary-value problem:

$$-D_r(r D_r z_e(r)) = r e(r), \quad r \in (0, 1), \quad \lim_{r \rightarrow 0^+} r D_r z_e(r) = 0, \quad z_e(1) = 0. \quad (2.6.8)$$

By (2.6.7) and the Lax–Milgram theorem, this has a unique weak solution  $z_e \in \mathbf{H}_{\tilde{w},0}^1(0, 1)$  satisfying

$$(z_e, v)_{\mathbf{H}_{\tilde{w},0}^1(0,1)} = (e, v)_{L_{\tilde{w}}^2(0,1)} \quad \forall v \in \mathbf{H}_{\tilde{w},0}^1(0, 1). \quad (2.6.9)$$

Also, by (2.6.7),

$$\|z_e\|_{\mathbf{H}_{\tilde{w},0}^1(0,1)}^2 \leq \frac{5}{16} \|e\|_{L_{\tilde{w}}^2(0,1)}^2.$$

We shall show that in fact  $D_r^2 z_e \in L_{\tilde{w}}^2(0, 1)$ , and thereby  $z_e \in \mathbf{H}_{\tilde{w},0}^2(0, 1)$ . To this end, note that

$$D_r z_e(r) = -\frac{1}{r} \int_0^r s e(s) \, ds, \quad r \in (0, 1].$$

Hence,  $D_r z_e \in C(0, 1]$  and, on recalling that  $e \in L^2_{\tilde{w}}(0, 1)$ , the Cauchy–Schwarz inequality yields

$$|D_r z_e(r)|^2 \leq \frac{1}{2} \int_0^r s |e(s)|^2 ds, \quad r \in (0, 1]. \quad (2.6.10)$$

This inequality implies that  $\lim_{r \rightarrow 0^+} D_r z_e(r) = 0$  and that, for any  $\varepsilon \in (0, 1)$ ,

$$\int_\varepsilon^1 \frac{1}{r} |D_r z_e(r)|^2 dr \leq \frac{1}{2\varepsilon} \int_0^1 s |e(s)|^2 ds.$$

Thus,  $\sqrt{r}(r^{-1}D_r z_e) \in L^2(\varepsilon, 1)$ ; hence, by (2.6.8),  $\sqrt{r}D_r^2 z_e = -\sqrt{r}(e + r^{-1}D_r z_e) \in L^2(\varepsilon, 1)$ . Multiplying this equality by  $\sqrt{r}D_r^2 z_e$  and integrating over the interval  $(\varepsilon, 1)$  gives

$$\int_\varepsilon^1 r |D_r^2 z_e(r)|^2 dr + \int_\varepsilon^1 D_r z_e(r) D_r^2 z_e(r) dr = - \int_\varepsilon^1 r e(r) D_r^2 z_e(r) dr.$$

Hence, by computing explicitly the second integral on the left-hand side and applying Cauchy's inequality  $|\alpha\beta| \leq \frac{1}{2}(\alpha^2 + \beta^2)$  on the right-hand side, we obtain

$$\int_\varepsilon^1 r |D_r^2 z_e(r)|^2 dr + |D_r z_e(1)|^2 \leq \int_\varepsilon^1 r |e(r)|^2 dr + |D_r z_e(\varepsilon)|^2.$$

Passing to the limit  $\varepsilon \rightarrow 0_+$  and omitting the second term on the left-hand side gives that  $D_r^2 z_e \in L^2_{\tilde{w}}(0, 1)$  and

$$\int_0^1 r |D_r^2 z_e(r)|^2 dr \leq \int_0^1 r |e(r)|^2 dr.$$

Combining this with our earlier bound from (2.6.9), we have that  $\|z_e\|_{\mathbb{H}_{\tilde{w}}^2(0,1)}^2 \leq \frac{21}{16} \|e\|_{L^2_{\tilde{w}}(0,1)}^2$ .

We are now ready to embark on the analysis of the projection error in the  $L^2_{\tilde{w}}(0, 1)$  norm. Recalling that  $e = \tilde{g} - P_N^J \tilde{g} \in \mathbb{H}_{\tilde{w},0}^1(0, 1)$ , we deduce from the weak formulation (2.6.9), the definition of the orthogonal projector  $P_N^J$ , the Cauchy–Schwarz inequality, (2.6.5) and the  $\mathbb{H}_{\tilde{w}}^2(0, 1)$  norm bound just derived that

$$\begin{aligned} \|\tilde{g} - P_N^J \tilde{g}\|_{L^2_{\tilde{w}}(0,1)}^2 &= (e, \tilde{g} - P_N^J \tilde{g})_{L^2_{\tilde{w}}(0,1)} = (z_e, \tilde{g} - P_N^J \tilde{g})_{\mathbb{H}_{\tilde{w},0}^1(0,1)} \\ &= (\tilde{g} - P_N^J \tilde{g}, z_e - P_N^J z_e)_{\mathbb{H}_{\tilde{w},0}^1(0,1)} \\ &\leq \|\tilde{g} - P_N^J \tilde{g}\|_{\mathbb{H}_{\tilde{w},0}^1(0,1)} \|z_e - P_N^J z_e\|_{\mathbb{H}_{\tilde{w},0}^1(0,1)} \\ &\leq cN^{1-k} \|\tilde{g}\|_{\mathbb{H}_{\tilde{w}}^k(0,1)} \cdot N^{-1} \|z_e\|_{\mathbb{H}_{\tilde{w}}^2(0,1)} \\ &\leq cN^{-k} \|\tilde{g}\|_{\mathbb{H}_{\tilde{w}}^k(0,1)} \|\tilde{g} - P_N^J \tilde{g}\|_{L^2_{\tilde{w}}(0,1)}, \quad k \geq 1. \end{aligned}$$

Dividing the left-most and the right-most term in this chain by  $\|\tilde{g} - P_N^J \tilde{g}\|_{L^2_{\tilde{w}}(0,1)}$  gives (2.6.6).  $\square$

Next, for  $\tilde{g} \in \tilde{\mathbb{H}}_{\tilde{w},0}^1(R)$ , with decomposition given in (2.6.4), we define the projection operator  $\tilde{\Pi}_N : \tilde{\mathbb{H}}_{\tilde{w},0}^1(R) \rightarrow \mathcal{P}_N(R)$  as:

$$(\tilde{\Pi}_N \tilde{g})(r, \theta) := (P_{N_\theta}^F \tilde{P}_{N_r}^J \tilde{g})(r, \theta) = (\tilde{P}_{N_r}^J P_{N_\theta}^F \tilde{g})(r, \theta),$$

where the finite-dimensional space  $\mathcal{P}_N(R)$  is defined as

$$\mathcal{P}_N(R) := \mathbb{P}_{N_r,0}(0,1) \oplus (r\mathbb{P}_{N_r,0}(0,1) \otimes \mathbb{S}_{N_\theta,0}(0,2\pi)).$$

The structure of this space reflects the decomposition (2.6.4). Note that the constant functions have been factored out of the space  $\mathbb{S}_{N_\theta}(0,2\pi)$  in the definition of  $\mathcal{P}_N(R)$ ; this is appropriate because, as observed above,  $(g_2(r, \cdot), 1)_{L^2(0,2\pi)} = 0$ . The lemma below establishes optimal order approximation results for this projector.

**Lemma 2.6.4** *Let  $\tilde{g} \in \tilde{\mathbb{H}}_{\tilde{w},0}^1(R)$ , with decomposition  $\tilde{g}(r, \theta) = \tilde{g}_1(r) + r\tilde{g}_2(r, \theta)$ , where  $\tilde{g}_1 = \frac{1}{2\pi}(\tilde{g}, 1)_{L^2(0,2\pi)} \in \mathbb{H}_{\tilde{w},0}^1(0,1)$ ,  $\tilde{g}_2 \in \mathbb{H}_{\tilde{w}}^{0,1}(R)$ ,  $\tilde{g}_2(1, \cdot) = 0$ , and assume, in addition, that  $\tilde{g}_2(\cdot, \theta) \in \mathbb{H}_{\tilde{w},0}^1(0,1)$  for a.e.  $\theta \in (0, 2\pi)$ . If  $\tilde{g}_1 \in \mathbb{H}_{\tilde{w}}^{k+1}(0,1)$  and  $\tilde{g}_2 \in \mathbb{H}_{\tilde{w}}^{k+1,0}(R) \cap \mathbb{H}_{\tilde{w}}^{k,1}(R) \cap \mathbb{H}_{\tilde{w}}^{0,l+1}(R) \cap \mathbb{H}_{\tilde{w}}^{1,l}(R)$  for some  $k, l \geq 1$ , then*

$$\begin{aligned} \|\tilde{g} - \tilde{\Pi}_N \tilde{g}\|_{\tilde{\mathbb{H}}_{\tilde{w}}^1(R)} &\leq C_1 N_r^{-k} \left( \|\tilde{g}_1\|_{\mathbb{H}_{\tilde{w}}^{k+1}(0,1)}^2 + \|\tilde{g}_2\|_{\mathbb{H}_{\tilde{w}}^{k+1,0}(R)}^2 + \|\tilde{g}_2\|_{\mathbb{H}_{\tilde{w}}^{k,1}(R)}^2 \right)^{\frac{1}{2}} \\ &\quad + C_2 N_\theta^{-l} \left( \|\tilde{g}_2\|_{\mathbb{H}_{\tilde{w}}^{0,l+1}(R)}^2 + \|\tilde{g}_2\|_{\mathbb{H}_{\tilde{w}}^{1,l}(R)}^2 \right)^{\frac{1}{2}}. \end{aligned} \quad (2.6.11)$$

If  $\tilde{g}_1 \in \mathbb{H}_{\tilde{w}}^k(0,1)$  and  $\tilde{g}_2 \in \mathbb{H}_{\tilde{w}}^{k,0}(R) \cap \mathbb{H}_{\tilde{w}}^{0,l}(R)$  for some  $k, l \geq 1$ , then

$$\|\tilde{g} - \tilde{\Pi}_N \tilde{g}\|_{L_{\tilde{w}}^2(R)} \leq C_1 N_r^{-k} \left( \|\tilde{g}_1\|_{\mathbb{H}_{\tilde{w}}^k(0,1)}^2 + \|\tilde{g}_2\|_{\mathbb{H}_{\tilde{w}}^{k,0}(R)}^2 \right)^{\frac{1}{2}} + C_2 N_\theta^{-l} \|\tilde{g}_2\|_{\mathbb{H}_{\tilde{w}}^{0,l}(R)}. \quad (2.6.12)$$

**Proof.** The left-hand side in (2.6.11) is given by:

$$\begin{aligned} \|\tilde{g} - \tilde{\Pi}_N \tilde{g}\|_{\tilde{\mathbb{H}}_{\tilde{w}}^1(R)}^2 &= \int_0^1 \tilde{w}(r) \int_0^{2\pi} \left\{ (\tilde{g} - \tilde{\Pi}_N \tilde{g})^2 + (D_r \tilde{g} - D_r(\tilde{\Pi}_N \tilde{g}))^2 \right\} d\theta dr \\ &\quad + \int_0^1 r^{-1} \int_0^{2\pi} (D_\theta \tilde{g} - D_\theta(\tilde{\Pi}_N \tilde{g}))^2 d\theta dr =: I + II. \end{aligned}$$

First consider term  $I$ . The two terms in the, inner,  $\theta$ -integral in  $I$  will be treated separately. Using the  $L^2$ -error bound for Fourier projection, as well as the fact that

$$\|P_{N_\theta}^F\|_{\mathcal{L}(L_p^2(0,2\pi), L_p^2(0,2\pi))} \leq 1,$$

it follows that

$$\begin{aligned} \|\tilde{g}(r, \cdot) - \tilde{\Pi}_N \tilde{g}(r, \cdot)\|_{L^2(0,2\pi)}^2 &\leq \left( \|\tilde{g}(r, \cdot) - P_{N_\theta}^F \tilde{g}(r, \cdot)\|_{L^2(0,2\pi)} + \|P_{N_\theta}^F(\tilde{g}(r, \cdot) - \tilde{P}_{N_r}^J \tilde{g}(r, \cdot))\|_{L^2(0,2\pi)} \right)^2 \\ &\leq \left( C_3 N_\theta^{-l} \|D_\theta^l \tilde{g}(r, \cdot)\|_{L^2(0,2\pi)} + \|\tilde{g}(r, \cdot) - \tilde{P}_{N_r}^J \tilde{g}(r, \cdot)\|_{L^2(0,2\pi)} \right)^2 \\ &\leq 2C_3^2 N_\theta^{-2l} \|D_\theta^l \tilde{g}(r, \cdot)\|_{L^2(0,2\pi)}^2 + 2\|\tilde{g}(r, \cdot) - \tilde{P}_{N_r}^J \tilde{g}(r, \cdot)\|_{L^2(0,2\pi)}^2, \end{aligned}$$

where  $D_\theta^l \tilde{g} = r D_\theta^l \tilde{g}_2$  and  $0 \leq r \leq 1$  have been used in the last line. Similarly,

$$\begin{aligned}
\|D_r \tilde{g}(r, \cdot) - D_r(\tilde{\Pi}_N \tilde{g}(r, \cdot))\|_{L^2(0, 2\pi)}^2 &\leq 2\|D_r \tilde{g} - P_{N_\theta}^F D_r \tilde{g}\|_{L^2(0, 2\pi)}^2 \\
&\quad + 2\|D_r P_{N_\theta}^F \tilde{g} - D_r P_{N_\theta}^F \tilde{P}_{N_r}^J \tilde{g}(r, \cdot)\|_{L^2(0, 2\pi)}^2 \\
&\leq 2C_3^2 N_\theta^{-2l} \|D_\theta^l D_r \tilde{g}(r, \cdot)\|_{L^2(0, 2\pi)}^2 \\
&\quad + 2\|D_r \tilde{g} - D_r \tilde{P}_{N_r}^J \tilde{g}(r, \cdot)\|_{L^2(0, 2\pi)}^2 \\
&\leq 4C_3^2 N_\theta^{-2l} \left( \|D_\theta^l \tilde{g}_2(r, \cdot)\|_{L^2(0, 2\pi)}^2 + \|D_r D_\theta^l \tilde{g}_2(r, \cdot)\|_{L^2(0, 2\pi)}^2 \right) \\
&\quad + 2\|D_r \tilde{g}(r, \cdot) - D_r \tilde{P}_{N_r}^J \tilde{g}(r, \cdot)\|_{L^2(0, 2\pi)}^2.
\end{aligned}$$

Therefore,

$$\begin{aligned}
I &\leq 6C_3^2 N_\theta^{-2l} \int_0^{2\pi} \left( \|D_\theta^l \tilde{g}_2(\cdot, \theta)\|_{L_{\tilde{w}}^2(0, 1)}^2 + \|D_r D_\theta^l \tilde{g}_2(\cdot, \theta)\|_{L_{\tilde{w}}^2(0, 1)}^2 \right) d\theta \\
&\quad + 2 \int_0^{2\pi} \|\tilde{g}(\cdot, \theta) - \tilde{P}_{N_r}^J \tilde{g}(\cdot, \theta)\|_{H_{\tilde{w}}^1(0, 1)}^2 d\theta.
\end{aligned}$$

The final term on the right-hand side of the last inequality can then be bounded using the univariate estimate (2.6.5):

$$\begin{aligned}
\|\tilde{g}(\cdot, \theta) - \tilde{P}_{N_r}^J \tilde{g}(\cdot, \theta)\|_{H_{\tilde{w}}^1(0, 1)}^2 &\leq 2\|\tilde{g}_1 - P_{N_r}^J \tilde{g}_1\|_{H_{\tilde{w}}^1(0, 1)}^2 + 2\|r(\tilde{g}_2(\cdot, \theta) - P_{N_r}^J \tilde{g}_2(\cdot, \theta))\|_{H_{\tilde{w}}^1(0, 1)}^2 \\
&\leq C^2 N_r^{-2k} \|\tilde{g}_1\|_{H_{\tilde{w}}^{k+1}(0, 1)}^2 \\
&\quad + 2 \int_0^1 \tilde{w}(r) \left\{ (2 + r^2)(\tilde{g}_2(r, \theta) - P_{N_r}^J \tilde{g}_2(r, \theta))^2 + 2r^2 (D_r(\tilde{g}_2(r, \theta) - P_{N_r}^J \tilde{g}_2(r, \theta)))^2 \right\} dr \\
&\leq C^2 N_r^{-2k} \|\tilde{g}_1\|_{H_{\tilde{w}}^{k+1}(0, 1)}^2 + 6\|\tilde{g}_2(\cdot, \theta) - P_{N_r}^J \tilde{g}_2(\cdot, \theta)\|_{H_{\tilde{w}}^1(0, 1)}^2 \\
&\leq C_4^2 N_r^{-2k} \left( \|\tilde{g}_1\|_{H_{\tilde{w}}^{k+1}(0, 1)}^2 + \|\tilde{g}_2(\cdot, \theta)\|_{H_{\tilde{w}}^{k+1}(0, 1)}^2 \right).
\end{aligned}$$

Therefore,

$$\begin{aligned}
I &\leq 6C_3^2 N_\theta^{-2l} \int_0^{2\pi} \left( \|D_\theta^l \tilde{g}_2(\cdot, \theta)\|_{L_{\tilde{w}}^2(0, 1)}^2 + \|D_r D_\theta^l \tilde{g}_2(\cdot, \theta)\|_{L_{\tilde{w}}^2(0, 1)}^2 \right) d\theta \\
&\quad + 2C_4^2 N_r^{-2k} \int_0^{2\pi} \left( \|\tilde{g}_1\|_{H_{\tilde{w}}^{k+1}(0, 1)}^2 + \|\tilde{g}_2(\cdot, \theta)\|_{H_{\tilde{w}}^{k+1}(0, 1)}^2 \right) d\theta, \tag{2.6.13}
\end{aligned}$$

which is an optimal-order bound on  $I$ .

Next, consider  $II$ . Since  $\theta$ -differentiation commutes with the projectors  $P_{N_r}^J$  and  $P_{N_\theta}^F$ , it follows that

$$\begin{aligned}
II &\leq 2 \int_0^1 r^{-1} \int_0^{2\pi} |D_\theta \tilde{g}(r, \theta) - P_{N_\theta}^F D_\theta \tilde{g}(r, \theta)|^2 d\theta dr \\
&\quad + 2 \int_0^1 r^{-1} \int_0^{2\pi} |P_{N_\theta}^F D_\theta \tilde{g}(r, \theta) - \tilde{P}_{N_r}^J (P_{N_\theta}^F D_\theta \tilde{g}(r, \theta))|^2 d\theta dr.
\end{aligned}$$

Therefore,

$$\begin{aligned}
II &\leq 2 \int_0^1 r^{-1} \int_0^{2\pi} |r D_\theta \tilde{g}_2(r, \theta) - r P_{N_\theta}^F D_\theta \tilde{g}_2(r, \theta)|^2 d\theta dr \\
&\quad + 2 \int_0^{2\pi} \int_0^1 r^{-1} |r P_{N_\theta}^F D_\theta \tilde{g}_2(r, \theta) - \tilde{P}_{N_r}^J (r P_{N_\theta}^F D_\theta \tilde{g}_2(r, \theta))|^2 dr d\theta \\
&\leq C_5^2 N_\theta^{-2l} \int_0^1 \tilde{w}(r) \int_0^{2\pi} |D_\theta^{l+1} \tilde{g}_2(r, \theta)|^2 d\theta dr \\
&\quad + 2 \int_0^{2\pi} \int_0^1 \tilde{w}(r) |P_{N_\theta}^F D_\theta \tilde{g}_2(r, \theta) - P_{N_r}^J (P_{N_\theta}^F D_\theta \tilde{g}_2(r, \theta))|^2 dr d\theta \\
&\leq C_5^2 N_\theta^{-2l} \int_0^{2\pi} \|D_\theta^{l+1} \tilde{g}_2(\cdot, \theta)\|_{L_w^2(0,1)}^2 d\theta + C_6^2 N_r^{-2k} \int_0^{2\pi} \|P_{N_\theta}^F D_\theta \tilde{g}_2(r, \theta)\|_{H_w^k(0,1)}^2 d\theta.
\end{aligned}$$

Where the  $L_w^2(0,1)$  norm error bound for  $P_{N_r}^J$  stated in (2.6.6), as well as the fact that  $\tilde{P}_{N_r}^J(r\tilde{g}_2) = r P_{N_r}^J(\tilde{g}_2)$  have been used in the argument above. For the second integral in the last line in the bound on  $II$ ,

$$\sum_{j=0}^k \int_0^1 \tilde{w}(r) \|P_{N_\theta}^F D_r^j D_\theta \tilde{g}_2(\cdot, r)\|_{L^2(0,2\pi)}^2 dr \leq \sum_{j=0}^k \int_0^1 \tilde{w}(r) \|D_r^j D_\theta \tilde{g}_2(\cdot, r)\|_{L^2(0,2\pi)}^2 dr.$$

Therefore,

$$II \leq C_5^2 N_\theta^{-2l} \int_0^{2\pi} \|D_\theta^{l+1} \tilde{g}_2(\cdot, \theta)\|_{L_w^2(0,1)}^2 d\theta + C_6^2 N_r^{-2k} \int_0^{2\pi} \|D_\theta \tilde{g}_2(\cdot, \theta)\|_{H_w^k(0,1)}^2 d\theta.$$

Combining the bounds for  $I$  and  $II$  with suitable constants  $C_1$  and  $C_2$ , gives

$$\begin{aligned}
\|\tilde{g} - P_{N_\theta}^F \tilde{P}_{N_r}^J \tilde{g}\|_{\tilde{H}_w^1(R)} &\leq C_1 N_r^{-k} \left\{ \int_0^{2\pi} (\|\tilde{g}_1\|_{H_w^{k+1}(0,1)}^2 + \|\tilde{g}_2\|_{H_w^{k+1}(0,1)}^2 + \|D_\theta \tilde{g}_2\|_{H_w^k(0,1)}^2) d\theta \right\}^{\frac{1}{2}} \\
&\quad + C_2 N_\theta^{-l} \left\{ \int_0^{2\pi} (\|D_\theta^{l+1} \tilde{g}_2\|_{L_w^2(0,1)}^2 + \|D_\theta^l \tilde{g}_2\|_{H_w^1(0,1)}^2) d\theta \right\}^{\frac{1}{2}}, \quad (2.6.14)
\end{aligned}$$

which is (2.6.11). The proof of the  $L_w^2(R)$  norm bound (2.6.12) is very similar: its main ingredients are, in fact, contained in the argument above. Therefore, for the sake of brevity, the details are omitted here.  $\square$

The bounds (2.6.11) and (2.6.12) can now be straightforwardly mapped from  $R$  to  $D_0$ . We define  $\mathcal{P}_N(D)$  as  $\mathcal{P}_N(R)$  mapped from  $R$  to  $D_0$  using the polar coordinate transformation (2.6.1), and we suppose that  $\hat{\psi} \in \mathcal{H}^{k+1, l+1}(D)$ , with  $k, l \geq 1$ , where

$$\begin{aligned}
\mathcal{H}^{k,l}(D) &:= \{g \in H_0^1(D) : \tilde{g} \in \tilde{H}_{w,0}^1(R) \text{ has a decomposition } \tilde{g}(r, \theta) = \tilde{g}_1(r) + r\tilde{g}_2(r, \theta), \\
&\quad \text{with } \tilde{g}_1 = \frac{1}{2\pi}(\tilde{g}, 1)_{L^2(0,2\pi)} \in H_w^k(0,1) \\
&\quad \text{and } \tilde{g}_2 \in H_{w,0}^{k,0}(R) \cap H_w^{k-1,1}(R) \cap H_w^{0,l}(R) \cap H_w^{1,l-1}(R)\},
\end{aligned}$$

equipped with the norm  $\|g\|_{\mathcal{H}^{k,l}(D)} := \left( \|g\|_{\mathcal{H}_r^k(D)}^2 + \|g\|_{\mathcal{H}_\theta^l(D)}^2 \right)^{\frac{1}{2}}$  where, for  $g \in \mathcal{H}^{k,l}(D)$  with  $\tilde{g}(r, \theta) = \tilde{g}_1(r) + r\tilde{g}_2(r, \theta)$ ,

$$\begin{aligned} \|g\|_{\mathcal{H}_r^k(D)} &:= \left( \|\tilde{g}_1\|_{\mathbb{H}_w^k(0,1)}^2 + \|\tilde{g}_2\|_{\mathbb{H}_w^{k,0}(R)}^2 + \|\tilde{g}_2\|_{\mathbb{H}_w^{k-1,1}(R)}^2 \right)^{\frac{1}{2}}, \\ \|g\|_{\mathcal{H}_\theta^l(D)} &:= \left( \|\tilde{g}_2\|_{\mathbb{H}_w^{0,l}(R)}^2 + \|\tilde{g}_2\|_{\mathbb{H}_w^{1,l-1}(R)}^2 \right)^{\frac{1}{2}}. \end{aligned}$$

We define

$$\hat{\Pi}_N : \mathcal{H}^{1,1}(D) \rightarrow \mathcal{P}_N(D) \quad \text{by} \quad (\hat{\Pi}_N g)(q_1, q_2) = (\tilde{\Pi}_N \tilde{g})(r, \theta), \quad g \in \mathcal{H}^{1,1}(D).$$

Thus, recalling (2.2.2) and noting that  $\mathcal{H}^{k,l}(D) \subset \mathbb{H}_0^1(D) \subset \mathbb{H}_0^1(D; M)$ ,  $k, l \geq 1$ , we deduce from (2.6.11) that

$$\|\hat{\psi} - \hat{\Pi}_N \hat{\psi}\|_{\mathbb{H}_0^1(D; M)} \leq C_1 N_r^{-k} \|\hat{\psi}\|_{\mathcal{H}_r^{k+1}(D)} + C_2 N_\theta^{-l} \|\hat{\psi}\|_{\mathcal{H}_\theta^{l+1}(D)} \quad (2.6.15)$$

for all  $\hat{\psi} \in \mathcal{H}^{k+1, l+1}(D)$ , with  $k, l \geq 1$ . Similarly, we obtain from (2.6.12) that

$$\|\hat{\psi} - \hat{\Pi}_N \hat{\psi}\|_{L^2(D)} \leq C_1 N_r^{-k} \|\hat{\psi}\|_{\mathcal{H}_r^k(D)} + C_2 N_\theta^{-l} \|\hat{\psi}\|_{\mathcal{H}_\theta^l(D)} \quad (2.6.16)$$

for all  $\hat{\psi} \in \mathcal{H}^{k,l}(D)$ , with  $k, l \geq 1$ .

## 2.7 Convergence analysis of the numerical method

In this section we use the two-dimensional approximation results derived in Section 2.6 to complete the convergence analysis of the fully-discrete numerical method (2.5.1), (2.5.2), based on the symmetrised form of the Fokker–Planck equation. We shall *assume* as much regularity as is needed in order to establish an optimal-order bound on the discretisation error. At the end of the section we shall comment on the extension of our results to a fully-discrete method that stems from the alternative semidiscretisation (2.4.3) in the case of the FENE model.

We see from (2.5.7) that in order to obtain bounds on the norms of  $\xi$  appearing on the left-hand side of (2.5.7) we need to bound the following terms:

$$\|\eta^0\|, \quad \|\eta\|_{\ell^2(0,T; \mathbb{H}_0^1(D; M))} \quad \text{and} \quad \left\| \frac{\partial \eta}{\partial t} \right\|_{L^2(0,T; L^2(D))}.$$

It follows from (2.6.16), (2.6.15) and the definition of  $\eta := \hat{\psi} - \hat{\Pi}_N \hat{\psi}$  that

$$\begin{aligned} \|\eta^0\| &\leq \|\hat{\psi}^0 - \hat{\Pi}_N \hat{\psi}^0\| \leq C_1 N_r^{-k} \|\hat{\psi}^0\|_{\mathcal{H}_r^k(D)} + C_2 N_\theta^{-l} \|\hat{\psi}^0\|_{\mathcal{H}_\theta^l(D)}, \\ \|\eta\|_{\ell^2(0,T; \mathbb{H}_0^1(D; M))} &\leq C_1 N_r^{-k} \|\hat{\psi}\|_{\ell^2(0,T; \mathcal{H}_r^{k+1}(D))} + C_2 N_\theta^{-l} \|\hat{\psi}\|_{\ell^2(0,T; \mathcal{H}_\theta^{l+1}(D))}, \\ \left\| \frac{\partial \eta}{\partial t} \right\|_{L^2(0,T; L^2(D))} &\leq C_1 N_r^{-k} \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0,T; \mathcal{H}_r^k(D))} + C_2 N_\theta^{-l} \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0,T; \mathcal{H}_\theta^l(D))}, \end{aligned}$$

with  $k, l \geq 1$ , provided that  $\hat{\psi}$  is such that the right-hand sides of these inequalities are finite. Substituting these three bounds into the right-hand side of (2.5.7) we deduce, with  $m\Delta t \leq T$ ,  $m = 0, 1, \dots, N_T$ , that

$$\begin{aligned} & \|\xi\|_{\ell^\infty(0,T;L^2(D))} + \|\nabla_M \xi\|_{\ell^2(0,T;L^2(D))} \\ & \leq C_1 N_r^{-k} \left( \|\hat{\psi}^0\|_{\mathcal{H}_r^k(D)} + \|\hat{\psi}\|_{\ell^2(0,T;\mathcal{H}_r^{k+1}(D))} + \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0,T;\mathcal{H}_r^k(D))} \right) \\ & \quad + C_2 N_\theta^{-l} \left( \|\hat{\psi}^0\|_{\mathcal{H}_\theta^l(D)} + \|\hat{\psi}\|_{\ell^2(0,T;\mathcal{H}_\theta^{l+1}(D))} + \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0,T;\mathcal{H}_\theta^l(D))} \right) \\ & \quad + C_3 \Delta t \left\| \frac{\partial^2 \hat{\psi}}{\partial t^2} \right\|_{L^2(0,T;L^2(D))}. \end{aligned} \quad (2.7.1)$$

Note, also, that

$$\|\eta\|_{\ell^\infty(0,T;L^2(D))} \leq C_1 N_r^{-k} \|\hat{\psi}\|_{\ell^\infty(0,T;\mathcal{H}_r^k(D))} + C_2 N_\theta^{-l} \|\hat{\psi}\|_{\ell^\infty(0,T;\mathcal{H}_\theta^l(D))}, \quad (2.7.2)$$

$$\|\nabla_M \eta\|_{\ell^2(0,T;L^2(D))} \leq C_1 N_r^{-k} \|\hat{\psi}\|_{\ell^2(0,T;\mathcal{H}_r^{k+1}(D))} + C_2 N_\theta^{-l} \|\hat{\psi}\|_{\ell^2(0,T;\mathcal{H}_\theta^{l+1}(D))}. \quad (2.7.3)$$

Now, by the triangle inequality,

$$\begin{aligned} & \|\hat{\psi} - \hat{\psi}_N\|_{\ell^\infty(0,T;L^2(D))} + \|\nabla_M(\hat{\psi} - \hat{\psi}_N)\|_{\ell^2(0,T;L^2(D))} \\ & \leq \|\xi\|_{\ell^\infty(0,T;L^2(D))} + \|\nabla_M \xi\|_{\ell^2(0,T;L^2(D))} \\ & \quad + \|\eta\|_{\ell^\infty(0,T;L^2(D))} + \|\nabla_M \eta\|_{\ell^2(0,T;L^2(D))}, \end{aligned}$$

whereby (2.7.1), (2.7.2) and (2.7.3) give

$$\begin{aligned} & \|\hat{\psi} - \hat{\psi}_N\|_{\ell^\infty(0,T;L^2(D))} + \|\nabla_M(\hat{\psi} - \hat{\psi}_N)\|_{\ell^2(0,T;L^2(D))} \\ & \leq C_1 N_r^{-k} \left( \|\hat{\psi}\|_{\ell^\infty(0,T;\mathcal{H}_r^k(D))} + \|\hat{\psi}\|_{\ell^2(0,T;\mathcal{H}_r^{k+1}(D))} + \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0,T;\mathcal{H}_r^k(D))} \right) \\ & \quad + C_2 N_\theta^{-l} \left( \|\hat{\psi}\|_{\ell^\infty(0,T;\mathcal{H}_\theta^l(D))} + \|\hat{\psi}\|_{\ell^2(0,T;\mathcal{H}_\theta^{l+1}(D))} + \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0,T;\mathcal{H}_\theta^l(D))} \right) \\ & \quad + C_3 \Delta t \left\| \frac{\partial^2 \hat{\psi}}{\partial t^2} \right\|_{L^2(0,T;L^2(D))}. \end{aligned}$$

We recall that  $\psi = \sqrt{M} \hat{\psi}$ , and we define  $\psi_N^R := \sqrt{M} \hat{\psi}_N^R$ . Consequently,

$$\begin{aligned} & \|\psi - \psi_N\|_{\ell^\infty(0,T;\mathfrak{H})} + \|\psi - \psi_N\|_{\ell^2(0,T;\mathfrak{K})} \\ & \leq C_1 N_r^{-k} \left( \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^\infty(0,T;\mathcal{H}_r^k(D))} + \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^2(0,T;\mathcal{H}_r^{k+1}(D))} + \left\| \frac{1}{\sqrt{M}} \frac{\partial \psi}{\partial t} \right\|_{L^2(0,T;\mathcal{H}_r^k(D))} \right) \\ & \quad + C_2 N_\theta^{-l} \left( \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^\infty(0,T;\mathcal{H}_\theta^l(D))} + \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^2(0,T;\mathcal{H}_\theta^{l+1}(D))} + \left\| \frac{1}{\sqrt{M}} \frac{\partial \psi}{\partial t} \right\|_{L^2(0,T;\mathcal{H}_\theta^l(D))} \right) \\ & \quad + C_3 \Delta t \left\| \frac{1}{\sqrt{M}} \frac{\partial^2 \psi}{\partial t^2} \right\|_{L^2(0,T;L^2(D))}, \end{aligned} \quad (2.7.4)$$

with  $k, l \geq 1$ , provided that  $\psi$  is such that right-hand side is finite.

That completes the convergence analysis of the method in the case of  $d = 2$ . For  $d = 3$  the argument is identical, and rests on a three-dimensional analogue of Lemma 2.6.2; this is discussed further in Section 2.8.3.

Starting from the second stability inequality stated in Lemma 2.3.6 and proceeding in an identical manner as above, one can derive analogous error bounds in the  $h^1(0, T; \mathfrak{H})$  and  $\ell^\infty(0, T; \mathfrak{K})$  norms.

**Remark 2.7.1** In the case of the FENE Maxwellian,  $\sqrt{M} \in \mathcal{P}_N(D)$  if, and only if, there exists a positive integer  $m$  such that  $b = 4m$  and  $N_r \geq 2m$ . In order to ensure that, more generally,  $\sqrt{M} \in \mathcal{P}_N(D)$  regardless of the specific choice of  $b$  and the value of  $N_r$ , one can simply enrich  $\mathcal{P}_N(D)$  by adding  $\sqrt{M}$  as an extra basis function. However, in general the polynomials in  $\mathcal{P}_N(D)$  approximate  $\sqrt{M}$  very closely, so this leads to a highly ill-conditioned basis. A better solution is to add the component of  $\sqrt{M}$  orthogonal to  $\mathcal{P}_N(D)$  (in the  $L^2(D)$  inner product, for example,) to the basis, rather than  $\sqrt{M}$  itself. This is implemented in Section 2.8 for a numerical example in which  $b$  is not divisible by 4 and is shown to work well in that case.  $\diamond$

**Remark 2.7.2** We make a second comment regarding the FENE model. Starting from the variant of the inequality (2.5.7) alluded to in Remark 2.4.1 in connection with the fully-discrete spectral method based on the semidiscretisation (2.4.3) with  $b \geq 4s^2/(2s - 1)$  and  $s > 1/2$ , one can derive an optimal-order error bound analogous to (2.7.4). The core of the argument is identical to the one above, and is therefore omitted.  $\diamond$

## 2.8 Numerical results

Numerical methods for solving the Fokker–Planck equation arising from the FENE dumbbell model for dilute polymeric fluids have been the focus of some attention recently; Du *et al.* [42] developed a finite difference scheme that preserved the unit integral property and the positivity of  $\psi$ , Chauvière & Lozinski [32, 33, 90, 91] developed a spectral method for this problem and Ammar *et al.* [3, 4] proposed a reduced-basis method for solving the Fokker–Planck equation for FENE dumbbell chains. For a survey of, alternative, stochastic techniques for the numerical simulation of polymeric liquids we refer to the monograph of Öttinger [101], the article of Jourdain, Lelièvre, and Le Bris [61] or the survey paper [80], for example. The computational results we present in this section are for the FENE potential only, although it would be straightforward to modify the numerical methods to apply to more general potentials that satisfy Hyptheses A and B.

In Section 2.8.1, we discuss the implementation of two spectral Galerkin methods for the case of  $d = 2$  based on the formulation (2.5.1), (2.5.2). We then present computational results for these schemes in order to illustrate their behaviour in practice, as well as to provide experimental support for the convergence theory developed in Section 2.7. Next, we compare the two spectral Galerkin methods based on the formulation (2.5.1), (2.5.2) with the method of Chauvière & Lozinski based on the ‘original’ form (2.1.4) of the Fokker–Planck equation (or, more precisely, its transformed version (2.4.1) resulting from the substitution (2.8.10), with  $s = 2$ ). Section 2.8.1 is concluded with a discussion of the convergence rate of the extra-stress tensor,  $\underline{\underline{\tau}}$ .

In Section 2.8.2, we present some numerical results for a semi-implicit temporal discretisation of the Fokker–Planck equation in order to compare its performance with the backward Euler scheme that has been discussed in this section. Finally, we consider the implementation of our spectral Galerkin method in three spatial dimensions in Section 2.8.3, and we show some computational results to demonstrate that the 3-dimensional scheme exhibits essentially the same behaviour as the schemes considered in the case of  $d = 2$  in Section 2.8.1.

### 2.8.1 Numerical methods in the two-dimensional case

With  $D := B(0, \sqrt{b}) \subset \mathbb{R}^2$ , we suppose that  $\hat{\psi} \in H_0^1(D)$  and hence,  $\tilde{\psi} \in \tilde{H}_{\tilde{w},0}^1(R)$ , where  $\tilde{\psi}(r, \theta) := \hat{\psi}(q_1, q_2)$  with  $q_1 = \sqrt{b} r \cos \theta$ ,  $q_2 = \sqrt{b} r \sin \theta$ . Using the decomposition (2.6.4),  $\tilde{\psi}$  can be written in polar coordinates as follows:

$$\tilde{\psi}(r, \theta) = \tilde{\psi}_1(r) + r\tilde{\psi}_2(r, \theta), \quad (r, \theta) \in R = (0, 1) \times (0, 2\pi), \quad (2.8.1)$$

where, as in Section 2.6,  $r$  has been scaled from  $(0, \sqrt{b})$  to  $(0, 1)$ , and  $\tilde{\psi}_1 := \frac{1}{2\pi}(\tilde{\psi}, 1)_{L^2(0, 2\pi)}$ . In the context of spectral methods in polar coordinates, (2.8.1) is referred to by Shen as the *essential pole condition* [114]. This condition is a ‘first-order’ form of the following full pole-condition [44]: in order that a function  $\tilde{\psi}$ , defined by

$$\tilde{\psi}(r, \theta) = \sum_{n \in \mathbb{Z}} \tilde{\gamma}_n(r) \tilde{E}_n(\theta), \quad \text{where } \tilde{E}_n(\theta) := \frac{1}{\sqrt{2\pi}} \exp(in\theta),$$

is infinitely differentiable when transformed from polar to cartesian coordinates, it is necessary that, for each  $n \in \mathbb{Z} \setminus \{0\}$ ,

$$\tilde{\gamma}_n(r) = \mathcal{O}(r^{|n|}) \quad \text{as } r \rightarrow 0_+. \quad (2.8.2)$$

That (2.8.1) is a ‘first-order’ form of the full pole condition is easily seen by writing  $\tilde{\gamma}_n(r) = r^{|n|} \tilde{G}_n(r)$ , with  $\tilde{G}_n(r) = \mathcal{O}(1)$  as  $r \rightarrow 0_+$ ; hence,

$$\tilde{\psi}(r, \theta) = \frac{1}{\sqrt{2\pi}} \tilde{\gamma}_0(r) + r \sum_{n \in \mathbb{Z} \setminus \{0\}} r^{|n|-1} \tilde{G}_n(r) \tilde{E}_n(\theta) =: \tilde{\psi}_1(r) + r\tilde{\psi}_2(r, \theta),$$

with  $\tilde{\psi}_1(r) = \tilde{\gamma}_0(r)/\sqrt{2\pi} = \frac{1}{2\pi}(\tilde{\psi}, 1)_{L^2(0, 2\pi)}$ , as required.

The full pole condition (2.8.2) is consistent with the result established in the proof of Lemma 2.6.2 stating that the expansion coefficients  $\tilde{\gamma}_n$ ,  $n \in \mathbb{Z} \setminus \{0\}$ , of a function in  $\tilde{H}_{\tilde{w},0}^1(R)$  satisfy  $\tilde{\gamma}_n(r) = o(1)$  as  $r \rightarrow 0_+$ , although the conditions (2.8.2) are clearly much more restrictive.

In order to fit into the framework of the numerical analysis in Sections 2.6 and 2.7, each element of  $\mathcal{P}_N(R)$  should satisfy (2.8.1) to ensure that  $\mathcal{P}_N(D)$  is contained in  $H_0^1(D)$ . The discrete space  $\mathcal{P}_N(R)$ , introduced in Section 2.6, satisfies this property. In this section we define a spectral Galerkin method for the Fokker–Planck equation based on a particular basis (denoted  $\mathcal{A}$ ) for  $\mathcal{P}_N(R)$  that satisfies the same decomposition.

For the purpose of comparison, we also introduce a second basis,  $\mathcal{B}$ , in which each function satisfies the full pole condition, (2.8.2). Thus, on mapping  $\mathcal{B}$  from  $R$  to  $D$  we obtain a basis for a finite-dimensional subspace of  $C^\infty(\bar{D}) \cap C_0(\bar{D}) \subset H_0^1(D)$ . The reason for considering this second basis is that typical solutions of the FENE Fokker–Planck equation are smooth on  $D$ , and therefore it is likely that in practice a Galerkin method based on  $\mathcal{B}$  will be more accurate

than a method based on  $\mathcal{A}$ : mapping the basis  $\mathcal{A}$  from  $R$  to  $D$  yields a finite-dimensional subspace of  $H_0^1(D)$  only, which contains functions that are not smooth at the origin in  $D$ . We note, however, that the span of  $\mathcal{B}$  does not coincide with  $\mathcal{P}_N(R)$ , and therefore the approximation properties of  $\mathcal{B}$  are not covered by the results in Section 2.6 that led to the error bounds in Section 2.7. Hence, the numerical results for basis  $\mathcal{A}$  are intended to verify the analysis developed in the previous sections, while basis  $\mathcal{B}$  is introduced to indicate the gain in performance that can be obtained by satisfying (2.8.2). By requiring more regularity from the basis than it being a finite-dimensional subspace of  $H_0^1(D)$  one could modify the arguments in Section 2.6 to derive convergence estimates based on a pole condition of higher order than (2.6.4), but this would make the derivation of the approximation results more laborious (*e.g.*, the projector  $\tilde{P}_N^J$  would have to obey (2.8.2) rather than (2.8.1)). Before introducing bases  $\mathcal{A}$  and  $\mathcal{B}$ , we make the following observation.

**Remark 2.8.1** Let  $\hat{\psi}$  be the weak solution of (2.1.6) corresponding to a given initial condition  $\hat{\psi}^0$ , define  $\hat{\psi}^*(\underline{q}, t) := \hat{\psi}(-\underline{q}, t)$  and suppose that  $\hat{\psi}^0$  is invariant under the change of independent variable  $\underline{q} \mapsto -\underline{q}$ , *i.e.*,  $\hat{\psi}^0(\underline{q}) = \hat{\psi}^0(-\underline{q})$  for a.e.  $\underline{q} \in D$ . On noting that  $M(\underline{q}) = M(-\underline{q})$ ,  $\underline{q} \in D$ , it follows that the weak formulation (2.1.6) is also invariant under this change of variable; hence  $\hat{\psi}$  and  $\hat{\psi}^*$  are weak solutions to the same initial-boundary-value problem. It follows by uniqueness of the weak solution established in Section 2.3 that  $\hat{\psi}(\underline{q}, t) \equiv \hat{\psi}^*(\underline{q}, t)$ , *i.e.*,  $\hat{\psi}(\underline{q}, t) = \hat{\psi}(-\underline{q}, t)$  for a.e.  $\underline{q} \in D$  and a.e.  $t \in [0, T]$ . This evenness of  $\hat{\psi}$  in the  $D$  domain with respect to  $\underline{q}$  translates into  $\pi$ -periodicity of  $\tilde{\psi}$  in the  $R$  domain with respect to  $\theta$ . An identical statement applies to the numerical solution  $(\hat{\psi}_N^n)_{n=0}^{N_T}$  defined by (2.5.1), (2.5.2), provided  $\mathcal{P}_N(D) \subset H_0^1(D)$  is such that whenever a function  $\underline{q} \mapsto v(\underline{q})$  belongs to  $\mathcal{P}_N(D)$  its even reflection  $\underline{q} \mapsto v(-\underline{q})$  also belongs to  $\mathcal{P}_N(D)$ : if  $\hat{\psi}^0(\underline{q}) = \hat{\psi}^0(-\underline{q})$  for a.e.  $\underline{q} \in D$ , uniqueness of the  $L^2(D)$  projection of  $\hat{\psi}^0$  onto  $\mathcal{P}_N(D)$  implies that  $\hat{\psi}_N^0(\underline{q}) = \hat{\psi}_N^0(-\underline{q})$  for a.e.  $\underline{q} \in D$ . Uniqueness of the numerical solution then yields  $\hat{\psi}_N^n(\underline{q}) = \hat{\psi}_N^n(-\underline{q})$  for a.e.  $\underline{q} \in D$  and all  $n = 0, \dots, N_T$ .  $\diamond$

The above remark demonstrates that (2.1.6) captures an important symmetry property of the dumbbell model for polymeric fluids: the configuration probability density function  $\psi$  is required to be symmetric about the origin in  $D$  because the beads of a dumbbell are indistinguishable. As long as  $\hat{\psi}^0$  and  $\mathcal{P}_N(D)$  are invariant under the change of independent variable  $\underline{q} \mapsto -\underline{q}$  described in Remark 2.8.1, the numerical solution will inherit the symmetry of the analytical solution implied by the symmetry of the initial condition. A consequence of this observation is that we should require the basis functions in  $\mathcal{A}$  and  $\mathcal{B}$  to obey the same symmetry condition; following [92], this is achieved in the definitions below by only including even trigonometric modes in  $\theta$ . Strictly speaking therefore  $\mathcal{A}$  is chosen to be a basis for the linear subspace of  $\mathcal{P}_N(R)$  consisting of all  $\pi$ -periodic functions. Note, however, that if the solution were  $2\pi$ -periodic, then one could simply include odd trigonometric modes as well. We are now ready to define the bases  $\mathcal{A}$  and  $\mathcal{B}$ .

**Basis  $\mathcal{A}$ :** Let  $\mathcal{A} := \mathcal{A}_1 \cup \mathcal{A}_2$  where:

$$\begin{aligned} \mathcal{A}_1 &:= \{(1-r)P_k(r) : k = 0, \dots, N_r - 1\}, \\ \mathcal{A}_2 &:= \{r(1-r)P_k(r)\Phi_{il}(\theta) : k = 0, \dots, N_r - 1; \quad i = 0, 1; \quad l = 1, \dots, N_\theta\}. \end{aligned}$$

$P_k$  is a polynomial of degree  $k$  in  $r \in [0, 1]$  and  $\Phi_{il}(\theta) = (1 - i) \cos(2l\theta) + i \sin(2l\theta)$ ,  $\theta \in [0, \pi]$ . We denote by  $P_k$  the  $k$ th Chebyshev polynomial scaled from  $[-1, 1]$  to  $[0, 1]$ . The numerical method is not particularly sensitive to this choice of polynomial, however, and other choices work well also. Notice that the polynomials in  $\mathcal{A}_1$  and  $\mathcal{A}_2$  both contain the factor  $(1 - r)$  in order to impose the homogeneous Dirichlet boundary condition on  $\partial D$ , and functions in  $\mathcal{A}_2$  contain an extra factor of  $r$  to enforce the essential pole condition. Basis  $\mathcal{A}$  is chosen so as to mimic the decomposition (2.8.1) of the analytical solution  $\tilde{\psi} \in \tilde{\mathbb{H}}_{\tilde{w},0}^1(R)$  in polar coordinates: the role of  $\text{span}(\mathcal{A}_1)$  is to approximate  $\tilde{\psi}_1$  while  $\text{span}(\mathcal{A}_2)$  is meant to approximate  $r\tilde{\psi}_2$ .

**Basis  $\mathcal{B}$ :** This is, effectively, the basis proposed by Matsushima and Marcus [97] and Verkley [121], except that, as above, we ensure that the functions are zero at  $r = 1$  and that they are  $\pi$ -periodic in  $\theta$ :

$$\mathcal{B} = \{W_{lk}(r)\Phi_{il}(\theta) : k = 0, \dots, N_r - 1; i = 0, 1; l = i, \dots, N_\theta\}, \quad (2.8.3)$$

where  $W_{lk}(r) = r^{2l}(1 - r^2)J_k^{(0,2l)}(2r^2 - 1)$  and  $J_k^{(\alpha,\beta)}(x)$  is the Jacobi polynomial on  $[-1, 1]$  of degree  $k$  with respect to the weight  $(1 - x)^\alpha(1 + x)^\beta$  ( $\Phi_{il}$  is the same as in  $\mathcal{A}$ ). Each element of  $\mathcal{B}$  satisfies (2.8.2).

$\mathcal{A}$  and  $\mathcal{B}$  both have cardinality  $N := N_r(2N_\theta + 1)$ . Expressing trial and test functions in terms of either  $\mathcal{A}$  or  $\mathcal{B}$ , it is now straightforward to determine the discretisation matrices corresponding to the integrals

$$\int_D \hat{\psi}_N^{n+1} \hat{\varphi} \, dq, \quad \int_D \nabla_M \hat{\psi}_N^{n+1} \cdot \nabla_M \hat{\varphi} \, dq, \quad \int_D (\underline{\kappa}^{n+1} \underline{q} \hat{\psi}_N^{n+1}) \cdot \nabla_M \hat{\varphi} \, dq \quad (2.8.4)$$

from (2.5.1). We label these matrices  $\mathbf{M}$ ,  $\mathbf{S}$  and  $\mathbf{C}^{n+1}$  for mass, stiffness and convection respectively.

Using the ansatz  $\tilde{\psi}_N^{n+1}(r, \theta) = \sum_{v=1}^N \tilde{\Psi}_v^{n+1} Y_v(r, \theta)$  for trial functions, where  $Y_v$  is a basis function (from either  $\mathcal{A}$  or  $\mathcal{B}$ ) for  $1 \leq v \leq N$ , denoting test functions as  $Y_u$  for  $1 \leq u \leq N$  and mapping (2.8.4) from  $D$  to  $R$  yields:

$$\mathbf{M}_{uv} = \int_0^1 \int_0^\pi b r Y_v(r, \theta) Y_u(r, \theta) \, dr \, d\theta, \quad (2.8.5)$$

$$\begin{aligned} \mathbf{S}_{uv} = \int_0^1 \int_0^\pi \left\{ r \frac{\partial Y_v}{\partial r} \frac{\partial Y_u}{\partial r} + \frac{1}{r} \frac{\partial Y_v}{\partial \theta} \frac{\partial Y_u}{\partial \theta} \right. \\ \left. + \frac{b}{2} \frac{r^2}{1 - r^2} \frac{\partial}{\partial r} (Y_u Y_v) + \frac{b^2}{4} \frac{r^3}{(1 - r^2)^2} Y_v Y_u \right\} \, dr \, d\theta, \end{aligned} \quad (2.8.6)$$

$$\begin{aligned} \mathbf{C}_{uv}^{n+1} = \int_0^1 \int_0^\pi b r Y_v \frac{\partial Y_u}{\partial \theta} (-\kappa_{11}^{n+1} \sin 2\theta - \kappa_{12}^{n+1} \sin^2 \theta + \kappa_{21} \cos^2 \theta) \, dr \, d\theta \\ + \int_0^1 \int_0^\pi \left( b r^2 Y_v \frac{\partial Y_u}{\partial r} + \frac{b^2}{2} \frac{r^3}{1 - r^2} Y_v Y_u \right) \times \\ \left( \kappa_{11}^{n+1} \cos 2\theta + \frac{1}{2} (\kappa_{12}^{n+1} + \kappa_{21}^{n+1}) \sin 2\theta \right) \, dr \, d\theta. \end{aligned} \quad (2.8.7)$$

Note that if the  $Y_u, Y_v$  do not satisfy (2.8.1), then the entries of  $\mathbf{S}$  may be undefined.

With these discretisation matrices in hand, the numerical solution is computed by solving the following linear system for the coefficient vector  $\tilde{\Psi}^{n+1} := (\tilde{\Psi}_1^{n+1}, \dots, \tilde{\Psi}_N^{n+1})^\top \in \mathbb{R}^N$ ,

$n = 0, 1, \dots, N_T - 1$ :

$$\left( \mathbf{M} + \Delta t \left( \frac{1}{2\text{Wi}} \mathbf{S} - \mathbf{C}^{n+1} \right) \right) \tilde{\Psi}^{n+1} = \mathbf{M} \tilde{\Psi}^n, \quad (2.8.8)$$

with  $\tilde{\Psi}^0$  defined by the initial datum. Then, the numerical approximation to the probability density function itself is obtained as  $\psi_N^{n+1}(q) = \sqrt{M(q)} \tilde{\psi}_N^{n+1}(r, \theta)$ , where  $r = |q|/\sqrt{b}$  and  $\tilde{\psi}_N^{n+1}(r, \theta) = \sum_{v=1}^N \tilde{\Psi}_v^{n+1} Y_v(r, \theta)$ .

For ease of evaluation, the integrals in (2.8.5), (2.8.6) and (2.8.7) can be factorised into products of 1-dimensional integrals over  $r$  and  $\theta$ . We evaluate the  $\theta$ -integrals exactly using trigonometric identities, and, noting that the  $r$ -integrands are all polynomials, we use Gauss quadrature to evaluate the  $r$ -integrals to machine precision.  $\mathbf{M}$  and  $\mathbf{S}$  are constant matrices, which can be pre-computed and reused, but if  $\underline{\kappa}$  is time-varying, we must reassemble  $\mathbf{C}^{n+1}$  at every time-step. However, it is straightforward to factor out the dependence of  $\mathbf{C}^{n+1}$  on  $\underline{\kappa}$  so that the integrals that determine  $\mathbf{C}^{n+1}$  need not be evaluated more than once. We use LU-decomposition to solve (2.8.8), which is appropriate because the spectral discretisation matrices are generally of moderate size.

We now present some numerical results. For simplicity, in the computations considered below we always use the normalised Maxwellian (which satisfies the symmetry property required in Remark 2.8.1 and also has unit volume) as the initial condition, so that  $\hat{\psi}^0(q) = \sqrt{M(q)}$ . Also, most of the results presented in this section are for computations in which  $b$  was chosen to be divisible by 4 so that the spaces  $\text{span}(\mathcal{A})$  and  $\text{span}(\mathcal{B})$  naturally contain  $\sqrt{M}$ , as in Remark 2.7.1. However, the basis enrichment technique described in Remark 2.7.1 was implemented to obtain the results in Table 2.3 (in which  $b = 10$ ) and, as discussed below, it worked well for that problem.

Henceforth, the two numerical methods that use basis  $\mathcal{A}$  and basis  $\mathcal{B}$ , respectively, will be referred to as method  $\mathcal{A}$  and method  $\mathcal{B}$ .

First of all we present results from solving the Fokker–Planck equation with parameters  $b = 16$ ,  $\text{Wi} = 1.2$  and  $\kappa_{11} = -\kappa_{22} = 1.1$ ,  $\kappa_{12} = 0.9$ ,  $\kappa_{21} = -0.6$  and with  $\Delta t = 0.05$ . These parameters were chosen somewhat arbitrarily, but the intention here is to visualise a typical evolution of  $\psi_N$  towards steady state, and to provide an initial qualitative comparison of methods  $\mathcal{A}$  and  $\mathcal{B}$  (quantitative convergence results will be presented below). By taking  $(N_r, N_\theta) = (26, 20)$  with basis  $\mathcal{A}$  and  $(N_r, N_\theta) = (21, 15)$  with basis  $\mathcal{B}$ , the solutions from the two methods were indistinguishable to the eye and appear to be fully resolved. As foreshadowed above,  $\mathcal{A}$  required more degrees-of-freedom than  $\mathcal{B}$  to resolve the solution to comparable accuracy in this case because, as can be seen in Figure 2.1,  $\psi_N$  is smooth at the origin in cartesian coordinates whereas the basis functions in  $\mathcal{A}$  are not necessarily smooth there. Nevertheless, a clear advantage of basis  $\mathcal{A}$  over basis  $\mathcal{B}$  is that it is built by relying on the essential pole condition only, as manifested by the decomposition in Lemma 2.6.2, which only requires the most basic smoothness hypothesis, that  $\tilde{\psi} \in \tilde{H}_{\tilde{w},0}^1(R)$  (implied by the assumption that the weak solution  $\hat{\psi} \in H_0^1(D; M)$  belongs to  $H_0^1(D)$ ).

Figure 2.1 shows snapshots of  $\psi_N$  at  $t = 0$ ,  $t = 1$ ,  $t = 2$  and  $t = 3$ , and  $\psi_N$  is close to steady state at  $t = 3$ .

To provide a quantitative study of the spatial accuracy of the numerical methods defined in this section, we use the fact that when  $\underline{\kappa}$  is a symmetric tensor the exact steady-state solution of the Fokker–Planck equation (2.1.1) with boundary condition (2.1.7), and unit

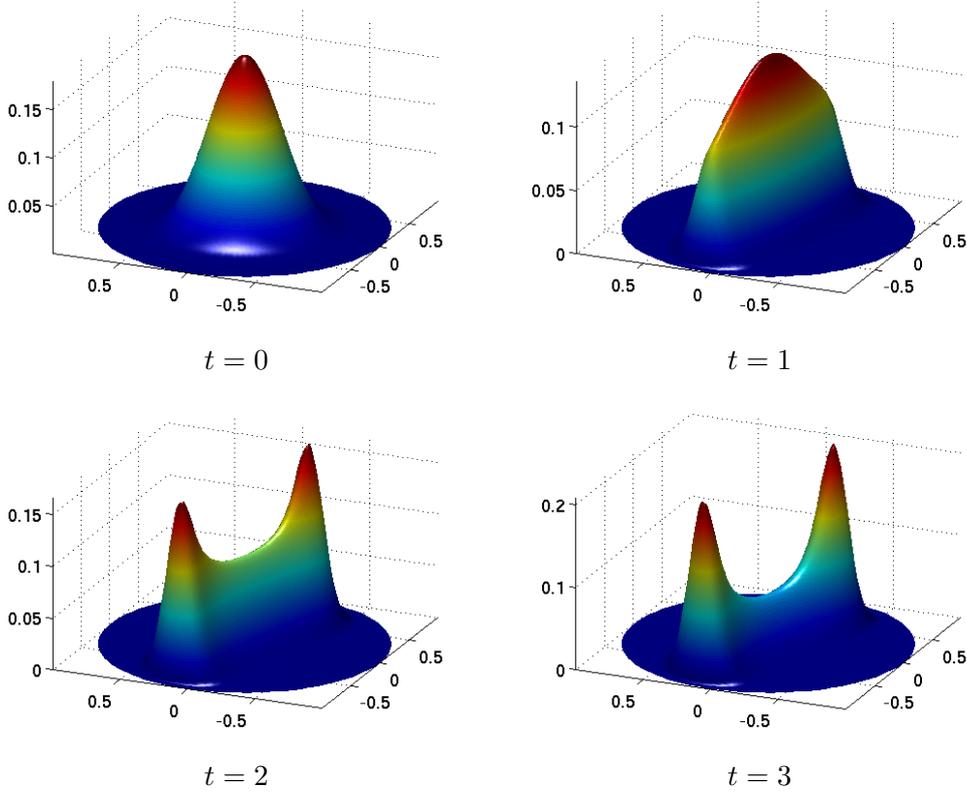


Figure 2.1: Snapshots of  $\psi_N$  at  $t = 0$ ,  $t = 1$ ,  $t = 2$  and  $t = 3$  illustrating evolution towards steady state. In this case, we have  $\Delta t = 0.05$ ,  $b = 16$ ,  $Wi = 1.2$  and  $\kappa_{11} = -\kappa_{22} = 1.1$ ,  $\kappa_{12} = 0.9$ ,  $\kappa_{21} = -0.6$ . This computation was performed using basis  $\mathcal{A}$  and basis  $\mathcal{B}$  with  $(N_r, N_\theta) = (26, 20)$  and  $(N_r, N_\theta) = (21, 15)$ , respectively. The solutions were fully resolved in each of these two cases.

volume, is given by

$$\psi_{\text{exact}}(q) := C M(q) \exp(Wi q^T \underline{\underline{\kappa}} q), \quad (2.8.9)$$

where  $C$  is a normalization constant chosen so that  $\int_D \psi_{\text{exact}}(q) dq = 1$ ; see, [23]. We now consider a particular case, referred to as *extensional flow*, in which  $\underline{\underline{\kappa}} = \text{diag}(\delta, -\delta)$ . This generally provides a good test case for numerical methods for the Fokker–Planck equation because it yields particularly sharp solution profiles that are challenging to resolve, and also the exact steady-state solution is available for comparison. In order to compare the convergence rates of methods  $\mathcal{A}$  and  $\mathcal{B}$ , we solved two distinct extensional flow problems for: (i)  $(b, Wi, \delta) = (12, 1, 1)$  and (ii)  $(b, Wi, \delta) = (20, 1, 2)$ , with a range of choices of  $(N_r, N_\theta)$ . In order to compare to the known exact steady-state solution, we took 2000 time-steps (with  $\Delta t = 0.05$  and  $T = 100$ ) in each case so that the final numerical solution is a very close approximation to the steady-state solution. This allows us to compare the spatial convergence rates of the two numerical methods without worrying about temporal discretisation error. Tables 2.1 and 2.2 show the relative errors (in the  $L^2(D)$  and  $H^1(D; M)$  norms) between the exact and the computed steady-state solutions for extensional flows (i) and (ii), respectively.

We can see from the data in the tables that methods  $\mathcal{A}$  and  $\mathcal{B}$  converge rapidly for

both problem (i) and problem (ii) and that for each choice of  $(N_r, N_\theta)$ , basis  $\mathcal{B}$  outperforms basis  $\mathcal{A}$  – again this is because the solution profiles are smooth at the origin in cartesian coordinates, see Figure 2.2. Nevertheless, the rapid convergence of method  $\mathcal{A}$  is consistent with the spectral error estimates established in Section 2.7 (recall that these error estimates do not apply to method  $\mathcal{B}$  because  $\text{span}(\mathcal{B})$  is not the same as  $\mathcal{P}_N(R)$  analysed in Section 2.6). It is also clear that problem (ii) is more challenging to resolve than problem (i); with both  $\mathcal{A}$  and  $\mathcal{B}$ , more basis functions are required to attain a given accuracy for problem (ii) than for problem (i). Note that the greater difficulty of resolving extensional flow (ii) is encoded in the convergence estimates in Section 2.7 because the constants in these estimates depend exponentially on  $b$ ,  $\delta$  (via  $\|\underline{\kappa}\|_{L^\infty(0,T)}$ ) and  $T$ . Moreover, the factor  $e^{2c_0m\Delta t}$  on the right-hand side in Lemma 2.3.1 permits exponential growth in time of the norm of  $\hat{\psi}_N$ , and this is reflected in the first row of Table 2.2 in which the solutions computed with  $(N_r, N_\theta) = (10, 10)$  for extensional flow (ii) resulted in numerical overflow.<sup>6</sup> Note that this overflow behaviour was only observed in the case of under-resolved computations that led to numerical solutions containing numerical oscillations *i.e.* it was not observed in rows 2, 3 and 4 of Table 2.2; note also that Chauvière & Lozinski’s method behaves in the same way for under-resolved solutions, as shown in Table 2.3.

$(N_r, N_\theta)$	Relative $L^2(D)$ error		Relative $H^1(D; M)$ error	
	Basis $\mathcal{A}$	Basis $\mathcal{B}$	Basis $\mathcal{A}$	Basis $\mathcal{B}$
(10,10)	$3.63 \times 10^{-2}$	$4.61 \times 10^{-3}$	$7.90 \times 10^{-2}$	$8.82 \times 10^{-3}$
(15,15)	$3.36 \times 10^{-3}$	$9.19 \times 10^{-6}$	$8.58 \times 10^{-3}$	$2.33 \times 10^{-5}$
(20,20)	$5.13 \times 10^{-5}$	$4.63 \times 10^{-9}$	$1.64 \times 10^{-4}$	$1.52 \times 10^{-8}$
(25,25)	$2.94 \times 10^{-7}$	$1.74 \times 10^{-12}$	$1.13 \times 10^{-6}$	$6.94 \times 10^{-12}$
(30,30)	$8.31 \times 10^{-10}$	$1.70 \times 10^{-13}$	$3.77 \times 10^{-9}$	$1.70 \times 10^{-13}$

Table 2.1: Relative errors in the  $L^2(D)$  and  $H^1(D; M)$  norms (*i.e.*  $\|\hat{\psi}_N - \hat{\psi}_{\text{exact}}\|/\|\hat{\psi}_{\text{exact}}\|$  and  $\|\hat{\psi}_N - \hat{\psi}_{\text{exact}}\|_{H^1(D;M)}/\|\hat{\psi}_{\text{exact}}\|_{H^1(D;M)}$ , respectively) for extensional flow (i) at steady-state, *i.e.*  $b = 12$ ,  $\text{Wi} = 1$  and  $\delta = 1$ .  $\hat{\psi}_N$  is an approximation to the steady-state solution obtained by taking 2000 time-steps with  $\Delta t = 0.05$ , and  $\hat{\psi}_{\text{exact}}$  is the exact steady-state solution, which is known in this case because  $\underline{\kappa}$  is symmetric.

The (fully resolved) solutions corresponding to extensional flow problems (i) and (ii) are shown in Figure 2.2, and in each case both  $\psi_N$  and  $\tilde{\psi}_N$  are plotted. It is clear that the solution profiles corresponding to (ii) are much more severe, and therefore it is not surprising that more modes were required in this case. The quantity of interest in these computations is  $\psi_N$ , but  $\tilde{\psi}_N$  is also plotted to emphasise the numerical difficulties that are encountered as  $b$  and  $\delta$  are increased. In the plots corresponding to (i), the peaks in  $\tilde{\psi}_N$  are higher than in  $\psi_N$ , but only by a factor of about 20. For (ii) on the other hand, the peaks in  $\tilde{\psi}_N$  are higher by a factor of roughly 5000. The causes of this behaviour are two-fold: with  $\delta = 2$  the flow has stronger extensional character and therefore the solution peaks are expected to be more concentrated and also, the larger value of  $b$  means that  $\sqrt{M}$  is more strongly degenerate near  $\partial D$  so that  $\hat{\psi}_N = \psi_N/\sqrt{M}$  takes larger values near the boundary. This second point can be

<sup>6</sup>When  $\underline{q}^T \underline{\kappa}(t) \underline{q} = 0$  for all  $t \in [0, T]$ , Lemma 2.3.1, with  $\mu = 0$  and  $\nu = 0$ , can be sharpened. The inequality holds with  $c_0 = 0$ , showing that the expression on the left-hand side of the inequality is bounded by  $\|\hat{\psi}^0\|^2$ , uniformly in  $T$ ,  $b$  and  $\|\underline{\kappa}\|_{L^\infty(0,T)}$ .

$(N_r, N_\theta)$	Relative $L^2(D)$ error		Relative $H^1(D; M)$ error	
	Basis $\mathcal{A}$	Basis $\mathcal{B}$	Basis $\mathcal{A}$	Basis $\mathcal{B}$
(10,10)	–	–	–	–
(15,15)	$2.47 \times 10^{-1}$	$9.57 \times 10^{-2}$	$1.79 \times 10^{-1}$	$9.53 \times 10^{-2}$
(20,20)	$3.91 \times 10^{-2}$	$1.72 \times 10^{-3}$	$4.88 \times 10^{-2}$	$2.54 \times 10^{-3}$
(25,25)	$9.07 \times 10^{-3}$	$1.71 \times 10^{-4}$	$9.77 \times 10^{-3}$	$2.37 \times 10^{-4}$
(30,30)	$1.50 \times 10^{-3}$	$2.97 \times 10^{-6}$	$2.61 \times 10^{-3}$	$4.49 \times 10^{-6}$
(35,35)	$3.37 \times 10^{-4}$	$2.14 \times 10^{-8}$	$5.60 \times 10^{-4}$	$3.66 \times 10^{-8}$
(40,40)	$2.54 \times 10^{-5}$	$5.97 \times 10^{-9}$	$4.55 \times 10^{-5}$	$5.94 \times 10^{-9}$

Table 2.2: Relative errors in the  $L^2(D)$  and  $H^1(D; M)$  norms for extensional flow (ii) at steady-state, *i.e.*  $(b, \text{Wi}, \delta) = (20, 1, 2)$ . The time-stepping strategy to compute the approximate steady-state solution was the same as in Table 2.1. The hyphens in the first row indicate that we obtained numerical overflow in those computations.

seen as a drawback, for  $b \gg 1$ , of the fully-discrete numerical method (2.5.1), (2.5.2), based on the symmetrised form of the Fokker–Planck equation. Presumably Chauvière & Lozinski [33] fixed their value of  $s$  ( $s = 2$  for  $d = 2$  and  $s = 2.5$  for  $d = 3$ ) in the transformation

$$\hat{\psi}(q) := \psi(q)/[M(q)]^{2s/b} = \psi(q)/(1 - |q|^2/b)^s \quad (2.8.10)$$

so as to avoid a similar effect; indeed, they presented some numerical results for  $b = 200$ . Values of  $b$  this large do not appear to be feasible with the fully-discrete method (2.5.1), (2.5.2), based on the substitution  $\hat{\psi}_N = \psi_N/\sqrt{M}$ .

As has been noted in Remark 2.5.2, there is in fact no difference between the stability properties of the method based on (2.5.1), (2.5.2) and of a Chauvière–Lozinski type method. However, if  $b \gg 1$ , for a typical  $\psi$  we have that  $\|\psi/\sqrt{M}\|_{L^\infty(D)} = \|\psi/(1 - |q|^2/b)^{b/4}\|_{L^\infty(D)} \gg \|\psi/(1 - |q|^2/b)^2\|_{L^\infty(D)}$ . Hence, compared to a Chauvière–Lozinski type method with the recommended choice of  $s = 2$  for  $d = 2$ , the maximum value of the numerical approximation  $\hat{\psi}_N$  to the function  $\hat{\psi}$  defined by the scheme (2.5.1), (2.5.2) can be much larger when  $b \gg 1$ , and can thereby require greater computational effort to resolve to a given accuracy. The computational results that we consider in this section are therefore restricted to moderate values of  $b$ . It has to be said, however, that when  $b \gg 1$  the FENE Maxwellian is very close to the Maxwellian of the Hookean model, uniformly in  $q$ ,  $|q| \leq \sqrt{b}$ ;<sup>7</sup> thus, instead of a FENE model with  $b \gg 1$ , one might as well use the, simpler, Hookean dumbbell model, which, as we shall see in Chapter 5, has an exact macroscopic closure (the Oldroyd-B) model. In the setting of the FENE model the practically relevant values of  $b$  are those of small to moderate size, and in this range the symmetrised method works well.

With these precursors, we now compare the accuracy of methods  $\mathcal{A}$  and  $\mathcal{B}$  to that of the spectral method of Chauvière & Lozinski discussed in [33]. In Table 2 of that paper, the authors presented convergence data for the (1,1)-component of the *polymeric extra-stress* tensor,  $\underline{\underline{\tau}} = (\tau_{ij})$ , computed for an extensional flow at steady state for the parameters

<sup>7</sup>On extending the FENE Maxwellian from  $B(0, \sqrt{b})$  to the whole of  $\mathbb{R}^d$  by 0 and denoting the resulting function by  $M_b$ , it is easily seen that  $M_b$  converges, as  $b \rightarrow \infty$ , to the Maxwellian of the Hookean model, uniformly on  $\mathbb{R}^d$ .

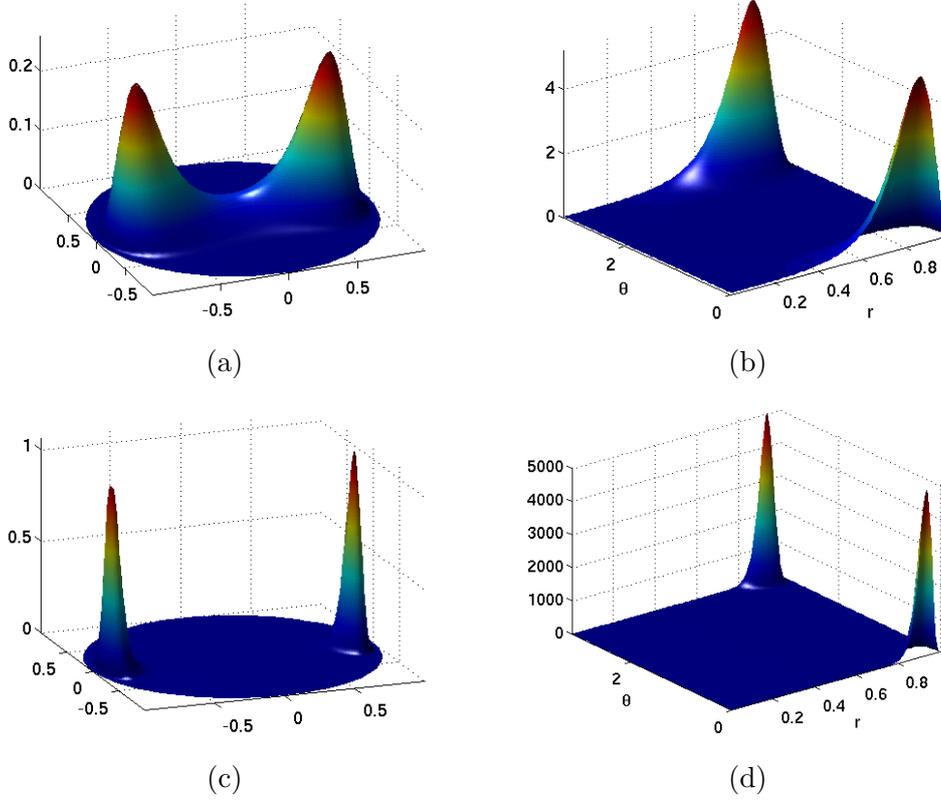


Figure 2.2: Numerical approximations to the steady state solution for extensional flow problems (i) and (ii) using  $(N_r, N_\theta) = (30, 30)$  and  $(N_r, N_\theta) = (40, 40)$ , respectively. Plots (a) and (b) show  $\psi_N$  and  $\psi_N$  respectively, at steady state for problem (i) and (c), (d) show  $\psi_N$  and  $\hat{\psi}_N$  for (ii). The purpose of plots (b) and (d) is to demonstrate that  $\hat{\psi}_N$  usually has a much steeper solution profile than  $\psi_N$  and this effect is amplified if either  $\delta$  or  $b$  (or both) are increased.

$(b, \lambda, \delta) = (10, 1, 5)$ . Note that when  $\psi$  is a function of  $q$  and  $t$  only,  $\underline{\tau}$  is defined as:

$$\underline{\tau}(t) := \int_D \underline{F} \otimes \underline{q} \psi(\underline{q}, t) d\underline{q} = \int_D \underline{F} \otimes \underline{q} \sqrt{M} \hat{\psi}(\underline{q}, t) d\underline{q}, \quad (2.8.11)$$

where  $\underline{F}$  is taken to be the FENE spring force here. Table 2.3 reproduces Chauvière & Lozinski's results and compares them to the corresponding results for methods  $\mathcal{A}$  and  $\mathcal{B}$ . Note that in this problem  $b$  is not divisible by 4. Therefore, in order to ensure that the volume of  $\psi_N$  is conserved with methods  $\mathcal{A}$  and  $\mathcal{B}$ , we added the component of  $\sqrt{M}$  orthogonal to  $\text{span}(\mathcal{A})$  (resp.  $\text{span}(\mathcal{B})$ ) to the bases to obtain an enriched discrete space that contains  $\sqrt{M}$  (cf. Remark 2.7.1).<sup>8</sup> This ensured that the volume of  $\psi_N$  was conserved to machine precision (except in the cases that rounding error polluted the results, these are indicated by hyphens in the table).

The data in Table 2.3 show that for this problem method  $\mathcal{B}$  converges at a comparable rate to the method of Chauvière & Lozinski, whereas  $\mathcal{A}$  appears to converge more slowly.

<sup>8</sup>Orthogonalisation was performed in the  $L^2(D)$  inner product.

Note that the reason why method  $\mathcal{B}$  and Chauvière & Lozinski’s method converge at a similar rate (at least in this case where  $b$  is relatively low) is that both methods involve ansatzes that impose extra regularity at the origin in cartesian coordinates; basis  $\mathcal{B}$  satisfies the pole condition (2.8.2), and Chauvière & Lozinski use a transformation that enforces  $\frac{\partial \psi}{\partial r} \Big|_{r=0} = 0$ , which, when combined with  $\pi$ -periodicity in  $\theta$ , has a similar effect.

$(N_r, N_\theta)$	Relative error of $\tau_{11}$		
	Basis $\mathcal{A}$	Basis $\mathcal{B}$	Chauvière & Lozinski
(11,5)	–	–	–
(13,6)	–	$4.8 \times 10^{-2}$	0.35
(21,10)	$1.8 \times 10^{-3}$	$2.0 \times 10^{-2}$	$2.0 \times 10^{-2}$
(31,15)	$2.1 \times 10^{-4}$	$1.4 \times 10^{-4}$	$1.4 \times 10^{-4}$
(41,20)	$1.3 \times 10^{-5}$	$8.7 \times 10^{-7}$	$2.1 \times 10^{-7}$

Table 2.3: Comparison of the relative errors in  $\tau_{11}$  for extensional flow with  $(b, \text{Wi}, \delta) = (10, 1, 5)$ . The three schemes compared are methods  $\mathcal{A}$  and  $\mathcal{B}$  and the spectral method of Chauvière & Lozinski. The data for the method of Chauvière & Lozinski is taken from Table 2 in [33].

In fact, as discussed in Section 1.3.3, in the context of deterministic multiscale computations for the micro-macro model, the primary reason for solving the Fokker–Planck equation is to obtain an approximation of  $\underline{\tau}$ . Therefore, the computational results in Table 2.3 are of great interest, and to shed further light on these results we now consider the convergence of  $\underline{\tau}$  from a theoretical point of view.

Let  $\tilde{\psi} \in \tilde{\mathbf{H}}_{\bar{w},0}^1(R)$  be the weak solution of (2.1.6) (transformed to polar coordinates). As in the proof of Lemma 2.6.2, we have

$$\tilde{\psi}(r, \theta, t) = \tilde{\psi}_1(r, t) + r \sum_{l=1}^{\infty} \left( \tilde{A}_l(r, t) \cos(2l\theta) + \tilde{B}_l(r, t) \sin(2l\theta) \right), \quad (2.8.12)$$

where we have only taken even modes in the sum (*cf.* Remark 2.8.1) and we use sin and cos functions in (2.8.12) rather than complex exponentials to match the structure of bases  $\mathcal{A}$  and  $\mathcal{B}$ . For simplicity, we shall restrict our attention to the component  $\tau_{11}$  of  $\underline{\tau}$ , although the other components can be treated in exactly the same way.

We consider  $\tau_{11}$  to be a functional defined on  $\hat{\psi} \in L^2(D)$  as follows:

$$\tau_{11}(\hat{\psi}) = \int_D F_1(\underline{q}) q_1 \sqrt{M(\underline{q})} \hat{\psi}(\underline{q}, t) d\underline{q}, \quad (2.8.13)$$

whereby,

$$\begin{aligned} |\tau_{11}(\hat{\psi})| &= \left| \int_D q_1^2 U'(\tfrac{1}{2}|q|^2) \sqrt{M(\underline{q})} \hat{\psi} d\underline{q} \right| \leq b \left( \int_D U'(\tfrac{1}{2}|q|^2)^2 M(\underline{q}) d\underline{q} \right)^{\frac{1}{2}} \|\hat{\psi}\| \\ &= \frac{b}{\sqrt{Z}} \left( \int_D \left(1 - |q|^2/b\right)^{\frac{b}{2}-2} d\underline{q} \right)^{\frac{1}{2}} \|\hat{\psi}\| = \frac{b}{\sqrt{Z}} \left( 2\pi b \int_0^1 (1-r^2)^{\frac{b}{2}-2} r dr \right)^{\frac{1}{2}} \|\hat{\psi}\| \\ &\leq \frac{b}{\sqrt{Z}} \left( 2\pi b \max\left(1, 2^{\frac{b}{2}-2}\right) \int_0^1 (1-r)^{\frac{b}{2}-2} dr \right)^{\frac{1}{2}} \|\hat{\psi}\|, \end{aligned} \quad (2.8.14)$$

where  $Z$  is the normalisation constant from (1.3.21). Hence, we require  $b > 2$  so that  $\tau_{11} \in L^2(D)' = L^2(D)$ ; this is the same condition that we assume for  $b$  throughout this work.

Applying  $\tau_{11}$  to (2.8.12) gives:

$$\begin{aligned}\tau_{11}(\hat{\psi}) &= \frac{b^2}{\sqrt{Z}} \int_0^1 \int_0^{2\pi} (1-r^2)^{\frac{b}{4}-1} r^3 \cos^2(\theta) \tilde{\psi}(r, \theta, t) dr d\theta \\ &= \frac{\pi b^2}{\sqrt{Z}} \int_0^1 r^3 (1-r^2)^{\frac{b}{4}-1} \left( \tilde{\psi}_1(r, t) + \frac{r}{2} \left( \tilde{A}_1(r, t) \right) \right) dr.\end{aligned}\quad (2.8.15)$$

This shows that, quite remarkably, due to orthogonality with  $\cos^2(\theta) = \frac{1}{2} + \frac{1}{2} \cos(2\theta)$  over  $\theta \in (0, 2\pi)$ , the functional  $\tau_{11}$  filters out all but two terms of the infinite series in (2.8.12). The same filtering occurs for Galerkin spectral methods that use trigonometric polynomials in  $\theta$ , such as method  $\mathcal{A}$ , method  $\mathcal{B}$  or the method of Chauvière & Lozinski. We consider method  $\mathcal{A}$  below, but the same approach could be applied to the other methods.

Suppose, using basis  $\mathcal{A}$ , that our numerical solution is defined as follows:

$$\tilde{\psi}_N(r, \theta) = (1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{0,k} P_k(r) + r(1-r) \sum_{i=0}^1 \sum_{l=1}^{N_\theta} \sum_{k=0}^{N_r-1} \tilde{\Psi}_{l,k}^i P_k(r) \Phi_{il}(\theta).$$

Then, assuming  $N_\theta \geq 1$ , we have

$$\tau_{11}(\hat{\psi}_N) = \frac{\pi b^2}{\sqrt{Z}} \int_0^1 r^3 (1-r^2)^{\frac{b}{4}-1} \left[ \left( (1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{0,k} P_k(r) \right) + \frac{r}{2} \left( (1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{1,k}^0 P_k(r) \right) \right] dr.$$

It follows that

$$\begin{aligned}\tau_{11}(\hat{\psi}(t^n)) - \tau_{11}(\hat{\psi}_N^n) &= \frac{\pi b^2}{\sqrt{Z}} \int_0^1 r^3 (1-r^2)^{\frac{b}{4}-1} \left[ \left( \tilde{\psi}_1(r, t^n) - (1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{0,k}^n P_k(r) \right) \right. \\ &\quad \left. + \frac{1}{2} \left( r \tilde{A}_1(r, t^n) - r(1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{1,k}^{0,n} P_k(r) \right) \right] dr.\end{aligned}\quad (2.8.16)$$

Applying the Cauchy-Schwarz inequality gives

$$\begin{aligned}|\tau_{11}(\hat{\psi}(t^n)) - \tau_{11}(\hat{\psi}_N^n)|^2 &\leq C_* \left\| \tilde{\psi}_1(r, t^n) - (1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{0,k}^n P_k(r) \right\|_{L_w^2(0,1)}^2 \\ &\quad + \frac{C_*}{4} \left\| r \tilde{A}_1(r, t^n) - r(1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{1,k}^{0,n} P_k(r) \right\|_{L_w^2(0,1)}^2,\end{aligned}\quad (2.8.17)$$

where,

$$C_* = \begin{cases} \frac{2\pi^2 b^4}{(b/2-1)Z}, & 2 < b < 4 \\ \frac{\pi^2 b^4}{3Z}, & b \geq 4 \end{cases}\quad (2.8.18)$$

and, as in Section 2.6,  $L_w^2(0,1)$  is the  $r$ -weighted  $L^2$  space.

On the other hand, using Parseval's identity, we have

$$\begin{aligned}
\|\hat{\psi}(\cdot, t^n) - \hat{\psi}_N^n(\cdot)\|_{L^2(D)}^2 &= b \int_0^1 \int_0^{2\pi} |\tilde{\psi}(r, \theta, t^n) - \tilde{\psi}_N^n(r, \theta)|^2 r \, dr \, d\theta \\
&= 2\pi b \left\| \tilde{\psi}_1(r, t^n) - (1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{0,k}^n P_k(r) \right\|_{L_w^2(0,1)}^2 \\
&\quad + \pi b \sum_{l=1}^{N_\theta} \left\| r \tilde{A}_l(r, t^n) - r(1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{l,k}^{0,n} P_k(r) \right\|_{L_w^2(0,1)}^2 \\
&\quad + \pi b \sum_{l=1}^{N_\theta} \left\| r \tilde{B}_l(r, t^n) - r(1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{l,k}^{1,n} P_k(r) \right\|_{L_w^2(0,1)}^2 \\
&\quad + \pi b \sum_{l=N_\theta+1}^{\infty} \left( \left\| r \tilde{A}_l(r, t^n) \right\|_{L_w^2(0,1)}^2 + \left\| r \tilde{B}_l(r, t^n) \right\|_{L_w^2(0,1)}^2 \right). \quad (2.8.19)
\end{aligned}$$

It follows that

$$\|\tau_{11}(\tilde{\psi}) - \tau_{11}(\tilde{\psi}_N)\|_{\ell^\infty(0,T)} \leq \sqrt{\frac{C_*}{2\pi b}} \|\hat{\psi} - \hat{\psi}_N\|_{\ell^\infty(0,T;L^2(D))}. \quad (2.8.20)$$

However, more importantly, we can see that the bound in (2.8.17) contains only two terms from the infinite sum in (2.8.19) (albeit with different constants) and therefore we expect that the error in  $\tau_{11}$  will typically be much smaller than the error in  $\hat{\psi}$ .

In practical computations, this manifests itself as superconvergence of  $\tau_{11}$ . We demonstrate this superconvergence here by comparing the  $L^2(D)$  convergence data for  $\hat{\psi}$  from Tables 2.1 and 2.2 with the corresponding errors in  $\tau_{11}$ . These results are plotted in Figure 2.3 and we can clearly see that, prior to stagnation due to rounding error,  $\tau_{11}$  converges at a faster rate than  $\hat{\psi}$ , and the error in  $\tau_{11}$  is typically orders of magnitude smaller than the  $\hat{\psi}$  error. This behaviour is extremely advantageous for micro-macro computations where the accuracy of  $\tau_{11}$  (rather than  $\hat{\psi}$ ) is crucial.

One interesting thing to note from Figure 2.3(b) is that the error in  $\tau_{11}$  appears to stagnate at around  $10^{-10}$  with both method  $\mathcal{A}$  and method  $\mathcal{B}$ , and in fact, the error increases to some extent when the number of spectral basis functions is increased further (*e.g.* compare  $N_r = N_\theta = 35$  to  $N_r = N_\theta = 40$  in the plot); this increase in error is due to the fact that the condition number of the linear system (2.8.8) increases with  $N_r$  and  $N_\theta$  and hence we can lose extra digits of accuracy for larger values of  $N_r$ ,  $N_\theta$ . Similarly, the condition number is larger for method  $\mathcal{A}$  than for method  $\mathcal{B}$  for the computations considered in Figure 2.3(a), which is why the error in  $\tau_{11}$  for method  $\mathcal{A}$  stagnates at around  $10^{-11}$ , whereas the error from method  $\mathcal{B}$  stagnates at  $10^{-13}$ .

**Remark 2.8.2** It was proved in Lemma 2.3.4 that the weak solution of the initial-boundary-value problem (2.1.1), (2.1.2), (2.1.7) is nonnegative a.e. on  $D$ . This property is not guaranteed to hold for the numerical solution. However, our numerical experiments consistently show that if there are sufficiently many modes in the approximation space to accurately resolve the solution then this nonnegativity property is preserved under discretisation. This

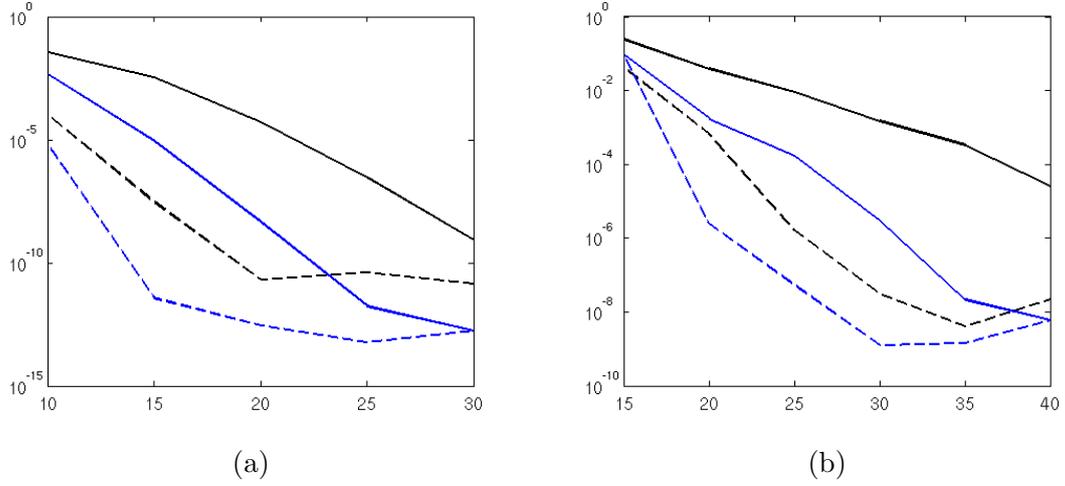


Figure 2.3: Comparison of convergence of  $\hat{\psi}$  and  $\tau_{11}$ . In both plots, the horizontal axis shows the value of  $N_r$  and  $N_\theta$  (chosen to be equal in these computations). Plot (a) shows data for the computations considered in Table 2.1 and plot (b) corresponds to Table 2.2. In both (a) and (b), the solid black line represents the relative  $L^2(D)$  error in  $\hat{\psi}$  for method  $\mathcal{A}$ , and the solid blue line represents the corresponding data for method  $\mathcal{B}$ . The dashed black line shows  $\tau_{11}$  errors arising from method  $\mathcal{A}$ , and the dashed blue line is analogous for method  $\mathcal{B}$ .

is illustrated in Figure 2.4 in which two cross-sections of the numerical solution for the  $(b, \text{Wi}, \delta) = (12, 1, 5)$  extensional flow are shown: the numerical solution on the left is fully resolved, while the one on the right is under-resolved. In the under-resolved case there are oscillations and clearly  $\psi_N \geq 0$  is not satisfied throughout  $D$ , whereas the nonnegativity property is accurately captured in the fully resolved case.  $\diamond$

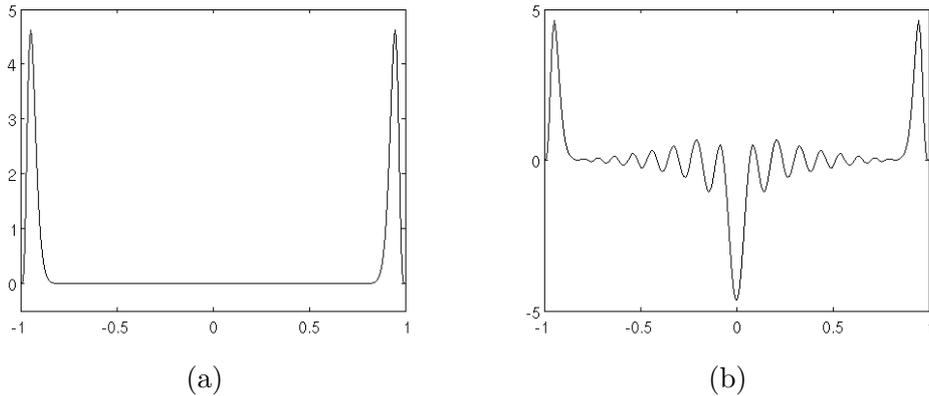


Figure 2.4: Cross-sections of the solution of the extensional flow problem with  $b = 12$ ,  $\text{Wi} = 1$  and  $\delta = 5$  at steady state, obtained using method  $\mathcal{B}$ . The fully-resolved solution in (a) was obtained using  $(N_r, N_\theta) = (41, 20)$ , and the under-resolved solution in (b) was obtained with  $(N_r, N_\theta) = (26, 20)$ .

### 2.8.2 The semi-implicit numerical method

Up until now we have confined our attention to the backward Euler temporal discretisation of the Fokker–Planck equation, as defined in (2.5.1). However, as we shall see, the semi-implicit discretisation, which is identical to (2.5.1) except that the term  $\int_D \underline{\underline{k}}_q \hat{\psi}_N \cdot \nabla_M \hat{\varphi} \, dq$  is treated explicitly in time, is important in Chapter 3. Therefore, as a precursor to the next section, we consider this semi-implicit scheme here.

It should be noted that all of the analytical results that we obtained for the backward Euler temporal discretisation (also referred to from now on as the fully-implicit discretisation) in this section also carry across to the semi-implicit scheme – we do not consider the details here, but, for example, in the process of proving Lemma 3.4.1 in the next section, we establish a stability result for the semi-implicit scheme that is almost identical to Lemma 2.3.1. However, although the  $L^2(D)$  stability estimates (and therefore also the asymptotic convergence results) are essentially identical for the fully-implicit and semi-implicit schemes, we show in this section that for practical computations, the fully-implicit discretisation tends to be much more stable in the sense that solutions obtained from the semi-implicit scheme are more likely to exhibit the exponential growth in time in the  $L^2(D)$  norm that is allowed due to the constant  $e^{2c_0 m \Delta t}$  in Lemma 2.3.1. This is not surprising; it is well-known that fully-implicit schemes are generally more stable than semi-implicit and explicit schemes for parabolic and hyperbolic PDEs.

All the details of the implementation of the semi-implicit method carry over from the discussion of the fully-implicit method above; the only difference is that instead of (2.8.8), the linear system in this case is:

$$\left( \mathbf{M} + \frac{\Delta t}{2\text{Wi}} \mathbf{S} \right) \tilde{\Psi}^{n+1} = (\mathbf{M} + \Delta t \mathbf{C}^n) \tilde{\Psi}^n. \quad (2.8.21)$$

We now present some numerical results that compare the fully-implicit and semi-implicit schemes. We only consider method  $\mathcal{B}$  here, with the understanding that the behaviour for method  $\mathcal{A}$  is essentially the same.

First of all, we repeated the computations in Table 2.1 (for an extensional flow with  $(b, \text{Wi}, \delta) = (12, 1, 1)$ ) using the semi-implicit scheme, and the results were identical to those reported in Table 2.1 for the backward Euler discretisation. However, on increasing the Weissenberg number from 1 to 5 we then observed significant differences between the two schemes. The results for the  $\text{Wi} = 5$  computations are summarised in Table 2.4.

We can see from the table that with  $(N_r, N_\theta) = (15, 15)$ , the semi-implicit scheme led to numerical solutions for which the  $L^2(D)$  norm error grew rapidly in time for all three time-step sizes,  $\Delta t = 0.1, 0.05$  and  $0.01$  (indicated by hyphens in the table), whereas the fully-implicit scheme had an  $\mathcal{O}(1)$  error in each of these cases. In the computations with  $(N_r, N_\theta) = (20, 20)$ , the fully-implicit method again performed better; we needed to take  $\Delta t = 0.01$  in order to get an accurate solution with the semi-implicit scheme, whereas the fully-implicit scheme was accurate with  $\Delta t = 0.1$ . Finally, for  $(N_r, N_\theta) = (25, 25)$  and  $(N_r, N_\theta) = (30, 30)$ , the two schemes behaved identically for  $\Delta t = 0.05$  and  $\Delta t = 0.01$ , but the fully-implicit scheme remained accurate for  $\Delta t = 0.1$  whereas the semi-implicit scheme did not.

These observations indicate that the fully-implicit scheme is reliable for coarser spatial discretisations and larger  $\Delta t$  than the semi-implicit scheme. This is especially noticeable

$(N_r, N_\theta)$	$\Delta t = 0.1, N_T = 250$		$\Delta t = 0.05, N_T = 500$		$\Delta t = 0.01, N_T = 2500$	
	Imp.	Semi-Imp.	Imp.	Semi-Imp.	Imp.	Semi-Imp.
(15,15)	1.49	–	1.49	–	1.49	–
(20,20)	$4.67 \times 10^{-2}$	–	$4.67 \times 10^{-2}$	$1.14 \times 10^{+2}$	$4.67 \times 10^{-2}$	$4.67 \times 10^{-2}$
(25,25)	$2.96 \times 10^{-3}$	–	$2.96 \times 10^{-3}$	$2.96 \times 10^{-3}$	$2.96 \times 10^{-3}$	$2.96 \times 10^{-3}$
(30,30)	$1.44 \times 10^{-4}$	–	$1.44 \times 10^{-4}$	$1.44 \times 10^{-4}$	$1.44 \times 10^{-4}$	$1.44 \times 10^{-4}$

Table 2.4: This table shows the relative  $L^2(D)$  error, with respect to the exact steady-state solution, for the implicit and semi-implicit schemes (using method  $\mathcal{B}$ ) applied to an extensional flow problem with  $(b, \text{Wi}, \delta) = (12, 5, 1)$ . Three different time-step sizes were tested, and the total number of time-steps,  $N_T$ , was varied in order to ensure that  $T = N_T \Delta t$  was the same in each case.

when the Weissenberg number is increased (recall that the two methods behaved identically for the extensional flow with  $\text{Wi} = 1$ ). Note also that scaling  $\underline{\kappa}$  has roughly the same effect as scaling  $\text{Wi}$ , *e.g.* the steady state solution (assuming it exists) depends on the product  $\text{Wi} \underline{\kappa}$  and not on  $\text{Wi}$  or  $\underline{\kappa}$  separately.<sup>9</sup> Hence, based on the results in Table 2.4, we conclude that it is preferable to use the fully-implicit temporal discretisation for problems in which  $\text{Wi}$  or  $|\underline{\kappa}|$ , or both, are large (compared to, say, 1).

### 2.8.3 Three-dimensional implementation of the spectral method

We now consider the implementation of the spectral method developed in this chapter for  $d = 3$ . This is closely related to the two-dimensional case, the primary differences being that we now use the spherical coordinate change of variables:

$$q = (\sqrt{br} \cos \theta \sin \phi, \sqrt{br} \sin \theta \sin \phi, \sqrt{br} \cos \phi), \quad (r, \theta, \phi) \in R := (0, 1) \times (0, 2\pi) \times (0, \pi),$$

instead of (2.6.1) and, following Chauvière & Lozinski [32], we choose each of our basis functions to be a product of a spherical harmonic in  $(\theta, \phi)$  and polynomial in  $r$ . Discretisations of this type have also been considered in the recent paper by Guo and Huang [57]. Note that in this section,  $\tilde{g}(r, \theta, \phi) := g(q_1, q_2, q_3)$ .

First of all, we redefine the space  $\tilde{H}^1(R)$  for the purposes of this section, in order to ensure that if  $g \in H^1(D)$ ,  $D \subset \mathbb{R}^3$  then  $\tilde{g} \in \tilde{H}^1(R)$ . Following the approach in the case of  $d = 2$ , we define  $\|\tilde{g}\|_{\tilde{H}^1(R)}^2$  by transforming  $\|g\|_{H^1(D)}^2$  from cartesian to spherical coordinates, and hence we have

$$\|\tilde{g}\|_{\tilde{H}^1(R)}^2 := \int_R r^2 \sin \phi \left( |\tilde{g}|^2 + \left| \frac{\partial \tilde{g}}{\partial r} \right|^2 + \frac{1}{r^2} \left| \frac{\partial \tilde{g}}{\partial \phi} \right|^2 + \frac{1}{r^2 \sin^2 \phi} \left| \frac{\partial \tilde{g}}{\partial \theta} \right|^2 \right) dr d\theta d\phi,$$

and,

$$\tilde{H}^1(R) := \{ \tilde{f} \in L_{\text{loc}}^2(R) \quad : \quad \tilde{f}(r, \cdot, \phi) \in H_p^1(0, 2\pi) \text{ for a.e. } (r, \phi) \in (0, 1) \times (0, \pi) \\ \text{and } \|\tilde{f}\|_{\tilde{H}^1(R)} < \infty \}.$$

<sup>9</sup>This can be seen by scaling  $\text{Wi}$  and  $\underline{\kappa}$  in (2.1.1) and noting that  $\frac{\partial \psi}{\partial t}$  vanishes at steady state.

We denote the spherical harmonics by  $S_{l,m} : (\theta, \phi) \mapsto S_{l,m}(\theta, \phi) \in \mathbb{R}$ . They are the solutions of the equation

$$\frac{1}{\sin \phi} \frac{\partial}{\partial \phi} \left( \sin \phi \frac{\partial}{\partial \phi} S_{l,m}(\theta, \phi) \right) + \frac{1}{\sin^2 \phi} \frac{\partial^2}{\partial \theta^2} S_{l,m}(\theta, \phi) + l(l+1) S_{l,m}(\theta, \phi) = 0, \quad (2.8.22)$$

for a.e.  $(\theta, \phi) \in (0, 2\pi) \times (0, \pi)$ , where (2.8.22) is the angular part of Laplace's equation in spherical coordinates. It can be shown, by separation of variables, that the solutions of (2.8.22) are of the form,

$$S_{l,m}(\theta, \phi) = C(l, m) P_l^m(\cos \phi) e^{im\theta}, \quad (2.8.23)$$

for  $l \in \mathbb{Z}_{\geq 0}$ ,  $|m| \leq l$ , where  $P_l^m$  denotes an associated Legendre function and  $C(l, m)$  is a normalisation constant. Also, the (appropriately normalised) spherical harmonics satisfy the following orthogonality property:

$$\int_0^{2\pi} \int_0^\pi S_{l_1, m_1}(\theta, \phi) \overline{S_{l_2, m_2}(\theta, \phi)} \sin \phi \, d\theta \, d\phi = \delta_{m_1, m_2} \delta_{l_1, l_2}, \quad (2.8.24)$$

where the overline notation denotes complex conjugation.

The next lemma will motivate our definition of a spectral basis in the case of  $d = 3$ .

**Lemma 2.8.3** *Let  $\tilde{g}(r, \theta) = \sum_{l=0}^{N_{\text{sph}}} \sum_{|m| \leq l} \tilde{\gamma}_l^m(r) S_{l,m}(\theta, \phi)$ ,  $N_{\text{sph}} \in \mathbb{Z}_{\geq 0}$ ,  $\tilde{\gamma}_0^0 \in H_{r^2}^1(0, 1)$  where  $H_{r^2}^1(0, 1)$  is the  $r^2$ -weighted  $H^1$ -space, and*

$$\tilde{\gamma}_l^m \in H^1(0, 1; 1, r^2) := \left\{ \tilde{f} \in H_{\text{loc}}^1(0, 1) : \int_0^1 \left( |\tilde{f}(r)|^2 + r^2 |\tilde{f}'(r)|^2 \right) dr < \infty \right\},$$

for  $l > 0$ ; then  $\tilde{g} \in \tilde{H}^1(R)$ .

**Proof.** Periodicity of  $\tilde{g}$  in  $\theta$  follows directly from the definition of the spherical harmonics, hence it only remains to verify that  $\|\tilde{g}\|_{\tilde{H}^1(R)} < \infty$ .

Integrating by parts in  $\theta$  and  $\phi$  (which is valid for spherical harmonics), we obtain:

$$\begin{aligned} \|\tilde{g}\|_{\tilde{H}^1(R)}^2 &= \int_R r^2 \sin \phi \left( |\tilde{g}|^2 + \left| \frac{\partial \tilde{g}}{\partial r} \right|^2 \right) dr \, d\theta \, d\phi \\ &\quad - \int_R \sin \phi \tilde{g} \left( \frac{1}{\sin \phi} \frac{\partial}{\partial \phi} \left( \sin \phi \frac{\partial \tilde{g}}{\partial \phi} \right) + \frac{1}{\sin^2 \phi} \frac{\partial^2 \tilde{g}}{\partial \theta^2} \right) dr \, d\theta \, d\phi, \end{aligned} \quad (2.8.25)$$

where the boundary conditions vanish due to periodicity. Substituting the series expression of  $\tilde{g}$  into (2.8.25) and using (2.8.22) and (2.8.24), we get:

$$\begin{aligned} \|\tilde{g}\|_{\tilde{H}^1(R)}^2 &= \sum_{l=0}^{N_{\text{sph}}} \sum_{|m| \leq l} \int_0^1 r^2 \left\{ |\tilde{\gamma}_l^m(r)|^2 dr + \left| \frac{d\tilde{\gamma}_l^m}{dr} \right|^2 \right\} dr \\ &\quad + \int_R \left\{ \sum_{l_1=0}^{N_{\text{sph}}} \sum_{|m_1| \leq l_1} \tilde{\gamma}_{l_1}^{m_1}(r) S_{l_1, m_1}(\theta, \phi) \right\} \left\{ \sum_{l_2=0}^{N_{\text{sph}}} \sum_{|m_2| \leq l_2} l_2(l_2+1) \tilde{\gamma}_{l_2}^{m_2}(r) \overline{S_{l_2, m_2}(\theta, \phi)} \right\} \sin \phi \, d\phi \, d\theta \, dr \\ &= \sum_{l=0}^{N_{\text{sph}}} \sum_{|m| \leq l} \int_0^1 \left\{ r^2 |\tilde{\gamma}_l^m(r)|^2 + r^2 \left| \frac{d}{dr} \tilde{\gamma}_l^m(r) \right|^2 + l(l+1) |\tilde{\gamma}_l^m(r)|^2 \right\} dr. \end{aligned} \quad (2.8.26)$$

By the hypotheses on the  $\tilde{\gamma}_l^m$ , it follows that  $\|\tilde{g}\|_{\tilde{H}^1(R)}$  is finite.  $\square$

Note that the  $\tilde{\gamma}_l^m$  in Lemma 2.8.3 need not be bounded on  $(0, 1)$  since, for example,  $r^{-1/4} \in H_{r,2}^1(0, 1) \cap H^1(0, 1; 1, r^2)$ .

It will be convenient from now on to use the real and imaginary parts of the spherical harmonics rather than the complex exponentials in (2.8.23), *i.e.*:

$$S_{l,m}^i(\theta, \phi) := C(l, m) P_l^m(\cos \phi) ((1 - i) \cos(m\theta) + i \sin(m\theta)), \quad (2.8.27)$$

where now  $0 \leq l \leq N_{\text{sph}}$ ,  $i \in \{0, 1\}$ , and  $i \leq m \leq l$ . In this section, we consider basis functions of the following form:

$$Y_{lm}^{ik}(r, \theta, \phi) := (1 - r) Q_k(r) S_{l,m}^i(\theta, \phi), \quad (2.8.28)$$

where  $(1 - r) Q_k \in \mathbb{P}_{N_r, 0}(0, 1)$  (as in the  $d = 2$  case,  $Q_k$  is taken to be a Chebyshev polynomial of degree  $k$ ,  $0 \leq k \leq N_r - 1$ , mapped from  $[-1, 1]$  to  $[0, 1]$ , although other polynomial choices could be considered also). Since  $\mathbb{P}_{N_r, 0}(0, 1) \subset H_{r,2}^1(0, 1) \cap H^1(0, 1; 1, r^2)$ , it follows from Lemma 2.8.3, that any finite linear combination of basis functions of the form (2.8.28) is contained in  $\tilde{H}_0^1(R)$ . This is a simpler situation than in two dimensions, since now we do not need to impose a specialised decomposition in order to guarantee inclusion in  $\tilde{H}^1(R)$ .

Below we shall introduce a basis on which our Galerkin spectral method in three dimensions will be based on. Before defining this basis, however, we first consider the symmetry property discussed in Remark 2.8.1 in the  $d = 3$  case. In fact, most of Remark 2.8.1 carries over to three dimensions unchanged; the only difference is that now the evenness of  $\hat{\psi}$  in the  $D$  domain with respect to  $q$  translates to requiring that we only use spherical harmonics in  $R$  for which  $l$  is an even number. This can be seen by the following argument. Suppose, using the change of variables to spherical coordinates, that  $q \mapsto (r, \theta, \phi)$ . Then also  $-q \mapsto (r, \theta + \pi, \pi - \phi)$ . Now, the symmetry condition we wish to impose is that for any basis function  $Y_{lm}^{ik}$  defined in (2.8.28), we have  $Y_{lm}^{ik}(r, \theta, \phi) = Y_{lm}^{ik}(r, \theta + \pi, \pi - \phi)$ . This, in turn, requires that  $S_{l,m}^i(\theta, \phi) = S_{l,m}^i(\theta + \pi, \pi - \phi)$ . Noting that,

$$S_{l,m}^i(\theta, \phi) = P_l^m(\cos \phi) ((1 - i) \cos(m\theta) + i \sin(m\theta)),$$

and

$$S_{l,m}^i(\theta + \pi, \pi - \phi) = (-1)^m P_l^m(-\cos \phi) ((1 - i) \cos(m\theta) + i \sin(m\theta)),$$

it follows that we can only use associated Legendre functions for which

$$P_l^m(x) = (-1)^m P_l^m(-x), \quad x \in [-1, 1].$$

Since the associated Legendre functions are defined as,

$$P_l^m(x) = (-1)^m (1 - x^2)^{m/2} \frac{d^m}{dx^m} (P_l(x)),$$

where  $P_l(x)$  is a Legendre polynomial of degree  $l$  (for which  $P_l(x) = (-1)^l P_l(-x)$ ), it follows that the required symmetry condition is satisfied if, and only if,  $l$  is an even number (for any  $m = 0, \dots, l$ ).

**Remark 2.8.4** In [32], Chauvière & Lozinski restricted their attention to two-dimensional macroscopic velocity fields, in which case a more restrictive symmetry condition was appropriate, *i.e.* that  $\psi(r, \theta, \phi) = \psi(r, \theta + \pi, \phi)$ , and hence they only considered spherical harmonics for which both  $l$  and  $m$  were even numbers. Compared to the more general symmetry condition considered above, the condition of Chauvière & Lozinski leads to a reduction in computational effort because for a given  $N_{\text{sph}}$ , fewer basis functions are used since the spherical harmonics with odd  $m$  are discarded, and also it is only necessary to consider  $\theta \in (0, \pi)$ . In this work, however, we are interested in treating the case in which the macroscopic velocity field can be three-dimensional, and therefore we require the symmetry condition for  $Y_{lm}^{ik}$  identified above.  $\diamond$

With the considerations discussed above in mind, we can now define a basis, denoted  $\mathcal{C}$ , as follows:

$$\mathcal{C} := \{Y_{lm}^{ik} : 0 \leq k \leq N_r - 1, i \in \{0, 1\}, l \in \{0, 2, 4, \dots, N_{\text{sph}}\} \text{ and } i \leq m \leq l\}.$$

From now on, the numerical method that uses basis  $\mathcal{C}$  will be referred to as method  $\mathcal{C}$ .

At this point, we could take a detour to consider three-dimensional approximation results for  $\text{span}(\mathcal{C}) \subset \tilde{\mathbb{H}}_0^1(R)$ , which would then allow us to extend our convergence results from Section 2.7 to the  $d = 3$  case. However, given that we have already considered approximation results in detail for  $d = 2$ , and given that the approach in the  $d = 3$  case would be completely analogous, for the sake of brevity, we omit discussion of approximation theory in three dimensions here. Note, however, that Guo & Huang [57] recently derived approximation results for a spectral method on the unit ball in  $\mathbb{R}^3$ , which could be applied to the convergence analysis of method  $\mathcal{C}$  (*e.g.* see Theorem 2.3 in that paper, which is similar to our approximation result (2.6.16)).

Below we shall test the performance of method  $\mathcal{C}$  on some model problems. First of all, however, we specify the spherical coordinate form of the discretisation matrices defined in (2.8.4). Using the same notation that we used for the discretisation matrices in polar coordinates, we let  $N$  denote the total number of basis functions, we set

$$\tilde{\psi}_N^{n+1}(r, \theta, \phi) = \sum_{v=1}^N \tilde{\Psi}_v^{n+1} Y_v(r, \theta, \phi),$$

where  $Y_v$  is a basis function from  $\mathcal{C}$  for  $1 \leq v \leq N$ , and we denote the test functions by  $Y_u$  for  $1 \leq u \leq N$ . Then,

$$(M_q)_{uv} = \int_0^1 \int_0^{2\pi} \int_0^\pi b^{3/2} Y_u Y_v r^2 \sin \phi \, dr \, d\theta \, d\phi, \quad (2.8.29)$$

$$(S_q)_{uv} = \int_0^1 \int_0^{2\pi} \int_0^\pi \left\{ b^{1/2} \frac{\partial Y_u}{\partial r} \frac{\partial Y_v}{\partial r} r^2 \sin \phi + b^{1/2} \frac{1}{\sin \phi} \frac{\partial Y_u}{\partial \theta} \frac{\partial Y_v}{\partial \theta} + b^{1/2} \frac{\partial Y_u}{\partial \phi} \frac{\partial Y_v}{\partial \phi} \sin \phi \right. \\ \left. + \frac{b^{3/2}}{2} r^3 \sin \phi (1 - r^2)^{-1} \left[ \frac{\partial Y_u}{\partial r} Y_v + Y_u \frac{\partial Y_v}{\partial r} \right] + \frac{b^{5/2}}{4} r^4 \sin \phi (1 - r^2)^{-2} Y_u Y_v \right\} dr \, d\theta \, d\phi, \quad (2.8.30)$$

$$(C_q^m)_{uv} = \int_0^1 \int_0^{2\pi} \int_0^\pi Y_v \left\{ k_r \left[ b^{3/2} r^3 \frac{\partial Y_u}{\partial r} + \frac{b^{5/2}}{2} r^4 (1 - r^2)^{-1} Y_u \right] \right. \\ \left. + k_\theta b^{3/2} r^2 \frac{\partial Y_u}{\partial \theta} + k_\phi b^{3/2} r^2 \sin \phi \frac{\partial Y_u}{\partial \phi} \right\} dr \, d\theta \, d\phi, \quad (2.8.31)$$

where  $k_r = (\underline{\kappa}(\underline{x}_m)\underline{e}_r) \cdot \underline{e}_r$ ,  $k_\theta = (\underline{\kappa}(\underline{x}_m)\underline{e}_r) \cdot \underline{e}_\theta$  and  $k_\phi = (\underline{\kappa}(\underline{x}_m)\underline{e}_r) \cdot \underline{e}_\phi$ , with  $\underline{e}_r$ ,  $\underline{e}_\theta$ ,  $\underline{e}_\phi$  the unit vectors in the  $r$ ,  $\theta$  and  $\phi$  directions:

$$\begin{aligned}\underline{e}_r &= (\cos \theta \sin \phi, \sin \theta \sin \phi, \cos \phi), \\ \underline{e}_\theta &= (-\sin \theta, \cos \theta, 0), \\ \underline{e}_\phi &= (\cos \theta \cos \phi, \sin \theta \cos \phi, -\sin \phi).\end{aligned}$$

Note that  $\underline{\kappa}q = \sqrt{b}r(k_r\underline{e}_r + k_\theta\underline{e}_\theta + k_\phi\underline{e}_\phi)$ , and  $(\underline{e}_r, \underline{e}_\theta, \underline{e}_\phi)$  is an orthonormal basis for  $\mathbb{R}^3$  for any  $(\theta, \phi) \in (0, 2\pi) \times (0, \pi)$ . We refer to Section 2.8.1 for the details of computing the discretisation matrices and the solution of the resulting linear system; the approach is completely analogous here.

Next, we present some computational results for method  $\mathcal{C}$ . We consider the backward Euler temporal discretisation of the FENE Fokker–Planck equation here, as opposed to the semi-implicit scheme considered in Section 2.8.2, and we restrict our attention to producing plots of the same type as in Figure 2.3 in order to visualise the convergence rates for  $\hat{\psi}$  and  $\underline{\tau}$ , and also to verify that we obtain superconvergence of  $\underline{\tau}$  in the  $d = 3$  case. Note that theoretical underpinning of the superconvergence of  $\underline{\tau}$  characterised in (2.8.17) and (2.8.19) can also be established in 3-dimensions; the reasoning is the same, except that we use Parseval’s identity based on spherical harmonics, as in Lemma 2.8.3.

As in the two-dimensional case, we know the exact steady state solution for problems in which  $\underline{\kappa}$  is a symmetric  $3 \times 3$  tensor (*cf.* (2.8.9)). We now consider two distinct problems; for each problem we have  $\underline{\kappa} = \underline{\kappa}^T$  so that we can compare the numerical solution with the exact steady state solution, and as in Tables 2.1 and 2.2 we take 2000 time-steps with  $\Delta t = 0.05$  to obtain an accurate approximation to the steady state solution.

The first problem we consider is a three-dimensional extensional flow with  $b = 12$ ,  $\text{Wi} = 1$  and  $\underline{\kappa}$  defined as follows:

$$\underline{\kappa} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & -1/2 & 0 \\ 0 & 0 & -1/2 \end{pmatrix}. \quad (2.8.32)$$

Figure 2.5(a) shows the convergence plots for  $\hat{\psi}$  and  $\tau_{11}$  for this problem. It is clear from the figure that we obtain spectral convergence of  $\hat{\psi}$ , and also, just as in Figure 2.3, we observe superconvergence of  $\tau_{11}$ .

Next, we consider a problem in which  $\underline{\kappa}$  is a full tensor:

$$\underline{\kappa} = \begin{pmatrix} 0.5 & 0.2 & 0.5 \\ 0.2 & -0.25 & -0.4 \\ 0.5 & -0.4 & -0.25 \end{pmatrix}, \quad (2.8.33)$$

and where  $b = 12$  and  $\text{Wi} = 1$  again. The convergence plot for this computation is shown in Figure 2.5(b), and the behaviour is much the same as in Figure 2.5(a).

## 2.9 Conclusions

The purpose of this chapter has been to develop a rigorous foundation for the numerical approximation of Fokker–Planck equations. We restricted our attention to the configuration

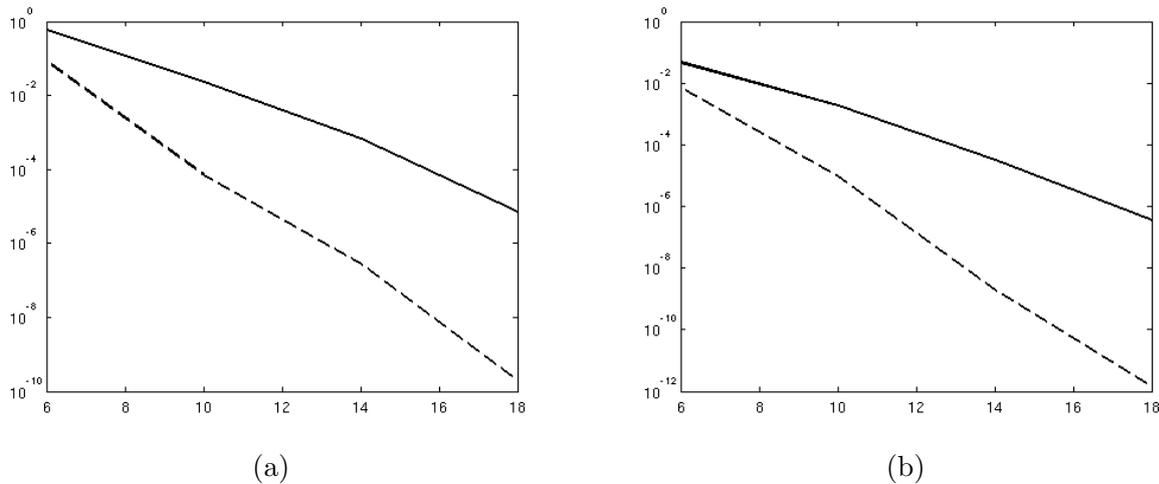


Figure 2.5: Comparison of convergence of  $\hat{\psi}$  and  $\tau_{11}$  for method  $\mathcal{C}$  for two different problems (we compared to the exact steady state solution, (2.8.9), by taking 2000 time-steps with  $\Delta t = 0.05$ ). Plot (a) corresponds to a three-dimensional extensional flow problem with  $b = 12$  and  $Wi = 1$  and with  $\kappa$  defined in (2.8.32). Plot (b) is analogous, except that in this case  $\kappa$  is as in (2.8.33). In both plots, the horizontal axis represents  $N_r$  and  $N_{\text{sph}}$  (chosen to be equal in these computations), and the solid and dashed lines show the relative  $L^2(D)$  error and relative  $\tau_{11}$  error, respectively.

space part of (1.3.36), but the work in this chapter will be built upon in subsequent chapters in order to develop numerical methods on  $\Omega \times D$ .

We focused on the symmetrised weak formulation of the Maxwellian-transformed equation, and we used the substitution  $\hat{\psi} = \psi/\sqrt{M}$ . The resulting formulation (2.1.6) facilitated the development of a number of analytical results in Sections 2.3 and 2.5. Using the approximation results derived in Section 2.6, optimal-order convergence of the fully-discrete spectral Galerkin method (2.5.1), (2.5.2) was established for the case of  $d = 2$ ; an analogous procedure could be carried out for  $d = 3$ . This analysis was performed for spring potentials that satisfy Hypotheses A and B; see Example 2.1.1.

In the case of the FENE model, we indicated the extension of our analysis to a class of numerical methods based on another change of variable, proposed by Chauvière & Lozinski; here a different transformation, (2.8.10), is applied to the Fokker–Planck equation. We showed that, at the analytical level at least, the two approaches lead to methods with very similar stability and accuracy properties.

Section 2.8 addressed issues related to the implementation of numerical methods for the FENE Fokker–Planck equation. In Section 2.8.1 we considered two distinct implementations, methods  $\mathcal{A}$  and  $\mathcal{B}$ , for the  $d = 2$  case, and these methods were also compared to the spectral method discussed in the paper of Chauvière & Lozinski [33] on the basis of numerical results reported therein. We showed that methods  $\mathcal{A}$  and  $\mathcal{B}$  work well for values of  $b$  up to about 20, and are comparable to the method formulated in [33] in terms of computational efficiency in this parameter range, with method  $\mathcal{B}$  being more accurate than method  $\mathcal{A}$ , and of a very similar accuracy as the method in [33]. Also, we demonstrated that the convergence of  $\tau_{11}$  tends to be much more rapid than the convergence of  $\hat{\psi}$  using our Galerkin spectral methods;

this is highly advantageous in the context of the micro-macro computations. In Section 2.8.3 we considered the implementation of the Galerkin spectral method, based on the symmetrised formulation, in three spatial dimensions. We constructed a  $\tilde{H}^1(R)$ -conforming spectral basis,  $\mathcal{C}$ , and demonstrated that the convergence properties of the spectral method based on  $\mathcal{C}$  are essentially the same as for the two-dimensional spectral methods considered in Section 2.8.1.

The numerical methods and analytical results developed in this chapter are built upon in Chapter 3, where we consider the Fokker–Planck equation on  $\Omega \times D$ .



## Chapter 3

# Alternating-direction methods for the full Fokker–Planck equation

### 3.1 Introduction

In this chapter, we develop numerical methods for the Maxwellian-transformed Fokker–Planck equation posed on  $\Omega \times D \times (0, T]$ :

$$\frac{\partial \psi}{\partial t} + \underline{u} \cdot \nabla_x \psi + \nabla_q \cdot (\underline{\kappa} \underline{q} \psi) = \frac{1}{2\mathbb{W}_i} \nabla_q \cdot \left( M \nabla_q \frac{\psi}{M} \right), \quad (\underline{x}, \underline{q}, t) \in \Omega \times D \times (0, T], \quad (3.1.1)$$

$$\psi(\underline{x}, \underline{q}, 0) = \psi^0(\underline{x}, \underline{q}), \quad (\underline{x}, \underline{q}) \in \Omega \times D. \quad (3.1.2)$$

Throughout this chapter we assume that  $\underline{u} : (\underline{x}, t) \in \Omega \times (0, T] \mapsto \underline{u}(\underline{x}, t) \in \mathbb{R}^d$  is an *a priori* defined vector field (hence  $\underline{\kappa} = \nabla_x \underline{u}$  is known *a priori* also). The precise hypotheses on  $\underline{u}$  and  $\underline{\kappa}$  shall be specified below.

The above equation will be referred to as the *full* Fokker–Planck equation, to distinguish it from the equation posed on  $D \times (0, T]$  only, that was studied in Chapter 2. From now on, we focus on the Maxwellian-transformed form of the Fokker–Planck equation given above (and its weak formulation in which the principal part of the differential operator is symmetric). However, it should be noted that the numerical methods developed and analysed in the forthcoming sections could just as well be based on the Chauvière–Lozinski-transformed equation that was studied in Section 2.4, and was also used to solve the full FENE Fokker–Planck equation in [32, 33, 91].

As discussed in Chapter 1, due to the cartesian product structure of the domain  $\Omega \times D$ , a natural approach to solving (3.1.1), (3.1.2) is to use an operator-splitting/alternating-direction approach, *cf.* (1.4.4), (1.4.5). This is the approach that we pursue in this chapter. The Galerkin spectral method on  $D$  that was developed in Chapter 2 will be used to solve (1.4.4), and a finite element method for (1.4.5) will also be introduced. A finite element method is convenient for the  $\underline{x}$ -direction solver because the physical space domain,  $\Omega$ , need not have simple geometry. As in Chapter 2, all of the analysis in this chapter is valid for any spring potential that satisfies Hypotheses A and B, but in the computational results section we consider the FENE model only.

We propose a fully-practical alternating-direction Galerkin method for (3.1.1). The approach is similar in spirit to the alternating-direction method used by Chauvière & Lozinski in [32, 33, 91]. However, there are some important theoretical questions related to applying

alternating-direction methods in this context, which have not previously been addressed in the literature, and we focus on these questions in this chapter. In particular, we consider the stability and convergence analysis of our alternating-direction scheme for (3.1.1) in Sections 3.4, 3.5, 3.6 and 3.7. It is not obvious *a priori* what effect applying a splitting of the form (1.4.4), (1.4.5) will have on a discretisation of (3.1.1), and therefore it is important to rigorously establish the stability and convergence properties of the alternating-direction numerical methods developed here.

The reader will note that the alternating-direction method under consideration here is nonstandard in the sense we consider  $d$ -dimensional cross-sections (instead of one-dimensional cross-sections) of  $\Omega \times D$ . This poses a formidable computational challenge because, as shall be seen in Section 3.3, we typically need to solve a large number problems posed in  $d$  spatial dimensions in each time-step. However, the method is extremely well suited to implementation on a parallel architecture since the  $q$ -direction solves are completely independent from one another, and similarly the  $x$ -direction solves are decoupled also. We discuss the parallel implementation of our alternating-direction scheme in Section 3.8, and our computational results in Section 3.9 were obtained using this parallel implementation.

The structure of this chapter is as follows. The weak formulation of the full Fokker–Planck equation is discussed in Section 3.2. We then introduce a quadrature-based alternating-direction procedure in Section 3.3 and derive stability results for this scheme in Section 3.4. Using the approximation results in Section 3.6, we then derive convergence estimates in Section 3.7. The implementation of the numerical method is described in Section 3.8, and in Section 3.9, numerical results for the FENE Fokker–Planck equation are presented in the simplified case that the macroscopic velocity,  $u$ , is taken to be a constant-in-time vector field.

## 3.2 Weak formulation and spatial discretisation

The full Fokker–Planck equation considered in this chapter depends on  $x \in \Omega$  as well as  $q \in D$ , and therefore we will require the use of slightly different function spaces than in Chapter 2. Let  $L^2(\Omega \times D)$  be defined in the obvious way, and let  $(\cdot, \cdot)$  and  $\|\cdot\|$  denote the  $L^2$  inner-product and norm over  $\Omega \times D$ :

$$(f, g) := \int_{\Omega \times D} f(x, q)g(x, q) \, dx \, dq \quad \text{and} \quad \|f\|^2 := (f, f).$$

We assume throughout this chapter that  $u$  is a divergence-free  $d$ -component vector function, *i.e.*

$$\nabla_x \cdot u(x, t) = 0 \quad \text{for a.e. } (x, t) \in \Omega \times (0, T]. \quad (3.2.1)$$

It would be straightforward to adapt the arguments in this chapter to the case where  $u$  is not divergence free, but this would make the analysis more messy and it would shed no further light on the properties of the numerical methods under consideration. Therefore in the interests of clarity and brevity, in this chapter we restrict our attention to the case when (3.2.1) is satisfied.

Also, we suppose that

$$u \in L^\infty(0, T; \mathbb{L}^\infty(\Omega)) \quad \text{and} \quad \nabla_x u = \kappa \in W^{1, \infty}(0, T; \mathbb{L}^\infty(\Omega)), \quad (3.2.2)$$

where, to simplify notation, we do not explicitly label the  $d$  or  $d \times d$  dimensionality of the function spaces for  $\underline{u}(\underline{x}, t) \in \mathbb{R}^d$  and  $\underline{\kappa}(\underline{x}, t) \in \mathbb{R}^{d \times d}$ . The assumption in (3.2.2) for  $\underline{\kappa}$  is stronger than the assumptions in Chapter 2; recall that in Lemma 2.3.1 and Theorem 2.3.2 we required  $\underline{\kappa} \in \underline{\mathbb{C}}[0, T]$  and in Lemma 2.4.1 we required  $\underline{\kappa} \in \underline{\mathbb{H}}^1(0, T)$ .

We shall also use the following space:

$$\mathcal{X} := \left\{ \varphi \in L^2(\Omega \times D) : \varphi \in L^2(\Omega; H_0^1(D; M)) \cap H^1(\Omega; L^2(D)) \right\},$$

equipped with the following norm:

$$\|\varphi\|_{\mathcal{X}} := \left\{ \int_{\Omega \times D} (|\varphi|^2 + |\nabla_M \varphi|^2) \, d\underline{x} \, d\underline{q} \right\}^{\frac{1}{2}}.$$

We note that the integrand in the definition of  $\|\varphi\|_{\mathcal{X}}$  does not include  $|\nabla_x \varphi|^2$ ; this is intentional.

Employing the substitution  $\hat{\psi} = \psi/\sqrt{M}$  that was used in Chapter 2, the weak formulation of (3.1.1) is as follows: Given  $\hat{\psi}^0 \in L^2(\Omega \times D)$ , find  $\hat{\psi} \in L^\infty(0, T; L^2(\Omega \times D)) \cap L^2(0, T; \mathcal{X})$  such that

$$\frac{d}{dt}(\hat{\psi}, \zeta) + (\underline{u} \cdot \nabla_x \hat{\psi}, \zeta) - (\underline{\kappa} \underline{q} \hat{\psi}, \nabla_M \zeta) + \frac{1}{2\text{Wi}} (\nabla_M \hat{\psi}, \nabla_M \zeta) = 0 \quad \forall \zeta \in \mathcal{X}, \quad (3.2.3)$$

$$\hat{\psi}(\underline{x}, \underline{q}, 0) = \hat{\psi}^0(\underline{x}, \underline{q}), \quad (\underline{x}, \underline{q}) \in \Omega \times D, \quad (3.2.4)$$

in the sense of distributions on  $(0, T)$ . Following Chapter 2, we weakly imposed the boundary condition (1.3.24) on  $\Omega \times \partial D$  for  $t \in (0, T]$ . For simplicity, we avoid boundary conditions on  $\partial\Omega \times D$  by assuming that the macroscopic velocity field is an *enclosed flow*, i.e. that

$$\underline{u} \cdot \underline{n}_{\partial\Omega} = 0 \text{ on } \partial\Omega, \quad (3.2.5)$$

where  $\underline{n}_{\partial\Omega} \in \mathbb{R}^d$  is the unit outward normal to  $\partial\Omega$ . Also, the initial condition (3.2.4) is understood to be imposed in a weak sense and, as in Chapter 2,  $\psi$  is recovered by multiplying  $\hat{\psi}$  by  $\sqrt{M}$ .

The term containing  $\underline{\kappa}$  in (3.2.3) will be of particular interest since, as we shall see, it is the most difficult term to treat using an alternating-direction method. We introduce the following bilinear form notation for this term, which will be convenient later on:

$$C(\underline{\kappa}; f, g) := (\underline{\kappa} \underline{q} f, \nabla_M g). \quad (3.2.6)$$

Next, we establish a statement analogous to Lemma 1.3.3 for the weak solution of (3.2.3). Recall that

$$\varrho(\underline{x}, t) := \int_D \psi(\underline{x}, \underline{q}, t) \, d\underline{q} = \int_D \sqrt{M(\underline{q})} \hat{\psi}(\underline{x}, \underline{q}, t) \, d\underline{q}.$$

Noting from Hypothesis B in Chapter 2 that  $\sqrt{M} \in H_0^1(D) \subset H_0^1(D; M)$ , we set  $\zeta = \sqrt{M}$  in (3.2.3), to obtain

$$\left( \frac{\partial \hat{\psi}}{\partial t} + \underline{u} \cdot \nabla_x \hat{\psi}, \sqrt{M} \right) = \left( \frac{\partial \psi}{\partial t} + \underline{u} \cdot \nabla_x \psi, 1 \right) = \int_{\Omega} \left( \frac{\partial \varrho}{\partial t} + \underline{u} \cdot \nabla_x \varrho \right) \, d\underline{x} = 0. \quad (3.2.7)$$

Due to (3.2.5), the material volume  $\Omega$  does not change with time and therefore applying the Reynolds transport theorem as in Lemma 1.3.3, we obtain,

$$\frac{d}{dt} \int_{\Omega} \varrho(\underline{x}, t) d\underline{x} = 0, \quad (3.2.8)$$

or equivalently,  $\int_{\Omega} \varrho(\underline{x}, t) d\underline{x} = \int_{\Omega} \varrho^0(\underline{x}) d\underline{x}$  for  $t \in (0, T]$ .

**Remark 3.2.1** By taking test functions of the form  $\zeta = \chi_S \sqrt{M}$ , where  $\chi_S$  is a mollified characteristic function for  $S \subset \Omega$ , one could extend the above result to arbitrary subsets of  $\Omega$  and therefore recover Lemma 1.3.3 in its full generality for the weak solution.  $\diamond$

We now introduce the spatial discretisation of (3.2.3), (3.2.4). Let  $V_h$  be a  $N_{\Omega}$ -dimensional  $H^1(\Omega)$ -conforming finite element space corresponding to a triangulation  $\mathcal{T}_h$  of  $\Omega$ . Also, as in Chapter 2, let  $\mathcal{P}_N(D) \subset H_0^1(D) \subset H_0^1(D; M)$  be an  $N_D$ -dimensional space spanned by a set of spectral basis functions on  $D$  (such as  $\mathcal{A}, \mathcal{B}$  or  $\mathcal{C}$  from Section 2.8). Noting that  $V_h \otimes \mathcal{P}_N(D) \subset \mathcal{X}$ , we obtain a spatially discrete formulation of the full Fokker–Planck equation as follows:

Let  $\hat{\psi}_{h,N}(\cdot, \cdot, 0) \in V_h \otimes \mathcal{P}_N(D)$  be the  $L^2(\Omega \times D)$  projection of  $\hat{\psi}^0$  onto  $V_h \otimes \mathcal{P}_N(D)$ . Find  $\hat{\psi}_{h,N}(\cdot, \cdot, t) \in V_h \otimes \mathcal{P}_N(D)$ ,  $t \in (0, T]$  satisfying (3.2.3) for all  $\zeta \in V_h \otimes \mathcal{P}_N(D)$  in the sense of distributions on  $(0, T)$ .

It would be possible to finite difference in time the spatially discrete formulation defined above in order to obtain a fully-discrete numerical method. However, this would be impractical in the present context because the discrete problem at each time-level would be posed on the domain  $\Omega \times D$ . As we have indicated, a more reasonable alternative is to use an alternating-direction method to split each  $2d$ -dimensional solve into a sequence of  $d$ -dimensional solves. This idea is considered in detail in the next section.

### 3.3 The alternating-direction numerical method

We begin this section by presenting a brief general overview of alternating-direction methods and we will then consider how to derive an alternating-direction method for (3.2.3), (3.2.4).

We concentrate on schemes that use a Galerkin spatial discretisation since this will allow us to use arguments analogous to those in Sections 2.3 and 2.5 in order to establish stability and convergence properties. The seminal work on alternating-direction methods of this type is by Douglas & Dupont [41]. In the example below, we illustrate the approach of Douglas & Dupont by considering a Galerkin-based alternating-direction method for the constant-coefficient heat equation in two spatial dimensions.

**Example 3.3.1** Suppose  $(x, y, t) \in (a_1, a_2) \times (b_1, b_2) \times (0, T) \mapsto u(x, y, t) \in \mathbb{R}$ , with  $u(\cdot, \cdot, 0) = u^0(\cdot, \cdot)$  and

$$\frac{\partial u}{\partial t} - \Delta u = 0, \quad \text{on } (x, y, t) \in (a_1, a_2) \times (b_1, b_2) \times (0, T),$$

with homogeneous Dirichlet boundary conditions in space. The corresponding weak formulation of this problem is:

Find  $u \in L^\infty(0, T; L^2((a_1, b_1) \times (a_2, b_2))) \cap L^2(0, T; H_0^1((a_1, b_1) \times (a_2, b_2)))$  satisfying

$$\int_{\Omega} \frac{\partial u}{\partial t} v \, dx \, dy + \int_{\Omega} \nabla_x u \cdot \nabla_x v \, dx \, dy = 0 \quad \forall v \in H_0^1((a_1, b_1) \times (a_2, b_2)), \quad (3.3.1)$$

$$u(x, y, 0) = u^0(x, y), \quad (x, y) \in (a_1, a_2) \times (b_1, b_2), \quad (3.3.2)$$

in the sense of distributions on  $(0, T)$ .

Suppose that  $X_h$  and  $Y_h$  are  $H_0^1(a_1, b_1)$ - and  $H_0^1(a_2, b_2)$ -conforming finite element spaces, respectively, with bases  $\{v_i \in X_h : 1 \leq i \leq N\}$  and  $\{w_i \in Y_h : 1 \leq i \leq N\}$  such that  $X_h = \text{span}(\{v_i\}_{1 \leq i \leq N})$  and  $Y_h = \text{span}(\{w_i\}_{1 \leq i \leq N})$ . Let  $X_h \otimes Y_h$  denote the following tensor product space:

$$X_h \otimes Y_h := \left\{ z : z = \sum_{i,j=1}^N \alpha_{ij} v_i w_j, \alpha_{ij} \in \mathbb{R} \text{ for each } 1 \leq i, j \leq N \right\}.$$

It follows that  $X_h \otimes Y_h \subset H_0^1(a_1, b_1; H_0^1(a_2, b_2)) \subset H_0^1((a_1, b_1) \times (a_2, b_2))$ . Using this tensor product finite element space we define a finite element scheme for this problem by replacing  $H_0^1((a_1, b_1) \times (a_2, b_2))$  with  $X_h \otimes Y_h$  in the weak formulation above. Also, supposing we employ Crank–Nicolson finite differencing to discretise (3.3.1) in time, then we obtain the following fully discrete problem (written in matrix form) at each time-step: Given  $u_h^n \in X_h \otimes Y_h$ , find  $u_h^{n+1} \in X_h \otimes Y_h$  satisfying

$$\begin{aligned} & \left( M_x \otimes M_y + \frac{\Delta t}{2} (S_x \otimes M_y + M_x \otimes S_y) \right) u_h^{n+1} \\ & = \left( M_x \otimes M_y - \frac{\Delta t}{2} (S_x \otimes M_y + M_x \otimes S_y) \right) u_h^n, \end{aligned} \quad (3.3.3)$$

where  $M_x$  and  $S_x$  (resp.  $M_y$  and  $S_y$ ) are the  $X_h$  (resp.  $Y_h$ ) mass and stiffness matrices, and the matrix tensor product<sup>1</sup> is defined as follows for matrices  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{p \times q}$ :

$$A \otimes B = \begin{bmatrix} a_{11}B & \dots & a_{1n}B \\ \vdots & \ddots & \vdots \\ a_{m1}B & \dots & a_{mn}B \end{bmatrix} \in \mathbb{R}^{mp \times nq}.$$

Since the matrices in (3.3.3) are tensor products of the  $x$ - and  $y$ -direction discretisation matrices, we can approximate (3.3.3) using the following two stage method:

$$\left( M_x + \frac{\Delta t}{2} S_x \right) \otimes I u_h^{n*} = \left( M_x - \frac{\Delta t}{2} S_x \right) \otimes I u_h^n \quad (3.3.4)$$

$$I \otimes \left( M_y + \frac{\Delta t}{2} S_y \right) u_h^{n+1} = I \otimes \left( M_y - \frac{\Delta t}{2} S_y \right) u_h^{n*}. \quad (3.3.5)$$

These equations define the fully discrete Galerkin alternating-direction method for this problem. We refer to (3.3.4) as the  $x$ -direction stage and to (3.3.5) as the  $y$ -direction stage.

<sup>1</sup>Also referred to as the Kronecker product.

By multiplying (3.3.4) by  $I \otimes (M_y - \Delta t/2S_y)$  and (3.3.5) by  $(M_x + \Delta t/2S_x) \otimes I$ , we see that the Galerkin alternating-direction method is equivalent to the following:

$$\begin{aligned} & \left( M_x \otimes M_y + \frac{\Delta t}{2} (S_x \otimes M_y + M_x \otimes S_y) + \frac{(\Delta t)^2}{4} S_x \otimes S_y \right) u_h^{n+1} \\ &= \left( M_x \otimes M_y - \frac{\Delta t}{2} (S_x \otimes M_y + M_x \otimes S_y) + \frac{(\Delta t)^2}{4} S_x \otimes S_y \right) u_h^n. \end{aligned} \quad (3.3.6)$$

This is referred to as the equivalent one-step method for (3.3.4), (3.3.5). We can see that the one-step method is identical to the Crank-Nicolson scheme, (3.3.3), except for the presence of the  $\frac{1}{4}(\Delta t)^2 S_x \otimes S_y$  perturbation terms in (3.3.6).

Using the approach of Douglas & Dupont, the next step is to rewrite (3.3.6) in inner product form as follows: Given  $u_h^n \in X_h \otimes Y_h$ , find  $u_h^{n+1} \in X_h \otimes Y_h$  satisfying

$$\begin{aligned} & \int_{\Omega} \frac{u_h^{n+1} - u_h^n}{\Delta t} v_h \, dx \, dy + \frac{1}{2} \int_{\Omega} \left\{ \nabla_x u_h^{n+1} \cdot \nabla_x v_h + \frac{\Delta t}{2} \left( \frac{\partial u_h^{n+1}}{\partial x} \frac{\partial v_h}{\partial y} + \frac{\partial u_h^{n+1}}{\partial y} \frac{\partial v_h}{\partial x} \right) \right\} dx \, dy \\ &= \frac{1}{2} \int_{\Omega} \left\{ -\nabla_x u_h^n \cdot \nabla_x v_h + \frac{\Delta t}{2} \left( \frac{\partial u_h^n}{\partial x} \frac{\partial v_h}{\partial y} + \frac{\partial u_h^n}{\partial y} \frac{\partial v_h}{\partial x} \right) \right\} dx \, dy \end{aligned} \quad (3.3.7)$$

for all  $v_h \in X_h \otimes Y_h$ . From here, one can use standard energy analysis to establish stability and convergence properties of (3.3.7), and therefore, equivalently, of (3.3.4), (3.3.5).

We now apply the approach described in Example 3.3.1 to the weak formulation, (3.2.3). First of all, define the bases

$$\{Y_k \in \mathcal{P}_N(D) : 1 \leq k \leq N_D\} \quad \text{and} \quad \{X_i \in V_h : 1 \leq i \leq N_{\Omega}\}, \quad (3.3.8)$$

such that  $\text{span}(\{Y_k\}_{1 \leq k \leq N_D}) = \mathcal{P}_N(D)$  and  $\text{span}(\{X_i\}_{1 \leq i \leq N_{\Omega}}) = V_h$ . Recalling (2.8.4), we define  $M_q, S_q \in \mathbb{R}^{N_D \times N_D}$  as

$$(M_q)_{lk} := \int_D Y_k(\underline{q}) Y_l(\underline{q}) \, d\underline{q}, \quad (3.3.9)$$

$$(S_q)_{lk} := \int_D \nabla_M Y_k(\underline{q}) \cdot \nabla_M Y_l(\underline{q}) \, d\underline{q}. \quad (3.3.10)$$

Similarly,  $M_x, T_x \in \mathbb{R}^{N_{\Omega} \times N_{\Omega}}$  are defined as follows:

$$(M_x)_{ij} := \int_{\Omega} X_i(\underline{x}) X_j(\underline{x}) \, d\underline{x}, \quad (3.3.11)$$

$$(T_x)_{ij} := \int_{\Omega} (\underline{y} \cdot \nabla_x X_j(\underline{x})) X_i(\underline{x}) \, d\underline{x}. \quad (3.3.12)$$

A fully discrete form of (3.2.3) using a backward-Euler time discretisation can be written as follows: Given  $\hat{\psi}_N^n = \sum_{jl} \gamma_{jl}^n X_j Y_l \in V_h \otimes \mathcal{P}_N(D)$ , find the vector  $\underline{\gamma}^{n+1} \in \mathbb{R}^{N_D N_{\Omega}}$ , defining a function  $\hat{\psi}_N^{n+1} = \sum_{jl} \gamma_{jl}^{n+1} X_j Y_l \in V_h \otimes \mathcal{P}_N(D)$ , such that

$$\begin{aligned} M_x \otimes M_q \left( \frac{\underline{\gamma}^{n+1} - \underline{\gamma}^n}{\Delta t} \right) + T_x \otimes M_q \underline{\gamma}^{n+1} + \frac{1}{2W_i} M_x \otimes S_q \underline{\gamma}^{n+1} \\ - C(\underline{\mathfrak{k}}^{n+1}; \hat{\psi}_N^{n+1}, \zeta_{ik}) = 0, \end{aligned} \quad (3.3.13)$$

where  $\zeta_{ik} = X_i \times Y_k \in V_h \otimes \mathcal{P}_N(D)$ . It is also possible to obtain a tensor product form discretisation matrix of  $C(\underline{\kappa}; \cdot, \cdot)$ , *i.e.* consider  $C(\underline{\kappa}; \zeta_{jl}, \zeta_{ik})$  as follows:

$$\begin{aligned} C(\underline{\kappa}; \zeta_{jl}, \zeta_{ik}) &:= \int_{\Omega \times D} \left( \underline{\kappa}^{n+1}(\underline{x}) \underline{q} X_j(\underline{x}) Y_l(\underline{q}) \right) \cdot \nabla_M (X_i(\underline{x}) Y_k(\underline{q})) \, d\underline{x} \, d\underline{q} \\ &= \sum_{s,t=1}^d \left( \int_{\Omega} \kappa_{st}^{n+1}(\underline{x}) X_i(\underline{x}) X_j(\underline{x}) \, d\underline{x} \right) \left( \int_D q_t Y_l(\underline{q}) \sqrt{M} \frac{\partial}{\partial q_s} \left( \frac{Y_k(\underline{q})}{\sqrt{M}} \right) \, d\underline{q} \right). \end{aligned}$$

Therefore, we define the matrices  $C_x^{st} \in \mathbb{R}^{N_\Omega \times N_\Omega}$  and  $C_q^{st} \in \mathbb{R}^{N_D \times N_D}$  for  $1 \leq s, t \leq d$  such that

$$(C_x^{st})_{ij} := \int_{\Omega} \kappa_{st}^{n+1}(\underline{x}) X_i(\underline{x}) X_j(\underline{x}) \, d\underline{x}, \quad (3.3.14)$$

$$(C_q^{st})_{kl} := \int_D q_t Y_l(\underline{q}) \sqrt{M} \frac{\partial}{\partial q_s} \left( \frac{Y_k(\underline{q})}{\sqrt{M}} \right) \, d\underline{q}. \quad (3.3.15)$$

Hence, we can rewrite the term on the final line of (3.3.13) as  $\sum_{s,t=1}^d C_x^{st} \otimes C_q^{st} \underline{\gamma}^{n+1}$ .

However, since this matrix expression for  $C(\underline{\kappa}; \cdot, \cdot)$  contains neither  $M_x$  nor  $M_q$ , we can no longer factorise the resulting equation in the same way as in (3.3.4), (3.3.5). That is, the term  $C(\underline{\kappa}; \cdot, \cdot)$  causes difficulties because its ‘coefficient’,  $\underline{\kappa}(\underline{x})\underline{q}$ , depends on both the  $\underline{x}$ - and  $\underline{q}$ -directions.

This issue has been considered a number of times in the literature. For example, in the context of collocation-based alternating-direction schemes Celia & Pinder [29, 30] and Bialecki & Fernandes [20] developed methods that could handle equations with general variable coefficients. However, as indicated earlier, our focus is on developing a Galerkin-based framework, and therefore, again, the work of Douglas & Dupont is the most relevant here. In [41], Douglas & Dupont developed a ‘‘Laplace modification’’ scheme for the heat equation with general coefficients which involved discretising the equation

$$\frac{\partial u}{\partial t} = \nabla_x \cdot (a(x, y, t, u) \nabla_x u) + f(x, y, t, u),$$

as follows,

$$\left( \frac{u^{n+1} - u^n}{\Delta t}, v \right) + (a^n(u^n) \nabla_x u^n, \nabla_x v) + \lambda (\nabla_x (u^{n+1} - u^n), \nabla_x v) = (f^n(u^n), v),$$

where  $\lambda$  is a constant scalar, which must satisfy a lower bound condition related to the supremum of  $|a|$  in order to ensure the stability of the numerical method. This discretisation then allows the use of a standard Galerkin alternating-direction method, as in Example 3.3.1, because the term containing  $a$  can be moved to the right-hand side and treated as a source term.

However, it is not obvious how to apply this kind of approach to (3.3.13), because our problematic term is a convection term rather than a diffusion term. The most natural idea in the spirit of Douglas & Dupont would be to move the  $C(\underline{\kappa}; \cdot, \cdot)$  term to the right-hand side of (3.3.13) and treat it explicitly in time. This idea is feasible, but for the purposes of practical computations, we would like to have the option of using a fully-implicit temporal

discretisation. Indeed, the numerical results in Section 2.8.2 demonstrated that the semi-implicit temporal discretisation of the Fokker–Planck equation in which the term  $C(\underline{\kappa}; \cdot, \cdot)$  was treated explicitly in time was less stable than the backward Euler discretisation, especially for problems in which the product  $\text{Wi} \|\underline{\kappa}\|_{L^\infty(0,T;L^\infty(\Omega))}$  is significantly larger than 1.

In order to circumvent this limitation, we develop a Galerkin alternating-direction approach that is an amalgamation of the Douglas & Dupont framework and a new quadrature-based method. Using this approach, we can define either a fully-implicit in time or a semi-implicit in time alternating-direction method for the Fokker–Planck equation. We shall consider both options in detail in this chapter.

### 3.3.1 The hybrid alternating-direction scheme

The first ingredient of this scheme is a quadrature rule on  $\Omega$ .

Let  $\{(\underline{x}_m, w_m), w_m > 0, \underline{x}_m \in \bar{\Omega}, m = 1, \dots, Q_\Omega\}$  define an element-based quadrature rule on the triangulation  $\mathcal{T}_h$ , where the  $\underline{x}_m$  are the quadrature points and the  $w_m$  are the corresponding weights. Therefore, for functions  $f, g \in C^0(\Omega)$ , the quadrature sum is evaluated element-wise as follows,

$$\sum_{m=1}^{Q_\Omega} w_m f(\underline{x}_m) g(\underline{x}_m) = \sum_{K \in \mathcal{T}_h} \sum_{l=1}^{Q_K} w_l^K f(\underline{x}_l^K) g(\underline{x}_l^K), \quad (3.3.16)$$

where  $Q_K$  is the number of quadrature points in element  $K$ . From now on, we will use the left-hand side of (3.3.16) as a shorthand for the right-hand side.

We now introduce two alternative hypotheses on the accuracy of the quadrature rule, Quadrature Hypothesis 1 (QH1) and Quadrature Hypothesis 2 (QH2).

**Quadrature Hypothesis 1 (QH1).** The quadrature rule satisfies

$$\sum_{m=1}^{Q_\Omega} w_m \kappa_{ij}(\underline{x}_m) f(\underline{x}_m) g(\underline{x}_m) = \int_{\Omega} \kappa_{ij}(\underline{x}) f(\underline{x}) g(\underline{x}) \, d\underline{x}, \quad (3.3.17)$$

for all  $f, g \in V_h$  and for each component  $\kappa_{ij}$  of  $\underline{\kappa}$ .  $\diamond$

As discussed in Chapter 4, in the context of the Navier–Stokes–Fokker–Planck system, we compute the macroscopic velocity field,  $\underline{u}$  by solving the Navier–Stokes equations using a finite element method on the triangulation  $\mathcal{T}_h$ , *i.e.* the same triangulation that is used for the alternating-direction method for the Fokker–Planck equation. As a result, it is reasonable to assume that the components of  $\underline{\kappa} = \nabla_{\underline{x}} \underline{u}$  are represented by piecewise polynomials on  $\mathcal{T}_h$  and in this case it is certainly possible to satisfy QH1 by choosing an appropriate element-based quadrature rule.

**Quadrature Hypothesis 2 (QH2).** The quadrature rule satisfies

$$\sum_{m=1}^{Q_\Omega} w_m f(\underline{x}_m) g(\underline{x}_m) = \int_{\Omega} f(\underline{x}) g(\underline{x}) \, d\underline{x}, \quad (3.3.18)$$

for all  $f, g \in V_h$ .  $\diamond$

QH1 is a stronger hypothesis than QH2, and therefore in general we will require a larger value of  $Q_\Omega$  in order to satisfy QH1. Some results in the following analysis will require

QH1, whereas for others, QH2 will suffice. Refer to Section 3.8 for a discussion of specific quadrature rules that we use to satisfy QH1 and QH2 in practice.

Next, let  $\hat{\psi}_{h,N} \in V_h \otimes \mathcal{P}_N(D)$  denote the numerical solution of the full Fokker–Planck equation. Recalling the bases from (3.3.8),  $\hat{\psi}_{h,N}$  can be written in terms of coefficients  $\{\hat{\psi}_{ik}\}$  as follows:

$$\hat{\psi}_{h,N} := \sum_{i=1}^{N_\Omega} \sum_{k=1}^{N_D} \hat{\psi}_{ik} X_i Y_k \in V_h \otimes \mathcal{P}_N(D). \quad (3.3.19)$$

Define the *line functions*,  $\hat{\psi}_k$ , for  $k = 1, \dots, N_D$ , as follows:

$$\hat{\psi}_k := \sum_{i=1}^{N_\Omega} \hat{\psi}_{ik} X_i \in V_h, \quad (3.3.20)$$

and note that (3.3.19) can be rewritten using (3.3.20) as follows:

$$\hat{\psi}_{h,N}(\underline{x}, \underline{q}) = \sum_{k=1}^{N_D} \hat{\psi}_k(\underline{x}) Y_k(\underline{q}). \quad (3.3.21)$$

The formula (3.3.21) shall be useful in the discussion of the alternating-direction methods below.

As discussed above, the term  $C(\underline{\kappa}; \cdot, \cdot)$  is the most problematic in terms of applying an alternating-direction method to the Fokker–Planck equation. Therefore we begin by considering how to use a quadrature-based scheme to derive an alternating-direction type of formulation of this term.

Suppose that QH1 is satisfied and that we have the line function decomposition (3.3.21) for  $\hat{\psi}_{h,N}$ , in which  $\hat{\psi}_k \in V_h$  for  $k = 1, \dots, N_D$ . Also, let  $\zeta = X \times Y \in V_h \otimes \mathcal{P}_N(D)$ . Then,

$$\begin{aligned} C(\underline{\kappa}; \hat{\psi}_{h,N}, \zeta) &= \int_{\Omega \times D} (\underline{\kappa} \underline{q} \hat{\psi}_{h,N}(\underline{x}, \underline{q})) \cdot \nabla_M \zeta(\underline{x}, \underline{q}) \, d\underline{q} \, d\underline{x} \\ &= \int_D \sum_{k=1}^{N_D} \int_{\Omega} [\underline{\kappa} \underline{q} \hat{\psi}_k(\underline{x}) Y_k(\underline{q})] \cdot \nabla_M (X(\underline{x}) Y(\underline{q})) \, d\underline{x} \, d\underline{q} \\ &= \int_D \sum_{k=1}^{N_D} \sum_{m=1}^{Q_\Omega} w_m [\underline{\kappa}(\underline{x}_m) \underline{q} \hat{\psi}_k(\underline{x}_m) Y_k(\underline{q})] \cdot \nabla_M (X(\underline{x}_m) Y(\underline{q})) \, d\underline{q} \\ &= \sum_{m=1}^{Q_\Omega} w_m X(\underline{x}_m) \left\{ \sum_{k=1}^{N_D} \hat{\psi}_k(\underline{x}_m) \left( \int_D (\underline{\kappa}(\underline{x}_m) \underline{q} Y_k(\underline{q})) \cdot \nabla_M Y(\underline{q}) \, d\underline{q} \right) \right\} \end{aligned} \quad (3.3.22)$$

This shows the equivalence between the Galerkin formulation of  $C(\underline{\kappa}; \cdot, \cdot)$  on  $\Omega \times D$  and the quadrature sum over  $m = 1, \dots, Q_\Omega$  of the term

$$\sum_{k=1}^{N_D} \hat{\psi}_k(\underline{x}_m) \left( \int_D (\underline{\kappa}(\underline{x}_m) \underline{q} Y_k(\underline{q})) \cdot \nabla_M Y(\underline{q}) \, d\underline{q} \right), \quad (3.3.23)$$

which is the  $q$ -direction discretisation of  $C(\underline{\kappa}; \cdot, \cdot)$ .

Note that (3.3.23) is exactly the discretisation of the  $q$ -convection term that was used in the spectral method in Chapter 2, except that now  $\underline{\kappa}$  depends on  $\underline{x} \in \Omega$ , and we sample

$\underline{\kappa}$  at the quadrature points  $\underline{x}_m$ . Also, the coefficient vector in (3.3.23) corresponding to the quadrature point  $\underline{x}_m$  is the set of sampled line functions,  $\hat{\psi}_k(\underline{x}_m)$ ,  $k = 1, \dots, N_D$ .

The preceding discussion relied on QH1, however we can use an analogous argument when only QH2 is assumed, in which case the quadrature rule is no longer exact for the  $\underline{\kappa}$ -weighted integral in (3.3.17) and therefore we do not have equality between the second and third lines of (3.3.22). Instead, a quadrature error,  $E$ , is introduced as follows:

$$\sum_{m=1}^{Q_\Omega} w_m \kappa_{ij}(\underline{x}_m) \hat{\psi}_k(\underline{x}_m) X(\underline{x}_m) = \int_{\Omega} \kappa_{ij}(\underline{x}) \hat{\psi}_k(\underline{x}) X(\underline{x}) d\underline{x} + E(\kappa_{ij}, \hat{\psi}_k, X). \quad (3.3.24)$$

Modifying (3.3.22) to include this error term, we obtain:

$$\begin{aligned} \sum_{m=1}^{Q_\Omega} w_m X(\underline{x}_m) \left\{ \sum_{k=1}^{N_D} \hat{\psi}_k(\underline{x}_m) \left( \int_D (\underline{\kappa}(\underline{x}_m) \underline{q} Y_k(\underline{q})) \cdot \nabla_M X(\underline{q}) d\underline{q} \right) \right\} \\ = C(\underline{\kappa}; \hat{\psi}_h, N, \zeta) + \sum_{k=1}^{N_D} \int_D \underline{E}(\underline{\kappa}, \hat{\psi}_k, X) \underline{q} Y_k(\underline{q}) \cdot \nabla_M Y(\underline{q}) d\underline{q}, \end{aligned} \quad (3.3.25)$$

where  $\left( \underline{E}(\underline{\kappa}, \hat{\psi}_k, X) \right)_{ij} := E(\kappa_{ij}, \hat{\psi}_k, X)$ . Of course, the precise nature of  $\underline{E}$  will depend on the choice of quadrature rule and the problem at hand. Nevertheless, if appropriate hypotheses on the rate of decay of  $\underline{E}$  are specified, it would be possible to consider the stability and convergence properties of an alternating-direction method that includes a quadrature error term of this form. However, for simplicity and brevity, we do not consider such quadrature error terms in the numerical analysis in this chapter. It is worth noting though that we develop a stability argument in Section 3.4 that only relies on QH2, and in which we do not need to consider quadrature error terms such as in (3.3.24).

It is clear from (3.3.22) that sampling functions at the quadrature points  $\{\underline{x}_m \in \bar{\Omega}, m = 1, \dots, Q_\Omega\}$  will play an important role in the alternating-direction methods we define below. We will also require a reconstruction operator, which maps from a set of values at the quadrature points to a function in  $V_h$ . We now introduce this operator. To simplify notation, we first define the following discrete inner product and norm over  $\Omega$  for  $\{f_m\}, \{g_m\} \in \mathbb{R}^{Q_\Omega}$ :

$$(\{f_m\}, \{g_m\})_{\ell^2(\Omega)} := \sum_{m=1}^{Q_\Omega} w_m f_m g_m, \quad \text{and} \quad \|\{f_m\}\|_{\ell^2(\Omega)} := (\{f_m\}, \{f_m\})_{\ell^2(\Omega)}^{\frac{1}{2}}. \quad (3.3.26)$$

Note that, by (3.3.17) or (3.3.18), for  $f, g \in V_h$ ,  $(\{f(\underline{x}_m)\}, \{g(\underline{x}_m)\})_{\ell^2(\Omega)} = (f, g)_{L^2(\Omega)}$ , where  $(\cdot, \cdot)_{L^2(\Omega)}$  is the standard  $L^2$  inner product on  $\Omega$ . Next we define the reconstruction operator  $\mathcal{R} : \{f_m\} \in \mathbb{R}^{Q_\Omega} \mapsto \mathcal{R}\{f_m\} \in V_h$  such that

$$(\mathcal{R}\{f_m\}, X)_{L^2(\Omega)} = (\{f_m\}, \{X(\underline{x}_m)\})_{\ell^2(\Omega)} \quad \forall X \in V_h. \quad (3.3.27)$$

**Remark 3.3.2** For any  $\mathcal{R}\{f_m\} \in V_h$ , there exist real numbers  $\gamma_1, \dots, \gamma_{N_\Omega}$  such that  $\mathcal{R}\{f_m\} = \sum_{j=1}^{N_\Omega} \gamma_j X_j$ . Letting  $X = X_i$ ,  $i = 1, \dots, N_\Omega$ , above it is clear that (3.3.27) is equivalent to the linear system  $M_x \underline{\gamma} = \underline{F}$  where  $M_x \in \mathbb{R}^{N_\Omega \times N_\Omega}$  is the  $V_h$  mass matrix,  $\underline{\gamma} = (\gamma_1, \dots, \gamma_{N_\Omega})^T$ , and  $\underline{F} \in \mathbb{R}^{N_\Omega}$  is such that  $F_i = (\{f_m\}, \{X_i(\underline{x}_m)\})_{\ell^2(\Omega)}$ . The matrix  $M_x$  is nonsingular, and therefore the reconstruction operator defined in (3.3.27) is well-defined.  $\diamond$

We are now in a position to discuss the alternating-direction Galerkin methods that are the focus of this chapter. We introduce two algorithms below, denoted method I and method II. Each method utilises a hybrid alternating-direction method, which combines the quadrature approach illustrated in (3.3.22) with a standard Douglas-Dupont type Galerkin alternating-direction method.

The distinction between method I and method II is that method I uses a semi-implicit spectral method in the  $q$ -direction (*i.e.* the term  $C(\underline{\kappa}; \cdot, \cdot)$  is treated explicitly in time) whereas method II uses a fully-implicit temporal discretisation.

### 3.3.2 Method I: Semi-implicit scheme

Method I is initialised by computing the  $L^2(\Omega \times D)$  projection,  $\hat{\psi}_{h,N}^0$ , of the initial datum  $\hat{\psi}^0 \in L^2(\Omega \times D)$  onto  $V_h \otimes \mathcal{P}_N(D)$ , so that  $\hat{\psi}_{h,N}^0 \in V_h \otimes \mathcal{P}_N(D)$ , satisfies

$$\left( \hat{\psi}^0, \zeta \right) = \left( \hat{\psi}_{h,N}^0, \zeta \right) \quad \text{for all } \zeta \in V_h \otimes \mathcal{P}_N(D). \quad (3.3.28)$$

Then, as in (1.4.4), (1.4.5), this alternating-direction method consists of two stages at each time-step: the  $q$ -direction stage and the  $x$ -direction stage. We begin with the  $q$ -direction stage, in which we essentially use the Galerkin spectral method in  $D$  from Chapter 2.

Suppose  $\hat{\psi}_{h,N}^n \in V_h \otimes \mathcal{P}_N(D)$ . Then in the  $q$ -direction stage we compute  $\hat{\psi}_{h,N}^{n*}(\underline{x}_m, \cdot) \in \mathcal{P}_N(D)$  for each  $m = 1, \dots, Q_\Omega$  satisfying

$$\begin{aligned} \int_D \frac{\hat{\psi}_{h,N}^{n*}(\underline{x}_m, \underline{q}) - \hat{\psi}_{h,N}^n(\underline{x}_m, \underline{q})}{\Delta t} Y_l(\underline{q}) \, d\underline{q} + \frac{1}{2\overline{\text{Wi}}} \int_D \nabla_M \hat{\psi}_{h,N}^{n*}(\underline{x}_m, \underline{q}) \cdot \nabla_M Y_l(\underline{q}) \, d\underline{q} \\ = \int_D (\underline{\kappa}^n(\underline{x}_m) \underline{q} \hat{\psi}_{h,N}^n(\underline{x}_m, \underline{q})) \cdot \nabla_M Y_l(\underline{q}) \, d\underline{q}, \end{aligned} \quad (3.3.29)$$

for  $l = 1, \dots, N_D$ . (3.3.29) defines an  $N_D \times N_D$  linear system at each quadrature point. In order to separate out the  $x$ - and  $q$ -direction dependencies more clearly, we rewrite this equation in terms of line functions using (3.3.20), *i.e.*:

$$\begin{aligned} \sum_{k=1}^{N_D} \hat{\psi}_k^{n*}(\underline{x}_m) \left( \int_D Y_k(\underline{q}) Y_l(\underline{q}) \, d\underline{q} + \frac{\Delta t}{2\overline{\text{Wi}}} \int_D \nabla_M Y_k(\underline{q}) \cdot \nabla_M Y_l(\underline{q}) \, d\underline{q} \right) \\ = \sum_{k=1}^{N_D} \hat{\psi}_k^n(\underline{x}_m) \left( \int_D Y_k(\underline{q}) Y_l(\underline{q}) \, d\underline{q} + \Delta t \int_D (\underline{\kappa}^n(\underline{x}_m) \underline{q} Y_k(\underline{q})) \cdot \nabla_M Y_l(\underline{q}) \, d\underline{q} \right), \end{aligned} \quad (3.3.30)$$

for  $l = 1, \dots, N_D$ . This system is solved at each quadrature point  $\underline{x}_m$ ,  $m = 1, \dots, Q_\Omega$ .

Equation (3.3.30) shows that in the  $q$ -direction stage, the sampled values of the line functions, *i.e.*  $\hat{\psi}_k^{n*}(\underline{x}_m)$ ,  $k = 1, \dots, N_D$ ,  $m = 1, \dots, Q_\Omega$ , are the coefficients to be computed. We determine these values by solving a different linear system at each quadrature point. Note that these linear systems are completely independent from one another. This independence enables parallel computation to be used very effectively in this context; this will be discussed in more detail later.

The  $q$ -direction stage is complete once the values  $\hat{\psi}_k^{n*}(\underline{x}_m)$ ,  $k = 1, \dots, N_D$ ,  $m = 1, \dots, Q_\Omega$  have been computed, and then we can begin solving in the  $x$ -direction. In the  $x$ -direction

stage, we use a finite element discretisation of the transport equation (1.4.5) to update the output data from the  $q$ -direction stage. That is, for a given  $k$ , we find  $\hat{\psi}_k^{n+1} \in V_h$ , satisfying:

$$\int_{\Omega} \hat{\psi}_k^{n+1} X_i \, d\tilde{x} + \Delta t \int_{\Omega} \left( y^{n+1} \cdot \nabla_x \hat{\psi}_k^{n+1} \right) X_i \, d\tilde{x} = \int_{\Omega} \mathcal{R}\{\hat{\psi}_k^{n*}(\tilde{x}_m)\} X_i \, d\tilde{x}, \quad (3.3.31)$$

for  $i = 1, \dots, N_{\Omega}$ .

Note, however, that based on (3.3.27), for the right-hand side in (3.3.31) we have:

$$\int_{\Omega} \mathcal{R}\{\hat{\psi}_k^{n*}(\tilde{x}_m)\} X_i \, d\tilde{x} = \sum_{m=1}^{Q_{\Omega}} w_m \hat{\psi}_k^{n*}(\tilde{x}_m) X_i(\tilde{x}_m) =: F_i. \quad (3.3.32)$$

Hence we do not actually have to explicitly compute  $\mathcal{R}\{\hat{\psi}_k^{n*}(\tilde{x}_m)\} \in V_h$  in order to solve (3.3.31), since it is equivalent to solve the following system:

$$\int_{\Omega} \hat{\psi}_k^{n+1} X_i \, d\tilde{x} + \Delta t \int_{\Omega} \left( y^{n+1} \cdot \nabla_x \hat{\psi}_k^{n+1} \right) X_i \, d\tilde{x} = F_i, \quad (3.3.33)$$

for  $i = 1, \dots, N_{\Omega}$ . We solve (3.3.33) for each  $k = 1, \dots, N_D$ , and, just as in the  $q$ -direction, these computations are decoupled from one another.

Once the  $\tilde{x}$ -direction computations are complete, we have the numerical solution at time level  $n + 1$ :

$$\hat{\psi}_{h,N}^{n+1} = \sum_{k=1}^{N_D} \hat{\psi}_k^{n+1} Y_k \in V_h \otimes \mathcal{P}_N(D).$$

Hence method I is defined by the initialisation (3.3.28), the  $q$ -direction spectral method (3.3.30) and the  $\tilde{x}$ -direction finite element method (3.3.33).

Before continuing further, we first verify that the  $q$ - and  $\tilde{x}$ -direction numerical methods are well-defined.

**Lemma 3.3.3** *Let  $A_q \in \mathbb{R}^{N_D \times N_D}$  denote the matrix appearing on the left-hand side of (3.3.30), i.e.*

$$A_q := M_q + \frac{\Delta t}{2W_i} S_q, \quad (3.3.34)$$

and let  $A_x \in \mathbb{R}^{N_{\Omega} \times N_{\Omega}}$  be the matrix from the left-hand side of (3.3.31),

$$A_x := M_x + \Delta t T_x. \quad (3.3.35)$$

The matrices  $A_q$  and  $A_x$  are nonsingular.

**Proof.** The result follows straightforwardly from the positive-definiteness of the bilinear forms,  $\mathfrak{B}_q(\cdot, \cdot) : \mathcal{P}_N(D) \times \mathcal{P}_N(D) \mapsto \mathbb{R}$ , and  $\mathfrak{B}_x(\cdot, \cdot) : V_h \times V_h \mapsto \mathbb{R}$ , defining  $A_q$  and  $A_x$  respectively.

Consider  $\mathfrak{B}_q(X, X)$  for any  $X \in \mathcal{P}_N(D) \setminus \{0\}$ :

$$\mathfrak{B}_q(X, X) = \|X\|_{L^2(D)}^2 + \frac{\Delta t}{2W_i} \|\nabla_M X\|_{L^2(D)}^2 \geq \|X\|_{L^2(D)}^2 > 0. \quad (3.3.36)$$

Similarly, integrating by parts and utilising the enclosed flow and divergence free assumptions for  $\mathfrak{B}_x(Y, Y)$  with  $Y \in V_h \setminus \{0\}$ , we have,

$$\mathfrak{B}_x(Y, Y) = \|Y\|_{L^2(\Omega)}^2 - \frac{\Delta t}{2} \int_{\Omega} (\nabla_x \cdot \underline{u}^{n+1}) Y^2 d\underline{x} = \|Y\|_{L^2(\Omega)}^2 > 0. \quad (3.3.37)$$

This completes the proof.  $\square$

In the next lemma we derive a Galerkin formulation posed on  $\Omega \times D$  for method I. This will allow us to apply arguments analogous to those in Chapter 2 to the numerical analysis of method I.

**Lemma 3.3.4** *Suppose the  $\underline{x}$ -direction quadrature rule satisfies QH1. Method I is equivalent to the following fully-discrete formulation:*

Given  $\hat{\psi}_{h,N}^0 \in V_h \otimes \mathcal{P}_N(D)$  defined as in (3.3.28), for each  $n = 0, \dots, N_T - 1$ ,  $\hat{\psi}_{h,N}^{n+1} \in V_h \otimes \mathcal{P}_N(D)$  satisfies

$$\begin{aligned} & \left( \frac{\hat{\psi}_{h,N}^{n+1} - \hat{\psi}_{h,N}^n}{\Delta t}, \zeta \right) + \left( \underline{u} \cdot \nabla_x \hat{\psi}_{h,N}^{n+1}, \zeta \right) + \frac{1}{2\text{Wi}} \left( \nabla_M \hat{\psi}_{h,N}^{n+1}, \nabla_M \zeta \right) \\ & + \frac{\Delta t}{2\text{Wi}} \left( \nabla_M \left( \underline{u} \cdot \nabla_x \hat{\psi}_{h,N}^{n+1} \right), \nabla_M \zeta \right) - \left( \underline{\kappa}^n \underline{q} \hat{\psi}_{h,N}^n, \nabla_M \zeta \right) = 0, \end{aligned} \quad (3.3.38)$$

for all  $\zeta \in V_h \otimes \mathcal{P}_N(D)$ .

**Proof.** Multiplying (3.3.30) through by  $X_i(\underline{x}_m)$ , where  $X_i \in V_h$ , and performing the weighted sum according to (3.3.16) gives,

$$\begin{aligned} & \sum_{k=1}^{N_D} (\{\hat{\psi}_k^{n*}(\underline{x}_m)\}, \{X_i(\underline{x}_m)\})_{\ell^2(\Omega)} \left( \int_D Y_k(\underline{q}) Y_l(\underline{q}) d\underline{q} + \frac{\Delta t}{2\text{Wi}} \int_D \nabla_M Y_k(\underline{q}) \cdot \nabla_M Y_l(\underline{q}) d\underline{q} \right) \\ & = \sum_{k=1}^{N_D} (\{\hat{\psi}_k^n(\underline{x}_m)\}, \{X_i(\underline{x}_m)\})_{\ell^2(\Omega)} \left( \int_D Y_k(\underline{q}) Y_l(\underline{q}) d\underline{q} \right) \\ & + \Delta t \sum_{m=1}^{Q_\Omega} w_m X_i(\underline{x}_m) \left\{ \sum_{k=1}^{N_D} \hat{\psi}_k^n(\underline{x}_m) \left( \int_D (\underline{\kappa}^n(\underline{x}_m) \underline{q} Y_k(\underline{q})) \cdot \nabla_M Y_l(\underline{q}) d\underline{q} \right) \right\}. \end{aligned} \quad (3.3.39)$$

Using the reconstruction operator, (3.3.27), with the  $\ell^2$  inner products and the argument of (3.3.22) on the term on the third line<sup>2</sup>, we obtain the following formulation for  $\mathcal{R}\hat{\psi}_{h,N}^{n*} \in V_h \otimes \mathcal{P}_N(D)$ ,

$$\begin{aligned} & \int_{\Omega \times D} \frac{\mathcal{R}\hat{\psi}_{h,N}^{n*}(\underline{x}, \underline{q}) - \hat{\psi}_{h,N}^n(\underline{x}, \underline{q})}{\Delta t} \zeta(\underline{x}, \underline{q}) d\underline{q} d\underline{x} + \frac{1}{2\text{Wi}} \int_{\Omega \times D} \nabla_M \mathcal{R}\hat{\psi}_{h,N}^{n*}(\underline{x}, \underline{q}) \cdot \nabla_M \zeta(\underline{x}, \underline{q}) d\underline{q} d\underline{x} \\ & = \int_{\Omega \times D} (\underline{\kappa}^n(\underline{x}) \underline{q} \hat{\psi}_{h,N}^n(\underline{x}, \underline{q})) \cdot \nabla_M \zeta(\underline{x}, \underline{q}) d\underline{q} d\underline{x}, \end{aligned} \quad (3.3.40)$$

<sup>2</sup>Note that  $\hat{\psi}_k$  in the term on the last line of (3.3.39) must be at time level  $n$  for the argument of (3.3.22) to apply since it relies on the values  $\{\hat{\psi}_k^n(\underline{x}_m)\}$  interpolating a function in  $V_h$ .

where  $\zeta = X_i \times Y_l$  is an element of  $V_h \otimes \mathcal{P}_N(D)$  and the numerical solution at the intermediate “time level”  $n^*$  is defined as:

$$\mathcal{R}\hat{\psi}_{h,N}^{n^*} := \sum_{k=1}^{N_D} \mathcal{R}\{\hat{\psi}_k^{n^*}(\underline{x}_m)\} Y_k \in V_h \otimes \mathcal{P}_N(D). \quad (3.3.41)$$

Equation (3.3.40) is the Galerkin formulation of (3.3.29) on  $\Omega \times D$  that is obtained by performing a quadrature sum over all  $Q_\Omega$  quadrature points in  $\Omega$ .

The  $\underline{x}$ -direction stage is more straightforward to deal with; we use the classical Douglas–Dupont Galerkin alternating-direction approach for (3.3.31), since it does not contain any  $q$ -dependent coefficients.

Let  $\mathcal{R}\{\hat{\psi}_k^{n^*}(\underline{x}_m)\} = \sum_{i=1}^{N_\Omega} \gamma_{ik}^{n^*} X_i$  so that according to (3.3.41), the vector

$$\underline{\gamma}^{n^*} = (\gamma_{11}^{n^*}, \dots, \gamma_{N_\Omega 1}^{n^*}, \gamma_{12}^{n^*}, \dots, \gamma_{N_\Omega N_D}^{n^*}) \in \mathbb{R}^{N_D N_\Omega}$$

defines  $\mathcal{R}\hat{\psi}_{h,N}^{n^*}$ . Similarly, denote the coefficient vector for  $\hat{\psi}_{h,N}^{n+1}$  as  $\underline{\gamma}^{n+1} \in \mathbb{R}^{N_D N_\Omega}$ , and since the vector entries are ordered in blocks according to the  $q$ -direction degrees-of-freedom, it follows that (3.3.31) can be written as a linear system where the matrices are in tensor product form, *i.e.*:

$$(I_q \otimes M_x + \Delta t I_q \otimes T_x) \underline{\gamma}^{n+1} = I_q \otimes M_x \underline{\gamma}^{n^*}, \quad (3.3.42)$$

where the discretisation matrices are as in (3.3.11) and (3.3.12), and  $I_q$  is the  $N_D \times N_D$  identity matrix.

Equation (3.3.40) can be written in tensor product matrix form also:

$$\left( M_q \otimes M_x + \frac{\Delta t}{2\text{Wi}} S_q \otimes M_x \right) \underline{\gamma}^{n^*} = M_q \otimes M_x \underline{\gamma}^n + \Delta t C(\underline{\kappa}^n; \hat{\psi}_{h,N}^n, \zeta_{il}), \quad (3.3.43)$$

where  $\zeta_{il} = X_i \times Y_l \in V_h \otimes \mathcal{P}_N(D)$ , for  $1 \leq i \leq N_\Omega$  and  $1 \leq l \leq N_D$ . Also,  $M_q$  and  $S_q$  are defined in (3.3.9), (3.3.10), respectively.

Multiplying (3.3.42) by  $(M_q \otimes I_x + \Delta t/(2\text{Wi}) S_q \otimes I_x)$ , where  $I_x$  is the  $N_\Omega \times N_\Omega$  identity matrix, yields

$$\begin{aligned} & \left( M_q \otimes M_x + \Delta t M_q \otimes T_x + \frac{\Delta t}{2\text{Wi}} S_q \otimes M_x + \frac{(\Delta t)^2}{2\text{Wi}} S_q \otimes T_x \right) \underline{\gamma}^{n+1} \\ & = \left( M_q \otimes M_x + \frac{\Delta t}{2\text{Wi}} S_q \otimes M_x \right) \underline{\gamma}^{n^*}. \end{aligned} \quad (3.3.44)$$

Equating the left-hand side of (3.3.43) with the right-hand side of (3.3.44) gives:

$$\begin{aligned} & \left( M_x \otimes M_q + \Delta t M_q \otimes T_x + \frac{\Delta t}{2\text{Wi}} S_q \otimes M_x + \frac{(\Delta t)^2}{2\text{Wi}} S_q \otimes T_x \right) \underline{\gamma}^{n+1} \\ & = M_q \otimes M_x \underline{\gamma}^n + \Delta t C(\underline{\kappa}^n; \hat{\psi}_{h,N}^n, \zeta_{il}). \end{aligned} \quad (3.3.45)$$

Equation (3.3.45) is equivalent to the inner product form in (3.3.38) and hence the proof is complete.  $\square$

Equation (3.3.38) will subsequently be referred to as the *equivalent one-step formulation* for method I. Note that (3.3.38) contains the cross-term,

$$\frac{\Delta t}{2\text{Wi}} \left( \nabla_M \left( \underline{u} \cdot \nabla_x \hat{\psi}_{h,N}^{n+1} \right), \nabla_M \zeta \right),$$

which is not present in the weak formulation (3.2.3). This is analogous to the alternating-direction formulation of the heat equation that was derived in Example 3.3.1, in which cross-terms of the form

$$\frac{\Delta t}{2} \left( \frac{\partial u_h^{n+1}}{\partial x} \frac{\partial v_h}{\partial y} + \frac{\partial u_h^{n+1}}{\partial y} \frac{\partial v_h}{\partial x} \right) \quad \text{and} \quad \frac{\Delta t}{2} \left( \frac{\partial u_h^n}{\partial x} \frac{\partial v_h}{\partial y} + \frac{\partial u_h^n}{\partial y} \frac{\partial v_h}{\partial x} \right),$$

were generated.

### 3.3.3 Method II: Fully-implicit scheme

Method II is very similar to method I, the sole difference being that the term  $C(\underline{\kappa}; \cdot, \cdot)$  is now treated implicitly in time, and therefore we refer to method II as a fully-implicit scheme. We do not discuss the initialisation step or the  $\underline{x}$ -direction scheme here because they are the same as in method I. Instead, we move immediately to discussing the  $\underline{q}$ -direction stage of method II.

Using the line function notation of (3.3.30), the  $\underline{q}$ -direction numerical method is defined as follows: Given the line functions  $\hat{\psi}_k^n \in V_h$ ,  $k = 1, \dots, N_D$ , determine the values  $\hat{\psi}_k^{n*}(\underline{x}_m)$  satisfying

$$\begin{aligned} \sum_{k=1}^{N_D} \hat{\psi}_k^{n*}(\underline{x}_m) \left( \int_D Y_k(\underline{q}) Y_l(\underline{q}) \, d\underline{q} + \frac{\Delta t}{2\text{Wi}} \int_D \nabla_M Y_k(\underline{q}) \cdot \nabla_M Y_l(\underline{q}) \, d\underline{q} \right. \\ \left. - \Delta t \int_D (\underline{\kappa}^{n+1}(\underline{x}_m) \underline{q} Y_k(\underline{q})) \cdot \nabla_M Y_l(\underline{q}) \, d\underline{q} \right) = \sum_{k=1}^{N_D} \hat{\psi}_k^n(\underline{x}_m) \int_D Y_k(\underline{q}) Y_l(\underline{q}) \, d\underline{q}, \end{aligned} \quad (3.3.46)$$

for all  $l = 1, \dots, N_D$ , and for each quadrature point  $\underline{x}_m$ ,  $m = 1, \dots, Q_\Omega$ .

Note that (3.3.46) is exactly the backward Euler Galerkin spectral method that was studied in Chapter 2. It follows as in Section 2.3 that for  $\Delta t$  sufficiently small the associated bilinear form is coercive, and therefore the linear system defined in (3.3.46) is nonsingular.

Unfortunately we cannot derive an equivalent one-step Galerkin formulation for method II using the same reasoning as in Lemma 3.3.4 because the proof of that lemma relied on the term  $C(\underline{\kappa}; \cdot, \cdot)$  being explicit-in-time (*cf.* footnote 2). In order to derive a one-step formulation for method II, we would need to recover an integral of  $\mathcal{R}\{\psi_k^{n*}(\underline{x}_m)\}$  over  $\Omega \times D$  by performing the quadrature sum of the discretisation of  $C(\underline{\kappa}; \cdot, \cdot)$  in (3.3.46). However, this is not possible because this would require a  $\underline{\kappa}$ -weighted reconstruction operator, as distinct from the unweighted reconstruction operator defined in (3.3.27).

Nevertheless, even without an equivalent one-step formulation, we are still able to prove that method II is stable. This is shown in the next section.

**Remark 3.3.5** It is possible to modify method II to obtain a Crank-Nicolson scheme, for example, by adding the term

$$-\frac{1}{2} \sum_{k=1}^{N_D} \hat{\psi}_k^n(\underline{x}_m) \left( \frac{\Delta t}{2\text{Wi}} \int_D \nabla_M Y_k(\underline{q}) \cdot \nabla_M Y_l(\underline{q}) \, d\underline{q} - \Delta t \int_D (\underline{\kappa}^n(\underline{x}_m) \underline{q} Y_k(\underline{q})) \cdot \nabla_M Y_l(\underline{q}) \, d\underline{q} \right)$$

to the right-hand side of (3.3.46), as well as adding the term

$$-\frac{1}{2} \int_{\Omega} \left( \underline{u}^n \cdot \nabla_x \hat{\psi}_k^n \right) X_i \, d\underline{x},$$

on the right-hand side of the  $\underline{x}$ -direction equation.

However, we are ultimately interested in solving the coupled Navier–Stokes–Fokker–Planck system and, as discussed in Chapter 4, the scheme we use for solving this coupled system introduces an  $\mathcal{O}(\Delta t)$  temporal discretisation error. Therefore, there will be no utility in using a Crank-Nicolson discretisation of the Fokker–Planck equation and hence we do not consider this idea any further.  $\diamond$

### 3.4 Stability of methods I and II

First of all, we consider the stability of method I. In this case, the availability of an equivalent one-step method allows the use of standard energy analysis as in the proof of Lemma 3.4.1 below.

Following Chapter 2, we introduce the following right-hand side forcing terms,

$$(\mu^{n+1}, \zeta), \quad (\underline{\nu}^{n+1}, \nabla_M \zeta), \quad (3.4.1)$$

where  $\mu \in L^2(\Omega \times D)$  and  $\underline{\nu} \in L^2(\Omega \times D)^d$ . Right-hand side terms of this form will be useful when we derive convergence estimates in Section 3.5.

**Lemma 3.4.1** *If QH1 holds, so that we have the equivalent one-step formulation for method I given in Lemma 3.3.4, then letting  $\Delta t = T/N_T$ ,  $N_T \geq 1$ ,  $\underline{\kappa} \in \mathcal{C}[0, T]$ ,  $\hat{\psi}_{h,N}^0 \in L^2(\Omega \times D)$ , for  $\hat{\psi}_{h,N}^s \in V_h \otimes \mathcal{P}_N(D)$  we have the following stability estimate:*

$$\begin{aligned} \|\hat{\psi}_{h,N}^s\|^2 + \sum_{n=0}^{s-1} \Delta t \left\| \frac{\hat{\psi}_{h,N}^{n+1} - \hat{\psi}_{h,N}^n}{\sqrt{\Delta t}} \right\|^2 + \sum_{n=0}^{s-1} \frac{\Delta t}{2\text{Wi}} \|\nabla_M \hat{\psi}_{h,N}^{n+1}\|^2 \\ \leq e^{Ks\Delta t} \left\{ \|\hat{\psi}_{h,N}^0\|^2 + \sum_{n=0}^{s-1} 2\Delta t (\|\mu^{n+1}\|^2 + 4\|\underline{\nu}^{n+1}\|^2) \right\}, \end{aligned} \quad (3.4.2)$$

for all  $s$  such that  $1 \leq s \leq N_T$ , where  $K := 2(1 + 4\text{Wi} b \|\underline{\kappa}\|_{L^\infty(0,T;L^\infty(\Omega))}^2)$ .

**Proof.** Consider (3.3.38) with the right-hand side terms of (3.4.1):

$$\begin{aligned} \left( \frac{\hat{\psi}_{h,N}^{n+1} - \hat{\psi}_{h,N}^n}{\Delta t}, \zeta \right) + \left( \underline{u} \cdot \nabla_x \hat{\psi}_{h,N}^{n+1}, \zeta \right) + \frac{1}{2\text{Wi}} \left( \nabla_M \hat{\psi}_{h,N}^{n+1}, \nabla_M \zeta \right) \\ + \frac{\Delta t}{2\text{Wi}} \left( \nabla_M \left( \underline{u} \cdot \nabla_x \hat{\psi}_{h,N}^{n+1} \right), \nabla_M \zeta \right) - \left( \underline{\kappa}^n \underline{q} \hat{\psi}_{h,N}^n, \nabla_M \zeta \right) \\ = (\mu^{n+1}, \zeta) + (\underline{\nu}^{n+1}, \nabla_M \zeta), \end{aligned} \quad (3.4.3)$$

for all  $\zeta \in V_h \otimes \mathcal{P}_N(D)$ . Set  $\zeta = \hat{\psi}_{h,N}^{n+1}$  in (3.4.3) to get

$$\begin{aligned} & \left( \frac{\hat{\psi}_{h,N}^{n+1} - \hat{\psi}_{h,N}^n}{\Delta t}, \hat{\psi}_{h,N}^{n+1} \right) + \left( \underline{y} \cdot \nabla_x \hat{\psi}_{h,N}^{n+1}, \hat{\psi}_{h,N}^{n+1} \right) + \frac{1}{2\text{Wi}} \|\nabla_M \hat{\psi}_{h,N}^{n+1}\|^2 \\ & + \frac{\Delta t}{2\text{Wi}} \left( \nabla_M \left( \underline{y} \cdot \nabla_x \hat{\psi}_{h,N}^{n+1} \right), \nabla_M \hat{\psi}_{h,N}^{n+1} \right) - \left( \underline{\kappa}^n \underline{q} \hat{\psi}_{h,N}^n, \nabla_M \hat{\psi}_{h,N}^{n+1} \right) \\ & = \left( \mu^{n+1}, \hat{\psi}_{h,N}^{n+1} \right) + \left( \underline{\nu}^{n+1}, \nabla_M \hat{\psi}_{h,N}^{n+1} \right). \end{aligned} \quad (3.4.4)$$

The  $\underline{x}$ -transport term vanishes because of (3.2.1) and (3.2.5). Similarly, the first term on the second line vanishes since

$$\begin{aligned} & \left( \nabla_M \left( \underline{y} \cdot \nabla_x \hat{\psi}_{h,N}^{n+1} \right), \nabla_M \hat{\psi}_{h,N}^{n+1} \right) = \int_{\Omega \times D} M \sum_{j=1}^d \left( \sum_{i=1}^d u_i \left( \frac{\partial}{\partial x_i} \frac{\partial}{\partial q_j} \frac{\hat{\psi}_{h,N}^{n+1}}{\sqrt{M}} \right) \left( \frac{\partial}{\partial q_j} \frac{\hat{\psi}_{h,N}^{n+1}}{\sqrt{M}} \right) \right) d\underline{x} dq \\ & = \frac{1}{2} \int_{\Omega \times D} M \sum_{j=1}^d \left( \sum_{i=1}^d u_i \frac{\partial}{\partial x_i} \left( \frac{\partial}{\partial q_j} \frac{\hat{\psi}_{h,N}^{n+1}}{\sqrt{M}} \right)^2 \right) d\underline{x} dq \\ & = -\frac{1}{2} \int_{\Omega \times D} M \sum_{j=1}^d \left( (\nabla_x \cdot \underline{y}) \left( \frac{\partial}{\partial q_j} \frac{\hat{\psi}_{h,N}^{n+1}}{\sqrt{M}} \right)^2 \right) d\underline{x} dq \\ & = -\frac{1}{2} \int_{\Omega \times D} (\nabla_x \cdot \underline{y}) |\nabla_M \hat{\psi}_{h,N}^{n+1}|^2 d\underline{x} dq = 0. \end{aligned}$$

Applying the identity  $2(a-b)a = a^2 - b^2 + (a-b)^2$  to the first term in (3.4.4), yields

$$\begin{aligned} & \|\hat{\psi}_{h,N}^{n+1}\|^2 + \left\| \hat{\psi}_{h,N}^{n+1} - \hat{\psi}_{h,N}^n \right\|^2 + \frac{\Delta t}{\text{Wi}} \|\nabla_M \hat{\psi}_{h,N}^{n+1}\|^2 = \|\hat{\psi}_{h,N}^n\|^2 \\ & + 2\Delta t \left( \underline{\kappa}^n \underline{q} \hat{\psi}_{h,N}^n, \nabla_M \hat{\psi}_{h,N}^{n+1} \right) + 2\Delta t \left( \mu^{n+1}, \hat{\psi}_{h,N}^{n+1} \right) + 2\Delta t \left( \underline{\nu}^{n+1}, \nabla_M \hat{\psi}_{h,N}^{n+1} \right) \\ & =: \|\hat{\psi}_{h,N}^n\|^2 + T_1 + T_2 + T_3. \end{aligned} \quad (3.4.5)$$

Treating  $T_1, T_2$  and  $T_3$  as in the proof of Lemma 2.3.1, we obtain:

$$\begin{aligned} & (1 - \Delta t) \|\hat{\psi}_{h,N}^{n+1}\|^2 + \Delta t \left\| \frac{\hat{\psi}_{h,N}^{n+1} - \hat{\psi}_{h,N}^n}{\sqrt{\Delta t}} \right\|^2 + \frac{\Delta t}{2\text{Wi}} \|\nabla_M \hat{\psi}_{h,N}^{n+1}\|^2 \\ & \leq (1 + C_0 \Delta t) \|\hat{\psi}_{h,N}^n\|^2 + \Delta t (\|\mu^{n+1}\|^2 + 4\|\underline{\nu}^{n+1}\|^2), \end{aligned} \quad (3.4.6)$$

where  $C_0 := 4\text{Wi} b \|\underline{\kappa}\|_{L^\infty(0,T;L^\infty(\Omega))}^2$ . Suppose that  $\Delta t \leq 0.5$ ; then

$$\begin{aligned} & \|\hat{\psi}_{h,N}^{n+1}\|^2 + \Delta t \left\| \frac{\hat{\psi}_{h,N}^{n+1} - \hat{\psi}_{h,N}^n}{\sqrt{\Delta t}} \right\|^2 + \frac{\Delta t}{2\text{Wi}} \|\nabla_M \hat{\psi}_{h,N}^{n+1}\|^2 \\ & \leq \frac{1 + C_0 \Delta t}{1 - \Delta t} \|\hat{\psi}_{h,N}^n\|^2 + 2\Delta t (\|\mu^{n+1}\|^2 + 4\|\underline{\nu}^{n+1}\|^2) \\ & \leq (1 + K \Delta t) \|\hat{\psi}_{h,N}^n\|^2 + 2\Delta t (\|\mu^{n+1}\|^2 + 4\|\underline{\nu}^{n+1}\|^2), \end{aligned}$$

where  $K := 2(1 + C_0) = 2(1 + 4\text{Wi} b \|\underline{\kappa}\|_{L^\infty(0,T;L^\infty(\Omega))}^2)$ .

Summing over  $n = 0, \dots, s-1$  gives,

$$\begin{aligned} & \|\hat{\psi}_{h,N}^s\|^2 + \sum_{n=0}^{s-1} \Delta t \left\| \frac{\hat{\psi}_{h,N}^{n+1} - \hat{\psi}_{h,N}^n}{\sqrt{\Delta t}} \right\|^2 + \sum_{n=0}^{s-1} \frac{\Delta t}{2\text{Wi}} \|\nabla_M \hat{\psi}_{h,N}^{n+1}\|^2 \\ & \leq \left\{ \|\hat{\psi}_{h,N}^0\|^2 + \sum_{n=0}^{s-1} 2\Delta t (\|\mu^{n+1}\|^2 + 4\|\mathcal{L}^{n+1}\|^2) \right\} + K \sum_{n=0}^{s-1} \Delta t \|\hat{\psi}_{h,N}^n\|^2, \end{aligned}$$

and applying a discrete Gronwall lemma yields (3.4.2).  $\square$

We cannot apply an analogous argument for method II due to the absence of an equivalent one-step method. However, by combining stability results for the  $q$ -direction and  $x$ -direction methods we can establish the stability of method II, as shown in the next lemma.

**Lemma 3.4.2** *Suppose QH2 is satisfied and let  $\Delta t = T/N_T$ ,  $N_T \geq 1$ . Then for  $\hat{\psi}_{h,N}^n \in V_h \otimes \mathcal{P}_N(D)$  computed using alternating-direction method II we have*

$$\|\hat{\psi}_{h,N}^n\| \leq e^{c_0 n \Delta t} \|\hat{\psi}_{h,N}^0\|. \quad (3.4.7)$$

for  $1 \leq n \leq N_T$ , where  $c_0 := 1 + 4\text{Wi}b \|\mathcal{L}\|_{\tilde{\mathcal{L}}_{L^\infty(0,T;L^\infty(\Omega))}}^2$ .

**Proof.** From the proof of Lemma 2.3.1, we have the following bound for (3.3.46) at a given quadrature point  $\mathfrak{x}_m \in \bar{\Omega}$ ,

$$\|\hat{\psi}^{n*}(\mathfrak{x}_m, \cdot)\|_{L^2(D)}^2 \leq (1 + 2c_0 \Delta t) \|\hat{\psi}^n(\mathfrak{x}_m, \cdot)\|_{L^2(D)}^2. \quad (3.4.8)$$

Rewriting (3.4.8) in terms of a basis  $\{Y_1, \dots, Y_{N_D}\}$  of  $\mathcal{P}_N(D)$ , which, without loss of generality may be assumed to be orthogonal in the  $L^2(D)$  inner product, we obtain:

$$\sum_{k=1}^{N_D} \hat{\psi}_k^{n*}(\mathfrak{x}_m)^2 \|Y_k\|_{L^2(D)}^2 \leq (1 + 2c_0 \Delta t) \sum_{k=1}^{N_D} \hat{\psi}_k^n(\mathfrak{x}_m)^2 \|Y_k\|_{L^2(D)}^2. \quad (3.4.9)$$

Using (3.3.16) to sum (3.4.9) for  $m = 1, \dots, Q_\Omega$ , and then employing (3.3.26), we have

$$\sum_{k=1}^{N_D} \|\{\hat{\psi}_k^{n*}(\mathfrak{x}_m)\}\|_{\ell^2(\Omega)}^2 \|Y_k\|_{L^2(D)}^2 \leq (1 + 2c_0 \Delta t) \sum_{k=1}^{N_D} \|\{\hat{\psi}_k^n(\mathfrak{x}_m)\}\|_{\ell^2(\Omega)}^2 \|Y_k\|_{L^2(D)}^2. \quad (3.4.10)$$

Since  $\hat{\psi}_{h,N}^n \in V_h \otimes \mathcal{P}_N(D)$ , it follows that  $\hat{\psi}_k^n \in V_h$ , and therefore (as observed below (3.3.26)) the discrete  $\ell^2(\Omega)$  norm on the right-hand side above is equal to the continuous  $L^2(\Omega)$  norm, so that

$$\begin{aligned} \sum_{k=1}^{N_D} \|\{\hat{\psi}_k^{n*}(\mathfrak{x}_m)\}\|_{\ell^2(\Omega)}^2 \|Y_k\|_{L^2(D)}^2 & \leq (1 + 2c_0 \Delta t) \sum_{k=1}^{N_D} \|\hat{\psi}_k^n\|_{L^2(\Omega)}^2 \|Y_k\|_{L^2(D)}^2 \\ & = (1 + 2c_0 \Delta t) \|\hat{\psi}_{h,N}^n\|^2. \end{aligned} \quad (3.4.11)$$

Also, by (3.2.1) and (3.2.5), it follows easily from (3.3.31) that:

$$\|\hat{\psi}_k^{n+1}\|_{L^2(\Omega)}^2 \leq \|\mathcal{R}\{\hat{\psi}_k^{n*}(\mathfrak{x}_m)\}\|_{L^2(\Omega)}^2, \quad (3.4.12)$$

for each  $k$ . Multiplying through by  $\|Y_k\|_{L^2(D)}^2$  in (3.4.12) and summing over  $k = 1, \dots, N_D$  gives

$$\|\hat{\psi}_{h,N}^{n+1}\|^2 = \sum_{k=1}^{N_D} \|\hat{\psi}_k^{n+1}\|_{L^2(\Omega)}^2 \|Y_k\|_{L^2(D)}^2 \leq \sum_{k=1}^{N_D} \|\mathcal{R}\{\hat{\psi}_k^{n*}(\mathfrak{x}_m)\}\|_{L^2(\Omega)}^2 \|Y_k\|_{L^2(D)}^2. \quad (3.4.13)$$

By taking  $\{f_m\} = \{\hat{\psi}_k^{n*}(\mathfrak{x}_m)\}$  and  $X = \mathcal{R}\{\hat{\psi}_k^{n*}(\mathfrak{x}_m)\} \in V_h$  in (3.3.27) and applying the Cauchy-Schwarz inequality in the  $\ell^2$  inner product, we have

$$\begin{aligned} \|\mathcal{R}\{\hat{\psi}_k^{n*}(\mathfrak{x}_m)\}\|_{L^2(\Omega)}^2 &= \left( \{\hat{\psi}_k^{n*}(\mathfrak{x}_m)\}, \{\mathcal{R}\{\hat{\psi}_k^{n*}(\mathfrak{x}_m)\}(\mathfrak{x}_m)\} \right)_{\ell^2(\Omega)} \\ &\leq \|\{\hat{\psi}_k^{n*}(\mathfrak{x}_m)\}\|_{\ell^2(\Omega)} \|\mathcal{R}\{\hat{\psi}_k^{n*}(\mathfrak{x}_m)\}\|_{\ell^2(\Omega)} \\ &= \|\{\hat{\psi}_k^{n*}(\mathfrak{x}_m)\}\|_{\ell^2(\Omega)} \|\mathcal{R}\{\hat{\psi}_k^{n*}(\mathfrak{x}_m)\}\|_{L^2(\Omega)}, \end{aligned}$$

and therefore,

$$\|\mathcal{R}\{\hat{\psi}_k^{n*}(\mathfrak{x}_m)\}\|_{L^2(\Omega)} \leq \|\{\hat{\psi}_k^{n*}(\mathfrak{x}_m)\}\|_{\ell^2(\Omega)}. \quad (3.4.14)$$

Combining (3.4.11), (3.4.13) and (3.4.14), gives,

$$\|\hat{\psi}_{h,N}^{n+1}\|^2 \leq (1 + 2c_0\Delta t) \|\hat{\psi}_{h,N}^n\|^2, \quad (3.4.15)$$

from which (3.4.7) follows easily on noting that  $1 + 2c_0\Delta t \leq e^{2c_0\Delta t}$ .  $\square$

**Remark 3.4.3** The argument in Lemma 3.4.2 can also be applied to method I and hence it follows that method I is stable when only QH2 is satisfied.  $\diamond$

### 3.5 Convergence analysis for method I: Part 1

In this section, the equivalent one-step scheme (3.3.38) and Lemma 3.4.1 are used to prove that the numerical solution obtained using method I converges to the weak solution of (3.2.3), (3.2.4). The convergence argument presented here is analogous to the approach in Section 2.5. Note that we need access to an equivalent one-step formulation to use this approach in the context of alternating-direction methods, and therefore we only consider the convergence analysis of method I. As in the previous chapter, we shall *assume* as much regularity as is needed in order to establish an optimal-order bound on the discretisation error.

Let  $\hat{\psi}(\cdot, \cdot, t)$  be the weak solution of (3.2.3), (3.2.4) at time  $t \in (0, T)$ . To simplify the notation, we write  $\hat{\psi}(t) := \hat{\psi}(\cdot, \cdot, t)$  throughout the rest of this section. As in Section 2.5, we define

$$e_{h,N}^n := \hat{\psi}(t^n) - \hat{\psi}_{h,N}^n = (\hat{\psi}(t^n) - \Pi\hat{\psi}(t^n)) + (\Pi\hat{\psi}(t^n) - \hat{\psi}_{h,N}^n) =: \eta^n + \xi^n,$$

where  $\Pi$  is a projection operator that projects onto  $V_h \otimes \mathcal{P}_N(D)$ .  $\Pi$  shall be defined later.

Noting that  $\xi^n \in V_h \otimes \mathcal{P}_N(D)$ , we apply the equivalent one-step formulation for method I, (3.3.38), to  $\xi^n = \hat{\psi}(t^n) - \hat{\psi}_{h,N}^n - \eta^n$  and set  $\zeta = \xi^{n+1}$ , to obtain:

$$\begin{aligned}
& \left( \frac{\xi^{n+1} - \xi^n}{\Delta t}, \xi^{n+1} \right) + (\underline{u} \cdot \nabla_x \xi^{n+1}, \xi^{n+1}) + \frac{1}{2\text{Wi}} \|\nabla_M \xi^{n+1}\|^2 \\
& \quad + \frac{\Delta t}{2\text{Wi}} (\nabla_M (\underline{u} \cdot \nabla_x \xi^{n+1}), \nabla_M \xi^{n+1}) - (\underline{\kappa}^n \underline{q} \xi^n, \nabla_M \xi^{n+1}) \\
& = \left( \frac{\hat{\psi}(t^{n+1}) - \hat{\psi}(t^n)}{\Delta t}, \xi^{n+1} \right) + (\underline{u} \cdot \nabla_x \hat{\psi}(t^{n+1}), \xi^{n+1}) + \frac{1}{2\text{Wi}} (\nabla_M \hat{\psi}(t^{n+1}), \nabla_M \xi^{n+1}) \\
& \quad + \frac{\Delta t}{2} (\nabla_M (\underline{u} \cdot \nabla_x \hat{\psi}(t^{n+1})), \nabla_M \xi^{n+1}) - (\underline{\kappa}^n \underline{q} \hat{\psi}(t^n), \nabla_M \xi^{n+1}) \\
& \quad - \left( \frac{\eta^{n+1} - \eta^n}{\Delta t}, \xi^{n+1} \right) - (\underline{u} \cdot \nabla_x \eta^{n+1}, \xi^{n+1}) - \frac{1}{2\text{Wi}} (\nabla_M \eta^{n+1}, \nabla_M \xi^{n+1}) \\
& \quad - \frac{\Delta t}{2\text{Wi}} (\nabla_M (\underline{u} \cdot \nabla_x \eta^{n+1}), \nabla_M \xi^{n+1}) + (\underline{\kappa}^n \underline{q} \eta^n, \nabla_M \xi^{n+1}), \tag{3.5.1}
\end{aligned}$$

where the terms containing  $\hat{\psi}_{h,N}^n$  and  $\hat{\psi}_{h,N}^{n+1}$  vanish since  $\hat{\psi}_{h,N}$  satisfies (3.3.38).

First of all we use the identities

$$\underline{\kappa}^n = \underline{\kappa}^{n+1} - \int_{t^n}^{t^{n+1}} \frac{\partial \underline{\kappa}}{\partial t} dt \quad \text{and} \quad \hat{\psi}^n = \hat{\psi}^{n+1} - \int_{t^n}^{t^{n+1}} \frac{\partial \hat{\psi}}{\partial t} dt,$$

to obtain:

$$\begin{aligned}
& (\underline{\kappa}^n \underline{q} \hat{\psi}(t^n), \nabla_M \xi^{n+1}) \\
& = (\underline{\kappa}^{n+1} \underline{q} \hat{\psi}(t^{n+1}), \nabla_M \xi^{n+1}) - \left( \left( \int_{t^n}^{t^{n+1}} \frac{\partial \underline{\kappa}}{\partial t} dt \right) \underline{q} \hat{\psi}(t^n), \nabla_M \xi^{n+1} \right) \\
& \quad - \left( \underline{\kappa}^{n+1} \underline{q} \left( \int_{t^n}^{t^{n+1}} \frac{\partial \hat{\psi}}{\partial t} dt \right), \nabla_M \xi^{n+1} \right) + \left( \left( \int_{t^n}^{t^{n+1}} \frac{\partial \underline{\kappa}}{\partial t} dt \right) \underline{q} \left( \int_{t^n}^{t^{n+1}} \frac{\partial \hat{\psi}}{\partial t} dt \right), \nabla_M \xi^{n+1} \right) \\
& =: (\underline{\kappa}^{n+1} \underline{q} \hat{\psi}(t^{n+1}), \nabla_M \xi^{n+1}) - (\underline{K}_1, \nabla_M \xi^{n+1}) - (\underline{K}_2, \nabla_M \xi^{n+1}) + (\underline{K}_3, \nabla_M \xi^{n+1}).
\end{aligned}$$

Now, considering only the terms containing  $\hat{\psi}$  on the right-hand side of (3.5.1), we have:

$$\begin{aligned}
& \left( \frac{\hat{\psi}(t^{n+1}) - \hat{\psi}(t^n)}{\Delta t}, \xi^{n+1} \right) + (\underline{u} \cdot \nabla_x \hat{\psi}(t^{n+1}), \xi^{n+1}) + \frac{1}{2\text{Wi}} (\nabla_M \hat{\psi}(t^{n+1}), \nabla_M \xi^{n+1}) \\
& \quad + \frac{\Delta t}{2\text{Wi}} (\nabla_M (\underline{u} \cdot \nabla_x \hat{\psi}(t^{n+1})), \nabla_M \xi^{n+1}) - (\underline{\kappa}^n \underline{q} \hat{\psi}(t^n), \nabla_M \xi^{n+1}) \\
& = \left( \frac{\hat{\psi}(t^{n+1}) - \hat{\psi}(t^n)}{\Delta t} - \frac{\partial \hat{\psi}}{\partial t}(t^{n+1}), \xi^{n+1} \right) + \frac{\Delta t}{2\text{Wi}} (\nabla_M (\underline{u} \cdot \nabla_x \hat{\psi}(t^{n+1})), \nabla_M \xi^{n+1}) \\
& \quad + (\underline{K}_1, \nabla_M \xi^{n+1}) + (\underline{K}_2, \nabla_M \xi^{n+1}) - (\underline{K}_3, \nabla_M \xi^{n+1}), \tag{3.5.2}
\end{aligned}$$

where the fact that  $\hat{\psi}$  satisfies (3.2.3), and the expansion of the term  $(\underline{\kappa}^n \underline{q} \hat{\psi}(t^n), \nabla_M \xi^{n+1})$  from above, have been used on the right-hand side. Using (3.5.2) on the right-hand side of

(3.5.1), we have:

$$\begin{aligned} & \left( \frac{\xi^{n+1} - \xi^n}{\Delta t}, \xi^{n+1} \right) + (\underline{u} \cdot \nabla_x \xi^{n+1}, \xi^{n+1}) + \frac{1}{2\mathbb{W}_i} \|\nabla_M \xi^{n+1}\|^2 \\ & \quad + \frac{\Delta t}{2} (\nabla_M (\underline{u} \cdot \nabla_x \xi^{n+1}), \nabla_M \xi^{n+1}) - \left( \tilde{\kappa}^n \underline{q} \xi^n, \nabla_M \xi^{n+1} \right) \\ & = (\mu^{n+1}, \xi^{n+1}) + (\nu^{n+1}, \nabla_M \xi^{n+1}), \end{aligned} \quad (3.5.3)$$

where

$$\mu^{n+1} := \frac{\hat{\psi}(t^{n+1}) - \hat{\psi}(t^n)}{\Delta t} - \frac{\partial \hat{\psi}}{\partial t}(t^{n+1}) - \frac{\eta^{n+1} - \eta^n}{\Delta t} - \underline{u} \cdot \nabla_x \eta^{n+1}, \quad (3.5.4)$$

$$\begin{aligned} \nu^{n+1} & := \frac{\Delta t}{2\mathbb{W}_i} \nabla_M (\underline{u} \cdot \nabla_x \hat{\psi}(t^{n+1})) + \tilde{K}_1 + \tilde{K}_2 - \tilde{K}_3 - \frac{1}{2\mathbb{W}_i} \nabla_M \eta^{n+1} \\ & \quad - \frac{\Delta t}{2\mathbb{W}_i} \nabla_M (\underline{u} \cdot \nabla_x \eta^{n+1}) + \tilde{\kappa}^n \underline{q} \eta^n. \end{aligned} \quad (3.5.5)$$

Therefore, applying the stability result (3.4.2) to (3.5.3) gives

$$\|\xi^n\|^2 + \sum_{m=0}^{n-1} \frac{\Delta t}{2\mathbb{W}_i} \|\nabla_M \xi^{m+1}\|^2 \leq e^{Kn\Delta t} \left\{ \|\xi^0\|^2 + \sum_{m=0}^{n-1} 2\Delta t (\|\mu^{m+1}\|^2 + 4\|\nu^{m+1}\|^2) \right\}. \quad (3.5.6)$$

The next step is to bound the right-hand side of (3.5.6) in terms of norms of  $\eta$  and  $\hat{\psi}$ .

First of all, just as in Section 2.5, we have that  $\|\xi^0\| \leq \|\eta^0\|$ . Next we consider  $\|\mu^{m+1}\|$ :

$$\begin{aligned} \|\mu^{m+1}\|^2 & \leq 3 \left\| \frac{\hat{\psi}(t^{m+1}) - \hat{\psi}(t^m)}{\Delta t} - \frac{\partial \hat{\psi}}{\partial t}(t^{m+1}) \right\|^2 + 3 \left\| \frac{\eta^{m+1} - \eta^m}{\Delta t} \right\|^2 + 3 \|\underline{u} \cdot \nabla_x \eta^{m+1}\|^2 \\ & =: 3(I + II + III). \end{aligned} \quad (3.5.7)$$

For term  $I$ , applying Taylor's theorem with integral remainder yields

$$I \leq \Delta t \int_{t^m}^{t^{m+1}} \left\| \frac{\partial^2 \hat{\psi}}{\partial t^2}(\cdot, \cdot, t) \right\|^2 dt,$$

and for term  $II$  we have the following bound:

$$II \leq \int_{\Omega \times D} \frac{1}{\Delta t} \int_{t^m}^{t^{m+1}} \left| \frac{\partial \eta}{\partial t}(\underline{x}, \underline{q}, t) \right|^2 dt d\underline{x} d\underline{q} = \frac{1}{\Delta t} \int_{t^m}^{t^{m+1}} \left\| \frac{\partial \eta}{\partial t}(\cdot, \cdot, t) \right\|^2 dt.$$

Term  $III$  is simple to bound by pulling out the supremum of  $\underline{u}$ , as follows:

$$III = \int_{\Omega \times D} (\underline{u} \cdot \nabla_x \eta^{m+1})^2 d\underline{x} d\underline{q} \leq \|\underline{u}\|_{L^\infty(0, T; L^\infty(\Omega))}^2 \|\nabla_x \eta^{m+1}\|^2. \quad (3.5.8)$$

Therefore,

$$\begin{aligned}
\sum_{m=0}^{n-1} 2\Delta t \|\mu^{m+1}\|^2 &\leq 6 \sum_{m=0}^{n-1} \Delta t^2 \int_{t^m}^{t^{m+1}} \left\| \frac{\partial^2 \hat{\psi}}{\partial t^2}(\cdot, \cdot, t) \right\|^2 dt + 6 \sum_{m=0}^{n-1} \int_{t^m}^{t^{m+1}} \left\| \frac{\partial \eta}{\partial t}(\cdot, \cdot, t) \right\|^2 dt \\
&\quad + 6 \|u\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \sum_{m=0}^{n-1} \Delta t \|\nabla_x \eta^{m+1}\|^2 \\
&= 6\Delta t^2 \left\| \frac{\partial^2 \hat{\psi}}{\partial t^2} \right\|_{L^2(0,t^m;L^2(\Omega \times D))}^2 + 6 \left\| \frac{\partial \eta}{\partial t} \right\|_{L^2(0,t^m;L^2(\Omega \times D))}^2 \\
&\quad + 6 \|u\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \|\nabla_x \eta\|_{L^2(0,t^m;L^2(\Omega \times D))}^2. \tag{3.5.9}
\end{aligned}$$

Next we derive upper bounds for the norms of the terms on the right-hand side of (3.5.5). First of all, we consider the cross-term,

$$\begin{aligned}
\|\nabla_M(u \cdot \nabla_x \hat{\psi}(t^{n+1}))\|^2 &= \int_{\Omega \times D} \left| \nabla_M \left( \sum_{i=1}^d u_i \frac{\partial}{\partial x_i} \hat{\psi}(t^{n+1}) \right) \right|^2 dx dq \\
&= \int_{\Omega \times D} \sum_{j=1}^d \left\{ \sqrt{M} \frac{\partial}{\partial q_j} \left( \sum_{i=1}^d u_i \frac{\partial}{\partial x_i} \left( \frac{\hat{\psi}(t^{n+1})}{\sqrt{M}} \right) \right) \right\}^2 dx dq \\
&= \int_{\Omega \times D} \sum_{j=1}^d \left\{ \sum_{i=1}^d u_i \frac{\partial}{\partial x_i} \left( \sqrt{M} \frac{\partial}{\partial q_j} \left( \frac{\hat{\psi}(t^{n+1})}{\sqrt{M}} \right) \right) \right\}^2 dx dq \\
&= \int_{\Omega \times D} \sum_{j=1}^d \left\{ u \cdot \nabla_x \left( \sqrt{M} \frac{\partial}{\partial q_j} \left( \frac{\hat{\psi}(t^{n+1})}{\sqrt{M}} \right) \right) \right\}^2 dx dq \\
&\leq \int_{\Omega \times D} \sum_{j=1}^d \left( |u|^2 \left| \nabla_x \left( \sqrt{M} \frac{\partial}{\partial q_j} \left( \frac{\hat{\psi}(t^{n+1})}{\sqrt{M}} \right) \right) \right|^2 \right) dx dq \\
&\leq \|u\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \|\nabla_x \nabla_M \psi(t^{n+1})\|^2. \tag{3.5.10}
\end{aligned}$$

By the same reasoning as in (3.5.10), it follows that:

$$\|\nabla_M(u \cdot \nabla_x \eta^{n+1})\|^2 \leq \|u\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \|\nabla_x \nabla_M \eta^{n+1}\|^2. \tag{3.5.11}$$

Also, we have

$$\|\kappa^n q \eta^n\|^2 \leq b \|\kappa\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \|\eta^n\|^2, \tag{3.5.12}$$

and finally it remains to bound the norms of  $\tilde{K}_1$ ,  $\tilde{K}_2$  and  $\tilde{K}_3$ , for which we have,

$$\|\tilde{K}_1\|^2 = \int_{\Omega \times D} \left\{ \left( \int_{t^n}^{t^{n+1}} \frac{\partial \kappa}{\partial t} q dt \right) \hat{\psi}(t^n) \right\}^2 dx dq \leq \Delta t^2 b \left\| \frac{\partial \kappa}{\partial t} \right\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \|\hat{\psi}(t^n)\|^2, \tag{3.5.13}$$

$$\|\tilde{K}_2\|^2 = \int_{\Omega \times D} \left\{ \kappa^{n+1} q \left( \int_{t^n}^{t^{n+1}} \frac{\partial \hat{\psi}}{\partial t} dt \right) \right\}^2 dx dq \leq \Delta t b \|\kappa\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \int_{t^n}^{t^{n+1}} \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|^2 dt, \tag{3.5.14}$$

and

$$\begin{aligned} \|\tilde{K}_3\|^2 &= \int_{\Omega \times D} \left\{ \left( \int_{t^n}^{t^{n+1}} \frac{\partial \tilde{\kappa}}{\partial t} g \, dt \right) \left( \int_{t^n}^{t^{n+1}} \frac{\partial \hat{\psi}}{\partial t} \, dt \right) \right\}^2 \, d\tilde{x} \, d\tilde{q} \\ &\leq \Delta t^3 b \left\| \frac{\partial \tilde{\kappa}}{\partial t} \right\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \int_{t^n}^{t^{n+1}} \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|^2 \, dt, \end{aligned} \quad (3.5.15)$$

and it is convenient to bound  $\tilde{K}_2$  and  $\tilde{K}_3$  together as follows:

$$\|\tilde{K}_2\|^2 + \|\tilde{K}_3\|^2 \leq b \Delta t \|\tilde{\kappa}\|_{W^{1,\infty}(0,T;L^\infty(\Omega))}^2 \int_{t^n}^{t^{n+1}} \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|^2 \, dt.$$

Therefore,

$$\begin{aligned} &\sum_{m=0}^{n-1} 8\Delta t \|\tilde{\nu}^{m+1}\|^2 \\ &\leq \sum_{m=0}^{n-1} 56\Delta t \left( \frac{\Delta t^2}{4W_1^2} \left\| \nabla_M(\underline{u} \cdot \nabla_x \hat{\psi}(t^{m+1})) \right\|^2 + \frac{\Delta t^2}{4W_1^2} \|\nabla_M(\underline{u} \cdot \nabla_x \eta^{m+1})\|^2 \right. \\ &\quad \left. + \frac{1}{4W_1^2} \|\nabla_M \eta^{m+1}\|^2 + \|\tilde{\kappa}^m g \eta^m\|^2 + \|\tilde{K}_1\|^2 + \|\tilde{K}_2\|^2 + \|\tilde{K}_3\|^2 \right) \\ &\leq \sum_{m=0}^{n-1} 56\Delta t \left( \frac{\Delta t^2}{4W_1^2} \|u\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \left( \|\nabla_x \nabla_M \hat{\psi}(t^{m+1})\|^2 + \|\nabla_x \nabla_M \eta^{m+1}\|^2 \right) \right. \\ &\quad \left. + \frac{1}{4W_1^2} \|\nabla_M \eta^{n+1}\|^2 + b \|\tilde{\kappa}\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \|\eta^n\|^2 \right. \\ &\quad \left. + \Delta t^2 b \left\| \frac{\partial \tilde{\kappa}}{\partial t} \right\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \left\| \hat{\psi}(t^m) \right\|^2 + \Delta t b \|\tilde{\kappa}\|_{W^{1,\infty}(0,T;L^\infty(\Omega))}^2 \int_{t^m}^{t^{m+1}} \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|^2 \, dt \right) \\ &= \frac{14}{W_1^2} \Delta t^2 \|u\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \left( \|\nabla_x \nabla_M \hat{\psi}\|_{\ell^2(0,t^n;L^2(\Omega \times D))}^2 + \|\nabla_x \nabla_M \eta\|_{\ell^2(0,t^n;L^2(\Omega \times D))}^2 \right) \\ &\quad + \frac{14}{W_1^2} \|\nabla_M \eta\|_{\ell^2(0,t^n;L^2(\Omega \times D))}^2 + 56 b \|\tilde{\kappa}\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \|\eta\|_{\ell^2(0,t^n;L^2(\Omega \times D))}^2 \\ &\quad + 56 b \Delta t^2 \left\| \frac{\partial \tilde{\kappa}}{\partial t} \right\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \left\| \hat{\psi} \right\|_{\ell^2(0,t^n;L^2(\Omega \times D))}^2 \\ &\quad + 56 \Delta t^2 b \|\tilde{\kappa}\|_{W^{1,\infty}(0,T;L^\infty(\Omega))}^2 \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0,t^n;L^2(\Omega \times D))}^2. \end{aligned} \quad (3.5.16)$$

We now combine the bounds in (3.5.6), (3.5.9) and (3.5.16) to get:

$$\begin{aligned}
& \|\xi^n\|^2 + \sum_{m=0}^{n-1} \frac{\Delta t}{2W_1} \|\nabla_M \xi^{m+1}\|^2 \\
& \leq e^{Kn\Delta t} \left\{ \|\eta^0\|^2 + 6\Delta t^2 \left\| \frac{\partial^2 \hat{\psi}}{\partial t^2} \right\|_{L^2(0,t^n;L^2(\Omega \times D))}^2 + 6 \left\| \frac{\partial \eta}{\partial t} \right\|_{L^2(0,t^n;L^2(\Omega \times D))}^2 \right. \\
& + 6\|y\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \|\nabla_x \eta\|_{L^2(0,t^n;L^2(\Omega \times D))}^2 \\
& + \frac{14}{W_1^2} \Delta t^2 \|u\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \left( \|\nabla_x \nabla_M \hat{\psi}\|_{L^2(0,t^n;L^2(\Omega \times D))}^2 + \|\nabla_x \nabla_M \eta\|_{L^2(0,t^n;L^2(\Omega \times D))}^2 \right) \\
& + \frac{14}{W_1^2} \|\nabla_M \eta\|_{L^2(0,t^n;L^2(\Omega \times D))}^2 + 56b \|\underline{\kappa}\|_{L^\infty(0,t^n;L^\infty(\Omega))}^2 \|\eta\|_{L^2(0,t^n;L^2(\Omega \times D))}^2 \\
& + 56b\Delta t^2 \left\| \frac{\partial \underline{\kappa}}{\partial t} \right\|_{L^\infty(0,T;L^\infty(\Omega))}^2 \|\hat{\psi}\|_{L^2(0,t^n;L^2(\Omega \times D))}^2 \\
& \left. + 56\Delta t^2 b \|\underline{\kappa}\|_{W^{1,\infty}(0,T;L^\infty(\Omega))}^2 \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0,t^n;L^2(\Omega \times D))}^2 \right\}. \tag{3.5.17}
\end{aligned}$$

Now, just as in Chapter 2, we need to bound the terms containing  $\eta$  in (3.5.17). This is considered in the next section.

### 3.6 Approximation results on $\Omega \times D$

In order to use the approximation results from Section 2.6, we restrict our attention to the  $d = 2$  case here although, of course, analogous results could be obtained for the  $d = 3$  case. We denote the projection operator considered in Section 2.6 (referred to there as  $\hat{\Pi}_N$ ) by  $\Pi_q : \mathcal{H}^{1,1}(D) \rightarrow \mathcal{P}_N(D)$ . Also, we consider a quasi-interpolation operator,  $\mathcal{I}_x : L^1(\Omega) \rightarrow V_h$ , which is a generalisation of the standard finite element interpolant such that the quasi-interpolant is well-defined for nonsmooth functions; we refer to Section 4.8 of [26] for the details of the definition of this operator (alternatively, see [35] or [112]).

We have the following result for  $\mathcal{I}_x$  (cf. Theorem (4.8.12) in [26]):

**Theorem 3.6.1** *Suppose that  $\mathcal{T}_h$  is nondegenerate in the sense that there exists  $\rho > 0$  such that for all  $K \in \mathcal{T}_h$ ,  $\text{diam}(B_K) \geq \rho \text{diam}(K)$ , where  $B_K$  is the largest ball contained in  $K$ . Suppose also that the set of shape functions for each element  $K \in \mathcal{T}_h$  contains all polynomials of degree less than  $m$ . Then, there exists a positive constant  $C$  such that*

$$\left( \sum_{K \in \mathcal{T}_h} h_K^{p(s-k)} \|v - \mathcal{I}_x v\|_{W^{s,p}(K)}^p \right)^{1/p} \leq C |v|_{W^{k,p}(\Omega)},$$

for all  $v \in W^{k,p}(\Omega)$ ,  $0 \leq k \leq m$ ,  $1 \leq p \leq \infty$ ,  $0 \leq s \leq k$ , where  $h_K := \text{diam}(K)$ .

**Corollary 3.6.2** (cf. Corollary 4.8.15 in [26]) *Setting  $s = k$  in Theorem 3.6.1, it follows that*

$$\|\mathcal{I}_x v\|_{W^{k,p}(\Omega)} \leq C |v|_{W^{k,p}(\Omega)} \quad \forall v \in W^{k,p}(\Omega), \tag{3.6.1}$$

for  $0 \leq s, k \leq m$ , where  $m$  is as in Theorem 3.6.1, and  $1 \leq p \leq \infty$ . Also, letting  $h = \max_{K \in \mathcal{T}_h} \text{diam}(K)$  in Theorem 3.6.1, we obtain

$$\|v - \mathcal{I}_x v\|_{\mathbb{W}^{s,p}(\Omega)} \leq Ch^{k-s} |v|_{\mathbb{W}^{k,p}(\Omega)}, \quad (3.6.2)$$

for  $0 \leq s \leq k$ ,  $0 \leq k \leq m$ , and  $m, p$  as in (3.6.1).

For the projection operator  $\Pi_q$ , recall from Section 2.6 that:

$$\|\hat{\psi} - \Pi_q \hat{\psi}\|_{\mathbb{H}_0^1(D;M)} \leq C_1 N_r^{-k} \|\hat{\psi}\|_{\mathcal{H}_r^{k+1}(D)} + C_2 N_\theta^{-l} \|\hat{\psi}\|_{\mathcal{H}_\theta^{l+1}(D)}, \quad (3.6.3)$$

and

$$\|\hat{\psi} - \Pi_q \hat{\psi}\|_{L^2(D)} \leq C_1 N_r^{-k} \|\hat{\psi}\|_{\mathcal{H}_r^k(D)} + C_2 N_\theta^{-l} \|\hat{\psi}\|_{\mathcal{H}_\theta^l(D)}. \quad (3.6.4)$$

Now, let the projection operator  $\Pi : L^1(\Omega; \mathcal{H}^{1,1}(D)) \rightarrow V_h \otimes \mathcal{P}_N(D)$  be defined as

$$\Pi := \mathcal{I}_x \Pi_q = \Pi_q \mathcal{I}_x,$$

so that  $\eta := \hat{\psi} - \Pi \hat{\psi}$ . We will use the approximation properties listed above for  $\Pi_q$  and  $\mathcal{I}_x$  to derive bounds for the terms  $\|\eta\|$ ,  $\|\nabla_x \eta\|$ ,  $\|\nabla_M \eta\|$  and  $\|\nabla_x \nabla_M \eta\|$  that appear on the right-hand side of (3.5.17).

First of all, consider  $\|\eta\|$ :

$$\|\eta\| = \|\hat{\psi} - \mathcal{I}_x \Pi_q \hat{\psi}\| \leq \|\hat{\psi} - \mathcal{I}_x \hat{\psi}\| + \|\mathcal{I}_x \hat{\psi} - \Pi_q \mathcal{I}_x \hat{\psi}\| =: I + II.$$

From (3.6.2), we have that

$$I = \left( \int_D \|\hat{\psi} - \mathcal{I}_x \hat{\psi}\|_{L^2(\Omega)}^2 dq \right)^{\frac{1}{2}} \leq Ch^s \left( \int_D |\hat{\psi}|_{\mathbb{H}^s(\Omega)}^2 dq \right)^{\frac{1}{2}}.$$

Also,

$$\begin{aligned} II &= \left( \int_\Omega \|\mathcal{I}_x \hat{\psi} - \Pi_q \mathcal{I}_x \hat{\psi}\|_{L^2(D)}^2 dx \right)^{\frac{1}{2}} \\ &\leq C_1 N_r^{-k} \left( \int_\Omega \|\mathcal{I}_x \hat{\psi}\|_{\mathcal{H}_r^k(D)}^2 dx \right)^{\frac{1}{2}} + C_2 N_\theta^{-l} \left( \int_\Omega \|\mathcal{I}_x \hat{\psi}\|_{\mathcal{H}_\theta^l(D)}^2 dx \right)^{\frac{1}{2}} \\ &\leq C_1 N_r^{-k} \left( \int_\Omega \|\hat{\psi}\|_{\mathcal{H}_r^k(D)}^2 dx \right)^{\frac{1}{2}} + C_2 N_\theta^{-l} \left( \int_\Omega \|\hat{\psi}\|_{\mathcal{H}_\theta^l(D)}^2 dx \right)^{\frac{1}{2}}, \end{aligned}$$

where we used (3.6.1) with  $k = 0$ ,  $p = 2$  to obtain the last line.

We treat  $\|\nabla_x \eta\|$  similarly:

$$\begin{aligned} \|\nabla_x \eta\| &\leq \|\nabla_x \hat{\psi} - \nabla_x \mathcal{I}_x \hat{\psi}\| + \|\nabla_x \mathcal{I}_x \hat{\psi} - \Pi_q \nabla_x \mathcal{I}_x \hat{\psi}\| \\ &\leq Ch^s \left( \int_D |\hat{\psi}|_{\mathbb{H}^{s+1}(\Omega)}^2 dq \right)^{\frac{1}{2}} \\ &\quad + C_1 N_r^{-k} \left( \int_\Omega \|\nabla_x \mathcal{I}_x \hat{\psi}\|_{\mathcal{H}_r^k(D)}^2 dx \right)^{\frac{1}{2}} + C_2 N_\theta^{-l} \left( \int_\Omega \|\nabla_x \mathcal{I}_x \hat{\psi}\|_{\mathcal{H}_\theta^l(D)}^2 dx \right)^{\frac{1}{2}} \\ &\leq Ch^s \left( \int_D |\hat{\psi}|_{\mathbb{H}^{s+1}(\Omega)}^2 dq \right)^{\frac{1}{2}} \\ &\quad + C_1 N_r^{-k} \left( \int_\Omega \|\nabla_x \hat{\psi}\|_{\mathcal{H}_r^k(D)}^2 dx \right)^{\frac{1}{2}} + C_2 N_\theta^{-l} \left( \int_\Omega \|\nabla_x \hat{\psi}\|_{\mathcal{H}_\theta^l(D)}^2 dx \right)^{\frac{1}{2}}. \end{aligned}$$

Next, we have

$$\begin{aligned}
\|\nabla_M \eta\| &\leq \|\nabla_M \hat{\psi} - \mathcal{I}_x \nabla_M \hat{\psi}\| + \|\nabla_M \mathcal{I}_x \hat{\psi} - \nabla_M \Pi_q \mathcal{I}_x \hat{\psi}\| \\
&\leq Ch^s \left( \int_D |\nabla_M \hat{\psi}|_{\mathbb{H}^s(\Omega)}^2 dq \right)^{\frac{1}{2}} \\
&\quad + C_1 N_r^{-k} \left( \int_{\Omega} \|\mathcal{I}_x \hat{\psi}\|_{\mathcal{H}_r^{k+1}(D)}^2 dx \right)^{\frac{1}{2}} + C_2 N_{\theta}^{-l} \left( \int_{\Omega} \|\mathcal{I}_x \hat{\psi}\|_{\mathcal{H}_{\theta}^{l+1}(D)}^2 dx \right)^{\frac{1}{2}} \\
&\leq Ch^s \left( \int_D |\nabla_M \hat{\psi}|_{\mathbb{H}^s(\Omega)}^2 dq \right)^{\frac{1}{2}} \\
&\quad + C_1 N_r^{-k} \left( \int_{\Omega} \|\hat{\psi}\|_{\mathcal{H}_r^{k+1}(D)}^2 dx \right)^{\frac{1}{2}} + C_2 N_{\theta}^{-l} \left( \int_{\Omega} \|\hat{\psi}\|_{\mathcal{H}_{\theta}^{l+1}(D)}^2 dx \right)^{\frac{1}{2}}.
\end{aligned}$$

Finally, we derive a bound for the cross-term

$$\|\nabla_x \nabla_M \eta\|$$

as follows:

$$\begin{aligned}
\|\nabla_x \nabla_M \eta\| &\leq \|\nabla_x \nabla_M \hat{\psi} - \nabla_x \mathcal{I}_x \nabla_M \hat{\psi}\| + \|\nabla_M \nabla_x \mathcal{I}_x \hat{\psi} - \nabla_M \Pi_q \nabla_x \mathcal{I}_x \hat{\psi}\| \\
&\leq Ch^s \left( \int_D |\nabla_M \hat{\psi}|_{\mathbb{H}^{s+1}(\Omega)}^2 dq \right)^{\frac{1}{2}} \\
&\quad + C_1 N_r^{-k} \left( \int_{\Omega} \|\nabla_x \mathcal{I}_x \hat{\psi}\|_{\mathcal{H}_r^{k+1}(D)}^2 dx \right)^{\frac{1}{2}} + C_2 N_{\theta}^{-l} \left( \int_{\Omega} \|\nabla_x \mathcal{I}_x \hat{\psi}\|_{\mathcal{H}_{\theta}^{l+1}(D)}^2 dx \right)^{\frac{1}{2}} \\
&\leq Ch^s \left( \int_D |\nabla_M \hat{\psi}|_{\mathbb{H}^{s+1}(\Omega)}^2 dq \right)^{\frac{1}{2}} \\
&\quad + C_1 N_r^{-k} \left( \int_{\Omega} \|\nabla_x \hat{\psi}\|_{\mathcal{H}_r^{k+1}(D)}^2 dx \right)^{\frac{1}{2}} + C_2 N_{\theta}^{-l} \left( \int_{\Omega} \|\nabla_x \hat{\psi}\|_{\mathcal{H}_{\theta}^{l+1}(D)}^2 dx \right)^{\frac{1}{2}}.
\end{aligned}$$

Therefore, we have the following optimal order bounds for the terms on the right-hand side of (3.5.17):

$$\|\eta^0\| \leq Ch^s \|\hat{\psi}^0\|_{\mathbb{H}^s(\Omega; L^2(D))} + C_1 N_r^{-k} \|\hat{\psi}^0\|_{L^2(\Omega; \mathcal{H}_r^k(D))} + C_2 N_{\theta}^{-l} \|\hat{\psi}^0\|_{L^2(\Omega; \mathcal{H}_{\theta}^l(D))},$$

$$\begin{aligned}
\|\eta\|_{\ell^2(0, t^n; L^2(\Omega \times D))} &\leq Ch^s \left\| \hat{\psi} \right\|_{\ell^2(0, t^n; \mathbb{H}^s(\Omega; L^2(D)))} + C_1 N_r^{-k} \left\| \hat{\psi} \right\|_{\ell^2(0, t^n; L^2(\Omega; \mathcal{H}_r^k(D)))} \\
&\quad + C_2 N_{\theta}^{-l} \left\| \hat{\psi} \right\|_{\ell^2(0, t^n; L^2(\Omega; \mathcal{H}_{\theta}^l(D)))},
\end{aligned}$$

$$\begin{aligned}
\left\| \frac{\partial \eta}{\partial t} \right\|_{L^2(0, t^n; L^2(\Omega \times D))} &\leq Ch^s \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0, t^n; \mathbb{H}^s(\Omega; L^2(D)))} + C_1 N_r^{-k} \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0, t^n; L^2(\Omega; \mathcal{H}_r^k(D)))} \\
&\quad + C_2 N_{\theta}^{-l} \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0, t^n; L^2(\Omega; \mathcal{H}_{\theta}^l(D)))},
\end{aligned}$$

and similarly,

$$\begin{aligned} \|\nabla_x \eta\|_{\ell^2(0,t^n;L^2(\Omega \times D))} &\leq Ch^s \|\hat{\psi}\|_{\ell^2(0,t^n;H^{s+1}(\Omega;L^2(D)))} \\ &\quad + C_1 N_r^{-k} \|\hat{\psi}\|_{\ell^2(0,t^n;H^1(\Omega;\mathcal{H}_r^k(D)))} + C_2 N_\theta^{-l} \|\hat{\psi}\|_{\ell^2(0,t^n;H^1(\Omega;\mathcal{H}_\theta^l(D)))}, \end{aligned}$$

$$\begin{aligned} \|\nabla_M \eta\|_{\ell^2(0,t^n;L^2(\Omega \times D))} &\leq Ch^s \|\hat{\psi}\|_{\ell^2(0,t^n;H^s(\Omega;H_0^1(D;M)))} \\ &\quad + C_1 N_r^{-k} \|\hat{\psi}\|_{\ell^2(0,t^n;L^2(\Omega;\mathcal{H}_r^{k+1}(D)))} + C_2 N_\theta^{-l} \|\hat{\psi}\|_{\ell^2(0,t^n;L^2(\Omega;\mathcal{H}_\theta^{l+1}(D)))}, \end{aligned}$$

and

$$\begin{aligned} \|\nabla_x \nabla_M \eta\|_{\ell^2(0,t^n;L^2(\Omega \times D))} &\leq Ch^s \|\hat{\psi}\|_{\ell^2(0,t^n;H^{s+1}(\Omega;H_0^1(D;M)))} \\ &\quad + C_1 N_r^{-k} \|\hat{\psi}\|_{\ell^2(0,t^n;H^1(\Omega;\mathcal{H}_r^{k+1}(D)))} + C_2 N_\theta^{-l} \|\hat{\psi}\|_{\ell^2(0,t^n;H^1(\Omega;\mathcal{H}_\theta^{l+1}(D)))}. \end{aligned}$$

### 3.7 Convergence analysis for method I: Part 2

Putting the estimates derived above into (3.5.17), with appropriate constants  $C_1, C_2, C_3$  and  $C_4$ , we obtain:

$$\begin{aligned} &\|\xi\|_{\ell^\infty(0,T;L^2(\Omega \times D))} + \|\nabla_M \xi\|_{\ell^2(0,T;L^2(\Omega \times D))} \\ &\leq C_1 h^s \left( \|\hat{\psi}^0\|_{H^s(\Omega;L^2(D))} + \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0,T;H^s(\Omega;L^2(D)))} + \|\hat{\psi}\|_{\ell^2(0,T;H^s(\Omega;H_0^1(D;M)))} \right. \\ &\quad \left. + \|\hat{\psi}\|_{\ell^2(0,T;H^{s+1}(\Omega;L^2(D)))} \right) \\ &\quad + C_2 N_r^{-k} \left( \|\hat{\psi}^0\|_{L^2(\Omega;\mathcal{H}_r^k(D))} + \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0,T;L^2(\Omega;\mathcal{H}_r^k(D)))} + \|\hat{\psi}\|_{\ell^2(0,T;H^1(\Omega;\mathcal{H}_r^k(D)))} \right. \\ &\quad \left. + \|\hat{\psi}\|_{\ell^2(0,T;L^2(\Omega;\mathcal{H}_r^{k+1}(D)))} \right) \\ &\quad + C_3 N_\theta^{-l} \left( \|\hat{\psi}^0\|_{L^2(\Omega;\mathcal{H}_\theta^l(D))} + \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0,T;L^2(\Omega;\mathcal{H}_\theta^l(D)))} + \|\hat{\psi}\|_{\ell^2(0,T;H^1(\Omega;\mathcal{H}_\theta^l(D)))} \right. \\ &\quad \left. + \|\hat{\psi}\|_{\ell^2(0,T;L^2(\Omega;\mathcal{H}_\theta^{l+1}(D)))} \right) \\ &\quad + C_4 \Delta t \left( \|\hat{\psi}\|_{\ell^2(0,T;L^2(\Omega \times D))} + \|\hat{\psi}\|_{H^2(0,T;L^2(\Omega \times D))} + \|\nabla_x \nabla_M \hat{\psi}\|_{\ell^2(0,T;L^2(\Omega \times D))} \right. \\ &\quad \left. + N_r^{-k} \|\hat{\psi}\|_{\ell^2(0,T;H^1(\Omega;\mathcal{H}_r^{k+1}(D)))} + N_\theta^{-l} \|\hat{\psi}\|_{\ell^2(0,T;H^1(\Omega;\mathcal{H}_\theta^{l+1}(D)))} \right). \end{aligned}$$

Hence, by the triangle inequality:

$$\begin{aligned}
& \|\hat{\psi} - \hat{\psi}_{h,N}\|_{\ell^\infty(0,T;L^2(\Omega \times D))} + \|\nabla_M(\hat{\psi} - \hat{\psi}_{h,N})\|_{\ell^2(0,T;L^2(\Omega \times D))} \\
& \leq \|\xi\|_{\ell^\infty(0,T;L^2(\Omega \times D))} + \|\nabla_M \xi\|_{\ell^2(0,T;L^2(\Omega \times D))} + \|\eta\|_{\ell^\infty(0,T;L^2(\Omega \times D))} + \|\nabla_M \eta\|_{\ell^2(0,T;L^2(\Omega \times D))} \\
& \leq C_1 h^s \left( \|\hat{\psi}\|_{\ell^\infty(0,T;H^s(\Omega;L^2(D)))} + \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0,T;H^s(\Omega;L^2(D)))} + \|\hat{\psi}\|_{\ell^2(0,T;H^s(\Omega;H_0^1(D;M)))} \right. \\
& \quad \left. + \|\hat{\psi}\|_{\ell^2(0,T;H^{s+1}(\Omega;L^2(D)))} \right) \\
& + C_2 N_r^{-k} \left( \|\hat{\psi}\|_{\ell^\infty(0,T;L^2(\Omega;H_r^k(D)))} + \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0,T;L^2(\Omega;H_r^k(D)))} + \|\hat{\psi}\|_{\ell^2(0,T;H^1(\Omega;H_r^k(D)))} \right. \\
& \quad \left. + \|\hat{\psi}\|_{\ell^2(0,T;L^2(\Omega;H_r^{k+1}(D)))} \right) \\
& + C_3 N_\theta^{-l} \left( \|\hat{\psi}\|_{\ell^\infty(0,T;L^2(\Omega;H_\theta^l(D)))} + \left\| \frac{\partial \hat{\psi}}{\partial t} \right\|_{L^2(0,T;L^2(\Omega;H_\theta^l(D)))} + \|\hat{\psi}\|_{\ell^2(0,T;H^1(\Omega;H_\theta^l(D)))} \right. \\
& \quad \left. + \|\hat{\psi}\|_{\ell^2(0,T;L^2(\Omega;H_\theta^{l+1}(D)))} \right) \\
& + C_4 \Delta t \left( \|\hat{\psi}\|_{\ell^2(0,T;L^2(\Omega \times D))} + \|\hat{\psi}\|_{H^2(0,T;L^2(\Omega \times D))} + \|\nabla_x \nabla_M \hat{\psi}\|_{\ell^2(0,T;L^2(\Omega \times D))} \right. \\
& \quad \left. + N_r^{-k} \|\hat{\psi}\|_{\ell^2(0,T;H^1(\Omega;H_r^{k+1}(D)))} + N_\theta^{-l} \|\hat{\psi}\|_{\ell^2(0,T;H^1(\Omega;H_\theta^{l+1}(D)))} \right). \tag{3.7.1}
\end{aligned}$$

Therefore, with  $\psi_{h,N} = \sqrt{M} \hat{\psi}_{h,N}$ , the estimate analogous to (2.7.4) for alternating-direction method I is the following:

$$\begin{aligned}
& \|\psi - \psi_{h,N}\|_{\ell^\infty(0,T;L^2(\Omega;S))} + \|\psi - \psi_{h,N}\|_{\ell^2(0,T;L^2(\Omega;R))} \\
& \leq C_1 h^s \left( \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^\infty(0,T;H^s(\Omega;L^2(D)))} + \left\| \frac{1}{\sqrt{M}} \frac{\partial \psi}{\partial t} \right\|_{L^2(0,T;H^s(\Omega;L^2(D)))} + \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^2(0,T;H^s(\Omega;H_0^1(D;M)))} \right. \\
& \quad \left. + \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^2(0,T;H^{s+1}(\Omega;L^2(D)))} \right) \\
& + C_2 N_r^{-k} \left( \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^\infty(0,T;L^2(\Omega;H_r^k(D)))} + \left\| \frac{1}{\sqrt{M}} \frac{\partial \psi}{\partial t} \right\|_{L^2(0,T;L^2(\Omega;H_r^k(D)))} + \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^2(0,T;H^1(\Omega;H_r^k(D)))} \right. \\
& \quad \left. + \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^2(0,T;L^2(\Omega;H_r^{k+1}(D)))} \right) \\
& + C_3 N_\theta^{-l} \left( \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^\infty(0,T;L^2(\Omega;H_\theta^l(D)))} + \left\| \frac{1}{\sqrt{M}} \frac{\partial \psi}{\partial t} \right\|_{L^2(0,T;L^2(\Omega;H_\theta^l(D)))} + \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^2(0,T;H^1(\Omega;H_\theta^l(D)))} \right. \\
& \quad \left. + \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^2(0,T;L^2(\Omega;H_\theta^{l+1}(D)))} \right) \\
& + C_4 \Delta t \left( \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^2(0,T;L^2(\Omega \times D))} + \left\| \frac{\psi}{\sqrt{M}} \right\|_{H^2(0,T;L^2(\Omega \times D))} + \left\| \nabla_x \nabla_M \frac{\psi}{\sqrt{M}} \right\|_{\ell^2(0,T;L^2(\Omega \times D))} \right. \\
& \quad \left. + N_r^{-k} \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^2(0,T;H^1(\Omega;H_r^{k+1}(D)))} + N_\theta^{-l} \left\| \frac{\psi}{\sqrt{M}} \right\|_{\ell^2(0,T;H^1(\Omega;H_\theta^{l+1}(D)))} \right), \tag{3.7.2}
\end{aligned}$$

for  $s, k, l \geq 1$ , provided that  $\psi$  is such that the right-hand side is finite. Note than an obvious

difference between (3.7.2) and (2.7.4) is that in (3.7.2) we require

$$\left\| \nabla_x \nabla_M \frac{\psi}{\sqrt{M}} \right\|_{\ell^2(0,T;L^2(\Omega \times D))} < \infty.$$

This regularity condition is necessitated by the presence of the cross term,

$$\left( \nabla_M \left( \underline{y} \cdot \nabla_x \hat{\psi}_{h,N}^{n+1} \right), \nabla_M \zeta \right),$$

in (3.3.38).

**Remark 3.7.1** Looking at (3.7.2), it could be argued that there is a mismatch between the convergence rates of the finite element method in  $\Omega$  and the spectral method in  $D$ , in the sense that the spectral method will generally be far more accurate. This is a reasonable point, but we believe that in practice the numerical method analysed here is appropriate. First of all, while in general a finite element scheme will have a low-order convergence rate, its flexibility is invaluable when it comes to meshing physical space domains that may be complicated. Moreover, we do not have a diffusion operator in the  $\underline{x}$ -direction, so it is not obvious that  $\psi$  will be highly smooth in  $\Omega$ .

Nevertheless, it is certainly also reasonable to use a higher-order method for solving the transport equation in physical space, for example, Chauvière & Lozinski used a spectral element method for this purpose in [32,33]. Note that the analysis in this section would carry over essentially unchanged if we replaced the finite element discretisation of (3.3.31) by a higher-order method.

On the other hand, the  $\underline{q}$ -direction is much better suited to the use of a high-order method since  $D$  is always a ball in  $\mathbb{R}^d$ , and, as seen in Section 2.8, at least for the FENE potential, the solution profiles in  $D$  are generally very smooth. Note that in practice the spectral convergence of the  $\underline{q}$ -direction numerical method means that the discrete space  $\mathcal{P}_N(D)$  need only have a rather low dimensionality. This is highly advantageous because (a) each  $\underline{q}$ -direction solve requires relatively modest computational resources and (b) a reduction in the dimensionality of  $\mathcal{P}_N(D)$  reduces the number of  $\underline{x}$ -direction solves that need to be performed each time-step (*cf.* (3.3.31)).  $\diamond$

**Remark 3.7.2** In the preceding argument, we made use of the (pointwise) divergence-free assumption, (3.2.1). This assumption was made to simplify the argument, but it is not essential. Note that it follows from (3.2.2) that  $\nabla_x \cdot \underline{y} \in L^\infty(\Omega)$ ; hence if we allowed  $\nabla_x \cdot \underline{y}$  to be nonzero and assumed the existence of a constant  $c_* \in (0, 1)$  such that

$$1 - \frac{1}{2} \Delta t \| [\nabla_x \cdot \underline{y}]_+ \|_{L^\infty(\Omega)} \geq c_*,$$

then the stability estimate (3.4.2) still holds. Here, for  $x \in \mathbb{R}$ , we used the notation  $[x]_+ := \max(0, x)$  for the positive part of  $x$ . For example, on taking  $c_* = 1/2$  we deduce that the stability estimate (3.4.2) holds provided that  $\Delta t \| [\nabla_x \cdot \underline{y}]_+ \|_{L^\infty(\Omega)} \leq 1$ . It is important to note that this restriction of  $\Delta t$  is completely independent of the spatial discretization parameters  $h$  and  $N$ .  $\diamond$

Now, following the discussion in Section 2.8.1, we consider the convergence of  $\tau_{\approx}$ . In order to coincide with Section 2.8.1, here we consider only the FENE spring force and the case in which  $d = 2$ .

Using Parseval's identity from Chapter 2, we write the weak solution  $\hat{\psi}(\underline{x}, \underline{q}, t) = \tilde{\psi}(\underline{x}, r, \theta, t)$  as follows:

$$\tilde{\psi}(\underline{x}, r, \theta, t) = \tilde{\psi}_1(\underline{x}, r, t) + r \sum_{l=1}^{\infty} \left( \tilde{A}_l(\underline{x}, r, t) \cos(2l\theta) + \tilde{B}_l(\underline{x}, r, t) \sin(2l\theta) \right), \quad (3.7.3)$$

and supposing we use basis  $\mathcal{A}$  in the  $\underline{q}$ -direction, we define the numerical solution as:

$$\tilde{\psi}_{h,N}(\underline{x}, r, \theta) = (1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{0,k}(\underline{x}) P_k(r) + r(1-r) \sum_{i=0}^1 \sum_{l=1}^{N_\theta} \sum_{k=0}^{N_r-1} \tilde{\Psi}_{l,k}^i(\underline{x}) P_k(r) \Phi_{il}(\theta),$$

where  $\tilde{\Psi}_{0,k}, \tilde{\Psi}_{l,k}^i \in V_h$  are line functions as in (3.3.20).

Therefore, proceeding as in Section 2.8, we obtain

$$\begin{aligned} & \|\tau_{11}(\hat{\psi}(t^n)) - \tau_{11}(\hat{\psi}_{h,N}^n)\|_{L^2(\Omega)}^2 \\ & \leq C_* \int_{\Omega} \left\| \tilde{\psi}_1(\underline{x}, r, t^n) - (1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{0,k}^n(\underline{x}) P_k(r) \right\|_{L_w^2(0,1)}^2 d\underline{x} \\ & \quad + \frac{C_*}{4} \int_{\Omega} \left\| r \tilde{A}_1(\underline{x}, r, t^n) - r(1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{1,k}^{0,n}(\underline{x}) P_k(r) \right\|_{L_w^2(0,1)}^2 d\underline{x}, \end{aligned} \quad (3.7.4)$$

where  $C_*$  is defined in (2.8.18).

Also, the analogue of (2.8.19) here is:

$$\begin{aligned} & \|\hat{\psi}(\cdot, \cdot, t^n) - \hat{\psi}_N^n(\cdot, \cdot)\|_{L^2(\Omega \times D)}^2 \\ & = 2\pi b \int_{\Omega} \left\| \tilde{\psi}_1(\underline{x}, r, t^n) - (1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{0,k}^n(\underline{x}) P_k(r) \right\|_{L_w^2(0,1)}^2 d\underline{x} \\ & \quad + \pi b \sum_{l=1}^{N_\theta} \int_{\Omega} \left\| r \tilde{A}_l(\underline{x}, r, t^n) - r(1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{l,k}^{0,n}(\underline{x}) P_k(r) \right\|_{L_w^2(0,1)}^2 d\underline{x} \\ & \quad + \pi b \sum_{l=1}^{N_\theta} \int_{\Omega} \left\| r \tilde{B}_l(\underline{x}, r, t^n) - r(1-r) \sum_{k=0}^{N_r-1} \tilde{\Psi}_{l,k}^{1,n}(\underline{x}) P_k(r) \right\|_{L_w^2(0,1)}^2 d\underline{x} \\ & \quad + \pi b \sum_{l=N_\theta+1}^{\infty} \int_{\Omega} \left( \left\| r \tilde{A}_l(\underline{x}, r, t^n) \right\|_{L_w^2(0,1)}^2 + \left\| r \tilde{B}_l(\underline{x}, r, t^n) \right\|_{L_w^2(0,1)}^2 \right) d\underline{x}, \end{aligned} \quad (3.7.5)$$

and hence, once again, the  $\tau_{11}$  error only contains two terms from the infinite series in (3.7.5), and as in (2.8.20), we have

$$\|\tau_{11}(\hat{\psi}) - \tau_{11}(\hat{\psi}_{h,N})\|_{\ell^\infty(0,T;L^2(\Omega))} \leq \sqrt{\frac{C_*}{2\pi b}} \|\hat{\psi} - \hat{\psi}_{h,N}\|_{\ell^\infty(0,T;L^2(\Omega \times D))}. \quad (3.7.6)$$

Note that since the line functions  $\tilde{\Psi}_{0,k}^n$  and  $\tilde{\Psi}_{1,k}^{0,n}$  in (3.7.4) are computed by solving (3.3.31) using the  $\underline{x}$ -direction finite element method, we expect an  $\mathcal{O}(h^s)$  error to dominate the spatial convergence rate of  $\tau_{\approx}$ , just as in (3.7.2). However, by comparing (3.7.4) and (3.7.5), we can see that only relatively few terms in the  $q$ -direction spectral expansion of  $\hat{\psi}_{h,N}$  contribute to the  $\tau_{11}$  error. Hence, this suggests that the accuracy of  $\tau_{\approx}$  will be less sensitive to the resolution of the  $q$ -direction spectral method than the accuracy of  $\hat{\psi}_{h,N}$ . In Section 3.9 we show that this is indeed the case in practice.

### 3.8 Implementation of methods I and II

In this section we consider the implementation of the  $q$ -direction spectral method and the  $\underline{x}$ -direction finite element method in Sections 3.8.1 and 3.8.2, respectively, and then in Section 3.8.3 we discuss the  $\underline{x}$ -direction quadrature rule used to integrate these two methods into a single alternating-direction algorithm. Finally, we consider the parallel implementation of the alternating-direction methods in Section 3.8.4.

#### 3.8.1 The $q$ -direction stage

We note first of all that from an implementational point of view method I and method II are almost identical; the only difference between the two methods is that method I uses a semi-implicit temporal discretisation whereas method II uses the backward Euler scheme.

Therefore, letting  $\hat{\psi}^{n*}(\underline{x}_m) \in \mathbb{R}^{N_D}$  be the vector with  $k^{\text{th}}$  entry equal to  $\hat{\psi}_k^{n*}(\underline{x}_m)$  and defining  $\hat{\psi}^n(\underline{x}_m)$  analogously, the set of  $q$ -direction linear systems to be solved at time-level  $n$  for method I is:

$$\left( M_q + \frac{\Delta t}{2\mathbb{W}_i} S_q \right) \hat{\psi}^{n*}(\underline{x}_m) = (M_q + \Delta t C_q^m) \hat{\psi}^n(\underline{x}_m), \quad (3.8.1)$$

for  $m = 1, \dots, Q_\Omega$ , whereas for method II we solve:

$$\left( M_q + \frac{\Delta t}{2\mathbb{W}_i} S_q - \Delta t C_q^m \right) \hat{\psi}^{n*}(\underline{x}_m) = M_q \hat{\psi}^n(\underline{x}_m), \quad (3.8.2)$$

for  $m = 1, \dots, Q_\Omega$ . The matrices  $M_q$ ,  $S_q$  and  $C_q^m$  in (3.8.1) and (3.8.2) are as defined in (2.8.4), where  $\underline{k}$  in  $C_q^m$  is sampled at  $\underline{x}_m$ . These matrices depend on the choice of basis of  $\mathcal{P}_N(D)$ ; refer to Section 2.8 for a discussion of the construction of bases  $\mathcal{A}$  and  $\mathcal{B}$  for the  $d = 2$  case, and basis  $\mathcal{C}$  in the case of  $d = 3$ .

It is clear that for both method I and method II, we must solve an  $N_D \times N_D$  linear system  $Q_\Omega$  times per time-step in the  $q$ -direction.  $Q_\Omega$  can be very large in practice. For example, in Section 3.9 we consider some computations for which  $Q_\Omega$  is on the order of  $10^4$ . The use of parallel computation can be very helpful in this situation because the  $q$ -direction linear solves are independent and therefore it is straightforward to perform them in parallel (we discuss this in detail in Section 3.8.4).

It is also interesting to note that method I requires significantly less computational effort in each time-step than method II because the matrix on the left-hand side in (3.8.1) is constant for all  $m$  and therefore we need only perform one LU-factorisation per time-step with method I, whereas the linear system in (3.8.2) must be reassembled and solved afresh at

each quadrature point  $\underline{x}_m$  since in general  $\underline{\kappa}(\underline{x}_m)$  varies from one quadrature point to the next. On the other hand, the numerical experiments in Section 2.8.2 indicate that the backward Euler temporal discretisation of the  $q$ -direction equation is more stable, and it allows one to take larger time-steps, especially for larger values of  $Wi$  or  $\|\underline{\kappa}\|_{L^\infty(\Omega)}$ . Hence, there is a familiar trade-off in efficiency: each time-step is faster with method I, but we can take larger time-steps with method II. Therefore the optimal choice of numerical method depends on the problem at hand.

**Remark 3.8.1** The alternating direction method used by Chauvière & Lozinski in [33] is similar to method II in that it treats the  $\underline{\kappa}$  convection term implicitly in time. In the follow-up papers [32,91] the same authors developed a fast solver approach in which the computational work required for each  $q$ -direction solve was significantly reduced. However, their fast solver was based on an assumption that  $\underline{\kappa}$  arises from a two-dimensional velocity field (*i.e.* that  $\Omega \subset \mathbb{R}^2$ ) whereas in this work we are interested in developing numerical methods that are suitable for  $\Omega \subset \mathbb{R}^3$ .  $\diamond$

The  $q$ -direction solvers for methods I and II were implemented in the C++ programming language and PETSc [8] was used to perform the linear algebra operations. PETSc was a natural choice in this context because it is designed for use on parallel architectures, which is a feature we made extensive use of.

### 3.8.2 The $x$ -direction stage

In the  $x$ -direction, methods I and II are identical: For each line function,  $\hat{\psi}_k^{n*}$ ,  $k = 1, \dots, N_D$ , we solve the transport equation (3.3.33). This involves solving an  $N_\Omega \times N_\Omega$  linear system  $N_D$  times, although the system matrix  $M_x + \Delta t T_x$  only needs to be assembled once per time-step.

In our implementation, we used an  $H^1(\Omega)$ -conforming finite element method with quadratic shape functions to perform the  $x$ -direction computations, and we used GMRES to solve the resulting linear systems. Hence, assuming sufficient regularity for  $\psi/\sqrt{M}$ , we can set  $s = 2$  in (3.7.2), which yields  $\mathcal{O}(h^2)$  terms in the error estimate. Note that in order to strengthen the norm in which the  $x$ -direction solver is stable, Chauvière & Lozinski used an SUPG scheme to discretise the transport equation in [33]. It would be straightforward to integrate such a scheme into our alternating-direction framework, but since the analysis in the preceding sections was performed for a standard Galerkin formulation in the  $x$ -direction, for consistency, we prefer to use the Galerkin method in practice also. Moreover, our numerical results in Section 3.9 and in Section 4 demonstrate that the standard Galerkin formulation performs well in practice.

This method was implemented using the free, open source C++ finite element library `libMesh` [68]. Note also that the  $x$ -direction computations are independent from one another, and hence parallel computation can again be used effectively.

### 3.8.3 The $x$ -direction quadrature rule

We have a great deal of freedom in the choice of the  $x$ -direction quadrature rule. From the analytical point of view, it is preferable to choose a quadrature rule that satisfies QH1, since then, at least with method I, we have access to the equivalent one-step formulation (3.3.38), which was the foundation of the convergence analysis of Section 3.7. However, Lemma (3.4.11)

also shows that only QH2 is required for the stability of method I and method II. In practice, the overall computation time depends very strongly on  $Q_\Omega$  and hence it is often desirable to only satisfy QH2 in order to keep  $Q_\Omega$  as low as possible.

We now discuss some quadrature rules with which we can satisfy either QH2 or both QH1 and QH2 (recall that QH1 is a stronger hypothesis than QH2). Of course, the quadrature rules depend on the element type and the dimension; we will consider triangles and quadrilaterals in two dimensions and tetrahedra and hexahedra in three dimensions. We discuss element-based quadrature rules only. By combining the quadrature rule on each element of  $\mathcal{T}_h$  we obtain a global formula as in (3.3.16).

We assume that each element  $K \in \mathcal{T}_h$  is an affine mapping of some canonical element  $\hat{K}$ . Hence we only need to consider quadrature rules on  $\hat{K}$ .

**Tensor product elements:** In this case, we consider  $\overline{\hat{K}}$  to be either the square  $[-1, 1]^2$  or the cube  $[-1, 1]^3$ . Let  $\{\hat{x}_1, \dots, \hat{x}_n\}$  and  $\{\hat{w}_1, \dots, \hat{w}_n\}$  define the points and weights of a Gaussian quadrature rule, such that  $\hat{x}_i \in (-1, 1)$  and  $\hat{w}_i > 0$  for each  $i$  (e.g. see Chapter 10 of [117]). It is well known that a Gaussian quadrature rule with  $n$  points in one dimension is optimal in the sense that it integrates polynomials of degree  $2n - 1$  on  $\hat{x} \in [-1, 1]$  exactly.

For tensor product finite elements defined on the reference square  $[-1, 1]^2$ , the natural choice of quadrature rule is a tensor product Gaussian rule. For example, following [129], we use the quadrature points:

$$\{(\hat{x}_1, \hat{x}_1), (\hat{x}_1, \hat{x}_2), \dots, (\hat{x}_1, \hat{x}_n), (\hat{x}_2, \hat{x}_1), \dots, (\hat{x}_n, \hat{x}_n)\},$$

and corresponding weights:

$$\{\hat{w}_1 \hat{w}_1, \hat{w}_1 \hat{w}_2, \dots, \hat{w}_1 \hat{w}_n, \hat{w}_2 \hat{w}_1, \dots, \hat{w}_n \hat{w}_n\}.$$

This quadrature rule involves  $Q_{\hat{K}} = n^2$  points and weights and exactly integrates polynomials on  $[-1, 1]^2$  of degree  $2n - 1$  in each direction. A three dimensional tensor product Gauss quadrature rule on  $[-1, 1]^3$  can be defined analogously.

It is clear from the discussion above that we can construct tensor product Gauss quadrature rules to exactly integrate polynomials of arbitrarily high degree on  $[-1, 1]^2$  or  $[-1, 1]^3$ . We now consider how many quadrature points we require to satisfy QH1 or QH2 on tensor product elements in two and three dimensions.

In the computations considered in Section 3.9 and in Chapter 4, we use tensor product quadratic shape functions on each element  $K \in \mathcal{T}_h$  for  $\hat{\psi}_{h,N}$  and for  $\underline{u}$ . Hence the components of  $\underline{\kappa} = \nabla_x \underline{u}$  can also be quadratic in each direction. Therefore, in order to satisfy QH1, we need to be able to exactly integrate polynomials of degree six, and for QH2 we need to integrate polynomials of degree four exactly. Let  $p$  denote the highest degree polynomial that can be exactly integrated by a quadrature rule. We use the following tensor product quadrature rules on the reference square and cube:

- QH1,  $p = 7$ :  $Q_{\hat{K}} = 16$  on  $\overline{\hat{K}} = [-1, 1]^2$ , and  $Q_{\hat{K}} = 64$  on  $\overline{\hat{K}} = [-1, 1]^3$ .
- QH2,  $p = 5$ :  $Q_{\hat{K}} = 9$  on  $\overline{\hat{K}} = [-1, 1]^2$ , and  $Q_{\hat{K}} = 27$  on  $\overline{\hat{K}} = [-1, 1]^3$ .

These quadrature rules are implemented in the `libMesh` software package.

**Simplices:** In this case we assume that  $\hat{K}$  is either a triangle in two dimensions or a tetrahedron in three dimensions. We again consider quadratic shape functions for  $\underline{y}$  and  $\hat{\psi}_{h,N}$ , but since we are no longer using tensor product finite elements, the components of  $\underline{\kappa} = \nabla_x \underline{y}$  are only linear functions in this case, so that in order to satisfy QH1 we need to exactly integrate fifth degree polynomials. To satisfy QH2, we need to exactly integrate degree four polynomials, as in the tensor product case.

In our computations, we used the following quadrature rules, which are implemented in the `libMesh` software package on triangles and tetrahedra:

- QH1 on triangles,  $p = 5$ :  $Q_{\hat{K}} = 7$  [123].
- QH2 on triangles,  $p = 4$ :  $Q_{\hat{K}} = 6$  [94].
- QH1 & QH2 on tetrahedra,  $p = 5$ :  $Q_{\hat{K}} = 14$  [123].

Note that there is a fourth order 11 point quadrature rule on tetrahedra from [64] that is implemented in `libMesh` also, but it contains a negative weight and therefore we cannot use it for our alternating-direction method since we need the quadrature rule to define an inner product, *cf.* (3.3.26). Therefore we use the same  $p = 5$  rule on tetrahedra for both QH1 and QH2.

### 3.8.4 Parallel implementation of the alternating-direction method

It is clear that the computational effort required to solve the high-dimensional Fokker–Planck equation can be very large, particularly in the case  $d = 3$ . Parallel computation is a key ingredient in the alternating-direction framework developed in this chapter, since it makes many problems tractable that would otherwise be well beyond our reach. As indicated above, methods I and II are very well suited to implementation on a parallel architecture; indeed these algorithms are “embarrassingly parallel” in the sense that they involve performing a large number of independent solves in each time-step.

More specifically, suppose we use  $N_{\text{proc}}$  processors ( $N_{\text{proc}} \geq 1$ ) to solve a problem (using either method I or II) with parameters  $N_D$ ,  $N_\Omega$  denoting the number of basis functions in the  $q$ -direction and  $x$ -direction, respectively, and  $Q_\Omega$  defining the number of quadrature points in  $\Omega$ , as in (3.3.16). At time-level  $n$ , we store a dense matrix  $D^n \in \mathbb{R}^{Q_\Omega \times N_D}$ , where  $(D^n)_{ij} = \hat{\psi}_j^n(\underline{x}_i)$ , and  $\hat{\psi}_j^n \in V_h$  is a line function as in (3.3.20). The entries of  $D^n$  uniquely determine  $\hat{\psi}_{h,N}^n \in V_h \otimes \mathcal{P}_N(D)$ . In practice  $D^n$  can be a very large matrix, so we partition it among the processors so that each processor stores a subset of the rows (for  $q$ -direction solves) or columns (for  $x$ -direction solves) of  $D^n$ . We would like these submatrices to be equally sized to obtain ideal load balancing between processors, but depending on  $Q_\Omega, N_D$  and  $N_{\text{proc}}$ , this is often not possible. However, to simplify the discussion here, we will assume for the remainder of this section that  $N_{\text{proc}}$  is a common divisor of  $Q_\Omega$  and  $N_D$  and hence that the submatrices are equally sized.

Now, let us consider the  $q$ -direction computations at time-level  $n$  (we do not distinguish between methods I and II here because, from the point of view of the current discussion, they are identical). We distribute  $D^n$  so that each processor stores  $Q_\Omega/N_{\text{proc}}$  rows of the matrix. Then, simultaneously, each processor solves the  $Q_\Omega/N_{\text{proc}}$   $q$ -direction problems corresponding to its rows in  $D^n$  and updates the data in the matrix. In this manner,  $D^n$  is updated to  $D^{n*}$  where  $(D^{n*})_{ij} = \hat{\psi}_j^{n*}(\underline{x}_i)$ .

Next, we perform the  $\tilde{x}$ -direction computations. First of all, however, we need to redistribute  $D^{n*}$  so that each processor stores  $N_D/N_{\text{proc}}$  columns of the matrix.<sup>3</sup> This involves a global communication operation between all of the processors, which can be time consuming. The time required to perform this parallel communication step depends on the problem size and the number of processors being used. We discuss this issue with regard to some practical computations in Section 3.9, where we show that by selecting  $N_{\text{proc}}$  appropriately it is generally possible to ensure that the matrix redistribution steps take only a small proportion of the overall computation time.

So, once this matrix redistribution is complete, the  $\tilde{x}$ -direction computations on each processor proceed in the same way as in the  $\tilde{q}$ -direction. That is, each processor works sequentially through its  $N_D/N_{\text{proc}}$  columns, first solving (3.3.33), and then sampling the resulting line function  $\hat{\psi}_k^{n+1}$  at  $\tilde{x}_m$  for  $m = 1, \dots, Q_\Omega$  and writing these values back into the matrix. This yields the updated matrix  $D^{n+1}$  on completion of all of the  $\tilde{x}$ -direction solves.

This process is performed for each time-step,  $n = 1, \dots, N_T$ . Note that for computations with the Navier–Stokes–Fokker–Planck system we will need to compute the extra-stress tensor  $\tilde{\tau}$  also. This can be easily included into the framework described above. Suppose we have just finished the  $\tilde{x}$ -direction solves so that  $D^{n+1}$  has been computed and is stored column-wise so that each processor holds  $N_D/N_{\text{proc}}$  columns of the matrix. Then to begin the next time-step, we redistribute  $D^{n+1}$  again so that each processor holds  $Q_\Omega/N_{\text{proc}}$  rows. Once the redistribution is complete and before we begin the  $\tilde{q}$ -direction solves, for each  $m = 1, \dots, Q_\Omega$  we compute and store the values  $\tilde{\tau}^{n+1}(\tilde{x}_m) \in \mathbb{R}^{d \times d}$  using (1.3.37) on the  $\tilde{q}$ -direction cross-section  $\hat{\psi}_{h,N}^{n+1}(\tilde{x}_m, \cdot) \in \mathcal{P}_N(D)$ ; this is again done row by row, and hence each processor only performs  $Q_\Omega/N_{\text{proc}}$  computations with Kramers expression. Using (3.3.27), we can reconstruct  $\mathcal{R}\{\tilde{\tau}^{n+1}(\tilde{x}_m)\} \in (V_h)^{d \times d}$ , which can be used in the right-hand side of (1.3.34).

### 3.9 Numerical results

In this section, we present some numerical results for the alternating-direction approach considered in this chapter applied to a model problem for the FENE Fokker–Planck equation in the  $d = 2$  case. We take  $\underline{u}$  to be the solution of the steady incompressible Navier–Stokes equations with  $\text{Re} = 1$ , and with forcing term  $f(x, y) = (5 \sin(2\pi y), -5 \sin(2\pi x))$ , in the domain  $\Omega = (0, 1)^2$ . In this case,  $\|\underline{\xi}\|_{L^\infty(\Omega)} \approx 2$ . We imposed the Dirichlet boundary condition  $\underline{u} = \underline{0}$  on  $\partial\Omega$ , which ensures that (3.2.5) is satisfied. The streamlines of  $\underline{u}$  are shown in Figure 3.1, and we take  $\underline{u}$  to be constant in time throughout  $t \in (0, T]$ . This velocity field was obtained by solving the Navier–Stokes equations using the Taylor–Hood finite element scheme with quadratic shape functions for  $\underline{u}$  and linear shape functions for the pressure (this numerical method is discussed in more detail in Section 4.2), and we use the same finite element mesh,  $\mathcal{T}_h$ , for the Navier–Stokes equations as for the alternating-direction method, and hence  $\underline{u} \in V_h$ . Note that in general the Taylor–Hood scheme for the Navier–Stokes equations does not yield a (pointwise) divergence-free velocity field, and hence the assumption (3.2.1) is not satisfied for the computational results in this section. However, as noted in Remark 3.7.2, the analysis developed in this chapter can be extended essentially unchanged to the case in which  $\underline{u}$  is not divergence-free.

<sup>3</sup>In our implementation, we performed this redistribution using PETSc’s transpose operation for parallel dense matrices.

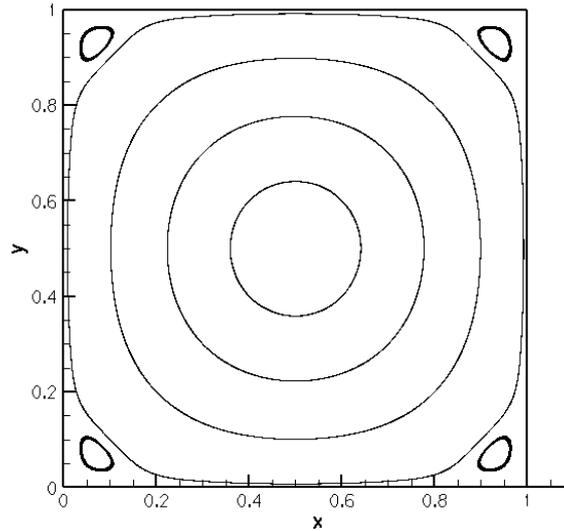


Figure 3.1: Streamlines of the macroscopic velocity field  $u$  driving the enclosed flow model problem. The velocity field is the solution of the steady Navier–Stokes equation with  $\text{Re} = 1$  on  $\Omega = (0, 1)^2$  with forcing  $f(x, y) = (5 \sin(2\pi y), -5 \sin(2\pi x))$ .

We now consider computations using methods I and II for the model problem described above, with the parameters  $\text{Wi} = 1$  and  $b = 12$ . Also, in each of the computations discussed below, we used the initial condition  $\hat{\psi}_{h,N}^0(\underline{x}, \underline{q}) = \sqrt{M(\underline{q})}$ , where  $M$  is the normalised Maxwellian and we ensured that  $N_r \geq 6$ , since according to Remark 2.7.1, that guarantees that  $\sqrt{M} \in \mathcal{P}_N(D)$  in this case. Our goal is to compare the performance of methods I and II, and to study the convergence of these methods under mesh refinement. All of the computations in this section were performed on the Lonestar parallel computer at the Texas Advanced Computing Center (TACC), <http://www.tacc.utexas.edu>, and we used the parallel implementation of the alternating direction method described in Section 3.8.

We do not know the exact solution of the Fokker–Planck equation with the velocity field in Figure 3.1 and therefore in order to obtain quantitative convergence results we first computed a “reference solution”,  $\hat{\psi}_{\text{ref}}$ , and corresponding polymeric extra-stress tensor,  $\underline{\tau}_{\text{ref}}$ , using method I with basis  $\mathcal{A}$  in the  $\underline{q}$ -direction and with a quadrature rule on  $\Omega$  that satisfied QH1. We obtained this reference solution using a highly refined discrete space,  $(V_h \otimes \mathcal{P}_N(D))_{\text{ref}}$ , for which  $\mathcal{T}_h$  was a  $40 \times 40$  uniform mesh of square finite elements and  $(N_r, N_\theta) = (14, 14)$ . In order to satisfy QH1 in this case we required  $Q_{\hat{K}} = 16$ , and hence  $Q_\Omega = 25600$  (*cf.* Section 3.8.3). We took 200 time-steps with  $\Delta t = 10^{-3}$  so that  $T = 0.2$ ; this value of  $\Delta t$  is sufficiently small so that temporal discretisation error does not contaminate the spatial convergence results presented below. The components of  $\underline{\tau}_{\text{ref}}$  at  $T = 0.2$  are shown in Figure 3.2.

In order to obtain convergence data, we then computed  $\hat{\psi}_{h,N}$  and the corresponding stress tensor  $\underline{\tau}$  for several coarser discrete spaces than  $(V_h \otimes \mathcal{P}_N(D))_{\text{ref}}$ . First of all we carried out this process using the same numerical method with which we obtained the reference solution, *i.e.* method I with basis  $\mathcal{A}$  and a quadrature rule that satisfied QH1. The solution data obtained from these computations are denoted  $\hat{\psi}_I$  and  $\underline{\tau}_I$  below. Then, we also computed a

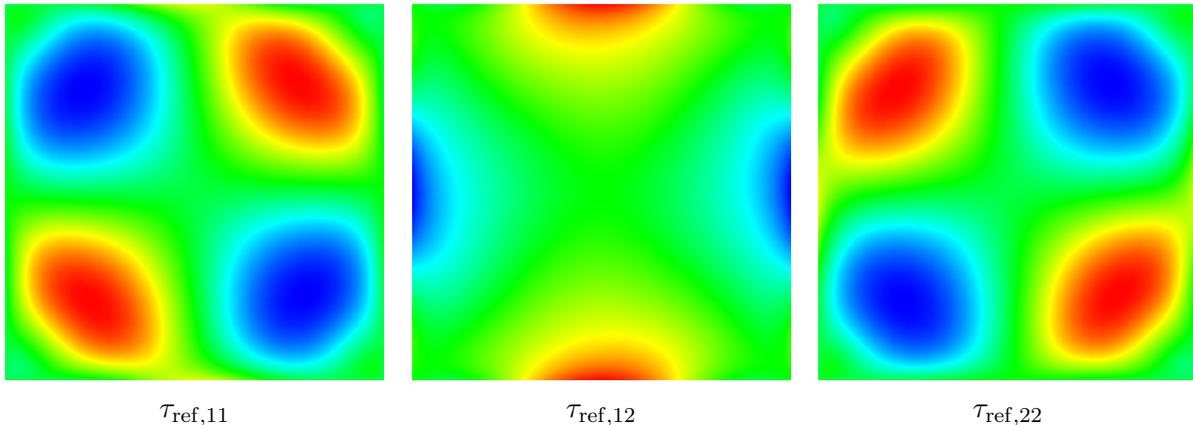


Figure 3.2: The components of  $\underline{\tau}_{\text{ref}}$  at  $T = 0.2$ . Note that we do not show  $\tau_{\text{ref},21}$  since it is identical to  $\tau_{\text{ref},12}$ . In the  $\tau_{\text{ref},11}$  and  $\tau_{\text{ref},22}$  plots, the values range from 0.882 (blue) to 1.15 (red), and in the  $\tau_{\text{ref},12}$  plot we have -0.229 (blue) to 0.229 (red).

corresponding set of numerical solutions on the same discrete spaces, but using method II with basis  $\mathcal{A}$  and a quadrature rule that only satisfied QH2.<sup>4</sup> We denote the solution data in this second case by  $\hat{\psi}_{\text{II}}$  and  $\underline{\tau}_{\text{II}}$ .

The numerical results for  $\hat{\psi}_{\text{I}}$  and  $\underline{\tau}_{\text{I}}$  were obtained using a numerical method that satisfies all of the hypotheses required by the convergence estimates in Section 3.7 (except the divergence-free assumption on  $\underline{u}$ , but, as mentioned above, this assumption is not essential; we only used it in order to simplify the analysis in this chapter). Hence, the  $\hat{\psi}_{\text{I}}$  and  $\underline{\tau}_{\text{I}}$  convergence data in the table allow us to compare the theoretical estimates with practical convergence results. Also, the numerical results enable us to compare the convergence behaviour of method I with QH1 to method II with QH2. These two methods are very similar to one another hence we expect to observe the same convergence behaviour in the two cases, but it is important to provide experimental evidence that these two methods converge to the same solution, and at the same rate, in practice because strictly speaking the convergence analysis in this chapter is only valid for method I with QH1.

The convergence estimates (3.7.2) and (3.7.6) indicate that if the error due to the  $q$ -direction spectral method is negligible compared to the error from the  $x$ -direction finite element method, we should obtain  $\mathcal{O}(h^2)$  convergence rates for both  $\hat{\psi}$  and  $\underline{\tau}$  as  $\mathcal{T}_h$  is refined. Table 3.1 gives the relative errors

$$\|\hat{\psi}_{\text{I}} - \hat{\psi}_{\text{ref}}\|_{\text{L}^2(\Omega \times D)} / \|\hat{\psi}_{\text{ref}}\|_{\text{L}^2(\Omega \times D)} \quad \text{and} \quad \|\hat{\psi}_{\text{II}} - \hat{\psi}_{\text{ref}}\|_{\text{L}^2(\Omega \times D)} / \|\hat{\psi}_{\text{ref}}\|_{\text{L}^2(\Omega \times D)}$$

as well as

$$\|\tau_{\text{I},11} - \tau_{\text{ref},11}\|_{\text{L}^2(\Omega)} / \|\tau_{\text{ref},11}\|_{\text{L}^2(\Omega)} \quad \text{and} \quad \|\tau_{\text{II},11} - \tau_{\text{ref},11}\|_{\text{L}^2(\Omega)} / \|\tau_{\text{ref},11}\|_{\text{L}^2(\Omega)},$$

at  $T = 0.2$ , for the discrete spaces that we considered.

In order to gain further insight into the convergence behaviour of the numerical methods, we plotted the data in Table 3.1 in Figures 3.3 and 3.4.

<sup>4</sup>Recall that we only require  $Q_{\hat{K}} = 9$  to satisfy QH2 on square finite elements.

$\mathcal{T}_h$	$(N_r, N_\theta)$	$\hat{\psi}_I$ error	$\tau_{I,11}$ error	$\hat{\psi}_{II}$ error	$\tau_{II,11}$ error
$5 \times 5$	(6, 6)	$2.07 \times 10^{-2}$	$1.63 \times 10^{-2}$	$2.08 \times 10^{-2}$	$1.63 \times 10^{-2}$
$5 \times 5$	(8, 8)	$2.05 \times 10^{-2}$	$1.63 \times 10^{-2}$	$2.06 \times 10^{-2}$	$1.63 \times 10^{-2}$
$5 \times 5$	(10, 10)	$2.05 \times 10^{-2}$	$1.63 \times 10^{-2}$	$2.06 \times 10^{-2}$	$1.63 \times 10^{-2}$
$10 \times 10$	(6, 6)	$6.25 \times 10^{-3}$	$4.22 \times 10^{-3}$	$6.30 \times 10^{-3}$	$4.24 \times 10^{-3}$
$10 \times 10$	(8, 8)	$5.62 \times 10^{-3}$	$4.22 \times 10^{-3}$	$5.65 \times 10^{-3}$	$4.23 \times 10^{-3}$
$10 \times 10$	(10, 10)	$5.54 \times 10^{-3}$	$4.22 \times 10^{-3}$	$5.58 \times 10^{-3}$	$4.23 \times 10^{-3}$
$20 \times 20$	(6, 6)	$3.29 \times 10^{-3}$	$9.95 \times 10^{-4}$	$3.40 \times 10^{-3}$	$1.07 \times 10^{-3}$
$20 \times 20$	(8, 8)	$1.80 \times 10^{-3}$	$9.90 \times 10^{-4}$	$1.89 \times 10^{-3}$	$1.04 \times 10^{-3}$
$20 \times 20$	(10, 10)	$1.52 \times 10^{-3}$	$9.90 \times 10^{-4}$	$1.67 \times 10^{-3}$	$1.04 \times 10^{-3}$

Table 3.1: Convergence of  $\hat{\psi}$  and  $\tau_{11}$  with respect to the reference solution  $\hat{\psi}_{\text{ref}}$  and reference polymeric stress tensor  $\tau_{\text{ref},11}$  for a series of increasingly refined discrete spaces. The errors are calculated in the  $L^2$  norm at  $T = 0.2$ , and are normalised by dividing by  $\|\hat{\psi}_{\text{ref}}(\cdot, \cdot, T)\|_{L^2(\Omega \times D)} = 0.31$  and  $\|\tau_{\text{ref},11}(\cdot, T)\|_{L^2(\Omega)} = 1.04$ .

In Figure 3.3, the convergence results for  $\hat{\psi}_I$  and  $\hat{\psi}_{II}$  with  $(N_r, N_\theta) = (6, 6)$  and  $(N_r, N_\theta) = (10, 10)$  are plotted on a log-log scale. We have also included a plot of  $h^2$  to show how the decay of the computed errors compare to the expected asymptotic rate. First of all, it is clear from the figure that the two numerical methods behave very similarly; the lines from  $\hat{\psi}_I$  and  $\hat{\psi}_{II}$  are almost indistinguishable. Also, Figure 3.3 shows that we obtain  $\mathcal{O}(h^2)$  convergence when  $(N_r, N_\theta) = (10, 10)$ . However, when  $(N_r, N_\theta) = (6, 6)$ , the plots plateau, which indicates that the error due to the spectral method dominates the  $\mathcal{O}(h^2)$  finite element error when  $\mathcal{T}_h$  is a  $20 \times 20$  mesh.

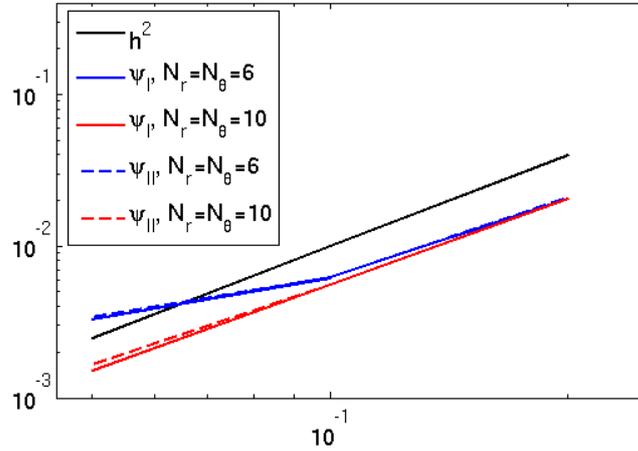


Figure 3.3: Plots of the  $\hat{\psi}_I$  and  $\hat{\psi}_{II}$  convergence data in Table 3.1. The black line shows the expected asymptotic decay rate,  $h^2$ , and the blue and red lines show the convergence of the two numerical methods when  $(N_r, N_\theta)$  is fixed at  $(6, 6)$  and  $(10, 10)$ , respectively.

The  $\tau_{I,11}$  and  $\tau_{II,11}$  convergence data is plotted in Figure 3.4. The data in Table 3.1 is almost identical for  $(N_r, N_\theta) = (6, 6), (8, 8)$  and  $(10, 10)$ , and therefore we only show the

$(N_r, N_\theta) = (6, 6)$  data in the figure. The plot shows that we obtained  $\mathcal{O}(h^2)$  convergence for both  $\tau_{I,11}$  and  $\tau_{II,11}$  as  $\mathcal{T}_h$  is refined from a  $5 \times 5$  mesh to  $20 \times 20$  mesh, when  $(N_r, N_\theta) = (6, 6)$ . This is markedly different from the convergence behaviour of  $\hat{\psi}_{h,N}$ , in which the  $q$ -direction spectral error for  $(N_r, N_\theta) = (6, 6)$  dominated the finite element error on the  $20 \times 20$   $\underline{x}$ -direction mesh. Therefore, this indicates that, just as in Section 2.8, the  $D$  domain spectral method exhibits superconvergence for  $\underline{\tau}$  compared to  $\hat{\psi}$ . This behaviour is dictated by (3.7.4), which indicates that only a small fraction of the terms in the expansion of  $\hat{\psi}_{h,N}$  in terms of spectral basis functions contribute to the error in  $\underline{\tau}$ . As has been noted earlier, the superconvergence of  $\underline{\tau}$  is extremely beneficial in the context of micro-macro computations for simulating dilute polymeric fluids because in that setting the error in  $\hat{\psi}$  is irrelevant; we are solely interested in the  $\underline{\tau}$  error.

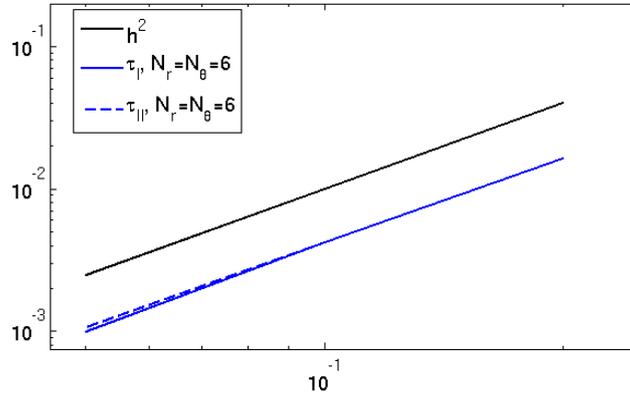


Figure 3.4: Plots of the  $\tau_{I,11}$  and  $\tau_{II,11}$  convergence data in Table 3.1. The black line shows the expected asymptotic decay rate,  $h^2$ , and the solid and dashed blue lines show, respectively, the  $\tau_{I,11}$  and  $\tau_{II,11}$  data for  $(N_r, N_\theta) = (6, 6)$ . The data for the other values of  $(N_r, N_\theta)$  are not plotted since the  $\tau_{11}$  convergence data in Table 3.1 is virtually unaffected by increasing the number of spectral basis functions.

Recall from the discussion in Section 3.8.1 that we expect method I to require significantly less computational work per time-step in the  $q$ -direction than method II. To demonstrate this in practice, we solved the same enclosed flow model problem using both method I and method II. We used a  $20 \times 20$  uniform mesh  $\mathcal{T}_h$  of square finite elements with  $Q_\Omega = 3600$  and basis  $\mathcal{B}$  with  $(N_r, N_\theta) = (15, 15)$  so that  $N_D = 465$ . With  $N_{\text{proc}} = 4$ , the total computation time per time-step for method I was 1.75 seconds, whereas for method II it was 3.42 seconds. This difference is due to the fact that method II took 2.37 seconds per time-step to perform the  $q$ -direction computations, whereas method I only took 0.70 seconds per time-step in the  $q$ -direction.

Nevertheless, for problems of physical interest, method II is often the preferred alternating-direction method. This is because the fully implicit temporal discretisation used by method II is more stable than the semi-implicit scheme in method I, especially for larger flow rates and Weissenberg numbers (*cf.* Section 2.8.2). Hence method I can require much smaller time-step sizes than method II, and this can often outweigh the reduced computational complexity per time-step of method I. Also, for large-scale problems we generally prefer to satisfy only QH2

rather than QH1 since with QH2 we can obtain a smaller value of  $Q_\Omega$ , which in turn reduces the computational work required in each time-step of the alternating-direction method.

We now move on to consider the scaling of the computation time as we increase the number of processors in the parallel implementation of the alternating-direction method. The enclosed-flow problem considered above provides a convenient test case with which we can quantify the parallel speedup for the alternating-direction method. We studied this speedup by, first of all, solving the enclosed flow problem on one node of the Lonestar parallel computer (each node contains 4 processors) to get the base computation time per time-step, which we denote  $T(1)$ . We then repeated the same computation, but using more computational nodes of the parallel computer and we recorded the computation time,  $T(N)$ , in each case, where  $N$  denotes the number of computational nodes that were used. We refer to the ratio  $T(1)/T(N)$  as the *parallel speedup*.

The parameters that have the most significant effect on the computation time of the parallel alternating-direction scheme are  $N_D$  and  $Q_\Omega$ , since these determine the number of  $\tilde{x}$ - and  $\tilde{q}$ -direction solves that need to be performed each time-step. Note that there are only two steps in the alternating-direction algorithm for which the computation time does not scale down proportionally to the number of processors being used: the matrix assembly for (3.3.33), which must be performed exactly once per time-step irrespective of  $N_{\text{proc}}$ , and also the dense matrix redistribution that precedes direction changes in the alternating-direction method. However, if the  $\tilde{x}$ - and  $\tilde{q}$ -direction solves dominate the overall computation time, then we can expect that the parallel speedup will scale linearly with the number of processors being used.

In order to examine the scaling of the parallel speedup in practice, we performed computations for two different discrete spaces, such that (i)  $N_D = 120$  and  $Q_\Omega = 3600$ , and (ii)  $N_D = 1800$  and  $Q_\Omega = 8100$ . We solved the enclosed flow problem for these spaces using a number of different choices of  $N_{\text{proc}}$ . We used method II with basis  $\mathcal{B}$  to obtain the data below, but the parallel speedup behaviour is essentially the same whether we use methods I or II or bases  $\mathcal{A}$  or  $\mathcal{B}$ . The base computation times were  $T(1) = 0.53$  seconds for the  $(N_D, Q_\Omega) = (120, 3600)$  computation, and  $T(1) = 157.0$  seconds for the  $(N_D, Q_\Omega) = (1800, 8100)$  case.

The parallel speedup of the alternating-direction method for the two discrete spaces discussed above is plotted in Figure 3.5. In the case that  $(N_D, Q_\Omega) = (1800, 8100)$ , we obtained a parallel speedup of 14.8 when  $N = 15$  (*i.e.*  $N_{\text{proc}} = 60$ ), whereas the speedup tailed off to less than 10 when  $N = 15$  for the computation with  $(N_D, Q_\Omega) = (120, 3600)$ . This difference in the scaling of the parallel speedup is primarily due to the fact that the overhead from the redistribution of  $D^n$  is much larger, as a proportion of the overall computation time, for the smaller problem. For example, for the  $(N_D, Q_\Omega) = (120, 3600)$  problem, matrix redistribution took 8.66% of the overall computation time when  $N = 1$ , but when  $N = 15$ , it increased to 30.4%. By contrast, in the larger problem with  $(N_D, Q_\Omega) = (1800, 8100)$ , more time is spent on the  $\tilde{q}$ - and  $\tilde{x}$ -direction solves in each time-step, so that only 0.89% of the computation time was taken for the matrix redistribution when  $N = 1$ , which increased to 2.25% when  $N = 15$ . Since 2.25% is still only a small proportion of the overall computation time, the matrix redistribution overhead does not significantly detract from the near optimal scaling of the parallel speedup shown in Figure 3.5 for the  $(N_D, Q_\Omega) = (1800, 8100)$  case. This indicates that as long as the values of  $N_D$  and  $Q_\Omega$  are large enough, the alternating-direction method can scale efficiently to a very large number of processors.

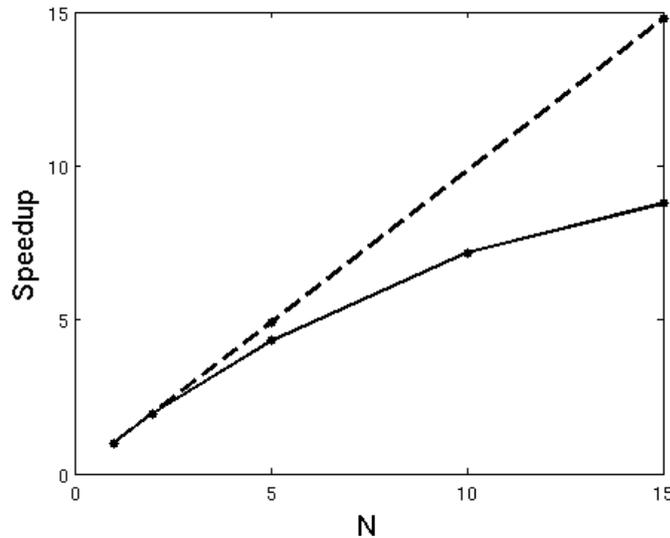


Figure 3.5: Plot of speedup, *i.e.*  $T(1)/T(N)$ , as the number of computational nodes is increased from 1 to 15. The speedup data for  $(N_D, Q_\Omega) = (120, 3600)$  is plotted as a solid line and the dashed line shows the data for  $(N_D, Q_\Omega) = (1800, 8100)$ . For each computation we chose the number of nodes so that  $N_{\text{proc}} (= 4N)$  was a common divisor of  $N_D$  and  $Q_\Omega$  in order to ensure optimal load balancing in each case so that the comparisons of computation time are fair.

### 3.10 Conclusions

In this chapter we developed an alternating-direction method for the Fokker–Planck equation, which is a hybrid of a classical Douglas–Dupont-type Galerkin alternating-direction scheme, and a new quadrature based scheme. We were able to derive a range of theoretical results for this scheme, including stability results in Section 3.4 and convergence estimates in Section 3.7. Much of this theory built upon the analysis of the Fokker–Planck equation in  $D$  that was considered in Section 2. We also put particular emphasis on practical computations in this chapter, and we discussed the implementation of the alternating-direction scheme in Section 3.8, and followed up in Section 3.9 by presenting a range of computational results for alternating-direction methods I and II applied to a model problem with a fixed velocity field,  $\underline{u}$ . We demonstrated that the convergence rates observed in practice for this model problem are accurately described by the theoretical results in Section 3.7. Moreover, we showed that, just as in Chapter 2, the  $q$ -direction spectral method yields a more accurate solution for  $\underline{\tau}$  than it does for  $\hat{\psi}$ , which means that if we are solely interested in the accuracy of  $\underline{\tau}$  – as is the case when we consider the Navier–Stokes–Fokker–Planck system – then we can take fewer spectral basis functions than we would need if  $\hat{\psi}$  were the quantity of primary interest. This leads to significant savings when we solve the Navier–Stokes–Fokker–Planck system, since the computational work required by the alternating-direction method for the Fokker–Planck equation depends strongly on  $N_D$ , the number of  $q$ -direction basis functions. In the next chapter we combine the numerical methods developed in this chapter for the Fokker–Planck equation with a finite element scheme for solving the Navier–Stokes equations to obtain an algorithm for solving the full micro-macro model for dilute polymeric fluids.



## Chapter 4

# The coupled Navier–Stokes–Fokker–Planck system

### 4.1 Introduction

In this chapter we develop an algorithm for solving the Navier–Stokes–Fokker–Planck system, (1.3.34)–(1.3.38), and we use this algorithm to obtain computational results for flow problems that are of physical interest. This chapter is relatively brief because the components of our algorithm are already well understood; we use a standard mixed finite element method for solving the Navier–Stokes equations and we couple this to the alternating-direction scheme for the Fokker–Planck equation that was considered in detail in Chapter 3. Our focus in this chapter is on obtaining practical computational results. The convergence analysis of a finite element approximation to the coupled Navier–Stokes–Fokker–Planck system will be carried out in Chapter 6; we note however that the the scheme studied there is based on a direct time-discretisation of the Fokker–Planck equation and does not including the alternating direction scheme developed in Chapter 3 and used herein.

The chapter is structured as follows. The numerical method for the Navier–Stokes–Fokker–Planck system is discussed in Section 4.2, and we present numerical results in Section 4.3. Note that throughout this chapter we consider the FENE potential only but, once again, the methodology would be the same for any spring potential that satisfies Hypotheses A and B.

### 4.2 Numerical method for the micro-macro model

The algorithm we use to couple the numerical methods for the Navier–Stokes equations and the Fokker–Planck equation is essentially the same as those used by Chauvière & Lozinski [32, 33, 91] and Helzel & Otto [55] for this purpose. We discuss this procedure below, but first we introduce numerical methods for the Navier–Stokes equations, and also for the Stokes equations.

Recall the nondimensionalised Navier–Stokes equations from Chapter 1, in which  $\nabla_x \cdot \underline{\tau}$

arises as a forcing term:

$$\frac{\partial \underline{u}}{\partial t} + (\underline{u} \cdot \nabla_x) \underline{u} + \nabla_x p = \frac{\gamma}{\text{Re}} \Delta_x \underline{u} + \frac{b+d+2}{b} \frac{1-\gamma}{\text{Re Wi}} \nabla_x \cdot \underline{\tau}, \quad (4.2.1)$$

$$\nabla_x \cdot \underline{u} = 0. \quad (4.2.2)$$

In this chapter we will also consider a Stokes–Fokker–Planck model, which is valid in the limit  $\text{Re} \rightarrow 0_+$ . In the Stokes equations the incompressibility condition (4.2.2) is unchanged, but we use the following momentum equation (in dimensional form):

$$\nabla_x p = \nu_s \Delta_x \underline{u} + \frac{1}{\rho} \nabla_x \cdot \underline{\tau}, \quad (4.2.3)$$

instead of (1.3.2). We nondimensionalise (4.2.3) by using (1.3.15) and the pressure rescaling  $p = (\nu U_0 / L_0) \hat{p}$ ,<sup>1</sup> to obtain:

$$\nabla_x p = \gamma \Delta_x \underline{u} + \frac{b+d+2}{b} \frac{1-\gamma}{\text{Wi}} \nabla_x \cdot \underline{\tau}. \quad (4.2.4)$$

Next, we introduce mixed finite element approximations of the incompressible Navier–Stokes and Stokes equations. The numerical analysis of these equations is well understood and therefore we discuss our approach only briefly; for further details see [46] or [51].

As in Chapter 3, let  $\mathcal{T}_h$  denote a finite element triangulation of  $\bar{\Omega}$ , and let  $V_h$  be the corresponding finite element space with quadratic shape functions that we used for the alternating-direction method for  $\hat{\psi}_{h,N}$  in Chapter 3. Also, let  $P_h$  denote the  $H^1(\Omega)$ -conforming finite element space based on  $\mathcal{T}_h$  that uses linear shape functions. Then  $\mathcal{V}_h := [V_h]^d$  and  $P_h$  are the Taylor–Hood finite element spaces for the Navier–Stokes equations (*cf.* Chapter 5 of [46]); these spaces are known to satisfy the inf-sup stability condition (*cf.* Section 12.6 of [26]). As noted in Chapter 3, in general the Taylor–Hood scheme does not yield a pointwise divergence free velocity field. In the context of the coupled Navier–Stokes–Fokker–Planck system, this may lead to undesirable effects, for example, related to the integral conservation property identified for the Fokker–Planck equation in (3.2.8). We did not examine the behaviour of this integral property in our numerical experiments presented in Section 4.3, but this is a question of interest for future research.

Using the discrete spaces introduced above, our numerical method for the Navier–Stokes system is defined as follows:

Suppose  $\underline{u}_h^0 \in \mathcal{V}_h$ ,  $p_h^0 \in P_h$  and  $\underline{\tau}_{h,N}^n \in \underline{\mathcal{L}}^2(\Omega) := (L^2(\Omega))^{d \times d}$  for  $n = 0, \dots, N_T - 1$  are given. Then, for  $n = 0, \dots, N_T - 1$ , find  $\underline{u}_h^{n+1} \in \mathcal{V}_h$  and  $p_h^{n+1} \in P_h$  satisfying:

$$\begin{aligned} & \int_{\Omega} \frac{\underline{u}_h^{n+1} - \underline{u}_h^n}{\Delta t} \cdot \underline{v}_h \, d\mathbf{x} + \int_{\Omega} ((\underline{u}_h^{n+1} \cdot \nabla_x) \underline{u}_h^{n+1}) \cdot \underline{v}_h \, d\mathbf{x} - \int_{\Omega} p_h^{n+1} \nabla_x \cdot \underline{v}_h \, d\mathbf{x} \\ & + \frac{\gamma}{\text{Re}} \int_{\Omega} \nabla_x \underline{u}_h^{n+1} : \nabla_x \underline{v}_h \, d\mathbf{x} + \frac{b+d+2}{b} \frac{1-\gamma}{\text{Re Wi}} \int_{\Omega} \underline{\tau}_{h,N}^n : \nabla_x \underline{v}_h \, d\mathbf{x} \\ & + \int_{\partial\Omega} \left[ \left( p_h^{n+1} \underline{I} - \frac{\gamma}{\text{Re}} \nabla_x \underline{u}_h^{n+1} - \frac{b+d+2}{b} \frac{1-\gamma}{\text{Re Wi}} \underline{\tau}_{h,N}^n \right) \underline{n}_{\partial\Omega} \right] \cdot \underline{v}_h \, ds = 0 \quad \forall \underline{v}_h \in \mathcal{V}_h, \end{aligned} \quad (4.2.5)$$

$$\int_{\Omega} q_h \nabla_x \cdot \underline{u}_h^{n+1} \, d\mathbf{x} = 0 \quad \forall q_h \in P_h. \quad (4.2.6)$$

<sup>1</sup>This pressure scaling is appropriate for creeping flow.

Here  $\underline{n}_{\partial\Omega}$  signifies the unit outward normal vector to  $\partial\Omega$ . For tensors  $\underline{A}$  and  $\underline{B}$ , the colon notation used above is defined as  $\underline{A} : \underline{B} := \sum a_{ij} b_{ij}$ .

In this section we consider channel flow problems in which we have an inflow boundary,  $\partial\Omega_{\text{in}}$ , an outflow boundary,  $\partial\Omega_{\text{out}}$  and channel wall boundaries  $\partial\Omega_0$ , such that  $\partial\Omega = \partial\Omega_{\text{in}} \cup \partial\Omega_{\text{out}} \cup \partial\Omega_0$ . We assume that the channel wall boundaries are stationary and we impose the no-slip boundary condition  $\underline{u}_h = \underline{0}$  on  $\partial\Omega_0$ . Also, we impose  $\underline{u}_h = \underline{u}_{\text{in}}$  on  $\partial\Omega_{\text{in}}$ , where  $\underline{u}_{\text{in}}$  is an inflow velocity profile corresponding to a fully-developed flow. In Section 4.3, the maximum of  $\underline{u}_{\text{in}}$  is denoted by  $U_{\text{max}}$ . As a result of these Dirichlet boundary conditions, we have  $\underline{u}_h = \underline{0}$  on  $\partial\Omega_{\text{in}} \cup \partial\Omega_0$ . Also, on  $\partial\Omega_{\text{out}} \times (0, T]$ , we impose

$$\left( p_{\approx}^I - \frac{\gamma}{\text{Re}_{\approx}} \nabla_x \underline{u} - \frac{b+d+2}{b} \frac{1-\gamma}{\text{Re Wi}_{\approx}} \underline{\tau} \right) \underline{n}_{\partial\Omega} = \underline{0}.$$

We approximate this boundary condition weakly, by omitting

$$\int_{\partial\Omega_{\text{out}}} \left[ \left( p_h^{n+1} \underline{I} - \frac{\gamma}{\text{Re}_{\approx}} \nabla_x \underline{u}_h^{n+1} - \frac{b+d+2}{b} \frac{1-\gamma}{\text{Re Wi}_{\approx}} \underline{\tau}_{h,N}^n \right) \underline{n}_{\partial\Omega} \right] \cdot \underline{v}_h \, ds$$

from (4.2.5). Hence in the boundary term in (4.2.5) vanishes on all of  $\partial\Omega$ . Note that the  $\underline{\tau}_{h,N}$  terms in (4.2.5) are at time-level  $n$  rather than  $n+1$ ; we shall see below that this enables us to couple the Fokker–Planck and Navier–Stokes equations in a convenient manner.

The momentum equation, (4.2.5), is nonlinear due to the term  $\int_{\Omega} ((\underline{u}_h^{n+1} \cdot \nabla_x) \underline{u}_h^{n+1}) \cdot \underline{v}_h \, dx$ . Hence, we use Newton’s method to solve the nonlinear system of equations arising from (4.2.5) and (4.2.6) at each time-level.

We now turn our attention to the Stokes equations, which we discretise in a very similar manner. The difference is that we replace (4.2.5) with the following equation:

$$\begin{aligned} & - \int_{\Omega} p_h^{n+1} \nabla_x \cdot \underline{v}_h \, dx + \gamma \int_{\Omega} \nabla_x \underline{u}_h^{n+1} : \nabla_x \underline{v}_h \, dx + \frac{b+d+2}{b} \frac{1-\gamma}{\text{Wi}} \int_{\Omega} \underline{\tau}_{h,N}^n : \nabla_x \underline{v}_h \, dx \\ & + \int_{\partial\Omega} \left[ \left( p_h^{n+1} \underline{I} - \gamma \nabla_x \underline{u}_h^{n+1} - \frac{b+d+2}{b} \frac{1-\gamma}{\text{Wi}} \underline{\tau}_{h,N}^n \right) \underline{n}_{\partial\Omega} \right] \cdot \underline{v}_h \, ds = 0 \quad \forall \underline{v}_h \in \underline{V}_h. \end{aligned} \quad (4.2.7)$$

We we apply the same boundary conditions as discussed above for the Navier–Stokes case, and therefore the boundary term in (4.2.7) vanishes also. Note that there is no time derivative in (4.2.4), and hence in this case the time-dependence comes only through  $\underline{\tau}_{h,N}^n$  and the boundary data. The Stokes equations are linear and therefore we do not require a Newton scheme in this case.

The mixed finite element methods described above for the Navier–Stokes and Stokes equations were implemented in the finite element library `libMesh` [68]. In both cases, we solve the linear systems that arise from the finite element discretisations using GMRES with incomplete LU factorisation as a preconditioner. In order to obtain faster convergence rates for the iterative solver one could apply more advanced preconditioning techniques, such as the techniques discussed in [46] that take advantage of the structure of the linear systems arising from the discretisation of Stokes or Navier–Stokes problems. However, there is little incentive for us to accelerate the convergence of our Navier–Stokes or Stokes solvers in this way because the overall computation time for computations with the Navier–Stokes–Fokker–Planck system is dominated by solving the Fokker–Planck equation on  $\Omega \times D$ .

In Chapter 3, we restricted our attention to enclosed flows to simplify the analysis in that chapter, but we are now interested in problems that have inflow and outflow boundaries.

Therefore, we need to define the boundary conditions for the Fokker–Planck equation on  $\partial\Omega_{\text{in}}$  and  $\partial\Omega_{\text{out}}$ .

In fact, since the Fokker–Planck equation on  $\Omega$  is a pure advection problem, we do not need to do anything different on  $\partial\Omega_{\text{out}}$  since by definition we have  $\underline{u} \cdot \underline{n}_{\partial\Omega} > 0$  there.<sup>2</sup> However, we do need to treat the inflow boundary differently. Suppose we set  $\underline{u}_h^n|_{\partial\Omega_{\text{in}}} = \underline{u}_{\text{in}}^n$  for the Stokes/Navier–Stokes system for  $n = 1, \dots, N_T$ . Then that boundary data also defines  $\underline{\kappa}_{\text{in}}^n = \nabla_x \underline{u}_{\text{in}}^n$  on  $\partial\Omega_{\text{in}}$ ,<sup>3</sup> and  $\underline{\kappa}_{\text{in}}$  in turn determines the inflow boundary data,  $\hat{\psi}_{\text{in}}$ , on  $\partial\Omega_{\text{in}} \times D$  for the Fokker–Planck equation. That is, for  $s \in \partial\Omega_{\text{in}}$ ,  $\hat{\psi}_{\text{in}}^n(s, \cdot) : \underline{q} \in D \mapsto \hat{\psi}_{\text{in}}^n(s, \underline{q}) \in \mathbb{R}$  for  $n = 1, \dots, N_T$  is determined by solving the  $\underline{q}$ -direction Fokker–Planck equation corresponding to  $\underline{\kappa}_{\text{in}}^n(s)$ , so that  $\hat{\psi}_{\text{in}}^n(s, \cdot) \in \mathcal{P}_N(D)$  for each  $n$ . Writing

$$\hat{\psi}_{\text{in}}(s, \underline{q}) = \sum_{k=1}^{N_D} \hat{\psi}_{\text{in},k}(s) Y_k(\underline{q}), \quad (s, \underline{q}) \in \partial\Omega_{\text{in}} \times D,$$

it then follows from (3.3.21) that  $\hat{\psi}_{\text{in},k}$  defines the inflow boundary data on  $\partial\Omega_{\text{in}}$  for  $\hat{\psi}_k$  in (3.3.33). In practice we only solve for  $\hat{\psi}_{\text{in}}$  at the nodes of  $\mathcal{T}_h$  on  $\partial\Omega_{\text{in}}$  so that we can impose the inflow boundary condition on the line function  $\hat{\psi}_k$  in an interpolatory sense. Notice also that we can compute the inflow boundary data for  $\hat{\psi}_{h,N}$  before we begin solving the Navier–Stokes–Fokker–Planck system, since  $\underline{u}_{\text{in}}$  and  $\underline{\kappa}_{\text{in}}$  are specified *a priori*.

We now define the algorithm for solving the Navier–Stokes–Fokker–Planck system. First of all, we initialise the system to the equilibrium state by setting  $\underline{u}_h^0 = \underline{0}$  on  $\Omega$ , and therefore  $\underline{\kappa}_{\text{in}}^0 = \nabla_x \underline{u}_h^0 = \underline{0}$  on  $\Omega$  also. Putting  $\underline{\kappa}_{\text{in}} = \underline{0}$  in (2.8.9), we can see that  $\psi = M$  is the corresponding equilibrium steady-state solution, and hence we set  $\hat{\psi}_{h,N}^0 = \sqrt{M} \in V_h \otimes \mathcal{P}_N(D)$  on  $\Omega \times D$ .<sup>4</sup> Also, for consistency with  $\hat{\psi}_{h,N}^0$ , we set  $\underline{\tau}_{h,N}^0 = \underline{I}$  on  $\Omega$ . Then, for  $n = 0, \dots, N_T - 1$ , we perform the following steps:

1. Compute  $\underline{u}_h^{n+1} \in \underline{V}_h$  and  $p_h^{n+1} \in P_h$  using the mixed finite element method discussed above for either the Navier–Stokes or Stokes system. We use the tensor  $\underline{\tau}_{h,N}^n$  in (4.2.5) or (4.2.7).
2. Use method I or method II to compute  $\hat{\psi}_{h,N}^{n+1} \in V_h \otimes \mathcal{P}_N(D)$  with  $\underline{\kappa}_{\text{in}}^n$  in (3.3.30) for method I or with  $\underline{\kappa}_{\text{in}}^{n+1}$  in (3.3.46) for method II, and  $\underline{u}_h^{n+1}$  in (3.3.33) for either method.
3. Using (1.3.37), compute  $\underline{\tau}_{h,N}^{n+1}$  on  $\Omega$  based on  $\hat{\psi}_{h,N}^{n+1} \in V_h \otimes \mathcal{P}_N(D)$ .
4. Return to 1. and continue marching in time.

Note that the  $\underline{\tau}_{h,N}$  terms in the momentum equations (4.2.5) or (4.2.7) are explicit in time. This allows the Stokes/Navier–Stokes equations to be coupled to the Fokker–Planck equation in a simple manner, but the drawback is that the algorithm defined in steps 1. to

<sup>2</sup>Strictly speaking, one has to be more careful: since the definition of the outflow boundary depends on  $\underline{u}$ , and  $\underline{u}$  is one of the unknowns, one cannot in general know *a priori* whether or not a specific portion of  $\partial\Omega$  is or isn't an outflow boundary. The numerical test problems considered here will, however, be simple enough to enable us to fix the location of  $\partial\Omega_{\text{out}}$  independently of  $\underline{u}$ .

<sup>3</sup>Since  $\underline{u}_{\text{in}}$  is a fully-developed flow, we assume that the velocity field upstream of  $\partial\Omega_{\text{in}}$  has the same profile  $\underline{u}_{\text{in}}$ ; this ensures that  $\nabla_x \underline{u}_{\text{in}}$  is well-defined on the inflow boundary.

<sup>4</sup>We assume here that  $\sqrt{M} \in \mathcal{P}_N(D)$ , which is reasonable according to Remark 2.7.1.

4. above is only conditionally stable. In Section 4.3 we use  $\Delta t = 0.01$  and this time-step size is sufficiently small to yield a reliable numerical method for the micro-macro problems that we consider.

### 4.3 Numerical results

In this section, we consider two distinct problems. The first is a planar contraction flow in the  $d = 2$  case, which we discuss in Section 4.3.1, and the second is a flow around a sphere in the  $d = 3$  case, considered in Section 4.3.2. For each of these two problems we present numerical results for one particular discrete space  $V_h \otimes \mathcal{P}_N(D)$ , but in each case we performed mesh refinement studies (*i.e.* we solved using a sequence of increasingly refined spaces) to ensure that the numerical results shown below are accurate.

#### 4.3.1 4-to-1 planar contraction flow

Contraction flows are standard benchmark problems in computational rheology because they are challenging from the numerical point of view and they also have practical relevance in industrial applications (for a detailed discussion of contraction flows see Chapter 8 of [103]). In this section we consider the coupled Navier–Stokes–Fokker–Planck model with  $\text{Re} = 1$  in a contracting domain, which is 10 units long, 4 units wide in the wider section and 1 unit wide in the narrow section. We set  $\partial\Omega_{\text{in}}$  and  $\partial\Omega_{\text{out}}$  to be the left-hand and right-hand boundaries of  $\Omega$ , respectively, and we let the top edge boundary be  $\partial\Omega_0$ . In this case, to save computational work we also imposed a symmetry boundary condition on the bottom boundary by setting the  $y$ -component of  $\underline{u}_h$  to zero there. We set  $\underline{u}_{\text{in}}$  to be a parabolic inflow profile, corresponding to steady Poiseuille flow in a channel, that vanishes at the top boundary and achieves its maximum value of  $U_{\text{max}} = 1$  at the symmetry boundary.

As specified in Chapter 3, we need  $\underline{\kappa} = \nabla_x \underline{u}_h \in \underline{\mathbb{L}}^\infty(\Omega)$  in order to use the alternating-direction methods I or II. Clearly, for any  $\underline{\mathbb{H}}^1(\Omega)$ -conforming finite element approximation,  $\underline{u}_h$ , this condition will be satisfied. Nevertheless, for the moment, let us consider the weak solution,  $\underline{u} \in \underline{\mathbb{H}}^k(\Omega)$  for some  $k > 0$ . In order to guarantee that  $\nabla_x \underline{u} \in \underline{\mathbb{L}}^\infty(\Omega)$ , we require the embedding  $\underline{\mathbb{H}}^{k-1}(\Omega) \subset \underline{\mathbb{L}}^\infty(\Omega)$  to hold; a sufficient condition for this embedding is that  $k - 1 > d/2$ , *i.e.* that  $k > 2$  when  $d = 2$ . However, contraction flows of polymeric fluids are typically simulated using ‘L-shaped’ domains and it is well known that the Stokes and Navier–Stokes equations exhibit a corner singularity on domains of this type so that in general  $\underline{u} \notin \underline{\mathbb{H}}^2(\Omega)$  (*cf.* Remark 5.10 in [46]). Therefore,  $\nabla_x \underline{u}$  will not, in general, belong to  $\underline{\mathbb{L}}^\infty(\Omega)$ , and hence the sequence  $\underline{\kappa}_h = \nabla_x \underline{u}_h$  will not be uniformly bounded in  $h$  as  $h \rightarrow 0_+$ . As a result, instead of an L-shaped domain, we use the physical space domain with a rounded corner shown in Figure 4.1. Also, in order to resolve the solution satisfactorily, the finite element mesh,  $\mathcal{T}_h$ , has been graded so that it is finer near the (rounded) corner.

We applied the algorithm defined in Section 4.2 for the coupled Navier–Stokes–Fokker–Planck system to the contraction flow problem described above. We set  $b = 12$ ,  $\text{Wi} = 0.8$ ,  $\gamma = 0.59$  and took 500 time-steps with  $\Delta t = 0.01$  so that  $T = 5$ . We used alternating-direction method II with basis  $\mathcal{A}$  and the  $p = 4$  quadrature rule on triangles for which  $Q_{\hat{K}} = 6$  (*cf.* Section 3.8.3) so that QH2 was satisfied. The mesh  $\mathcal{T}_h$  contained 905 triangular finite elements and therefore  $Q_\Omega = 5430$ .<sup>5</sup> Also, we used  $(N_r, N_\theta) = (20, 20)$  for the  $q$ -direction

<sup>5</sup>6335 quadrature points would have been required to satisfy QH1; hence we obtain a significant reduction



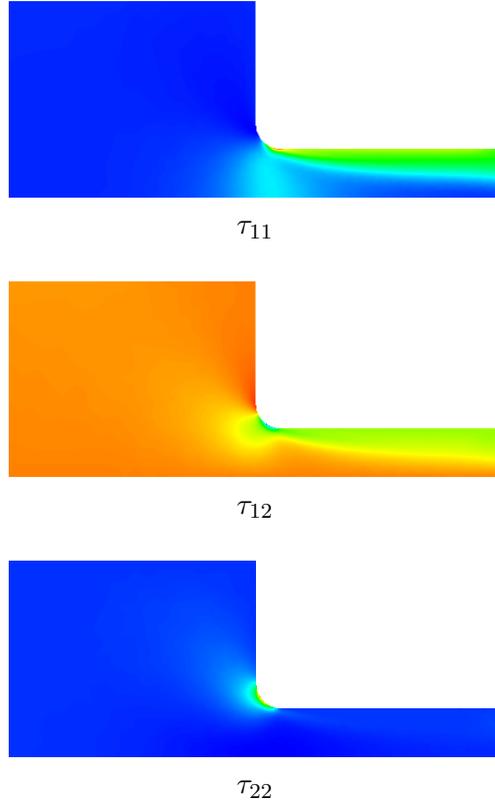


Figure 4.2: The components of  $\underline{\underline{\tau}}_{h,N}$  at  $T = 5$ . In the  $\tau_{11}$  plot, values range from 0.45 (blue) to 15.7 (red), in the  $\tau_{12}$  ( $= \tau_{21}$ ) plot we have -9.75 (blue) to 1.41 (red) and in the  $\tau_{22}$  plot, 0.46 (blue) to 11.5 (red). The polymeric extra-stress is largest in the region near the rounded corner.

the  $q$ -direction spectral method we used basis  $\mathcal{C}$  with  $(N_r, N_{\text{sph}}) = (12, 12)$ , so that  $N_D = 1092$ . Therefore, in each time-step, 72100 three-dimensional  $q$ -direction solves and 1092 three-dimensional  $x$ -direction solves were performed. We took 100 time-steps with  $\Delta t = 0.01$  to reach  $T = 1$ . Plots of the  $x$ -component of  $\underline{u}_h$  and of  $p_h$  at  $T = 1$  are shown in Figure 4.3. Also, the components of the polymeric extra-stress tensor at  $T = 1$  are shown in Figure 4.4. This computation was performed with  $N_{\text{proc}} = 128$  and it took 38.7 seconds to evaluate each time-step of the coupled Stokes-Fokker-Planck system.

## 4.4 Conclusions

In this chapter we introduced a deterministic multiscale algorithm for the micro-macro model of dilute polymeric fluids. This algorithm couples the alternating-direction scheme from Chapter 3 to a finite element method (for Stokes or Navier-Stokes) for computing the macroscopic velocity field. We used this algorithm to simulate two channel flows; a 4-to-1 contraction (with a rounded reentrant corner to avoid a singularity in  $\underline{u}$ ) in Section 4.3.1, and a flow around a spherical obstacle in a channel with square cross-section in Section 4.3.2.

We made extensive use of parallel computation in order to obtain the computational

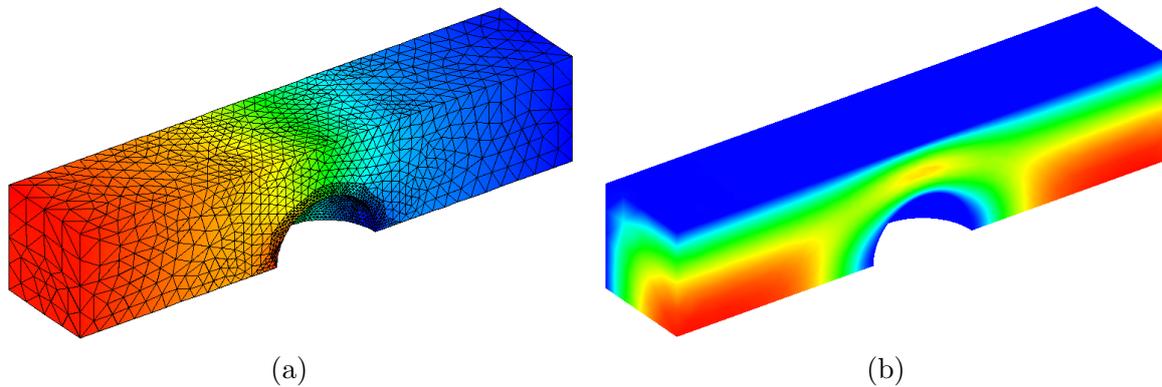


Figure 4.3: (a) Plot of the pressure,  $p_h \in P_h$ , at  $T = 1$ , with values ranging from 0.5 (blue) to 14.4 (red). Also, this plot shows the mesh  $\mathcal{T}_h$ . Note that the mesh is very fine in the vicinity of the spherical obstacle in order to resolve the solution structure in that region. (b) The  $x$ -component of the macroscopic velocity field at  $T = 1$ ; values range from 0 (blue) to 1 (red).

results in Section 4.3. In particular, to the best of our knowledge the micro-macro model has not previously been used in the case that  $\Omega \times D \in \mathbb{R}^6$  and this was only made feasible in Section 4.3.2 through the use of large-scale parallel computation.

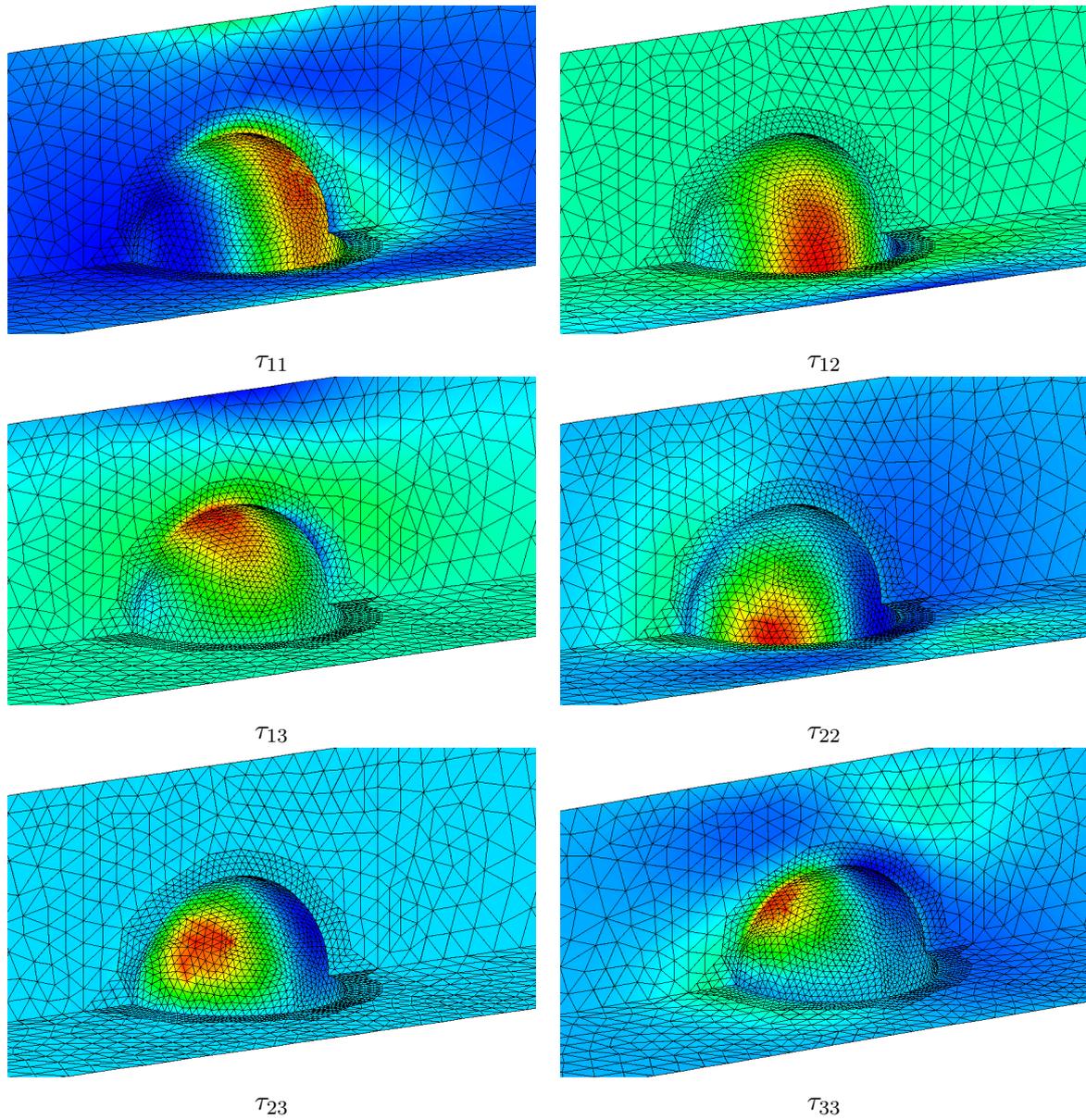


Figure 4.4: Plots of the components of the polymeric extra-stress tensor,  $\underline{\tau}_{h,N}$ , at  $T = 1$  for the channel flow around a spherical obstacle. The minimum (blue) and maximum (red) values in each plot are as follows;  $\tau_{11}$ : 0.53 to 6.25,  $\tau_{12}$ :  $-1.25$  to 2.41,  $\tau_{13}$ :  $-1.21$  to 2.5,  $\tau_{22}$ : 0.48 to 3.35,  $\tau_{23}$ :  $-0.33$  to 1.15 and  $\tau_{33}$ : 0.47 to 3.46.



## Chapter 5

# Existence of global weak solutions to Navier–Stokes–Fokker–Planck systems

### 5.1 Introduction

This chapter is concerned with the question of existence of global weak solutions to a system of nonlinear partial differential equations that arises from the kinetic theory of dilute polymer solutions. The solvent is an incompressible, viscous, isothermal Newtonian fluid confined to a bounded open set  $\Omega \subset \mathbb{R}^d$ ,  $d = 2$  or  $3$ , with boundary  $\partial\Omega$ . For the sake of simplicity of presentation, we shall suppose that  $\Omega$  has solid boundary  $\partial\Omega$ ; the velocity field  $\underline{u}$  will then satisfy the no-slip boundary condition  $\underline{u} = \underline{0}$  on  $\partial\Omega$ . The polymer chains, which are suspended in the solvent, are assumed not to interact with each other. The conservation of momentum and mass equations for the solvent then have the form of the incompressible Navier–Stokes equations in which the elastic *extra-stress* tensor  $\underline{\tau}$  (i.e., the polymeric part of the Cauchy stress tensor) appears as a source term:

Given  $T \in \mathbb{R}_{>0}$ , find  $\underline{u} : (\underline{x}, t) \in \bar{\Omega} \times [0, T] \mapsto \underline{u}(\underline{x}, t) \in \mathbb{R}^d$  and  $p : (\underline{x}, t) \in \Omega \times (0, T] \mapsto p(\underline{x}, t) \in \mathbb{R}$  such that

$$\frac{\partial \underline{u}}{\partial t} + (\underline{u} \cdot \nabla_x) \underline{u} - \nu \Delta_x \underline{u} + \nabla_x p = \underline{f} + \nabla_x \cdot \underline{\tau} \quad \text{in } \Omega \times (0, T], \quad (5.1.1a)$$

$$\nabla_x \cdot \underline{u} = 0 \quad \text{in } \Omega \times (0, T], \quad (5.1.1b)$$

$$\underline{u} = \underline{0} \quad \text{on } \partial\Omega \times (0, T], \quad (5.1.1c)$$

$$\underline{u}(\underline{x}, 0) = \underline{u}_0(\underline{x}) \quad \forall \underline{x} \in \Omega, \quad (5.1.1d)$$

where  $\underline{u}$  is the velocity field,  $p$  is the pressure of the fluid,  $\nu \in \mathbb{R}_{>0}$  is the viscosity of the solvent, and  $\underline{f}$  is the density of body forces acting on the fluid.

In the kinetic models under consideration here the extra-stress tensor  $\underline{\tau}$  is defined as the weighted mean of  $\psi$ , the probability density function of the (random) conformation vector of the polymer molecules (cf. (5.1.6) below). The Kolmogorov equation satisfied by  $\psi$  is a Fokker–Planck-type second-order parabolic equation whose transport coefficients depend on the velocity field  $\underline{u}$ .

Polymer solutions exhibit a range of non-Newtonian flow properties: in particular, the stress endured by a fluid element depends upon the history of deformations experienced by that element. Thereby, rheological properties of non-Newtonian fluids are governed by the flow-induced evolution of their internal microstructure. Following Keunings [66], a relevant feature of the microstructure is the *conformation* of the macromolecules, i.e., their orientation and the degree of stretching they experience. From the macroscopic viewpoint it is only the statistical distribution of conformations that matters: the macroscopic stress carried by each fluid element is governed by the distribution of polymer conformations within that element. Motivated by this observation, kinetic theories of polymeric fluids ignore quantum mechanical and atomistic effects and focus on “coarse-grained” models of the polymeric conformations. Depending on the level of coarse-graining, one may arrive at a hierarchy of kinetic models. For example, a dilute solution of linear polymers in a Newtonian solvent can be described in some detail by the freely jointed bead-rod *Kramers chain*, which comprises a number of beads (of the order of 100) connected by rigid linear segments. A coarser model of the same polymer is the freely jointed bead-spring chain, a *Rouse chain*, consisting of a smaller number of beads (of the order of 10) connected linearly by entropic springs. A coarser model still is the dumbbell model, which involves two beads connected by a spring; cf. Bird, Curtiss, Armstrong, and Hassager [23]. As has been emphasized by Keunings [66], such coarse-grained models are not meant to capture the detailed structure of the polymer. Rather, they are intended to describe, in more or less detail, the evolution of polymer conformations in a macroscopic flow.

Many of the interesting properties of dilute polymer solutions can be understood by modelling them as suspensions of simple coarse-grained objects (viz. dumbbells) in a Newtonian fluid. This chapter is devoted to the mathematical analysis of dumbbell models that are nonlinearly coupled Navier–Stokes–Fokker–Planck systems of partial differential equations: from the technical viewpoint these relatively simple models already exhibit many of the analytical difficulties encountered in the study of more complex models.

Suppose that the domain of admissible conformations (orientation vectors of polymer chains)  $D \subset \mathbb{R}^d$  is a balanced convex open set in  $\mathbb{R}^d$ ; the term *balanced* means that  $\underline{q} \in D$  if, and only if,  $-\underline{q} \in D$ . Hence, in particular,  $\underline{0} \in D$ . Typically,  $D$  is the whole of  $\mathbb{R}^d$  or a bounded open  $d$ -dimensional ball centred at the origin  $\underline{0} \in \mathbb{R}^d$ .

Let  $\mathcal{O} \subset [0, \infty)$  denote the image of  $D$  under the mapping  $\underline{q} \mapsto \frac{1}{2}|\underline{q}|^2$ , and consider the *spring-potential*  $U \in C^\infty(\mathcal{O}; \mathbb{R}_{\geq 0})$ . Clearly,  $0 \in \mathcal{O}$ . We shall suppose that  $U(0) = 0$  and that  $U$  is monotonic increasing and unbounded on  $\mathcal{O}$ . The elastic spring-force  $\underline{F} : D \subseteq \mathbb{R}^d \rightarrow \mathbb{R}^d$  is then defined by

$$\underline{F}(\underline{q}) = U'(\tfrac{1}{2}|\underline{q}|^2) \underline{q}. \quad (5.1.2)$$

**Example 5.1.1** In the Hookean dumbbell model, the spring force is defined by  $\underline{F}(\underline{q}) = \underline{q}$ , with  $\underline{q} \in D = \mathbb{R}^d$ , corresponding to  $U(s) = s$ ,  $s \in \mathcal{O} = [0, \infty)$ . Unfortunately, this simple model is physically unrealistic as it admits arbitrarily large extensions. We shall therefore assume in what follows that  $D$  is a bounded open ball in  $\mathbb{R}^d$  centred at the origin  $\underline{0} \in \mathbb{R}^d$ .  $\diamond$

We shall further suppose that there exist constants  $c_i > 0$ ,  $i = 1, 2, 3, 4$ , and  $\gamma > 1$  such that the (normalized) Maxwellian  $M$ , defined by

$$M(\underline{q}) = \frac{e^{-U(\frac{1}{2}|\underline{q}|^2)}}{\int_D e^{-U(\frac{1}{2}|\underline{q}|^2)} d\underline{q}},$$

and the associated potential  $U$  satisfy

$$c_1 [\text{dist}(\underline{q}, \partial D)]^\gamma \leq M(\underline{q}) \leq c_2 [\text{dist}(\underline{q}, \partial D)]^\gamma \quad \forall \underline{q} \in D, \quad (5.1.3a)$$

$$c_3 \leq [\text{dist}(\underline{q}, \partial D)] U'(\frac{1}{2}|\underline{q}|^2) \leq c_4 \quad \forall \underline{q} \in D. \quad (5.1.3b)$$

Observe that

$$M(\underline{q}) \nabla_{\underline{q}} [M(\underline{q})]^{-1} = -[M(\underline{q})]^{-1} \nabla_{\underline{q}} M(\underline{q}) = \nabla_{\underline{q}} U(\frac{1}{2}|\underline{q}|^2) = U'(\frac{1}{2}|\underline{q}|^2) \underline{q}. \quad (5.1.4)$$

Since  $[U(\frac{1}{2}|\underline{q}|^2)]^2 = (-\ln M(\underline{q}) + \text{Const.})^2$ , it follows from (5.1.3a,b) that (if  $\gamma > 1$ , as has been assumed here, then)

$$\int_D \left[ 1 + [U(\frac{1}{2}|\underline{q}|^2)]^2 + [U'(\frac{1}{2}|\underline{q}|^2)]^2 \right] M(\underline{q}) \, d\underline{q} < \infty. \quad (5.1.5)$$

**Example 5.1.2** In the FENE (finitely extensible nonlinear elastic) dumbbell model the spring force is given by

$$\underline{F}(\underline{q}) = \frac{1}{1 - |\underline{q}|^2/b} \underline{q}, \quad \underline{q} \in D = B(\underline{0}, b^{\frac{1}{2}}),$$

corresponding to  $U(s) = -\frac{b}{2} \ln(1 - \frac{2s}{b})$ ,  $s \in \mathcal{O} = [0, \frac{b}{2})$ . Here  $B(\underline{0}, b^{\frac{1}{2}})$  is a bounded open ball in  $\mathbb{R}^d$  centred at the origin  $\underline{0} \in \mathbb{R}^d$  and of fixed radius  $b^{\frac{1}{2}}$ , with  $b > 0$ . Direct calculations show that the Maxwellian  $M$  and the elastic potential  $U$  of the FENE model satisfy conditions (5.1.3a,b) with  $\gamma = \frac{b}{2}$  provided that  $b > 2$ . Thereby (5.1.5) also holds for  $b > 2$ .

It is interesting to note that in the (equivalent) stochastic version of the FENE model a solution to the system of stochastic differential equations associated with the Fokker–Planck equation exists and has trajectorial uniqueness if, and only if,  $b > 2$  (cf. Jourdain, Lelièvre, and Le Bris [62] for details). Thus, the assumption  $\gamma > 1$  can be seen as the weakest reasonable requirement on the decay-rate of  $M$  in (5.1.3a) as  $\text{dist}(\underline{q}, \partial D) \rightarrow 0$ .  $\diamond$

Due to the flow-induced thermal agitation, polymer molecules are subjected to Brownian forces. Let  $(\underline{x}, \underline{q}, t) \mapsto \psi(\underline{x}, \underline{q}, t)$  denote the probability density function corresponding to the vector-valued stochastic process  $(\underline{X}(t), \underline{Q}(t))$ , where  $\underline{X}(t) \in \Omega$  is the position vector of the centre of mass of the dumbbell at time  $t \geq 0$ , and  $\underline{Q}(t) \in D$  is the conformation (or end-to-end) vector of the dumbbell at time  $t \geq 0$ . Roughly speaking,  $\psi(\underline{x}, \underline{q}, t)$  represents the probability at time  $t$  of finding the centre of mass of a dumbbell at  $\underline{x}$  and having elongation vector  $\underline{q}$ .

The governing equations of the coupled Navier–Stokes–Fokker–Planck model are (5.1.1a–d), where the extra-stress tensor  $\underline{\tau}$  is defined by

$$\underline{\tau}(\underline{x}, t) = k_B \mathcal{T} \left( \int_D \underline{q} \underline{q}^T U'(\frac{1}{2}|\underline{q}|^2) \psi(\underline{x}, \underline{q}, t) \, d\underline{q} - \rho(\underline{x}, t) \underline{\underline{I}} \right), \quad (5.1.6)$$

with the density of polymer chains located at  $\underline{x}$  at time  $t$  given by

$$\rho(\underline{x}, t) = \int_D \psi(\underline{x}, \underline{q}, t) \, d\underline{q}. \quad (5.1.7)$$

The probability density function  $\psi$  is a solution of the Fokker–Planck equation

$$\frac{\partial \psi}{\partial t} + (\underline{u} \cdot \nabla_x) \psi + \nabla_q \cdot (\underline{\sigma}(\underline{u}) \underline{q} \psi) = \frac{1}{2\lambda} \nabla_q \cdot (\nabla_q \psi + U'(\frac{1}{2}|\underline{q}|^2) \underline{q} \psi) + \varepsilon \Delta_x \psi, \quad (5.1.8)$$

with  $\underline{\sigma}(\underline{v}) \equiv \nabla_x \underline{v}$ , where  $(\nabla_x \underline{v})(\underline{x}, t) \in \mathbb{R}^{d \times d}$  and  $\{\nabla_x \underline{v}\}_{ij} = \frac{\partial v_i}{\partial x_j}$  (cf. Barrett and Süli [11]). Here,  $\varepsilon = \ell_0^2/(8\lambda)$  is the centre-of-mass diffusion coefficient of the dumbbells,  $\ell_0 \ll \text{diam}(\Omega)$  is the characteristic microscopic length-scale (i.e. the characteristic dumbbell size) and  $\lambda = \zeta/4H$ . The parameter  $\lambda \in \mathbb{R}_{>0}$  characterizes the elastic relaxation property of the fluid,  $\zeta > 0$  is a friction coefficient,  $H > 0$  is a spring-constant,  $k_B > 0$  is the Boltzmann constant and  $\mathcal{T} > 0$  is the absolute temperature.

A noteworthy feature of (5.1.11) compared to classical Fokker–Planck equations for bead-spring models in the literature is the presence of the  $\underline{x}$ -dissipative centre-of-mass diffusion term  $\varepsilon \Delta_x \psi \equiv (\ell_0^2/8\lambda) \Delta_x \psi$  on the right-hand side of the Fokker–Planck equation (5.1.8). We refer to Barrett and Süli [11] for the derivation of (5.1.8) and the mathematical justification of the presence of the centre-of-mass diffusion term  $\varepsilon \Delta_x \psi$ ; see also the recent article by Schieber [108] concerning generalized dumbbell models with centre-of-mass diffusion. In standard derivations of bead-spring models the centre-of-mass diffusion term is routinely omitted, on the grounds that it is several orders of magnitude smaller than the other terms in the equation. Indeed, when  $L \approx 1$  is a characteristic macroscopic length-scale (such as, for example,  $\text{diam}(\Omega)$ ), Bhave, Armstrong, and Brown [19] estimate the ratio  $\ell_0^2/L^2$  to be in the range of about  $10^{-9}$  to  $10^{-7}$ . However, the omission of the term  $\varepsilon \Delta_x \psi$  from (5.1.8) in the case of a heterogeneous solvent velocity  $\underline{u}(\underline{x}, t)$  is a mathematically counterproductive model reduction. When  $\varepsilon \Delta_x \psi$  is absent, (5.1.8) becomes a degenerate parabolic equation exhibiting hyperbolic behaviour with respect to  $(\underline{x}, t)$ . Since the study of weak solutions to the coupled problem requires one to work with velocity fields  $\underline{u}$  that have very limited Sobolev regularity (typically  $\underline{u} \in L^\infty(0, T; \underline{L}^2(\Omega)) \cap L^2(0, T; \underline{H}_0^1(\Omega))$ ), one is then forced into the technically unpleasant framework of hyperbolically degenerate parabolic equations with rough transport coefficients (cf. Ambrosio [2] and DiPerna and Lions [40]). The resulting difficulties are further exacerbated by the fact that, when  $D$  is bounded, a typical spring force  $\underline{F}(\underline{q})$  for a finitely extensible model (such as FENE) explodes as  $\underline{q}$  approaches  $\partial D$ ; see Example 5.1.2 above. For these reasons, here we shall retain the centre-of-mass diffusion term in (5.1.8). At the macroscopic level, centre-of-mass diffusion can be seen as stress diffusion: in the case of the Hookean model with centre-of-mass diffusion, the corresponding macroscopic model is Oldroyd-B with stress diffusion. For a careful numerical study of the Oldroyd-B model with stress diffusion, we refer to the paper of Sureshkumar and Beris [118]; see also the paper of Bhave, Armstrong and Brown [19].

We conclude this introduction with a brief survey of recent developments on the analysis of classical bead-spring models; with the exception of Barrett and Süli [11] mentioned above and Bhave, Armstrong and Brown [19] and El-Kareh and Leal [45], all articles cited consider models that correspond to formally letting  $\varepsilon = 0$  in (5.1.8), i.e., omitting the centre-of-mass diffusion term.

An early contribution to the existence and uniqueness of local-in-time solutions to a family of bead-spring type polymeric flow models is due to Renardy [106]. While the class of potentials  $\underline{F}(\underline{q})$  considered by Renardy [106] (cf. hypotheses (F) and (F') on pp. 314–315) does include the case of Hookean dumbbells, it excludes the practically relevant case of the FENE model (see Example 5.1.2 above). More recently, E, Li, and Zhang [43] and Li,

Zhang, and Zhang [81] have revisited the question of local existence of solutions for dumbbell models. A further development in this direction is the work of Zhang and Zhang [127], where the local existence of regular solutions to FENE-type models has been shown. All of these papers require high regularity of the initial data. More recently, Lin, Zhang and Zhang [85] have shown the global existence of smooth solution to the two-dimensional FENE dumbbell model.

Constantin [37] has considered the Navier–Stokes equations coupled to nonlinear Fokker–Planck equations describing the evolution of the probability distribution of the particles interacting with the fluid. He described, in the case when  $D$  is a Riemannian manifold, relations determining the coefficients of the stresses added in the fluid by the particles; these relations link the extra stresses to the kinematic effect of the fluid velocity on the particles and to the interparticle interaction potential. In equations (of Type 1, in the terminology of Constantin [37]) where the extra stresses depend linearly on the particle distribution density, as is the case in the present chapter, the energy balance requires a response potential. In equations (of Type 2) where the added stresses depend quadratically on the particle distribution, it is shown that energy balance can be achieved without a dynamic response potential, and global existence of smooth solutions is shown if inertial effects are neglected. The necessary relationship (eq. (2.14) in Constantin [37]) for the existence of a Lyapunov function in the sense of Theorem 2.2 of Constantin [37] does not hold for the polymer models considered in the present chapter.

Otto and Tzavaras [102] have investigated the Doi model (which is similar to a Hookean model (cf. Example 5.1.1 above), except that  $D = S^2$ ) for suspensions of rod-like molecules in the dilute regime. For certain parameter values, the velocity gradient vs. stress relation defined by the stationary and homogeneous flow is not rank-one monotone. They considered the evolution of possibly large perturbations of stationary flows and proved that, even in the absence of a microscopic cutoff, discontinuities in the velocity gradient cannot occur in finite time.

Jourdain, Lelièvre, and Le Bris [62] studied the existence of solutions to the FENE model in the case of a simple Couette flow. By using tools from the theory of stochastic differential equations, they established the existence of a unique local-in-time solution to the FENE model in two space dimensions ( $d = 2$ ) when the velocity field  $\underline{u}$  is unidirectional and of the particular form  $\underline{u}(x_1, x_2) = (u_1(x_2), 0)^T$ . The notion of solution for which existence is proved in the paper of Jourdain, Lelièvre, and Le Bris [62] is mixed *deterministic-stochastic* in the sense that it is deterministic in the “macroscopic” variable  $\underline{x}$  but stochastic in the “microscopic” variable  $\underline{q}$ . In contrast, our notion of solution (cf. Section 5.3 below) is deterministic both macroscopically and microscopically, since the microscales are modelled here by the probability density function  $\psi(\underline{x}, \underline{q}, t)$ . The choice between these different notions of solution has far-reaching consequences on computational simulation: mixed deterministic-stochastic notions of solution necessitate the use of Monte Carlo-type algorithms for the numerical approximation of polymer configurations, as proposed in the monograph of Öttinger [101] and, for example, in the paper of Jourdain, Lelièvre, and Le Bris [61]; whereas weak solutions in the sense considered herein can be approximated by entirely deterministic (e.g., Galerkin-type) schemes, as was done, for example, in Lozinski, Chauvière, Fang, and Owens [92] and Lozinski, Owens, and Fang [93]—at the cost of solving a Fokker–Planck equation in  $2d$  spatial dimensions.

In the case of Hookean dumbbells, and assuming  $\varepsilon = 0$ , the coupled microscopic-macroscopic

model described above yields, formally, taking the second moment of  $q \mapsto \psi(q, \underline{x}, t)$ , the fully macroscopic, Oldroyd-B model of viscoelastic flow. Lions and Masmoudi [86] have shown the existence of global-in-time weak solutions to the Oldroyd-B model in a simplified corotational setting (i.e. with  $\sigma(\underline{y}) = \nabla_{\underline{x}} \underline{y}$  replaced by  $\frac{1}{2}(\nabla_{\underline{x}} \underline{y} - (\nabla_{\underline{x}} \underline{y})^T)$ ). The argument of Lions and Masmoudi [86] is based on exploiting the propagation in time of the compactness of the solution (i.e. the property that if one takes a sequence of weak solutions which converges weakly and such that the corresponding sequence of initial data converges strongly, then the weak limit is also a solution) and the DiPerna–Lions [40] theory of renormalized solutions to linear hyperbolic equations with nonsmooth transport coefficients. It is not known if an identical global existence result for the Oldroyd-B model also holds in the absence of the crucial assumption that the drag term is corotational. We note in passing that, assuming  $\varepsilon > 0$ , the coupled microscopic-macroscopic model above yields, taking the appropriate moments in the case of Hookean dumbbells, a dissipative version of the Oldroyd-B model. In this sense, the Hookean dumbbell model has a macroscopic closure: it is the Oldroyd-B model when  $\varepsilon = 0$ , and a dissipative version of Oldroyd-B when  $\varepsilon > 0$  (cf. Barrett and Süli [11]). In contrast, the FENE model is not known to have an exact closure at the macroscopic level, though Du, Liu and Yu [42] and Yu, Du, and Liu [125] have recently considered the analysis of approximate closures of the FENE model. Lions and Masmoudi [87] proved the global existence of weak solutions for the corotational FENE dumbbell model, once again corresponding to the case of  $\varepsilon = 0$ , and the Doi model, also called the rod model. As in Lions and Masmoudi [86], the proof is based on propagation of compactness; see also the related paper of Masmoudi [96].

Previously, El-Kareh and Leal [45] had proposed a macroscopic model, with added dissipation in the equation that governs the evolution of the conformation tensor  $\underline{\underline{A}}(\underline{x}, t) := \int_D \underline{q} \underline{q}^T U'(\frac{1}{2}|\underline{q}|^2) \psi(\underline{x}, \underline{q}, t) d\underline{q}$ , in order to account for Brownian motion across streamlines; the model can be thought of as an approximate macroscopic closure of a FENE-type micro-macro model with centre-of-mass diffusion.

Barrett, Schwab, and Süli [10] established the existence of, global in time, weak solutions to the coupled microscopic-macroscopic model (5.1.1a–d) and (5.1.8) with  $\varepsilon = 0$ , an  $\underline{x}$ -mollified velocity gradient in the Fokker–Planck equation and an  $\underline{x}$ -mollified probability density function  $\psi$  in the Kramers expression—admitting a large class of potentials  $U$  (including the Hookean dumbbell model as well as general FENE-type models); in addition to these mollifications,  $\underline{y}$  in the  $\underline{x}$ -convective term  $(\underline{y} \cdot \nabla_{\underline{x}})\psi$  in the Fokker–Planck equation was also mollified. Unlike Lions and Masmoudi [86], the arguments in Barrett, Schwab, and Süli [10] did not require the assumption that the drag term was corotational in the FENE case. The mollification  $S_\alpha$  of the velocity field  $\underline{y}$  that was considered in Barrett, Schwab and Süli [10] was stimulated by the Leray- $\alpha$  model of the incompressible Navier–Stokes equations (the viscous Camassa–Holm equations), proposed by Foias, Holm, and Titi [50], with the mollified velocity field  $S_\alpha \underline{y}$  defined as the solution of a Helmholtz–Stokes problem, thus ensuring that the mollified velocity field  $S_\alpha \underline{y}$  is still divergence-free and satisfies the same boundary condition as  $\underline{y}$ .

In Barrett and Süli [11], we derived the coupled Navier–Stokes–Fokker–Planck model with centre-of-mass diffusion stated above. The anisotropic Friedrichs mollifiers, which naturally arise in the derivation of the model in the Kramers expression for the extra stress tensor and in the drag term in the Fokker–Planck equation, were replaced by isotropic Friedrichs mollifiers. We established the existence of global-in-time weak solutions to the model for a general class of spring-force-potentials including in particular the FENE potential. We justified also,

through a rigorous limiting process, certain classical reductions of this model appearing in the literature that exclude the centre-of-mass diffusion term from the Fokker-Planck equation on the grounds that the diffusion coefficient is small relative to other coefficients featuring in the equation. In the case of a corotational drag term we performed a rigorous passage to the limit as the Friedrichs mollifiers in the Kramers expression and the drag term converge to identity operators.

In the present chapter neither the probability density function  $\psi$  in the Kramers expression (5.1.6) nor the velocity field  $\underline{u}$  in the drag term

$$\nabla_q \cdot (\underline{\sigma}(\underline{u}) \underline{q} \psi) = \nabla_q \cdot \left[ \underline{\sigma}(\underline{u}) \underline{q} M \left( \frac{\psi}{M} \right) \right] \quad (5.1.9)$$

appearing in (5.1.8) will be mollified. Instead, motivated by recent papers of Jourdain, Lelièvre, Le Bris, and Otto [63] and Lin, Liu, and Zhang [84] (see also Arnold, Markowich, Toscani, and Unterreiter [6], and Desvillettes and Villani [39]) concerning the convergence of the probability density function  $\psi$  to its equilibrium value  $\psi_\infty(\underline{x}, q) := M(q)$  (corresponding to the equilibrium value  $\underline{u}_\infty(\underline{x}) := \underline{0}$  of the velocity field) in the absence of body forces  $\underline{f}$ , we observe that if  $\psi/M$  is bounded above then, for  $L \in \mathbb{R}_{>0}$  sufficiently large, the drag term (5.1.9) is equal to

$$\nabla_q \cdot \left[ \underline{\sigma}(\underline{u}) \underline{q} M \beta^L \left( \frac{\psi}{M} \right) \right],$$

where  $\beta^L \in C(\mathbb{R})$  is a cut-off function defined as

$$\beta^L(s) := \begin{cases} s & \text{for } s \leq L, \\ L & \text{for } L \leq s. \end{cases} \quad (5.1.10)$$

It follows that, for  $L \gg 1$ , any solution  $\psi$  of (5.1.8), such that  $\psi/M$  is bounded above, also satisfies

$$\begin{aligned} \frac{\partial \psi}{\partial t} + (\underline{u} \cdot \nabla_x) \psi + \nabla_q \cdot \left[ \underline{\sigma}(\underline{u}) \underline{q} M \beta^L \left( \frac{\psi}{M} \right) \right] \\ = \frac{1}{2\lambda} \nabla_q \cdot \left( M \nabla_q \left( \frac{\psi}{M} \right) \right) + \varepsilon \Delta_x \psi \quad \text{in } \Omega \times D \times (0, T]. \end{aligned} \quad (5.1.11)$$

We impose the following boundary and initial conditions:

$$\left[ \frac{M}{2\lambda} \nabla_q \left( \frac{\psi}{M} \right) - \underline{\sigma}(\underline{u}) \underline{q} M \beta^L \left( \frac{\psi}{M} \right) \right] \cdot \frac{\underline{q}}{|\underline{q}|} = 0 \quad \text{on } \Omega \times \partial D \times (0, T], \quad (5.1.12a)$$

$$\varepsilon \nabla_x \psi \cdot \underline{n}_{\partial\Omega} = 0 \quad \text{on } \partial\Omega \times D \times (0, T], \quad (5.1.12b)$$

$$\psi(\underline{x}, q, 0) = \psi_0(\underline{x}, q) \geq 0 \quad \forall (\underline{x}, q) \in \Omega \times D; \quad (5.1.12c)$$

where  $\underline{q}$  is normal to  $\partial D$ , as  $D$  is a bounded ball centred at the origin, and  $\underline{n}_{\partial\Omega}$  is the unit outward normal to  $\partial\Omega$ . Here  $\int_D \psi_0(\underline{x}, q) dq = 1$  for a.e.  $\underline{x} \in \Omega$ .

The coupled problem (5.1.1a–d), (5.1.6), (5.1.7), (5.1.11), (5.1.12a–c) will be referred to as a *dumbbell model with microscopic cut-off*. In order to highlight the dependence on  $\varepsilon$  and

$L$ , in subsequent sections the solution to (5.1.11), (5.1.12a–c) will be labelled  $\psi_{\varepsilon,L}$ . Due to the coupling of (5.1.11) to (5.1.1a) through (5.1.6), the velocity and the pressure will also depend on  $\varepsilon$  and  $L$  and we shall therefore denote them in subsequent sections by  $u_{\varepsilon,L}$  and  $p_{\varepsilon,L}$ .

A detailed argument for introducing cut-off, albeit of a very different nature, was put forward in El-Kareh and Leal [45] (cf. (3.10a,b)); the authors used a nonnegative function  $q \in D \mapsto g(|q|)$  that is compactly supported in  $D$ , in both the right-hand side of the momentum equation and in the macroscopic counterpart of the Fokker–Planck equation, in order to truncate the unbounded function  $q \in D \mapsto U'(\frac{1}{2}|q|^2) = 1/(1 - |q|^2/b)$ ,  $|q|^2 < b$ , to a bounded compactly supported function  $q \in D \mapsto g(|q|)U'(\frac{1}{2}|q|^2)$ .

The cut-off  $\beta^L$  proposed here has several attractive properties. We observe that the couple  $\{u_\infty, \psi_\infty\}$ , defined by  $u_\infty(x) := 0$  and  $\psi_\infty(x, q) := M(q)$ , is still an equilibrium solution of (5.1.1a–d) with  $f = 0$ , (5.1.6), (5.1.7), (5.1.11), (5.1.12a–c) for all  $L > 0$ . Thus, unlike the truncation of the (unbounded) potential proposed in El-Kareh and Leal [45], the introduction of the cut-off function  $\beta^L$  into the Fokker–Planck equation (5.1.8) does not alter the equilibrium solution  $(u_\infty, \psi_\infty)$  of the original Navier–Stokes–Fokker–Planck system. In addition, the boundary conditions for  $\psi$  on  $\partial\Omega \times D \times (0, T]$  and  $\Omega \times \partial D \times (0, T]$  ensure that

$$\frac{1}{|\Omega|} \int_{\Omega \times D} \psi(x, q, t) dq dx = \frac{1}{|\Omega|} \int_{\Omega \times D} \psi_0(x, q) dq dx = 1 \quad \forall t \in \mathbb{R}_{\geq 0}.$$

Our objective is to establish the existence of, global in time, weak solutions to the the dumbbell model with microscopic cut-off. The chapter is structured as follows. We begin, in Section 5.2, by stating the weak formulation of the coupled Navier–Stokes–Fokker–Planck system with centre-of-mass diffusion and microscopic cut-off, for the general class of potentials  $U$  under consideration. In particular, the FENE model fits into the general setting. In Section 5.3 we embark on the proof of existence of weak solutions to our model. We introduce a family of weighted Sobolev spaces that provide the natural functional-analytic framework for the problem: the weight of the space is the Maxwellian induced by the potential  $U$  appearing in the Fokker–Planck equation. Our proof requires a special compact embedding result in these Maxwellian-weighted Sobolev spaces, which is proved in the Appendix to this chapter by combining compact embedding theorems by Antoci [5] and Shakhmurov [113]. The proof of existence of global weak solutions to the coupled Navier–Stokes–Fokker–Planck system (5.1.1a–d), (5.1.6), (5.1.7), (5.1.11), (5.1.12a–c) then rests on a weak-convergence argument. A key ingredient, resulting in sufficiently strong a-priori bounds, is a special testing procedure based on the convex entropy function

$$s \in \mathbb{R}_{\geq 0} \mapsto \mathcal{F}(s) := s(\ln s - 1) + 1 \in \mathbb{R}_{\geq 0}$$

in the weak formulation of the Fokker–Planck equation. This leads to a fortuitous cancellation of the extra stress term on the right-hand side of the Navier–Stokes equation with the drag term in the Fokker–Planck equation and results in an  $L^\infty(0, T; L^1(\Omega))$  bound on the relative entropy  $\mathcal{E}_M(\psi)$  of  $\psi$  with respect to the equilibrium solution  $\psi_\infty = M$ , where

$$\mathcal{E}_M(\psi) := \int_D \mathcal{F}\left(\frac{\psi}{M}\right) M(q) dq.$$

The choice of the entropy function  $\mathcal{F}$  in the present context has been motivated by the papers Arnold, Markowich, Toscani, and Unterreiter [6], Desvillettes and Villani [39],

Jourdain, Lelièvre, Le Bris, and Otto [63] and Lin, Liu, and Zhang [84] cited above. It is important to note that the cut-off function  $\beta^L$  and the entropy function  $\mathcal{F}$  are closely related, viz.  $\beta^L(s) = \min(1/\mathcal{F}''(s), L)$ , and this connection will play a crucial role in our argument. Due to the fact that  $\mathcal{F}''(s)$  is unbounded at  $s = 0$ , in Section 5.3 the strictly convex entropy function  $\mathcal{F}$  will be replaced by a strictly convex regularization  $\mathcal{F}_\delta^L$  whose second derivative is bounded above by  $1/\delta$  and bounded below by  $1/L$ ,  $\delta \in (0, 1)$ ,  $L > 1$ ; at the same time the cut-off function  $\beta^L$  will be replaced by a strictly positive cut-off function  $\beta_\delta^L$  defined by  $\beta_\delta^L(s) = 1/[\mathcal{F}_\delta^L]''(s)$ . The existence of global weak solutions to the regularized cut-off problem is shown in Section 5.3.1. In Section 5.3.2 we then pass to the limit  $\delta \rightarrow 0_+$  with the regularization parameter  $\delta$ , to deduce the existence of a global weak solution to the coupled Navier–Stokes–Fokker–Planck system (5.1.1a–d), (5.1.6), (5.1.7), (5.1.11), (5.1.12a–c) with microscopic cut-off. Ideally, one would like to replace  $\beta^L(s) = \min(s, L)$  by  $\beta(s) = s$  in the Fokker–Planck equation. However, our current proof of existence in the general noncorotational case requires the presence of the microscopic cut-off function  $\beta^L$  on the drag term. Nevertheless, in the case of a corotational drag term at least passage to the limit  $L \rightarrow \infty$  recovers the Fokker–Planck equation (5.1.8), without cut-off (see Remark 5.3.9).

The convergence analysis of a general class of Galerkin-type approximations to the coupled corotational Navier–Stokes–Fokker–Planck model, which is mentioned above and was formulated in Barrett and Süli [11], was considered in Barrett and Süli [13]; for the convergence analysis of finite element approximations to the general noncorotational model with cut-off, considered herein, we refer to the discussion in the next chapter.

## 5.2 The polymer model

We term polymer models, under consideration here, microscopic–macroscopic-type models, since the continuum mechanical *macroscopic* equations of incompressible fluid flow are coupled to a *microscopic* model: the Fokker–Planck equation describing the statistical properties of particles in the continuum. We first present these equations and collect assumptions on the parameters in the model.

Let  $\Omega \subset \mathbb{R}^d$  be a bounded open set with a Lipschitz-continuous boundary  $\partial\Omega$ , and suppose that the set  $D$  of admissible elongation vectors  $q$  in (5.1.8) is a bounded open ball in  $\mathbb{R}^d$ ,  $d = 2$  or  $3$ , centred at the origin.

Gathering (5.1.1a–d), (5.1.6), and (5.1.8) together, we then consider the following initial-boundary-value problem dependent on the parameters  $\varepsilon \ll 1$  and  $L \gg 1$ :

( $P_{\varepsilon,L}$ ) Find  $\underline{u}_{\varepsilon,L} : (\underline{x}, t) \in \mathbb{R}^{d+1} \mapsto \underline{u}_{\varepsilon,L}(\underline{x}, t) \in \mathbb{R}^d$  and  $p_{\varepsilon,L} : (\underline{x}, t) \in \mathbb{R}^{d+1} \mapsto p_{\varepsilon,L}(\underline{x}, t) \in \mathbb{R}$  such that

$$\frac{\partial \underline{u}_{\varepsilon,L}}{\partial t} + (\underline{u}_{\varepsilon,L} \cdot \nabla_x) \underline{u}_{\varepsilon,L} - \nu \Delta_x \underline{u}_{\varepsilon,L} + \nabla_x p_{\varepsilon,L} = \underline{f} + \nabla_x \cdot \underline{\tau}(\psi_{\varepsilon,L}) \quad (5.2.1a)$$

$$\text{in } \Omega \times (0, T],$$

$$\nabla_x \cdot \underline{u}_{\varepsilon,L} = 0 \quad \text{in } \Omega \times (0, T], \quad (5.2.1b)$$

$$\underline{u}_{\varepsilon,L} = \mathbf{0} \quad \text{on } \partial\Omega \times (0, T], \quad (5.2.1c)$$

$$\underline{u}_{\varepsilon,L}(\cdot, 0) = \underline{u}_0 \quad \text{on } \Omega, \quad (5.2.1d)$$

where  $\nu \in \mathbb{R}_{>0}$  is the given viscosity,  $f(\underline{x}, t)$  is the given body force and  $\underline{\tau}(\psi_{\epsilon, L}) : (\underline{x}, t) \in \mathbb{R}^{d+1} \mapsto \underline{\tau}(\psi_{\epsilon, L})(\underline{x}, t) \in \mathbb{R}^{d \times d}$  is the symmetric extra-stress tensor, dependent on a probability density function  $\psi_{\epsilon, L} : (\underline{x}, \underline{q}, t) \in \mathbb{R}^{2d+1} \mapsto \psi_{\epsilon, L}(\underline{x}, \underline{q}, t) \in \mathbb{R}$ , defined as

$$\underline{\tau}(\psi_{\epsilon, L}) = k_B \mathcal{T} (\underline{C}(\psi_{\epsilon, L}) - \rho(\psi_{\epsilon, L}) \underline{I}). \quad (5.2.2)$$

Here  $k_B, \mathcal{T} \in \mathbb{R}_{>0}$  are, respectively, the Boltzmann constant and the absolute temperature,  $\underline{I}$  is the unit  $d \times d$  tensor, and

$$\underline{C}(\psi_{\epsilon, L})(\underline{x}, t) = \int_D \psi_{\epsilon, L}(\underline{x}, \underline{q}, t) U'(\frac{1}{2}|\underline{q}|^2) \underline{q} \underline{q}^T d\underline{q} \quad (5.2.3a)$$

and

$$\rho(\psi_{\epsilon, L})(\underline{x}, t) = \int_D \psi_{\epsilon, L}(\underline{x}, \underline{q}, t) d\underline{q}. \quad (5.2.3b)$$

The Fokker–Planck equation with microscopic cut-off satisfied by  $\psi_{\epsilon, L}$  is:

$$\begin{aligned} \frac{\partial \psi_{\epsilon, L}}{\partial t} + (\underline{u}_{\epsilon, L} \cdot \underline{\nabla}_x) \psi_{\epsilon, L} + \underline{\nabla}_q \cdot \left[ \underline{g}(\underline{u}_{\epsilon, L}) \underline{q} M \beta^L \left( \frac{\psi_{\epsilon, L}}{M} \right) \right] \\ = \frac{1}{2\lambda} \underline{\nabla}_q \cdot \left( M \underline{\nabla}_q \left( \frac{\psi_{\epsilon, L}}{M} \right) \right) + \varepsilon \Delta_x \psi_{\epsilon, L} \quad \text{in } \Omega \times D \times (0, T]. \end{aligned} \quad (5.2.4)$$

Here,  $\underline{g}(\underline{u}) \equiv \underline{\nabla}_x \underline{u}$  and, for a given  $L \gg 1$ ,  $\beta^L \in C(\mathbb{R})$  is defined by (5.1.10).

We impose the following boundary and initial conditions:

$$\left[ \frac{M}{2\lambda} \underline{\nabla}_q \left( \frac{\psi_{\epsilon, L}}{M} \right) - \underline{g}(\underline{u}_{\epsilon, L}) \underline{q} M \beta^L \left( \frac{\psi_{\epsilon, L}}{M} \right) \right] \cdot \frac{\underline{q}}{|\underline{q}|} = 0 \quad \text{on } \Omega \times \partial D \times (0, T], \quad (5.2.5a)$$

$$\varepsilon \underline{\nabla}_x \psi_{\epsilon, L} \cdot \underline{n}_{\partial\Omega} = 0 \quad \text{on } \partial\Omega \times D \times (0, T], \quad (5.2.5b)$$

$$\psi_{\epsilon, L}(\underline{x}, \underline{q}, 0) = \psi_0(\underline{x}, \underline{q}) \geq 0 \quad \forall (\underline{x}, \underline{q}) \in \Omega \times D, \quad (5.2.5c)$$

where  $\underline{n}_{\partial\Omega}$  is the unit outward normal to  $\partial\Omega$ . Here  $\int_D \psi_0(\underline{x}, \underline{q}) d\underline{q} = 1$  for a.e.  $\underline{x} \in \Omega$ . The boundary conditions for  $\psi_{\epsilon, L}$  on  $\partial\Omega \times D \times (0, T]$  and  $\Omega \times \partial D \times (0, T]$  have been chosen so as to ensure that  $\int_{\Omega \times D} \psi_{\epsilon, L}(\underline{x}, \underline{q}, t) d\underline{q} d\underline{x} = \int_{\Omega \times D} \psi_0(\underline{x}, \underline{q}) d\underline{q} d\underline{x} = |\Omega|$  for all  $t \geq 0$ .

### 5.3 Existence of global weak solutions

Let

$$\underline{\mathbb{H}} := \{w \in \underline{\mathbb{L}}^2(\Omega) : \underline{\nabla}_x \cdot w = 0\} \quad \text{and} \quad \underline{\mathbb{V}} := \{w \in \underline{\mathbb{H}}_0^1(\Omega) : \underline{\nabla}_x \cdot w = 0\}, \quad (5.3.1)$$

where the divergence operator  $\underline{\nabla}_x \cdot$  is to be understood in the sense of vector-valued distributions on  $\Omega$ . Let  $\underline{\mathbb{V}}'$  be the dual of  $\underline{\mathbb{V}}$ . Let  $\underline{\mathcal{S}} : \underline{\mathbb{V}}' \rightarrow \underline{\mathbb{V}}$  be such that  $\underline{\mathcal{S}} \underline{y}$  is the unique solution to the Helmholtz–Stokes problem

$$\int_{\Omega} \underline{\mathcal{S}} \underline{y} \cdot w d\underline{x} + \int_{\Omega} \underline{\nabla}_x (\underline{\mathcal{S}} \underline{y}) : \underline{\nabla}_x w d\underline{x} = \langle \underline{y}, w \rangle_V \quad \forall w \in \underline{\mathbb{V}}, \quad (5.3.2)$$

where  $\langle \cdot, \cdot \rangle_V$  denotes the duality pairing between  $\mathcal{V}'$  and  $\mathcal{V}$ . We note that

$$\langle v, \mathcal{S} v \rangle_V = \|\mathcal{S} v\|_{\mathbb{H}^1(\Omega)}^2 \quad \forall v \in \mathcal{V}' \supset (\mathbb{H}_0^1(\Omega))', \quad (5.3.3)$$

and  $\|\mathcal{S} \cdot\|_{\mathbb{H}^1(\Omega)}$  is a norm on  $\mathcal{V}'$ .

For later purposes, we recall the following well-known Gagliardo–Nirenberg inequality. Let  $r \in [2, \infty)$  if  $d = 2$ , and  $r \in [2, 6]$  if  $d = 3$  and  $\theta = d \left(\frac{1}{2} - \frac{1}{r}\right)$ . Then, there is a constant  $C$ , depending only on  $\Omega$ ,  $r$  and  $d$ , such that the following inequality holds for all  $\eta \in \mathbb{H}^1(\Omega)$ :

$$\|\eta\|_{L^r(\Omega)} \leq C \|\eta\|_{L^2(\Omega)}^{1-\theta} \|\eta\|_{\mathbb{H}^1(\Omega)}^\theta. \quad (5.3.4)$$

Our aim here is to prove existence of a (global-in-time) solution of a weak formulation of the problem  $(P_{\varepsilon,L})$  for any fixed parameters  $\varepsilon \in (0, 1]$  and  $L > 1$  under the following assumptions on the data:

$$\begin{aligned} \partial\Omega \in C^{0,1}, \quad u_0 \in \mathbb{H}, \quad \hat{\psi}_0 := M^{-1} \psi_0 \in L_M^2(\Omega \times D) \text{ with } \hat{\psi}_0 \geq 0 \text{ a.e. in } \Omega \times D, \\ \gamma > 1 \text{ in (5.1.3a,b),} \quad \text{and} \quad \underset{\sim}{f} \in L^2(0, T; \mathcal{V}'). \end{aligned} \quad (5.3.5)$$

Here  $L_M^2(\Omega \times D)$  is the Maxwellian-weighted  $L^2$  space over  $\Omega \times D$  with norm

$$\|\hat{\varphi}\|_{L_M^2(\Omega \times D)} := \left\{ \int_{\Omega \times D} M |\hat{\varphi}|^2 d\underset{\sim}{q} d\underset{\sim}{x} \right\}^{\frac{1}{2}}.$$

Similarly, we introduce  $L_M^2(D)$ , the Maxwellian-weighted  $L^2$  space over  $D$ .

On introducing

$$\|\hat{\varphi}\|_{\mathbb{H}_M^1(\Omega \times D)} := \left\{ \int_{\Omega \times D} M \left[ |\hat{\varphi}|^2 + |\nabla_x \hat{\varphi}|^2 + |\nabla_q \hat{\varphi}|^2 \right] d\underset{\sim}{q} d\underset{\sim}{x} \right\}^{\frac{1}{2}}, \quad (5.3.6)$$

we then set

$$\hat{X} \equiv \mathbb{H}_M^1(\Omega \times D) := \left\{ \hat{\varphi} \in L_{\text{loc}}^1(\Omega \times D) : \|\hat{\varphi}\|_{\mathbb{H}_M^1(\Omega \times D)} < \infty \right\}. \quad (5.3.7)$$

It follows that

$$C^\infty(\overline{\Omega \times D}) \text{ is dense in } \hat{X}. \quad (5.3.8)$$

This can be shown, for example, by a simple adaptation of Lemma 3.1 in Barrett, Schwab, and Süli [10], which appeals to fundamental results on weighted Sobolev spaces in Triebel [120] and Kufner [77]. We have from Sobolev embedding that

$$L^s(\Omega; L_M^2(D)) \hookrightarrow \mathbb{H}^1(\Omega; L_M^2(D)), \quad (5.3.9)$$

where  $s \in [1, \infty)$  if  $d = 2$  or  $s \in [1, 6]$  if  $d = 3$ . Similarly to (5.3.4) we have, with  $r$  and  $\theta$  as defined there, that there exists a constant  $C$ , depending only on  $\Omega$ ,  $r$  and  $d$ , such that

$$\|\hat{\varphi}\|_{L^r(\Omega; L_M^2(D))} \leq C \|\hat{\varphi}\|_{L^2(\Omega; L_M^2(D))}^{1-\theta} \|\hat{\varphi}\|_{\mathbb{H}^1(\Omega; L_M^2(D))}^\theta \quad \forall \hat{\varphi} \in \mathbb{H}^1(\Omega; L_M^2(D)). \quad (5.3.10)$$

In addition, we note that the embeddings

$$L_M^2(D) \hookrightarrow H_M^1(D), \quad (5.3.11a)$$

$$L_M^2(\Omega \times D) \equiv L^2(\Omega; L_M^2(D)) \hookrightarrow H_M^1(\Omega \times D) \equiv L^2(\Omega; H_M^1(D)) \cap H^1(\Omega; L_M^2(D)) \quad (5.3.11b)$$

are compact if  $\gamma \geq 1$  in (5.1.3a,b); see 5.4.

Let  $\hat{X}'$  be the dual space of  $\hat{X}$  with  $L_M^2(\Omega \times D)$  being the pivot space. Then, similarly to (5.3.2), let  $\mathcal{G} : \hat{X}' \rightarrow \hat{X}$  be such that  $\mathcal{G} \hat{\eta}$  is the unique solution of

$$\begin{aligned} \int_{\Omega \times D} M \left[ (\mathcal{G} \hat{\eta}) \hat{\varphi} + \nabla_q (\mathcal{G} \hat{\eta}) \cdot \nabla_q \hat{\varphi} + \nabla_x (\mathcal{G} \hat{\eta}) \cdot \nabla_x \hat{\varphi} \right] dq dx \\ = \langle M \hat{\eta}, \hat{\varphi} \rangle_{\hat{X}} \quad \forall \hat{\varphi} \in \hat{X}, \end{aligned} \quad (5.3.12)$$

where  $\langle M \cdot, \cdot \rangle_{\hat{X}}$  denotes the duality pairing between  $\hat{X}'$  and  $\hat{X}$ . Then, similarly to (5.3.3), we have that

$$\langle M \hat{\eta}, \mathcal{G} \hat{\eta} \rangle_{\hat{X}} = \|\mathcal{G} \hat{\eta}\|_{\hat{X}}^2 \quad \forall \hat{\eta} \in \hat{X}', \quad (5.3.13)$$

and  $\|\mathcal{G} \cdot\|_{\hat{X}}$  is a norm on  $\hat{X}'$ .

We recall the following compactness result, see, e.g., Temam [119] and Simon [115]. Let  $\mathcal{B}_0$ ,  $\mathcal{B}$  and  $\mathcal{B}_1$  be Banach spaces,  $\mathcal{B}_i$ ,  $i = 0, 1$ , reflexive, with a compact embedding  $\mathcal{B}_0 \hookrightarrow \mathcal{B}$  and a continuous embedding  $\mathcal{B} \hookrightarrow \mathcal{B}_1$ . Then, for  $\alpha_i > 1$ ,  $i = 0, 1$ , the embedding

$$\left\{ \eta \in L^{\alpha_0}(0, T; \mathcal{B}_0) : \frac{\partial \eta}{\partial t} \in L^{\alpha_1}(0, T; \mathcal{B}_1) \right\} \hookrightarrow L^{\alpha_0}(0, T; \mathcal{B}) \quad (5.3.14)$$

is compact. We note here that the proof of Theorem 2.3.2 developed in [71] for the Fokker–Planck equation does not rely on the compact embedding of  $H_0^1(D; M)$  into  $L^2(D)$ . However that argument does not work for the coupled Fokker–Planck–Navier–Stokes system considered in this chapter since due to the presence of nonlinearities; thus we shall have to rely on the compact embedding of  $H_M^1(D)$  into  $L_M^2(D)$  and the compact embedding (5.3.14) with suitable choices of  $\mathcal{B}_0$ ,  $\mathcal{B}$ ,  $\mathcal{B}_1$ ,  $\alpha_0$  and  $\alpha_1$ .

Throughout we will assume that (5.3.5) hold, so that (5.1.5) and (5.3.11a,b) hold. We note for future reference that (5.2.3a) and (5.1.5) yield that, for  $\hat{\varphi} \in L_M^2(\Omega \times D)$ ,

$$\begin{aligned} \int_{\Omega} |C(M \hat{\varphi})|_{\approx}^2 dx &= \int_{\Omega} \sum_{i=1}^d \sum_{j=1}^d \left( \int_D M \hat{\varphi} U' q_i q_j dq \right)_{\approx}^2 dx \\ &\leq d \left( \int_D M (U')^2 |q|_{\approx}^4 dq \right) \left( \int_{\Omega \times D} M |\hat{\varphi}|_{\approx}^2 dq dx \right) \\ &\leq C \left( \int_{\Omega \times D} M |\hat{\varphi}|_{\approx}^2 dq dx \right), \end{aligned} \quad (5.3.15)$$

where  $C = C(d)$  is a positive constant.

In order to prove existence of weak solutions to  $(P_{\varepsilon, L})$ , we require a further regularization. Let  $\mathcal{F} \in C(\mathbb{R}_{>0})$  be defined by

$$\mathcal{F}(s) := s(\ln s - 1) + 1, \quad s > 0. \quad (5.3.16)$$

As  $\lim_{s \rightarrow 0^+} \mathcal{F}(s) = 1$ , the function  $\mathcal{F}$  can be considered to be defined and continuous on  $[0, \infty)$ , where it is a nonnegative, strictly convex function with  $\mathcal{F}(1) = 0$ .

We then introduce the following convex regularization  $\mathcal{F}_\delta^L \in C^{2,1}(\mathbb{R})$  of  $\mathcal{F}$  defined, for any  $\delta \in (0, 1)$  and  $L > 1$ , by

$$\mathcal{F}_\delta^L(s) := \begin{cases} \frac{s^2 - \delta^2}{2\delta} + s(\ln \delta - 1) + 1 & \text{for } s \leq \delta, \\ \mathcal{F}(s) \equiv s(\ln s - 1) + 1 & \text{for } \delta \leq s \leq L, \\ \frac{s^2 - L^2}{2L} + s(\ln L - 1) + 1 & \text{for } L \leq s. \end{cases} \quad (5.3.17)$$

Hence,

$$[\mathcal{F}_\delta^L]'(s) = \begin{cases} \frac{s}{\delta} + \ln \delta - 1 & \text{for } s \leq \delta, \\ \ln s & \text{for } \delta \leq s \leq L, \\ \frac{s}{L} + \ln L - 1 & \text{for } L \leq s, \end{cases} \quad (5.3.18a)$$

$$[\mathcal{F}_\delta^L]''(s) = \begin{cases} \delta^{-1} & \text{for } s \leq \delta, \\ s^{-1} & \text{for } \delta \leq s \leq L, \\ L^{-1} & \text{for } L \leq s. \end{cases} \quad (5.3.18b)$$

We note that

$$\mathcal{F}_\delta^L(s) \geq \begin{cases} \frac{s^2}{2\delta} & \text{for } s \leq 0, \\ \frac{s^2}{4L} - C(L) & \text{for } s \geq 0; \end{cases} \quad (5.3.19)$$

and that  $[\mathcal{F}_\delta^L]''(s)$  is bounded below by  $1/L$  for all  $s \in \mathbb{R}$ . Finally, we set

$$\beta_\delta^L(s) := ([\mathcal{F}_\delta^L]'')^{-1}(s) = \max\{\beta^L(s), \delta\}, \quad (5.3.20)$$

and observe that  $\beta_\delta^L(s)$  is bounded above by  $L$  for all  $s \in \mathbb{R}$ .

### 5.3.1 Existence for $(\mathbf{P}_{\varepsilon, L, \delta})$

$(\mathbf{P}_{\varepsilon, L, \delta})$ , with solution  $\{u_{\varepsilon, L, \delta}, \psi_{\varepsilon, L, \delta}\}$ , will denote problem  $(\mathbf{P}_{\varepsilon, L})$ , where  $\beta^L(\cdot)$  in (5.2.4) and (5.2.5a) is replaced by  $\beta_\delta^L(\cdot)$ ; recall (5.1.10) and (5.3.20). In this section we will prove existence of a solution to the following weak formulation of  $(\mathbf{P}_{\varepsilon, L, \delta})$  for given parameters  $\varepsilon, \delta \in (0, 1]$  and  $L > 1$  with  $\hat{\psi}_{\varepsilon, L, \delta} = \psi_{\varepsilon, L, \delta}/M$ :

$(\mathbf{P}_{\varepsilon, L, \delta})$  Find  $u_{\varepsilon, L, \delta} \in L^\infty(0, T; \mathbb{L}^2(\Omega)) \cap L^2(0, T; \mathbb{Y}) \cap W^{1, \frac{4}{d}}(0, T; \mathbb{Y}')$  as well as  $\hat{\psi}_{\varepsilon, L, \delta} \in L^\infty(0, T; \mathbb{L}_M^2(\Omega \times D)) \cap L^2(0, T; \hat{\mathbb{X}}) \cap W^{1, \frac{4}{d}}(0, T; \hat{\mathbb{X}}')$ , with  $\underline{\mathcal{C}}(M \hat{\psi}_{\varepsilon, L, \delta}) \in L^\infty(0, T; \underline{\mathbb{L}}^2(\Omega))$ ,

such that  $u_{\epsilon,L,\delta}(\cdot, 0) = u_0(\cdot)$ ,  $\hat{\psi}_{\epsilon,L,\delta}(\cdot, \cdot, 0) = \hat{\psi}_0(\cdot, \cdot)$  and

$$\begin{aligned} & \int_0^T \left\langle \frac{\partial u_{\epsilon,L,\delta}}{\partial t}, w \right\rangle_V dt \\ & \quad + \int_0^T \int_{\Omega} \left[ (u_{\epsilon,L,\delta} \cdot \nabla_x) u_{\epsilon,L,\delta} \right] \cdot w + \nu \nabla_x u_{\epsilon,L,\delta} : \nabla_x w \, dx dt \\ & = \int_0^T \langle f, w \rangle_V dt - k_B \mathcal{T} \int_0^T \int_{\Omega} C(M \hat{\psi}_{\epsilon,L,\delta}) : \nabla_x w \, dx dt \quad \forall w \in L^{\frac{4}{4-d}}(0, T; \mathbb{V}); \end{aligned} \quad (5.3.21a)$$

$$\begin{aligned} & \int_0^T \left\langle M \frac{\partial \hat{\psi}_{\epsilon,L,\delta}}{\partial t}, \hat{\varphi} \right\rangle_{\hat{X}} dt \\ & \quad + \int_0^T \int_{\Omega \times D} M \left[ \frac{1}{2\lambda} \nabla_q \hat{\psi}_{\epsilon,L,\delta} - [\sigma(u_{\epsilon,L,\delta}) q] \beta_\delta^L(\hat{\psi}_{\epsilon,L,\delta}) \right] \cdot \nabla_q \hat{\varphi} \, dq \, dx dt \\ & \quad + \int_0^T \int_{\Omega \times D} M \left[ \varepsilon \nabla_x \hat{\psi}_{\epsilon,L,\delta} - u_{\epsilon,L,\delta} \hat{\psi}_{\epsilon,L,\delta} \right] \cdot \nabla_x \hat{\varphi} \, dq \, dx dt = 0 \quad \forall \hat{\varphi} \in L^{\frac{4}{4-d}}(0, T; \hat{X}). \end{aligned} \quad (5.3.21b)$$

**Remark 5.3.1** If  $d = 2$ , then  $u_{\epsilon,L,\delta} \in C([0, T]; \mathbb{H})$  (cf. Lemma 1.2 on p. 176 of Temam [119]), whereas if  $d = 3$ , then  $u_{\epsilon,L,\delta}$  is weakly continuous only as a mapping from  $[0, T]$  into  $\mathbb{H}$  (similarly as in Theorem 3.1 on p. 191 in Temam [119]). It is in the latter, weaker sense that the imposition of the initial condition to the  $u_{\epsilon,L,\delta}$ -equation will be understood for  $d = 2, 3$ : that is,  $\lim_{t \rightarrow 0^+} \int_{\Omega} (u_{\epsilon,L,\delta}(x, t) - u_0(x)) \cdot v(x) \, dx = 0$  for all  $v \in \mathbb{H}$ . Similarly, for the initial conditions of the  $\hat{\psi}_{\epsilon,L,\delta}$ -equation for  $d = 2, 3$ :  $\lim_{t \rightarrow 0^+} \int_{\Omega \times D} M (\hat{\psi}_{\epsilon,L,\delta}(x, q, t) - \hat{\psi}_0(x, q)) \hat{\varphi}(x, q) \, dq \, dx = 0$  for all  $\hat{\varphi} \in L_M^2(\Omega \times D)$ .  $\diamond$

**Remark 5.3.2** We note that the change of variable  $\hat{\psi} := \psi/M$  considered here differs from the change of variable  $\hat{\psi} := \psi/\sqrt{M}$  considered in the previous three chapters. One can, however, easily adapt our earlier analysis to this alternative change of variable.  $\diamond$

In order to prove existence of a weak solution to  $(P_{\epsilon,L,\delta})$ , we discretize in time; and so for any  $T > 0$ , let  $N \Delta t = T$  and  $t^n = n \Delta t$ ,  $n = 0 \rightarrow N$ . To prove existence of weak solutions under minimal smoothness requirements on the initial data, recall (5.3.5), we introduce  $u^0 \in \mathbb{V}$  such that

$$\int_{\Omega} \left[ u^0 \cdot v + \Delta t \nabla_x u^0 : \nabla_x v \right] dx = \int_{\Omega} u_0 \cdot v \, dx \quad \forall v \in \mathbb{V}; \quad (5.3.22)$$

and so

$$\int_{\Omega} [ |u^0|^2 + \Delta t |\nabla_x u^0|^2 ] dx \leq \int_{\Omega} |u_0|^2 dx \leq C. \quad (5.3.23)$$

In addition, we have that  $u^0$  converges to  $u_0$  weakly in  $\mathbb{H}$  in the limit of  $\Delta t \rightarrow 0_+$ .

Let  $u_{\epsilon,L,\delta}^0 = u^0$  and  $\hat{\psi}_{\epsilon,L,\delta}^0 = \hat{\psi}_0$ . Then, for  $n = 1 \rightarrow N$ , given  $\{u_{\epsilon,L,\delta}^{n-1}, \hat{\psi}_{\epsilon,L,\delta}^{n-1}\} \in \mathbb{V} \times L_M^2(\Omega \times D)$ , find  $\{u_{\epsilon,L,\delta}^n, \hat{\psi}_{\epsilon,L,\delta}^n\} \in \mathbb{V} \times \hat{X}$  such that

$$\begin{aligned} & \int_{\Omega} \left[ \frac{u_{\epsilon,L,\delta}^n - u_{\epsilon,L,\delta}^{n-1}}{\Delta t} + (u_{\epsilon,L,\delta}^{n-1} \cdot \nabla_x) u_{\epsilon,L,\delta}^n \right] \cdot w \, dx + \nu \int_{\Omega} \nabla_x u_{\epsilon,L,\delta}^n : \nabla_x w \, dx \\ & = \int_{\Omega} f^n \cdot w \, dx - k_B \mathcal{T} \int_{\Omega} C(M \hat{\psi}_{\epsilon,L,\delta}^n) : \nabla_x w \, dx \quad \forall w \in \mathbb{V}, \end{aligned} \quad (5.3.24a)$$

$$\begin{aligned}
& \int_{\Omega \times D} M \frac{\hat{\psi}_{\epsilon, L, \delta}^n - \hat{\psi}_{\epsilon, L, \delta}^{n-1}}{\Delta t} \hat{\varphi} \, dq \, dx \\
& + \int_{\Omega \times D} M \left[ \frac{1}{2\lambda} \nabla_q \hat{\psi}_{\epsilon, L, \delta}^n - [\sigma(u_{\epsilon, L, \delta}^n) q] \beta_\delta^L(\hat{\psi}_{\epsilon, L, \delta}^n) \right] \cdot \nabla_q \hat{\varphi} \, dq \, dx \\
& + \int_{\Omega \times D} M \left[ \varepsilon \nabla_x \hat{\psi}_{\epsilon, L, \delta}^n - u_{\epsilon, L, \delta}^{n-1} \hat{\psi}_{\epsilon, L, \delta}^n \right] \cdot \nabla_x \hat{\varphi} \, dq \, dx = 0 \quad \forall \hat{\varphi} \in \hat{X}; \quad (5.3.24b)
\end{aligned}$$

where

$$f^n(\cdot) := \frac{1}{\Delta t} \int_{t^{n-1}}^{t^n} f(\cdot, t) \, dt \in V'. \quad (5.3.25)$$

Now, letting  $f^{\Delta t, +}(\cdot, t) := f^n(\cdot)$  for  $t \in (t^{n-1}, t^n]$ ,  $n = 1 \rightarrow N$ , (5.3.5) and (5.3.25) imply that

$$f^{\Delta t, +} \rightarrow f \quad \text{strongly in } L^{\frac{4}{d}}(0, T; V') \text{ as } \Delta t \rightarrow 0_+, \quad (5.3.26)$$

It is convenient to rewrite (5.3.24a) as

$$b(u_{\epsilon, L, \delta}^n, w) = \ell_b(\hat{\psi}_{\epsilon, L, \delta}^n)(w) \quad \forall w \in Y; \quad (5.3.27)$$

where for all  $w_i \in \mathbb{H}_0^1(\Omega)$ ,  $i = 1, 2$ ,

$$b(w_1, w_2) := \int_{\Omega} [w_1 + \Delta t (u_{\epsilon, L, \delta}^{n-1} \cdot \nabla_x) w_1] \cdot w_2 \, dx + \Delta t \nu \int_{\Omega} \nabla_x w_1 : \nabla_x w_2 \, dx, \quad (5.3.28a)$$

and for all  $w \in \mathbb{H}_0^1(\Omega)$  and  $\hat{\varphi} \in L_M^2(\Omega \times D)$

$$\ell_b(\hat{\varphi})(w) := \Delta t \langle f^n, w \rangle_V + \int_{\Omega} [u_{\epsilon, L, \delta}^{n-1} \cdot w - \Delta t k_B T C(M \hat{\varphi}) : \nabla_x w] \, dx. \quad (5.3.28b)$$

We note that

$$\begin{aligned}
& \int_{\Omega} [(v \cdot \nabla_x) w_1] \cdot w_2 \, dx \\
& = - \int_{\Omega} [(v \cdot \nabla_x) w_2] \cdot w_1 \, dx \quad \forall v \in Y, \quad \forall w_1, w_2 \in \mathbb{H}_0^1(\Omega), \quad (5.3.29)
\end{aligned}$$

and hence  $b(\cdot, \cdot)$  is a continuous nonsymmetric coercive bilinear functional on  $\mathbb{H}_0^1(\Omega) \times \mathbb{H}_0^1(\Omega)$ . In addition,  $\ell_b(\hat{\varphi})(\cdot)$  is a continuous linear functional on  $Y$  for any  $\hat{\varphi} \in L_M^2(\Omega \times D)$ .

For  $r > d$ , let

$$Y^r := \left\{ v \in L^r(\Omega) : \int_{\Omega} v \cdot \nabla_x w \, dx = 0 \quad \forall w \in W^{1, \frac{r}{r-1}}(\Omega) \right\}. \quad (5.3.30)$$

It is also convenient to rewrite (5.3.24b) as

$$a(\hat{\psi}_{\epsilon, L, \delta}^n, \hat{\varphi}) = \ell_a(u_{\epsilon, L, \delta}^n, \hat{\psi}_{\epsilon, L, \delta}^n)(\hat{\varphi}) \quad \forall \hat{\varphi} \in \hat{X}, \quad (5.3.31)$$

where, for all  $\hat{\varphi}_1, \hat{\varphi}_2 \in \hat{X}$ ,

$$a(\hat{\varphi}_1, \hat{\varphi}_2) := \int_{\Omega \times D} M \left( \hat{\varphi}_1 \hat{\varphi}_2 + \Delta t \left[ \varepsilon \nabla_x \hat{\varphi}_1 - u_{\varepsilon, L, \delta}^{n-1} \hat{\varphi}_1 \right] \cdot \nabla_x \hat{\varphi}_2 + \frac{\Delta t}{2\lambda} \nabla_q \hat{\varphi}_1 \cdot \nabla_q \hat{\varphi}_2 \right) dq dx, \quad (5.3.32a)$$

and, for all  $v \in \mathbb{H}^1(\Omega)$ ,  $\hat{\eta} \in L_M^2(\Omega \times D)$  and  $\hat{\varphi} \in \hat{X}$ ,

$$\ell_a(v, \hat{\eta})(\hat{\varphi}) := \int_{\Omega \times D} M \left[ \hat{\psi}_{\varepsilon, L, \delta}^{n-1} \hat{\varphi} + \Delta t [\sigma(v) q] \beta_\delta^L(\hat{\eta}) \cdot \nabla_q \hat{\varphi} \right] dq dx, \quad (5.3.32b)$$

It follows from (5.3.30) and (5.3.9) that for  $r > d$

$$\int_{\Omega \times D} M v \hat{\varphi} \cdot \nabla_x \hat{\varphi} dq dx = 0 \quad \forall v \in \mathbb{Y}^r, \quad \forall \hat{\varphi} \in \hat{X}; \quad (5.3.33)$$

and hence that  $a(\cdot, \cdot)$  is a continuous nonsymmetric coercive bilinear functional on  $\hat{X} \times \hat{X}$ . In addition,  $\ell_a(v, \hat{\eta})(\cdot)$  is a linear functional on  $\hat{X}$  for all  $v \in \mathbb{H}^1(\Omega)$  and  $\hat{\eta} \in L_M^2(\Omega \times D)$ .

In order to prove existence of a solution to (5.3.24a,b), we consider a fixed-point argument. Given  $\hat{\psi} \in L_M^2(\Omega \times D)$  let  $\{u^*, \hat{\psi}^*\} \in \mathbb{Y} \times \hat{X}$  be such that

$$b(u^*, w) = \ell_b(\hat{\psi})(w) \quad \forall w \in \mathbb{V}, \quad (5.3.34a)$$

$$a(\hat{\psi}^*, \hat{\varphi}) = \ell_a(u^*, \hat{\psi})(\hat{\varphi}) \quad \forall \hat{\varphi} \in \hat{X}. \quad (5.3.34b)$$

The Lax–Milgram theorem yields the existence of a unique solution to (5.3.34a,b), and so the overall procedure (5.3.34a,b) is well defined.

**Lemma 5.3.3** *Let  $G : L_M^2(\Omega \times D) \rightarrow \hat{X} \subset L_M^2(\Omega \times D)$  denote the nonlinear map that takes  $\hat{\psi}$  to  $\hat{\psi}^* = G(\hat{\psi})$  via the procedure (5.3.34a,b). Then  $G$  has a fixed point. Hence there exists a solution  $\{u_{\varepsilon, L, \delta}^n, \hat{\psi}_{\varepsilon, L, \delta}^n\} \in \mathbb{Y} \times \hat{X}$  to (5.3.24a,b).*

**Proof.** Clearly, a fixed point of  $G$  yields a solution of (5.3.24a,b). In order to show that  $G$  has a fixed point, we apply Schauder’s fixed-point theorem; that is, we need to show that (i)  $G : L_M^2(\Omega \times D) \rightarrow L_M^2(\Omega \times D)$  is continuous, that (ii) it is compact, and that (iii) there exists a  $C_\star \in \mathbb{R}_{>0}$  such that

$$\|\hat{\psi}\|_{L_M^2(\Omega \times D)} \leq C_\star \quad (5.3.35)$$

for every  $\hat{\psi} \in L_M^2(\Omega \times D)$  and  $\kappa \in (0, 1]$  satisfying  $\hat{\psi} = \kappa G(\hat{\psi})$ .

Let  $\{\hat{\psi}^{(i)}\}_{i \geq 0}$  be such that

$$\hat{\psi}^{(i)} \rightarrow \hat{\psi} \quad \text{strongly in } L_M^2(\Omega \times D) \quad \text{as } i \rightarrow \infty. \quad (5.3.36)$$

It follows immediately from (5.3.20) and (5.3.15) that

$$M^{\frac{1}{2}} \beta_\delta^L(\hat{\psi}^{(i)}) \rightarrow M^{\frac{1}{2}} \beta_\delta^L(\hat{\psi}) \quad \text{strongly in } L^\infty(\Omega \times D) \quad \text{as } i \rightarrow \infty, \quad (5.3.37a)$$

$$\underset{\approx}{C}(M \hat{\psi}^{(i)}) \rightarrow \underset{\approx}{C}(M \hat{\psi}) \quad \text{strongly in } L^2(\Omega) \quad \text{as } i \rightarrow \infty. \quad (5.3.37b)$$

We need to show that

$$\hat{\eta}^{(i)} := G(\hat{\psi}^{(i)}) \rightarrow G(\hat{\psi}) \quad \text{strongly in } L_M^2(\Omega \times D) \quad \text{as } i \rightarrow \infty, \quad (5.3.38)$$

in order to prove (i) above. We have from the definition of  $G$ , see (5.3.34a,b), that, for all  $i \geq 0$ ,

$$a(\hat{\eta}^{(i)}, \hat{\varphi}) = \ell_a(v^{(i)}, \hat{\psi}^{(i)})(\hat{\varphi}) \quad \forall \hat{\varphi} \in \hat{X}, \quad (5.3.39a)$$

where  $v^{(i)} \in \mathcal{V}$  satisfies

$$b(v^{(i)}, w) = \ell_b(\hat{\psi}^{(i)})(w) \quad \forall w \in \mathcal{V}. \quad (5.3.39b)$$

Choosing  $\hat{\varphi} = \hat{\eta}^{(i)}$  in (5.3.39a) yields, on noting the simple identity

$$2(s_1 - s_2)s_1 = s_1^2 + (s_1 - s_2)^2 - s_2^2 \quad \forall s_1, s_2 \in \mathbb{R}, \quad (5.3.40)$$

(5.3.33) and (5.3.20) that, for all  $i \geq 0$ ,

$$\begin{aligned} & \int_{\Omega \times D} M \left[ |\hat{\eta}^{(i)}|^2 + |\hat{\eta}^{(i)} - \hat{\psi}_{\epsilon, L, \delta}^{n-1}|^2 + \frac{\Delta t}{2\lambda} |\nabla_q \hat{\eta}^{(i)}|^2 + 2\varepsilon \Delta t |\nabla_x \eta^{(i)}|^2 \right] dq dx \\ & \leq \int_{\Omega \times D} M |\hat{\psi}_{\epsilon, L, \delta}^{n-1}|^2 dq dx + C(L, \lambda) \Delta t \int_{\Omega} |\nabla_x v^{(i)}|^2 dx. \end{aligned} \quad (5.3.41)$$

Choosing  $w \equiv v^{(i)}$  in (5.3.39b), and noting (5.3.40), (5.3.29), (5.3.15), (5.3.2), a Poincaré inequality and (5.3.36) yields, for all  $i \geq 0$ , that

$$\begin{aligned} & \int_{\Omega} \left[ |v^{(i)}|^2 + |v^{(i)} - u_{\epsilon, L, \delta}^{n-1}|^2 \right] dx + \Delta t \nu \int_{\Omega} |\nabla_x v^{(i)}|^2 dx \\ & \leq \int_{\Omega} |u_{\epsilon, L, \delta}^{n-1}|^2 dx + C \Delta t \|S f^n\|_{H^1(\Omega)}^2 + C \Delta t \int_{\Omega \times D} M |\hat{\psi}^{(i)}|^2 dq dx \leq C. \end{aligned} \quad (5.3.42)$$

Combining (5.3.41) and (5.3.42), we have for all  $i \geq 0$  that

$$\|\hat{\eta}^{(i)}\|_{\hat{X}} + \|v^{(i)}\|_{H^1(\Omega)} \leq C(L, (\Delta t)^{-1}). \quad (5.3.43)$$

It follows from (5.3.43), (5.3.9) and the compactness of the embedding (5.3.11b) that there exists a subsequence  $\{\hat{\eta}^{(i_k)}, v^{(i_k)}\}_{i_k \geq 0}$  and functions  $\hat{\eta} \in \hat{X}$  and  $v \in \mathcal{V}$  such that, as  $i_k \rightarrow \infty$ ,

$$\hat{\eta}^{(i_k)} \rightarrow \hat{\eta} \quad \text{weakly in } L^s(\Omega; L_M^2(D)), \quad (5.3.44a)$$

$$M^{\frac{1}{2}} \nabla_x \hat{\eta}^{(i_k)} \rightarrow M^{\frac{1}{2}} \nabla_x \hat{\eta} \quad \text{weakly in } L^2(\Omega \times D), \quad (5.3.44b)$$

$$M^{\frac{1}{2}} \nabla_q \hat{\eta}^{(i_k)} \rightarrow M^{\frac{1}{2}} \nabla_q \hat{\eta} \quad \text{weakly in } L^2(\Omega \times D), \quad (5.3.44c)$$

$$\hat{\eta}^{(i_k)} \rightarrow \hat{\eta} \quad \text{strongly in } L_M^2(\Omega \times D), \quad (5.3.44d)$$

$$v^{(i_k)} \rightarrow v \quad \text{weakly in } H^1(\Omega); \quad (5.3.44e)$$

where  $s \in [1, \infty)$  if  $d = 2$  or  $s \in [1, 6]$  if  $d = 3$ . It follows from (5.3.39b), (5.3.28a,b), (5.3.44e) and (5.3.37b) that  $\underline{v} \in \underline{\mathbb{V}}$  and  $\hat{\psi} \in \hat{\mathbb{X}}$  satisfy

$$b(\underline{v}, \underline{w}) = \ell_b(\hat{\psi})(\underline{w}) \quad \forall \underline{w} \in \underline{\mathbb{V}}. \quad (5.3.45)$$

It follows from (5.3.39a), (5.3.32a,b), (5.3.44a–e) and (5.3.37a) that  $\hat{\eta}, \hat{\psi} \in \hat{\mathbb{X}}$  and  $\underline{v} \in \underline{\mathbb{V}}$ , satisfy

$$a(\hat{\eta}, \hat{\varphi}) = \ell_a(\underline{v}, \hat{\psi})(\hat{\varphi}) \quad \forall \hat{\varphi} \in \hat{\mathbb{X}}. \quad (5.3.46)$$

Combining (5.3.46) and (5.3.45), we have that  $\hat{\eta} = G(\hat{\psi}) \in \hat{\mathbb{X}}$ . Therefore the whole sequence  $\hat{\eta}^{(i)} \equiv G(\hat{\psi}^{(i)}) \rightarrow G(\hat{\psi})$  strongly in  $L_M^2(\Omega \times D)$  as  $i \rightarrow \infty$ , and so (i) holds.

As the embedding  $\hat{\mathbb{X}} \hookrightarrow L_M^2(\Omega \times D)$  is compact, it follows that (ii) holds.

As regards (iii),  $\hat{\psi} = \kappa G(\hat{\psi})$  implies that  $\{\underline{v}, \hat{\psi}\} \in \underline{\mathbb{V}} \times \hat{\mathbb{X}}$  satisfies

$$b(\underline{v}, \underline{w}) = \ell_b(\hat{\psi})(\underline{w}) \quad \forall \underline{w} \in \underline{\mathbb{V}}, \quad (5.3.47a)$$

$$a(\hat{\psi}, \hat{\varphi}) = \kappa \ell_a(\underline{v}, \hat{\psi})(\hat{\varphi}) \quad \forall \hat{\varphi} \in \hat{\mathbb{X}}. \quad (5.3.47b)$$

Choosing  $\underline{w} \equiv \hat{\psi}$  in (5.3.47a) yields, similarly to (5.3.42), that

$$\begin{aligned} & \frac{1}{2} \int_{\Omega} \left[ |\underline{v}|^2 + |\underline{v} - \underline{u}_{\epsilon, L, \delta}^{n-1}|^2 - |\underline{u}_{\epsilon, L, \delta}^{n-1}|^2 \right] d\underline{x} + \Delta t \nu \int_{\Omega} |\underline{\nabla}_x \underline{v}|^2 d\underline{x} \\ & = \Delta t \left[ \langle \underline{f}^n, \underline{v} \rangle_V - k_B \mathcal{T} \int_{\Omega} \underline{C}(M \hat{\psi}) : \underline{\nabla}_x \underline{v} d\underline{x} \right]. \end{aligned} \quad (5.3.48)$$

Choosing  $\hat{\varphi} = [\mathcal{F}_{\delta}^L]'(\hat{\psi})$  in (5.3.47b) and noting (5.3.18a), (5.3.20), (5.3.8), (5.1.4), (5.2.3a) and that  $\underline{v}$  is divergence-free yield

$$\begin{aligned} & \int_{\Omega \times D} M \left[ \mathcal{F}_{\delta}^L(\hat{\psi}) - \mathcal{F}_{\delta}^L(\kappa \hat{\psi}_{\epsilon, L, \delta}^{n-1}) \right] d\underline{q} d\underline{x} \\ & + \Delta t \int_{\Omega \times D} M \left[ \varepsilon \underline{\nabla}_x \hat{\psi} \cdot \underline{\nabla}_x ([\mathcal{F}_{\delta}^L]'(\hat{\psi})) + \frac{1}{2\lambda} \underline{\nabla}_q \hat{\psi} \cdot \underline{\nabla}_q ([\mathcal{F}_{\delta}^L]'(\hat{\psi})) \right] d\underline{q} d\underline{x} \\ & \leq \kappa \Delta t \int_{\Omega \times D} M \underline{\sigma}(\underline{v}) \underline{q} \cdot \underline{\nabla}_q \hat{\psi} d\underline{q} d\underline{x} \\ & = \kappa \Delta t \int_{\Omega} \underline{C}(M \hat{\psi}) : \underline{\sigma}(\underline{v}) d\underline{x}. \end{aligned} \quad (5.3.49)$$

Combining (5.3.48) and (5.3.49), and noting (5.3.2) and a Poincaré inequality yields that

$$\begin{aligned} & \frac{\kappa}{2} \int_{\Omega} \left[ |\underline{v}|^2 + |\underline{v} - \underline{u}_{\epsilon, L, \delta}^{n-1}|^2 \right] d\underline{x} + \kappa \Delta t \nu \int_{\Omega} |\underline{\nabla}_x \underline{v}|^2 d\underline{x} + k_B \mathcal{T} \int_{\Omega \times D} M \mathcal{F}_{\delta}^L(\hat{\psi}) d\underline{q} d\underline{x} \\ & + k_B \mathcal{T} \Delta t \int_{\Omega \times D} M \left[ \varepsilon \underline{\nabla}_x \hat{\psi} \cdot \underline{\nabla}_x ([\mathcal{F}_{\delta}^L]'(\hat{\psi})) + \frac{1}{2\lambda} \underline{\nabla}_q \hat{\psi} \cdot \underline{\nabla}_q ([\mathcal{F}_{\delta}^L]'(\hat{\psi})) \right] d\underline{q} d\underline{x} \\ & \leq \kappa \Delta t \langle \underline{f}^n, \underline{v} \rangle_V + \frac{\kappa}{2} \int_{\Omega} |\underline{u}_{\epsilon, L, \delta}^{n-1}|^2 d\underline{x} + k_B \mathcal{T} \int_{\Omega \times D} M \mathcal{F}_{\delta}^L(\kappa \hat{\psi}_{\epsilon, L, \delta}^{n-1}) d\underline{q} d\underline{x} \\ & \leq \frac{\kappa}{2} \Delta t \nu \int_{\Omega} |\underline{\nabla}_x \underline{v}|^2 d\underline{x} + \kappa \Delta t C(\nu^{-1}) \| \underline{S} \underline{f}^n \|_{H^1(\Omega)}^2 \\ & \quad + \frac{\kappa}{2} \int_{\Omega} |\underline{u}_{\epsilon, L, \delta}^{n-1}|^2 d\underline{x} + k_B \mathcal{T} \int_{\Omega \times D} M \mathcal{F}_{\delta}^L(\kappa \hat{\psi}_{\epsilon, L, \delta}^{n-1}) d\underline{q} d\underline{x}. \end{aligned} \quad (5.3.50)$$

It is easy to show that  $\mathcal{F}_\delta^L(s)$  is nonnegative for all  $s \in \mathbb{R}$ , with  $\mathcal{F}_\delta^L(1) = 0$ . Furthermore, for any  $\kappa \in (0, 1]$ ,

$$\begin{aligned} \mathcal{F}_\delta^L(\kappa s) &\leq \mathcal{F}_\delta^L(s) && \text{if } s < 0 \text{ or } 1 \leq \kappa s, \\ \mathcal{F}_\delta^L(\kappa s) &\leq \mathcal{F}_\delta^L(0) \leq 1 && \text{if } 0 \leq \kappa s \leq 1. \end{aligned}$$

Thus we deduce that

$$\mathcal{F}_\delta^L(\kappa s) \leq \mathcal{F}_\delta^L(s) + 1 \quad \forall s \in \mathbb{R}, \quad \forall \kappa \in (0, 1]. \quad (5.3.51)$$

Hence, the bounds (5.3.50) and (5.3.51), on noting (5.3.19) and (5.3.18b), which implies that  $[\mathcal{F}_\delta^L(s)]'' \geq L^{-1}$  for all  $s \in \mathbb{R}$ , give rise to the desired bound (5.3.35) with  $C_*$  dependent only on  $L$ ,  $k_B$ ,  $\mathcal{T}$  and  $\hat{\psi}_{\epsilon, L, \delta}^{n-1}$ . Hence (iii) holds, and so  $G$  has a fixed point. Thus we have proved existence of a solution to (5.3.24a,b).  $\square$

Choosing  $w \equiv u_{\epsilon, L, \delta}^n$  in (5.3.27) and  $\hat{\varphi} \equiv [\mathcal{F}_\delta^L]'(\hat{\psi}_{\epsilon, L, \delta}^n)$ , and combining, then yields, similarly to (5.3.50), that

$$\begin{aligned} &\frac{1}{2} \int_{\Omega} \left[ |u_{\epsilon, L, \delta}^n|^2 + |u_{\epsilon, L, \delta}^n - u_{\epsilon, L, \delta}^{n-1}|^2 \right] dx + k_B \mathcal{T} \int_{\Omega \times D} M \mathcal{F}_\delta^L(\hat{\psi}_{\epsilon, L, \delta}^n) dq dx \\ &+ \Delta t \left[ \frac{\nu}{2} \int_{\Omega} |\nabla_x u_{\epsilon, L, \delta}^n|^2 dx + k_B \mathcal{T} \varepsilon \int_{\Omega \times D} M \nabla_x \hat{\psi}_{\epsilon, L, \delta}^n \cdot \nabla_x ([\mathcal{F}_\delta^L]'(\hat{\psi}_{\epsilon, L, \delta}^n)) dq dx \right. \\ &\quad \left. + \frac{k_B \mathcal{T}}{2\lambda} \int_{\Omega \times D} M \nabla_q \hat{\psi}_{\epsilon, L, \delta}^n \cdot \nabla_q ([\mathcal{F}_\delta^L]'(\hat{\psi}_{\epsilon, L, \delta}^n)) dq dx \right] \\ &\leq \Delta t C(\nu^{-1}) \|S f^n\|_{\mathbb{H}^1(\Omega)}^2 + \frac{1}{2} \int_{\Omega} |u_{\epsilon, L, \delta}^{n-1}|^2 dx \\ &\quad + k_B \mathcal{T} \int_{\Omega \times D} M \mathcal{F}_\delta^L(\hat{\psi}_{\epsilon, L, \delta}^{n-1}) dq dx. \end{aligned} \quad (5.3.52)$$

Summing (5.3.52) from  $n = 1 \rightarrow m$ , with  $1 \leq m \leq N$ , yields that

$$\begin{aligned} &\frac{1}{2} \left[ \int_{\Omega} |u_{\epsilon, L, \delta}^m|^2 dx + \sum_{n=1}^m \int_{\Omega} |u_{\epsilon, L, \delta}^n - u_{\epsilon, L, \delta}^{n-1}|^2 dx \right] + k_B \mathcal{T} \int_{\Omega \times D} M \mathcal{F}_\delta^L(\hat{\psi}_{\epsilon, L, \delta}^m) dq dx \\ &+ \sum_{n=1}^m \Delta t \left[ \frac{\nu}{2} \int_{\Omega} |\nabla_x u_{\epsilon, L, \delta}^n|^2 dx + k_B \mathcal{T} \varepsilon \int_{\Omega \times D} M \nabla_x \hat{\psi}_{\epsilon, L, \delta}^n \cdot \nabla_x ([\mathcal{F}_\delta^L]'(\hat{\psi}_{\epsilon, L, \delta}^n)) dq dx \right. \\ &\quad \left. + \frac{k_B \mathcal{T}}{2\lambda} \int_{\Omega \times D} \nabla_q \hat{\psi}_{\epsilon, L, \delta}^n \cdot \nabla_q ([\mathcal{F}_\delta^L]'(\hat{\psi}_{\epsilon, L, \delta}^n)) dq dx \right] \\ &\leq \frac{1}{2} \int_{\Omega} |u^0|^2 dx + k_B \mathcal{T} \int_{\Omega \times D} M \mathcal{F}_\delta^L(\hat{\psi}_0) dq dx + C(\nu^{-1}) \sum_{n=1}^m \Delta t \|S f^n\|_{\mathbb{H}^1(\Omega)}^2 \\ &\leq \frac{1}{2} \int_{\Omega} |u^0|^2 dx + k_B \mathcal{T} \int_{\Omega \times D} M \mathcal{F}_\delta^L(\hat{\psi}_0) dq dx + C(\nu^{-1}) \int_0^{tm} \|S f\|_{\mathbb{H}^1(\Omega)}^2 dt \leq C; \end{aligned} \quad (5.3.53)$$

where  $C$  is independent of  $\delta$ ,  $L$  and  $\Delta t$ , on assuming that  $L$  is chosen so that

$$0 \leq \hat{\psi}_0 \leq L \quad \text{a.e. in } \Omega \times D. \quad (5.3.54)$$

Choosing  $\hat{\varphi} = \hat{\psi}_{\epsilon,L,\delta}^n$  in (5.3.31), and noting (5.3.40), (5.3.33), (5.3.20) and (5.1.5), yields that

$$\begin{aligned}
& \int_{\Omega \times D} M \left[ |\hat{\psi}_{\epsilon,L,\delta}^n|^2 + |\hat{\psi}_{\epsilon,L,\delta}^n - \hat{\psi}_{\epsilon,L,\delta}^{n-1}|^2 \right] dq dx \\
& + \Delta t \int_{\Omega \times D} M \left[ 2\varepsilon \left| \nabla_x \hat{\psi}_{\epsilon,L,\delta}^n \right|^2 + \frac{1}{\lambda} \left| \nabla_q \hat{\psi}_{\epsilon,L,\delta}^n \right|^2 \right] dq dx \\
& = \int_{\Omega \times D} M \left[ |\hat{\psi}_{\epsilon,L,\delta}^{n-1}|^2 + \Delta t \left[ \sigma(\underline{u}_{\epsilon,L,\delta}^n) q \right] \beta_\delta^L(\hat{\psi}_{\epsilon,L,\delta}^n) \cdot \nabla_q \hat{\psi}_{\epsilon,L,\delta}^n \right] dq dx \\
& \leq \int_{\Omega \times D} M |\hat{\psi}_{\epsilon,L,\delta}^{n-1}|^2 dq dx + \frac{\Delta t}{2\lambda} \int_{\Omega \times D} M |\nabla_q \hat{\psi}_{\epsilon,L,\delta}^n|^2 dq dx \\
& \quad + C(L, \lambda) \Delta t \int_{\Omega} |\nabla_x \underline{u}_{\epsilon,L,\delta}^n|^2 dx. \tag{5.3.55}
\end{aligned}$$

Summing (5.3.55) from  $n = 1 \rightarrow m$ , with  $1 \leq m \leq N$ , yields, on noting (5.3.53), that

$$\begin{aligned}
& \int_{\Omega \times D} M |\hat{\psi}_{\epsilon,L,\delta}^m|^2 dq dx + \sum_{n=1}^m \int_{\Omega \times D} M |\hat{\psi}_{\epsilon,L,\delta}^n - \hat{\psi}_{\epsilon,L,\delta}^{n-1}|^2 dq dx \\
& + \sum_{n=1}^m \Delta t \int_{\Omega \times D} M \left[ 2\varepsilon \left| \nabla_x \hat{\psi}_{\epsilon,L,\delta}^n \right|^2 + \frac{1}{2\lambda} \left| \nabla_q \hat{\psi}_{\epsilon,L,\delta}^n \right|^2 \right] dq dx \\
& \leq \int_{\Omega \times D} M |\hat{\psi}_0|^2 dq dx + C(L) \sum_{n=1}^m \Delta t \int_{\Omega} |\nabla_x \underline{u}_{\epsilon,L,\delta}^n|^2 dx \leq C(L). \tag{5.3.56}
\end{aligned}$$

Choosing  $w \equiv \underline{S} \left( \frac{\underline{u}_{\epsilon,L,\delta}^n - \underline{u}_{\epsilon,L,\delta}^{n-1}}{\Delta t} \right) \in \mathbb{V}$  in (5.3.27) yields, on noting (5.3.2), (5.3.3) and (5.3.29), that

$$\begin{aligned}
& \int_{\Omega} \left[ \left| \nabla_x \left[ \underline{S} \left( \frac{\underline{u}_{\epsilon,L,\delta}^n - \underline{u}_{\epsilon,L,\delta}^{n-1}}{\Delta t} \right) \right] \right|^2 + \left| \underline{S} \left( \frac{\underline{u}_{\epsilon,L,\delta}^n - \underline{u}_{\epsilon,L,\delta}^{n-1}}{\Delta t} \right) \right|^2 \right] dx \\
& = \int_{\Omega} \left[ -\nu \nabla_x \underline{u}_{\epsilon,L,\delta}^n \cdot \nabla_x \left[ \underline{S} \left( \frac{\underline{u}_{\epsilon,L,\delta}^n - \underline{u}_{\epsilon,L,\delta}^{n-1}}{\Delta t} \right) \right] - k_B T C(M \hat{\psi}_{\epsilon,L,\delta}^n) : \nabla_x \left[ \underline{S} \left( \frac{\underline{u}_{\epsilon,L,\delta}^n - \underline{u}_{\epsilon,L,\delta}^{n-1}}{\Delta t} \right) \right] \right] dx \\
& + \int_{\Omega} \underline{u}_{\epsilon,L,\delta}^n \cdot \left[ (\underline{u}_{\epsilon,L,\delta}^{n-1} \cdot \nabla_x) \left[ \underline{S} \left( \frac{\underline{u}_{\epsilon,L,\delta}^n - \underline{u}_{\epsilon,L,\delta}^{n-1}}{\Delta t} \right) \right] \right] dx + \left\langle f^n, \underline{S} \left( \frac{\underline{u}_{\epsilon,L,\delta}^n - \underline{u}_{\epsilon,L,\delta}^{n-1}}{\Delta t} \right) \right\rangle_V \\
& \leq C \left[ \| \underline{S} f^n \|_{\mathbb{H}^1(\Omega)}^2 + \int_{\Omega} \left[ |C(M \hat{\psi}_{\epsilon,L,\delta}^n)|^2 + |\nabla_x \underline{u}_{\epsilon,L,\delta}^n|^2 + |\underline{u}_{\epsilon,L,\delta}^{n-1}|^2 |\underline{u}_{\epsilon,L,\delta}^n|^2 \right] dx \right]. \tag{5.3.57}
\end{aligned}$$

Applying the Cauchy–Schwarz inequality, the algebraic-geometric mean inequality, (5.3.4),

and a Poincaré inequality yields that

$$\begin{aligned}
\int_{\Omega} |u_{\sim, \epsilon, L, \delta}^{n-1}|^2 |u_{\sim, \epsilon, L, \delta}^n|^2 dx &\leq \left( \int_{\Omega} |u_{\sim, \epsilon, L, \delta}^{n-1}|^4 dx \right)^{\frac{1}{2}} \left( \int_{\Omega} |u_{\sim, \epsilon, L, \delta}^n|^4 dx \right)^{\frac{1}{2}} \\
&\leq \frac{1}{2} \sum_{m=n-1}^n \int_{\Omega} |u_{\sim, \epsilon, L, \delta}^m|^4 dx \\
&\leq C \sum_{m=n-1}^n \left[ \left( \int_{\Omega} |u_{\sim, \epsilon, L, \delta}^m|^2 dx \right)^{2-\frac{d}{2}} \left( \int_{\Omega} |\nabla_x u_{\sim, \epsilon, L, \delta}^m|^2 dx \right)^{\frac{d}{2}} \right]. \quad (5.3.58)
\end{aligned}$$

Taking the  $\frac{2}{d}$  power of both sides of (5.3.57), summing from  $n = 1 \rightarrow N$ , and noting (5.3.58), (5.3.15), (5.3.56), (5.3.53) and (5.3.23) yields that

$$\begin{aligned}
&\sum_{n=1}^N \Delta t \left( \int_{\Omega} \left[ \left| \nabla_x \left[ S \left( \frac{u_{\sim, \epsilon, L, \delta}^n - u_{\sim, \epsilon, L, \delta}^{n-1}}{\Delta t} \right) \right] \right|^2 + \left| S \left( \frac{u_{\sim, \epsilon, L, \delta}^n - u_{\sim, \epsilon, L, \delta}^{n-1}}{\Delta t} \right) \right|^2 \right] dx \right)^{\frac{2}{d}} \\
&\leq C \left[ \sum_{n=1}^N \Delta t \left( \int_{\Omega} |C(M \hat{\psi}_{\sim, \epsilon, L, \delta}^n)|^2 dx \right)^{\frac{2}{d}} \right] + C(T) \left[ \sum_{n=1}^N \Delta t \int_{\Omega} |\nabla_x u_{\sim, \epsilon, L, \delta}^n|^2 dx \right]^{\frac{2}{d}} \\
&\quad + C(T) \left[ \max_{n=0 \rightarrow N} \left( \int_{\Omega} |u_{\sim, \epsilon, L, \delta}^n|^2 dx \right)^{\frac{4}{d}-1} \right] \left[ \sum_{n=0}^N \Delta t \int_{\Omega} |\nabla_x u_{\sim, \epsilon, L, \delta}^n|^2 dx \right] \\
&\quad + \sum_{n=1}^N \Delta t \|Sf_{\sim}^n\|_{\mathbb{H}^1(\Omega)}^{\frac{4}{d}} \\
&\leq C(L, T) + C \int_0^T \|Sf_{\sim}\|_{\mathbb{H}^1(\Omega)}^{\frac{4}{d}} dt \leq C(L, T). \quad (5.3.59)
\end{aligned}$$

Choosing  $\hat{\varphi} \equiv \mathcal{G} \left( \frac{\hat{\psi}_{\sim, \epsilon, L, \delta}^n - \hat{\psi}_{\sim, \epsilon, L, \delta}^{n-1}}{\Delta t} \right) \in \hat{X}$  in (5.3.31) yields, on noting (5.3.12), (5.3.13), (5.3.20) and (5.1.5), that

$$\begin{aligned}
&\left\| \mathcal{G} \left( \frac{\hat{\psi}_{\sim, \epsilon, L, \delta}^n - \hat{\psi}_{\sim, \epsilon, L, \delta}^{n-1}}{\Delta t} \right) \right\|_{\hat{X}}^2 \\
&\leq C \left[ \|\hat{\psi}_{\sim, \epsilon, L, \delta}^n\|_{\hat{X}}^2 + \|u_{\sim, \epsilon, L, \delta}^n\|_{\mathbb{H}^1(\Omega)}^2 + \int_{\Omega \times D} M |u_{\sim, \epsilon, L, \delta}^{n-1}|^2 |\hat{\psi}_{\sim, \epsilon, L, \delta}^n|^2 dq dx \right]. \quad (5.3.60)
\end{aligned}$$

Similarly to (5.3.58), on noting (5.3.4) and (5.3.10), we have that

$$\begin{aligned}
\int_{\Omega \times D} M |u_{\sim, \epsilon, L, \delta}^{n-1}|^2 |\hat{\psi}_{\sim, \epsilon, L, \delta}^n|^2 dq dx &\leq \|u_{\sim, \epsilon, L, \delta}^{n-1}\|_{L^4(\Omega)}^2 \|\hat{\psi}_{\sim, \epsilon, L, \delta}^n\|_{L^4(\Omega; L_M^2(D))}^2 \\
&\leq C \left[ \|u_{\sim, \epsilon, L, \delta}^{n-1}\|_{L^2(\Omega)}^{4-d} \|\nabla_x u_{\sim, \epsilon, L, \delta}^{n-1}\|_{L^2(\Omega)}^d \right. \\
&\quad \left. + \|\hat{\psi}_{\sim, \epsilon, L, \delta}^n\|_{L^2(\Omega; L_M^2(D))}^{4-d} \|\hat{\psi}_{\sim, \epsilon, L, \delta}^n\|_{\mathbb{H}^1(\Omega; L_M^2(D))}^d \right]. \quad (5.3.61)
\end{aligned}$$

Taking the  $\frac{2}{d}$  power of both sides of (5.3.60), summing from  $n = 1 \rightarrow N$ , and noting (5.3.61), (5.3.56) and (5.3.23) yields, similarly to (5.3.59), that

$$\sum_{n=1}^N \Delta t \left\| \mathcal{G} \left( \frac{\hat{\psi}_{\epsilon,L,\delta}^n - \hat{\psi}_{\epsilon,L,\delta}^{n-1}}{\Delta t} \right) \right\|_{\hat{X}}^{\frac{4}{d}} \leq C(L, T). \quad (5.3.62)$$

Now we introduce some definitions prior to passing to the limit  $\Delta t \rightarrow 0_+$ . Let

$$u_{\epsilon,L,\delta}^{\Delta t}(\cdot, t) := \frac{t - t^{n-1}}{\Delta t} u_{\epsilon,L,\delta}^n(\cdot) + \frac{t^n - t}{\Delta t} u_{\epsilon,L,\delta}^{n-1}(\cdot), \quad t \in [t^{n-1}, t^n], \quad n \geq 1, \quad (5.3.63a)$$

and

$$u_{\epsilon,L,\delta}^{\Delta t,+}(\cdot, t) := u^n(\cdot), \quad u_{\epsilon,L,\delta}^{\Delta t,-}(\cdot, t) := u^{n-1}(\cdot), \quad t \in (t^{n-1}, t^n], \quad n \geq 1. \quad (5.3.63b)$$

We note for future reference that

$$u_{\epsilon,L,\delta}^{\Delta t} - u_{\epsilon,L,\delta}^{\Delta t,\pm} = (t - t^{n,\pm}) \frac{\partial u_{\epsilon,L,\delta}^{\Delta t}}{\partial t}, \quad t \in (t^{n-1}, t^n), \quad n \geq 1, \quad (5.3.64)$$

where  $t^{n,+} := t^n$  and  $t^{n,-} := t^{n-1}$ . Using the above notation, and introducing analogous notation for  $\{\hat{\psi}_{\epsilon,L,\delta}^n, f_{\sim}^n\}_{n=0}^N$ , (5.3.27) summed for  $n = 1 \rightarrow N$  can be restated as

$$\begin{aligned} & \int_0^T \left\langle \frac{\partial u_{\epsilon,L,\delta}^{\Delta t}}{\partial t}, w \right\rangle_V dt \\ & \quad + \int_0^T \int_{\Omega} \left[ \left[ (u_{\epsilon,L,\delta}^{\Delta t,-} \cdot \nabla_x) u_{\epsilon,L,\delta}^{\Delta t,+} \right] \cdot w + \nu \nabla_x u_{\epsilon,L,\delta}^{\Delta t,+} : \nabla_x w \right] dx dt \\ & = \int_0^T \left[ \langle f^{\Delta t,+}, w \rangle_V - k_B \mathcal{T} \int_{\Omega} C(M \hat{\psi}_{\epsilon,L,\delta}^{\Delta t,+}) : \nabla_x w dx \right] dt \quad \forall w \in L^{\frac{4}{4-d}}(0, T; V). \end{aligned} \quad (5.3.65)$$

Similarly, (5.3.31) summed for  $n = 1 \rightarrow N$  can be restated as

$$\begin{aligned} & \int_0^T \left\langle M \frac{\partial \hat{\psi}_{\epsilon,L,\delta}^{\Delta t}}{\partial t}, \hat{\varphi} \right\rangle_{\hat{X}} dq dx dt \\ & \quad + \int_0^T \int_{\Omega \times D} M \left[ \frac{1}{2\lambda} \nabla_q \hat{\psi}_{\epsilon,L,\delta}^{\Delta t,+} - [\sigma(u_{\epsilon,L,\delta}^{\Delta t,+}) q] \beta_{\delta}^L(\hat{\psi}_{\epsilon,L,\delta}^{\Delta t,+}) \right] \cdot \nabla_q \hat{\varphi} dq dx dt \\ & \quad + \int_0^T \int_{\Omega \times D} M \left[ \varepsilon \nabla_x \hat{\psi}_{\epsilon,L,\delta}^{\Delta t,+} - u_{\epsilon,L,\delta}^{\Delta t,-} \hat{\psi}_{\epsilon,L,\delta}^{\Delta t,+} \right] \cdot \nabla_x \hat{\varphi} dq dx dt = 0 \\ & \quad \forall \hat{\varphi} \in L^{\frac{4}{4-d}}(0, T; \hat{X}). \end{aligned} \quad (5.3.66)$$

We have from (5.3.53) and (5.3.63a,b), on noting (5.3.18b), that

$$\begin{aligned} & \sup_{t \in (0, T)} \left[ \int_{\Omega} |u_{\epsilon,L,\delta}^{\Delta t(\pm)}|^2 dx \right] + \int_0^T \int_{\Omega} \frac{|u_{\epsilon,L,\delta}^{\Delta t,+} - u_{\epsilon,L,\delta}^{\Delta t,-}|^2}{\Delta t} dx dt \\ & \quad + \nu \int_0^T \int_{\Omega} |\nabla_x u_{\epsilon,L,\delta}^{\Delta t(\pm)}|^2 dx dt \leq C(T). \end{aligned} \quad (5.3.67)$$

In the above, the notation  $\underline{u}_{\epsilon,L,\delta}^{\Delta t(\pm)}$  means  $\underline{u}_{\epsilon,L,\delta}^{\Delta t}$  with or without the superscripts  $\pm$ . Similarly, we have from (5.3.56), (5.3.53), (5.3.19), (5.3.15), (5.3.59), (5.3.62) and (5.3.63a,b) that

$$\begin{aligned}
& \sup_{t \in (0,T)} \left[ \int_{\Omega \times D} M |\hat{\psi}_{\epsilon,L,\delta}^{\Delta t(\pm)}|^2 dq dx \right] + \frac{1}{\delta} \sup_{t \in (0,T)} \left[ \int_{\Omega \times D} M [\hat{\psi}_{\epsilon,L,\delta}^{\Delta t(\pm)}]_-^2 dq dx \right] \\
& + \frac{1}{\lambda} \int_0^T \int_{\Omega \times D} M \left| \nabla_q \hat{\psi}_{\epsilon,L,\delta}^{\Delta t,+} \right|^2 dq dx dt + \varepsilon \int_0^T \int_{\Omega \times D} M \left| \nabla_x \hat{\psi}_{\epsilon,L,\delta}^{\Delta t,+} \right|^2 dq dx dt \\
& + \sup_{t \in (0,T)} \left[ \int_{\Omega} |C(M \hat{\psi}_{\epsilon,L,\delta}^{\Delta t(\pm)})|^2 dx \right] + \int_0^T \int_{\Omega \times D} M \frac{|\hat{\psi}_{\epsilon,L,\delta}^{\Delta t,+} - \hat{\psi}_{\epsilon,L,\delta}^{\Delta t,-}|^2}{\Delta t} dq dx dt \\
& + \int_0^T \int_{\Omega \times D} M \left[ \left\| \underline{S} \frac{\partial \underline{u}_{\epsilon,L,\delta}^{\Delta t}}{\partial t} \right\|_{\mathbf{H}^1(\Omega)}^{\frac{4}{d}} + \left\| \underline{G} \frac{\partial \hat{\psi}_{\epsilon,L,\delta}^{\Delta t}}{\partial t} \right\|_{\hat{\mathbf{X}}}^{\frac{4}{d}} \right] dq dx dt \leq C(L, T). \quad (5.3.68)
\end{aligned}$$

We are now in a position to prove the following convergence result.

**Lemma 5.3.4** *There exists a subsequence of  $\{\underline{u}_{\epsilon,L,\delta}^{\Delta t}, \hat{\psi}_{\epsilon,L,\delta}^{\Delta t}\}_{\Delta t > 0}$ , and functions  $\underline{u}_{\epsilon,L,\delta} \in L^\infty(0, T; \underline{L}^2(\Omega)) \cap L^2(0, T; \underline{Y}) \cap W^{1, \frac{4}{d}}(0, T; \underline{Y}')$  and  $\hat{\psi}_{\epsilon,L,\delta} \in L^\infty(0, T; L_M^2(\Omega \times D)) \cap L^2(0, T; \hat{\mathbf{X}}) \cap W^{1, \frac{4}{d}}(0, T; \hat{\mathbf{X}}')$  such that, as  $\Delta t \rightarrow 0_+$ ,*

$$\underline{u}_{\epsilon,L,\delta}^{\Delta t(\pm)} \rightarrow \underline{u}_{\epsilon,L,\delta} \quad \text{weak}^* \text{ in } L^\infty(0, T; \underline{L}^2(\Omega)), \quad (5.3.69a)$$

$$\underline{u}_{\epsilon,L,\delta}^{\Delta t(\pm)} \rightarrow \underline{u}_{\epsilon,L,\delta} \quad \text{weakly in } L^2(0, T; \underline{Y}), \quad (5.3.69b)$$

$$\underline{S} \frac{\partial \underline{u}_{\epsilon,L,\delta}^{\Delta t}}{\partial t} \rightarrow \underline{S} \frac{\partial \underline{u}_{\epsilon,L,\delta}}{\partial t} \quad \text{weakly in } L^{\frac{4}{d}}(0, T; \underline{Y}), \quad (5.3.69c)$$

$$\underline{u}_{\epsilon,L,\delta}^{\Delta t(\pm)} \rightarrow \underline{u}_{\epsilon,L,\delta} \quad \text{strongly in } L^2(0, T; \underline{L}^r(\Omega)), \quad (5.3.69d)$$

where  $r \in [1, \infty)$  if  $d = 2$  and  $r \in [1, 6)$  if  $d = 3$ ; and

$$M^{\frac{1}{2}} \hat{\psi}_{\epsilon,L,\delta}^{\Delta t(\pm)} \rightarrow M^{\frac{1}{2}} \hat{\psi}_{\epsilon,L,\delta} \quad \text{weak}^* \text{ in } L^\infty(0, T; L^2(\Omega \times D)), \quad (5.3.70a)$$

$$M^{\frac{1}{2}} \nabla_q \hat{\psi}_{\epsilon,L,\delta}^{\Delta t,+} \rightarrow M^{\frac{1}{2}} \nabla_q \hat{\psi}_{\epsilon,L,\delta} \quad \text{weakly in } L^2(0, T; L^2(\Omega \times D)), \quad (5.3.70b)$$

$$M^{\frac{1}{2}} \nabla_x \hat{\psi}_{\epsilon,L,\delta}^{\Delta t,+} \rightarrow M^{\frac{1}{2}} \nabla_x \hat{\psi}_{\epsilon,L,\delta} \quad \text{weakly in } L^2(0, T; L^2(\Omega \times D)), \quad (5.3.70c)$$

$$\underline{G} \frac{\partial \hat{\psi}_{\epsilon,L,\delta}^{\Delta t}}{\partial t} \rightarrow \underline{G} \frac{\partial \hat{\psi}_{\epsilon,L,\delta}}{\partial t} \quad \text{weakly in } L^{\frac{4}{d}}(0, T; \hat{\mathbf{X}}), \quad (5.3.70d)$$

$$M^{\frac{1}{2}} \hat{\psi}_{\epsilon,L,\delta}^{\Delta t(\pm)} \rightarrow M^{\frac{1}{2}} \hat{\psi}_{\epsilon,L,\delta} \quad \text{strongly in } L^2(0, T; L^2(\Omega \times D)), \quad (5.3.70e)$$

$$M^{\frac{1}{2}} \beta_\delta^L(\hat{\psi}_{\epsilon,L,\delta}^{\Delta t(\pm)}) \rightarrow M^{\frac{1}{2}} \beta_\delta^L(\hat{\psi}_{\epsilon,L,\delta}) \quad \text{strongly in } L^\infty(0, T; L^\infty(\Omega \times D)), \quad (5.3.70f)$$

$$\underline{C}(M \hat{\psi}_{\epsilon,L,\delta}^{\Delta t(\pm)}) \rightarrow \underline{C}(M \hat{\psi}_{\epsilon,L,\delta}) \quad \text{strongly in } L^2(0, T; L^2(\Omega)). \quad (5.3.70g)$$

**Proof.** The results (5.3.69a–c) follow immediately from the bounds (5.3.67) and the bound on  $\underline{u}_{\epsilon,L,\delta}^{\Delta t}$  in (5.3.68). The strong convergence result (5.3.69d) for  $\underline{u}_{\epsilon,L,\delta}^{\Delta t}$  follows immediately from (5.3.69a–c), (5.3.3) and (5.3.14), on noting that  $\mathbb{Y} \subset \mathbb{H}_0^1(\Omega)$  is compactly embedded in  $\mathbb{L}^r(\Omega)$  for the stated values of  $r$ . We now prove (5.3.69d) for  $\underline{u}_{\epsilon,L,\delta}^{\Delta t,\pm}$ . First, we obtain from the bound on the second term on the left-hand side of (5.3.67) and from (5.3.64) that

$$\|\underline{u}_{\epsilon,L,\delta}^{\Delta t} - \underline{u}_{\epsilon,L,\delta}^{\Delta t,\pm}\|_{\mathbb{L}^2(0,T;\mathbb{L}^2(\Omega))}^2 \leq C \Delta t. \quad (5.3.71)$$

Second, we note from (5.3.4) that, for all  $\eta \in \mathbb{L}^2(0,T;\mathbb{H}^1(\Omega))$ ,

$$\|\eta\|_{\mathbb{L}^2(0,T;\mathbb{L}^r(\Omega))} \leq C \|\eta\|_{\mathbb{L}^2(0,T;\mathbb{L}^2(\Omega))}^{1-\theta} \|\eta\|_{\mathbb{L}^2(0,T;\mathbb{H}^1(\Omega))}^{\theta} \quad (5.3.72)$$

for any  $r \in [2, \infty)$  if  $d = 2$  or any  $r \in [2, 6)$  if  $d = 3$ , where  $\theta = d(\frac{1}{2} - \frac{1}{r}) \in [0, 1)$ . Hence, combining (5.3.71), (5.3.72), and (5.3.69d) for  $\underline{u}_{\epsilon,L,\delta}^{\Delta t}$  yields (5.3.69d) for  $\underline{u}_{\epsilon,L,\delta}^{\Delta t,\pm}$ .

The result (5.3.70a) follows immediately from the bounds on the first and sixth terms on the left-hand side of (5.3.68). It follows immediately from the bound on the third term on the left-hand side of (5.3.68) that (5.3.70b) holds for some limit  $\underline{g} \in \mathbb{L}^2(0,T;\mathbb{L}^2(\Omega \times D))$ , which we need to identify. However, for any  $\eta \in \mathbb{L}^2(0,T;\mathbb{C}_0^\infty(\Omega \times D))$ , it follows from (5.1.4) and the compact support of  $\eta$  on  $D$  that  $[\nabla_q \cdot (M^{\frac{1}{2}} \eta)]/M^{\frac{1}{2}} \in \mathbb{L}^2(0,T;\mathbb{L}^2(\Omega \times D))$ , and hence the above convergence implies, noting (5.3.70a), that

$$\begin{aligned} \int_0^T \int_{\Omega \times D} \underline{g} \cdot \underline{\eta} \, d\underline{q} \, d\underline{x} \, dt &\leftarrow - \int_0^T \int_{\Omega \times D} M^{\frac{1}{2}} \hat{\psi}_{\epsilon,L,\delta}^{\Delta t,+} \frac{\nabla_q \cdot (M^{\frac{1}{2}} \underline{\eta})}{M^{\frac{1}{2}}} \, d\underline{q} \, d\underline{x} \, dt \\ &\rightarrow - \int_0^T \int_{\Omega \times D} M^{\frac{1}{2}} \hat{\psi}_{\epsilon,L,\delta} \frac{\nabla_q \cdot (M^{\frac{1}{2}} \underline{\eta})}{M^{\frac{1}{2}}} \, d\underline{q} \, d\underline{x} \, dt \end{aligned} \quad (5.3.73)$$

as  $\Delta t \rightarrow 0_+$ . Hence the desired result (5.3.70b) follows from (5.3.73), noting the denseness of  $\mathbb{C}_0^\infty(\Omega \times D)$  in  $\mathbb{L}^2(\Omega \times D)$ . Similar arguments prove (5.3.70c,d) on noting (5.3.70a), and the fourth and seventh bounds in (5.3.68). The strong convergence result (5.3.70e) for  $\hat{\psi}_{\epsilon,L,\delta}^{\Delta t}$  follows immediately from (5.3.70a–c), (5.3.13), (5.3.14) and (5.3.11b). Similarly to (5.3.71), the sixth bound in (5.3.68) then yields that (5.3.70e) holds for  $\hat{\psi}_{\epsilon,L,\delta}^{\Delta t,\pm}$ . Finally, the desired results (5.3.70f,g) follow immediately from (5.3.70e), (5.3.20), (5.2.3a) and (5.3.15).  $\square$

Similarly to (5.3.72), we have, for any  $r \in [2, \infty)$  if  $d = 2$  or any  $r \in [2, 6]$  if  $d = 3$ , that

$$\|\eta\|_{\mathbb{L}^{\frac{2}{\theta}}(0,T;\mathbb{L}^r(\Omega))} \leq C \|\eta\|_{\mathbb{L}^2(0,T;\mathbb{H}^1(\Omega))} \quad \text{if } \eta \in \mathbb{L}^\infty(0,T;\mathbb{L}^2(\Omega)), \quad (5.3.74a)$$

$$\|\hat{\varphi}\|_{\mathbb{L}^{\frac{2}{\theta}}(0,T;\mathbb{L}^r(\Omega;\mathbb{L}_M^2(D)))} \leq C \|\hat{\varphi}\|_{\mathbb{L}^2(0,T;\mathbb{H}^1(\Omega;\mathbb{L}_M^2(D)))} \quad \text{if } \hat{\varphi} \in \mathbb{L}^\infty(0,T;\mathbb{L}^2(\Omega;\mathbb{L}_M^2(D))); \quad (5.3.74b)$$

where  $\theta = d(\frac{1}{2} - \frac{1}{r}) \in [0, 1]$ . It follows from (5.3.69a–d), (5.3.70g), (5.3.29), (5.3.74a), (5.3.2) and (5.3.26) that we may pass to the limit,  $\Delta t \rightarrow 0_+$ , in (5.3.65) to obtain that  $\underline{u}_{\epsilon,L,\delta} \in \mathbb{L}^\infty(0,T;\mathbb{L}^2(\Omega)) \cap \mathbb{L}^2(0,T;\mathbb{Y}) \cap \mathbb{W}^{1,\frac{4}{d}}(0,T;\mathbb{Y}')$  and  $\underline{C}(M \hat{\psi}_{\epsilon,L,\delta}) \in \mathbb{L}^\infty(0,T;\mathbb{L}^2(\Omega))$  satisfy (5.3.21a). It also follows from (5.3.22) that  $\underline{u}_{\epsilon,L,\delta}(\cdot, 0) = \underline{u}_0(\cdot)$  in the required sense, recall Remark 5.3.1.

It follows from (5.3.70a–f), (5.3.69b,d), (5.3.74b) and (5.3.8) that we may pass to the limit  $\Delta t \rightarrow 0_+$  in (5.3.66) to obtain that  $\hat{\psi}_{\epsilon,L,\delta} \in \mathbb{L}^\infty(0,T;\mathbb{L}_M^2(\Omega \times D)) \cap \mathbb{L}^2(0,T;\hat{\mathbb{X}}) \cap \mathbb{W}^{1,\frac{4}{d}}(0,T;\hat{\mathbb{X}}')$  and  $\underline{u}_{\epsilon,L,\delta} \in \mathbb{L}^2(0,T;\mathbb{Y})$  satisfy (5.3.21b).

Hence we have proved existence of a global weak solution to  $(P_{\varepsilon,L,\delta})$ , (5.3.21a,b). Moreover, it follows from (5.3.67), (5.3.68), (5.3.69a–c) and (5.3.70a–g) that

$$\sup_{t \in (0,T)} \left[ \int_{\Omega} |u_{\varepsilon,L,\delta}|^2 dx \right] + \nu \int_0^T \int_{\Omega} |\nabla_x u_{\varepsilon,L,\delta}|^2 dx dt \leq C(T), \quad (5.3.75a)$$

$$\begin{aligned} & \sup_{t \in (0,T)} \left[ \int_{\Omega \times D} M |\hat{\psi}_{\varepsilon,L,\delta}|^2 dq dx \right] + \frac{1}{\delta} \sup_{t \in (0,T)} \left[ \int_{\Omega \times D} M [\hat{\psi}_{\varepsilon,L,\delta}]^2 dq dx \right] \\ & + \int_0^T \int_{\Omega \times D} M \left[ \frac{1}{\lambda} |\nabla_q \hat{\psi}_{\varepsilon,L,\delta}|^2 + \varepsilon |\nabla_x \hat{\psi}_{\varepsilon,L,\delta}|^2 \right] dq dx dt \\ & + \sup_{t \in (0,T)} \left[ \int_{\Omega} |C(M \hat{\psi}_{\varepsilon,L,\delta})|^2 dx \right] \\ & + \int_0^T \left[ \left\| \mathcal{S} \frac{\partial u_{\varepsilon,L,\delta}}{\partial t} \right\|_{\mathbb{H}^1(\Omega)}^{\frac{4}{d}} + \left\| \mathcal{G} \frac{\partial \hat{\psi}_{\varepsilon,L,\delta}}{\partial t} \right\|_{\hat{X}}^{\frac{4}{d}} \right] dt \leq C(L, T). \end{aligned} \quad (5.3.75b)$$

**Remark 5.3.5** Since the test functions in  $\mathbb{Y}$  are divergence-free, the pressure has been eliminated in (5.3.21a,b); it can be recovered in a very weak sense following the same procedure as for the incompressible Navier–Stokes equations discussed on p. 208 in Temam [119]; i.e., one obtains that  $\int_0^t p_{\varepsilon,L,\delta}(\cdot, s) ds \in C([0, T]; L^2(\Omega))$ .  $\diamond$

### 5.3.2 Existence for $(P_{\varepsilon,L})$

As the bounds (5.3.75a,b) are independent of the parameter  $\delta$ , it follows immediately, similarly to (5.3.69a–d), (5.3.70a–g), and (5.3.75a,b), that the following lemma holds.

**Lemma 5.3.6** *There exists a subsequence of  $\{u_{\varepsilon,L,\delta}, \hat{\psi}_{\varepsilon,L,\delta}\}_{\delta>0}$ , and functions*

$$u_{\varepsilon,L} \in L^\infty(0, T; \mathbb{L}^2(\Omega)) \cap L^2(0, T; \mathbb{Y}) \cap W^{1, \frac{4}{d}}(0, T; \mathbb{Y}')$$

and

$$\hat{\psi}_{\varepsilon,L} \in L^\infty(0, T; L^2_M(\Omega \times D)) \cap L^2(0, T; \hat{X}) \cap W^{1, \frac{4}{d}}(0, T; \hat{X}'),$$

with  $\hat{\psi}_{\varepsilon,L} \geq 0$  a.e. in  $\Omega \times D \times (0, T)$ , such that, as  $\delta \rightarrow 0_+$ ,

$$u_{\varepsilon,L,\delta} \rightharpoonup u_{\varepsilon,L} \quad \text{weak* in } L^\infty(0, T; \mathbb{L}^2(\Omega)), \quad (5.3.76a)$$

$$u_{\varepsilon,L,\delta} \rightharpoonup u_{\varepsilon,L} \quad \text{weakly in } L^2(0, T; \mathbb{Y}), \quad (5.3.76b)$$

$$\mathcal{S} \frac{\partial u_{\varepsilon,L,\delta}}{\partial t} \rightharpoonup \mathcal{S} \frac{\partial u_{\varepsilon,L}}{\partial t} \quad \text{weakly in } L^{\frac{4}{d}}(0, T; \mathbb{Y}), \quad (5.3.76c)$$

$$u_{\varepsilon,L,\delta} \rightarrow u_{\varepsilon,L} \quad \text{strongly in } L^2(0, T; \mathbb{L}^r(\Omega)), \quad (5.3.76d)$$

where  $r \in [1, \infty)$  if  $d = 2$  and  $r \in [1, 6)$  if  $d = 3$ ; and

$$M^{\frac{1}{2}} \hat{\psi}_{\epsilon, L, \delta} \rightarrow M^{\frac{1}{2}} \hat{\psi}_{\epsilon, L} \quad \text{weak* in } L^\infty(0, T; L^2(\Omega \times D)), \quad (5.3.77a)$$

$$M^{\frac{1}{2}} \nabla_{\tilde{q}} \hat{\psi}_{\epsilon, L, \delta} \rightarrow M^{\frac{1}{2}} \nabla_{\tilde{q}} \hat{\psi}_{\epsilon, L} \quad \text{weakly in } L^2(0, T; L^2(\Omega \times D)), \quad (5.3.77b)$$

$$M^{\frac{1}{2}} \nabla_x \hat{\psi}_{\epsilon, L, \delta} \rightarrow M^{\frac{1}{2}} \nabla_x \hat{\psi}_{\epsilon, L} \quad \text{weakly in } L^2(0, T; L^2(\Omega \times D)), \quad (5.3.77c)$$

$$\mathcal{G} \frac{\partial \hat{\psi}_{\epsilon, L, \delta}}{\partial t} \rightarrow \mathcal{G} \frac{\partial \hat{\psi}_{\epsilon, L}}{\partial t} \quad \text{weakly in } L^{\frac{4}{d}}(0, T; \hat{X}), \quad (5.3.77d)$$

$$M^{\frac{1}{2}} \hat{\psi}_{\epsilon, L, \delta} \rightarrow M^{\frac{1}{2}} \hat{\psi}_{\epsilon, L} \quad \text{strongly in } L^2(0, T; L^2(\Omega \times D)), \quad (5.3.77e)$$

$$M^{\frac{1}{2}} \beta_\delta^L(\hat{\psi}_{\epsilon, L, \delta}) \rightarrow M^{\frac{1}{2}} \beta^L(\hat{\psi}_{\epsilon, L}) \quad \text{strongly in } L^\infty(0, T; L^\infty(\Omega \times D)), \quad (5.3.77f)$$

$$\underset{\approx}{C}(M \hat{\psi}_{\epsilon, L, \delta}) \rightarrow \underset{\approx}{C}(M \hat{\psi}_{\epsilon, L}) \quad \text{strongly in } L^2(0, T; L^2(\Omega)). \quad (5.3.77g)$$

In addition, we have that

$$\sup_{t \in (0, T)} \left[ \int_{\Omega} |u_{\epsilon, L}|^2 dx \right] + \nu \int_0^T \int_{\Omega} |\nabla_x u_{\epsilon, L}|^2 dx dt \leq C(T), \quad (5.3.78a)$$

$$\begin{aligned} & \sup_{t \in (0, T)} \left[ \int_{\Omega \times D} M^{\frac{1}{2}} |\hat{\psi}_{\epsilon, L}|^2 dq dx \right] + \sup_{t \in (0, T)} \left[ \int_{\Omega} |C(M \hat{\psi}_{\epsilon, L})|^2 dx \right] \\ & + \int_0^T \int_{\Omega \times D} M \left[ \frac{1}{\lambda} \left| \nabla_q \hat{\psi}_{\epsilon, L} \right|^2 + \varepsilon \left| \nabla_x \hat{\psi}_{\epsilon, L} \right|^2 \right] dq dx dt \\ & + \int_0^T \left[ \left\| S \frac{\partial u_{\epsilon, L}}{\partial t} \right\|_{H^1(\Omega)}^{\frac{4}{d}} + \left\| \mathcal{G} \frac{\partial \hat{\psi}_{\epsilon, L}}{\partial t} \right\|_{\hat{X}}^{\frac{4}{d}} \right] dt \leq C(L, T). \end{aligned} \quad (5.3.78b)$$

In particular, the nonnegativity of  $\hat{\psi}_{\epsilon, L}$  in the above lemma follows from the second bound in (5.3.75b). Therefore we can then pass to limit  $\delta \rightarrow 0_+$  in  $(P_{\epsilon, L, \delta})$  to obtain global existence of a weak solution to the following problem for given  $\varepsilon \in (0, 1]$  and  $L > 1$ :

$(P_{\varepsilon, L})$  Find functions  $u_{\epsilon, L} \in L^\infty(0, T; L^2(\Omega)) \cap L^2(0, T; \mathcal{Y}) \cap W^{1, \frac{4}{d}}(0, T; \mathcal{Y}')$  and  $\hat{\psi}_{\epsilon, L} \in L^\infty(0, T; L_M^2(\Omega \times D)) \cap L^2(0, T; \hat{X}) \cap W^{1, \frac{4}{d}}(0, T; \hat{X}')$ , with  $\underset{\approx}{C}(M \hat{\psi}_{\epsilon, L}) \in L^\infty(0, T; L^2(\Omega))$ , such

that  $u_{\varepsilon,L}(\cdot, 0) = u_0(\cdot)$ ,  $\hat{\psi}_{\varepsilon,L}(\cdot, 0) = \hat{\psi}_0(\cdot)$  and

$$\begin{aligned} & \int_0^T \left\langle \frac{\partial u_{\varepsilon,L}}{\partial t}, w \right\rangle_V dt \\ & \quad + \int_0^T \int_{\Omega} \left[ (u_{\varepsilon,L} \cdot \nabla_x) u_{\varepsilon,L} \right] \cdot w + \nu \nabla_x u_{\varepsilon,L} : \nabla_x w \, dx dt \\ & = \int_0^T \langle f, w \rangle_V dt - k_B \mathcal{T} \int_0^T \int_{\Omega} C(M \hat{\psi}_{\varepsilon,L}) : \nabla_x w \, dx dt \quad \forall w \in L^{\frac{4}{4-d}}(0, T; V), \end{aligned} \quad (5.3.79a)$$

$$\begin{aligned} & \int_0^T \left\langle M \frac{\partial \hat{\psi}_{\varepsilon,L}}{\partial t}, \hat{\varphi} \right\rangle_{\hat{X}} dt \\ & \quad + \int_0^T \int_{\Omega \times D} M \left[ \frac{1}{2\lambda} \nabla_q \hat{\psi}_{\varepsilon,L} - [\sigma(u_{\varepsilon,L}) q] \beta^L(\hat{\psi}_{\varepsilon,L}) \right] \cdot \nabla_q \hat{\varphi} \, dq \, dx \, dt \\ & \quad + \int_0^T \int_{\Omega \times D} M \left[ \varepsilon \nabla_x \hat{\psi}_{\varepsilon,L} - u_{\varepsilon,L} \hat{\psi}_{\varepsilon,L} \right] \cdot \nabla_x \hat{\varphi} \, dq \, dx \, dt = 0 \quad \forall \hat{\varphi} \in L^{\frac{4}{4-d}}(0, T; \hat{X}). \end{aligned} \quad (5.3.79b)$$

**Remark 5.3.7** Although we have introduced  $x$ -diffusion and a cut-off above to  $\hat{\psi} = \psi/M$  in the drag term in the Fokker–Planck equation through the parameters  $\varepsilon \in (0, 1]$  and  $L > 1$  in the model  $(P_{\varepsilon,L})$  compared to the standard polymer model, (P); we wish to stress that the bounds on  $u_{\varepsilon,L}$ , the variable of real physical interest, in (5.3.78a) are independent of these parameters  $\varepsilon$  and  $L$ .  $\diamond$

**Remark 5.3.8** We also note that, for any  $s \in (0, T)$  and  $\Delta t$  sufficiently small such that  $0 < \Delta t < s$ , we can choose  $\hat{\varphi}(x, q, t) = \frac{1}{\Delta t} \{[s-t]_+ - [s-\Delta t-t]_+\}$  in  $(P_{\varepsilon,L})$  to yield that

$$\frac{1}{\Delta t} \int_{s-\Delta t}^s \int_{\Omega \times D} M \hat{\psi}_{\varepsilon,L}(x, q, t) \, dq \, dx \, dt = \int_{\Omega \times D} M \hat{\psi}_0(x, q) \, dq \, dx.$$

Passing to the limit  $\Delta t \rightarrow 0_+$ , we deduce that

$$\int_{\Omega \times D} M \hat{\psi}_{\varepsilon,L}(x, q, s) \, dq \, dx = \int_{\Omega \times D} M \psi_0(x, q) \, dq \, dx \quad \forall s \in (0, T).$$

An identical statement can be made about  $\hat{\psi}_{\varepsilon,L,\delta}$  in  $(P_{\varepsilon,L,\delta})$ .  $\diamond$

**Remark 5.3.9** In the case of a corotational model (i.e. with  $\underline{\sigma}(y) = \underline{\nabla}_x y$  replaced by  $\underline{\sigma}_{\text{corot}}(y) := \frac{1}{2} (\underline{\nabla}_x y - (\underline{\nabla}_x y)^T)$  in the drag term in the Fokker–Planck equation), the right-hand sides in the estimates (5.3.55) and (5.3.56) become independent of  $L$ , as one can exploit additional cancellations due to the skew-symmetry of  $\underline{\sigma}_{\text{corot}}(y)$ . Hence, (5.3.59) is then also independent of  $L$ . This raises the question whether in the case of a corotational model one can pass to the limit  $L \rightarrow \infty$  to recover the Fokker–Planck equation, *without* cut-off. The answer to this question is positive, however some modifications are required in the arguments

above in order to show this. For the sake of brevity, we omit the details and only highlight the key changes needed.

In our discussion above, because of the cut-off, we also control the time derivative of  $\hat{\psi}_{\varepsilon,L,\delta}$ ; without cut-off this does not appear to be possible. In addition, one should avoid (5.3.62) as the right-hand side of this inequality remains  $L$ -dependent regardless of whether or not the drag term is corotational. It is possible to get around these technical difficulties by proceeding as in Barrett and Süli [11]. Firstly, the time derivative has to be transferred from  $\hat{\psi}_{\varepsilon,L,\delta}$  to the (time-dependent) test function in the weak formulation of the Fokker–Planck equation. Secondly, as we will no longer have strong convergence of a subsequence of  $\{\hat{\psi}_{\varepsilon,L,\delta}\}_{\delta>0}$  to  $\hat{\psi}_{\varepsilon,L}$  as  $\delta \rightarrow 0_+$ , and of  $\{\hat{\psi}_{\varepsilon,L}\}_{L>1}$  to  $\hat{\psi}_{\varepsilon}$  as  $L \rightarrow \infty$ , the drag term has to be rewritten using the fact that for all  $\varrho \in \mathbb{H}_0^1(\Omega)$  and  $\hat{\varrho} \in \mathbb{H}^1(\Omega; L_M^2(D))$

$$\int_{\Omega \times D} M [\underline{g}_{\text{corot}}(\varrho) \underline{q}] \cdot \hat{\varrho} \, d\mathbf{q} \, d\mathbf{x} = \frac{1}{2} \int_{\Omega \times D} M \left[ (\varrho \cdot \underline{q}) (\nabla_x \cdot \hat{\varrho}) - [(\nabla_x \hat{\varrho}) \underline{q}] \cdot \varrho \right] \, d\mathbf{q} \, d\mathbf{x}.$$

One can then pass to the simultaneous limit  $\delta \rightarrow 0_+$  and  $L \rightarrow \infty$  in a very similar manner as we did in the final section of Barrett and Süli [11].  $\diamond$

## 5.4 Appendix: Compact embedding of Maxwellian-weighted spaces

Let us suppose that  $D$  is a bounded open ball in  $\mathbb{R}^d$  centred at  $\mathfrak{Q} \in \mathbb{R}^d$ , and let  $U$  and  $M$  be as in Section 5.1. Our aim is to prove that the embedding of the Maxwellian-weighted Sobolev space  $\mathbb{H}_M^1(\Omega \times D)$  into the Maxwellian-weighted Lebesgue space  $L_M^2(\Omega \times D)$  is compact. The proof proceeds in three steps.

### 5.4.1 Step 1: Compact embedding of $\mathbb{H}_M^1(D)$ into $L_M^2(D)$ , completeness, separability

We may suppose, with no loss of generality, that  $D = B(\mathfrak{Q}, b^{\frac{1}{2}})$ , with  $b > 0$ , as in the case of the FENE model, whereby  $\mathcal{O} = [0, \frac{b}{2})$ . As in Section 5.1, we shall assume that  $U \in C^\infty(\mathcal{O}; \mathbb{R}_{\geq 0})$ ,  $U(0) = 0$ ,  $U$  is monotonic increasing with  $\lim_{s \rightarrow (b/2)_-} U(s) = +\infty$ , and  $U$  and the associated Maxwellian  $M$  satisfy (5.1.3a,b) with  $\gamma \geq 1$ . Elsewhere in this section we require  $\gamma > 1$  (cf. (5.3.5)).

Let

$$f(r) := \lambda e^{-U(\frac{1}{2}(b^{\frac{1}{2}}-r)^2)} \quad \forall r \in (0, b^{\frac{1}{2}}],$$

where  $\lambda \in \mathbb{R}_{>0}$ . Clearly,  $\lim_{r \rightarrow 0_+} f(r) = 0$ ,  $\lim_{r \rightarrow 0_+} f'(r) = 0$ ,  $f$  is positive and increasing on  $(0, b^{\frac{1}{2}}]$ ,  $f \in C^1[0, b^{\frac{1}{2}}]$ , and

$$M(\underline{q}) = f(b^{\frac{1}{2}} - |\underline{q}|) \quad \forall \underline{q} : |\underline{q}| < b^{\frac{1}{2}}.$$

With this choice of  $D$  and  $f$ , the compactness of the embedding of the Maxwellian-weighted Sobolev space  $\mathbb{H}_M^1(D)$  into the Maxwellian-weighted Lebesgue space  $L_M^2(D)$  follows from Lemma 5.2 in Antoci [5], while Theorem 2.3 in Antoci [5] implies, with  $p = 2$ , that  $\mathbb{H}_M^1(D)$  and  $L_M^2(D)$  are Hilbert spaces.

As a matter of fact,  $H_M^1(D)$  and  $L_M^2(D)$  are separable Hilbert spaces. This, as we shall prove below, follows on noting that  $C^1(\overline{D})$  is a separable Banach space (e.g. the set  $\mathcal{P}$  of all polynomials with rational coefficients is a countable dense subset of  $C^1(\overline{D})$ ) and that, by Theorem 3.2.2(c) in Triebel [120],  $C^\infty(\overline{D})$  is dense in both  $H_M^1(D)$  and  $L_M^2(D)$ .

Indeed, given  $v \in H_M^1(D)$  and any  $\varepsilon > 0$  there exists  $\varphi \in C^\infty(\overline{D})$  such that

$$\|v - \varphi\|_{H_M^1(D)} < \frac{1}{2} \varepsilon.$$

Since  $C^1(\overline{D})$  is separable, there exists a countable dense set  $\mathcal{P} \subset C^1(\overline{D})$ ; hence, given  $\varepsilon > 0$  there exists  $p \in \mathcal{P}$  such that

$$\|\varphi - p\|_{C^1(\overline{D})} < \frac{1}{2} \left( \int_D M(\underline{q}) \, d\underline{q} \right)^{-1/2} \varepsilon.$$

Clearly,  $C^1(\overline{D}) \subset H_M^1(D)$  and therefore  $\mathcal{P} \subset H_M^1(D)$ . Thus,

$$\begin{aligned} \|v - p\|_{H_M^1(D)} &\leq \|v - \varphi\|_{H_M^1(D)} + \|\varphi - p\|_{H_M^1(D)} \\ &< \frac{1}{2} \varepsilon + \|\varphi - p\|_{C^1(\overline{D})} \left( \int_D M(\underline{q}) \, d\underline{q} \right)^{1/2} < \varepsilon. \end{aligned}$$

This shows that the countable set  $\mathcal{P} \subset H_M^1(D)$  is dense in  $H_M^1(D)$ . Therefore  $H_M^1(D)$  is separable. By an identical argument,  $L_M^2(D)$  is separable.

### 5.4.2 Step 2: Isometric isomorphisms

Let  $\Omega$  be a bounded open Lipschitz domain in  $\mathbb{R}^d$ . We now show the isometric isomorphism of the following pairs of spaces, respectively:  $L_M^2(\Omega \times D)$  and  $L^2(\Omega; L_M^2(D))$ ;  $H_M^{0,1}(\Omega \times D)$  and  $L^2(\Omega; H_M^1(D))$ ;  $H_M^{1,0}(\Omega \times D)$  and  $H^1(\Omega; L_M^2(D))$ . For a precise definition of  $H_M^{0,1}(\Omega \times D)$  and  $H_M^{1,0}(\Omega \times D)$ , see below.

**Isometric isomorphism of  $L_M^2(\Omega \times D)$  and  $L^2(\Omega; L_M^2(D))$ .** Let

$$L^2(\Omega; L_M^2(D)) := \{v \in \mathcal{M}_w(\Omega, L_M^2(D)) : \int_\Omega \|v(\underline{x})\|_{L_M^2(D)}^2 \, d\underline{x} < \infty\},$$

where

$$\mathcal{M}_w(\Omega, L_M^2(D)) := \{v : \Omega \rightarrow L_M^2(D) : v \text{ is weakly measurable on } \Omega\}.$$

Let  $\{\varphi_j\}_{j=1}^\infty$  be a complete orthonormal system in the (separable) Hilbert space  $L_M^2(D)$  with respect to the inner product  $(\cdot, \cdot)$  of  $L_M^2(D)$ . For  $v \in L^2(\Omega; L_M^2(D))$ , we define the function

$$V_N(\underline{x}, \underline{q}) := \sum_{j=1}^N (v(\underline{x}), \varphi_j) \varphi_j(\underline{q}).$$

As  $v$  is weakly measurable on  $\Omega$ , each of the functions  $\underline{x} \mapsto (v(\underline{x}), \varphi_j)$ ,  $j = 1, 2, \dots$ , is measurable on  $\Omega$ ; therefore  $(\underline{x}, \underline{q}) \mapsto (v(\underline{x}), \varphi_j)$  is measurable on  $\Omega \times D$  for all  $j = 1, 2, \dots$ . Similarly,  $\underline{q} \mapsto \varphi_j(\underline{q})$  is measurable on  $D$  for each  $j = 1, 2, \dots$ , and therefore  $(\underline{x}, \underline{q}) \mapsto \varphi_j(\underline{q})$  is measurable on  $\Omega \times D$ . Hence, also  $V_N$  is a measurable function on  $\Omega \times D$ . Now,

$$|V_N(\underline{x}, \underline{q})|^2 = \sum_{j=1}^N \sum_{k=1}^N (v(\underline{x}), \varphi_j) (v(\underline{x}), \varphi_k) \varphi_j(\underline{q}) \varphi_k(\underline{q}).$$

By the Cauchy–Schwarz inequality  $M\varphi_j\varphi_k = M^{\frac{1}{2}}\varphi_j \cdot M^{\frac{1}{2}}\varphi_k \in L^1(D)$  for all  $j, k \geq 1$ ; hence also  $M(\cdot)|V_N(\underline{x}, \cdot)|^2 \in L^1(D)$  for a.e.  $\underline{x} \in \Omega$ . Thus, by the orthonormality of the  $\varphi_j$ ,  $j = 1, 2, \dots$ , in  $L_M^2(D)$ ,

$$\int_D M(\underline{q}) |V_N(\underline{x}, \underline{q})|^2 d\underline{q} = \sum_{j=1}^N |(v(\underline{x}), \varphi_j)|^2, \quad \text{a.e. } \underline{x} \in \Omega.$$

By Bessel's inequality in  $L_M^2(D)$ , the right-hand side of this last equality is bounded by  $\|v(\underline{x})\|_{L_M^2(D)}^2$  for a.e.  $\underline{x} \in \Omega$ , and, by hypothesis,  $\underline{x} \mapsto v(\underline{x}) \in L^2(\Omega)$ ; therefore, by Fubini's theorem,  $M|V_N|^2 \in L^1(\Omega \times D)$ . Upon integrating both sides over  $\Omega$ , and using Fubini's theorem on the left-hand side to write the multiple integral over  $\Omega$  and  $D$  as an integral over  $\Omega \times D$ , we have

$$\|V_N\|_{L_M^2(\Omega \times D)}^2 := \int_{\Omega \times D} M(\underline{q}) |V_N(\underline{x}, \underline{q})|^2 d\underline{q} d\underline{x} = \sum_{j=1}^N \int_{\Omega} |(v(\underline{x}), \varphi_j)|^2 d\underline{x}. \quad (5.4.1)$$

Now, let

$$y_N(\underline{x}) := \sum_{j=1}^N |(v(\underline{x}), \varphi_j)|^2, \quad \underline{x} \in \Omega.$$

The sequence  $\{y_N(\underline{x})\}_{N=1}^{\infty}$  is monotonic increasing for almost all  $\underline{x} \in \Omega$ ; also, according to Bessel's inequality in  $L_M^2(D)$  we have that

$$0 \leq y_N(\underline{x}) \leq \|v(\underline{x})\|_{L_M^2(D)}^2 \quad \forall N \geq 1, \quad \text{a.e. } \underline{x} \in \Omega.$$

Thus  $\{y_N(\underline{x})\}_{N=1}^{\infty}$  is a bounded sequence of real numbers, for a.e.  $\underline{x} \in \Omega$ . Therefore, the sequence  $\{y_N(\underline{x})\}_{N=1}^{\infty}$  converges in  $\mathbb{R}$  for a.e.  $\underline{x} \in \mathbb{R}$ , with

$$y(\underline{x}) = \lim_{N \rightarrow \infty} y_N(\underline{x}) = \sum_{j=1}^{\infty} |(v(\underline{x}), \varphi_j)|^2, \quad \text{a.e. } \underline{x} \in \Omega.$$

By the monotone convergence theorem,

$$\begin{aligned} \lim_{N \rightarrow \infty} \sum_{j=1}^N \int_{\Omega} |(v(\underline{x}), \varphi_j)|^2 d\underline{x} &= \lim_{N \rightarrow \infty} \int_{\Omega} y_N(\underline{x}) d\underline{x} \\ &= \int_{\Omega} y(\underline{x}) d\underline{x} = \int_{\Omega} \sum_{j=1}^{\infty} |(v(\underline{x}), \varphi_j)|^2 d\underline{x}. \end{aligned} \quad (5.4.2)$$

This implies that

$$\left\{ \sum_{j=1}^N \int_{\Omega} |(v(\underline{x}), \varphi_j)|^2 d\underline{x} \right\}_{N=1}^{\infty}$$

is a convergent sequence of real numbers. Hence, it is also a Cauchy sequence in  $\mathbb{R}$ .

Since, for any  $N > L \geq 1$ ,

$$\int_{\Omega \times D} |V_N(\underline{x}, \underline{q}) - V_L(\underline{x}, \underline{q})|^2 d\underline{q} d\underline{x} = \sum_{j=L+1}^N \int_D |(v(\underline{x}), \varphi_j)|^2 d\underline{x},$$

it follows that  $\{V_N\}_{N=1}^\infty$  is a Cauchy sequence in  $L_M^2(\Omega \times D)$ . Since  $L_M^2(\Omega \times D)$  is a Hilbert space, there exists a unique  $V \in L_M^2(\Omega \times D)$  such that

$$V = \lim_{N \rightarrow \infty} V_N \quad \text{in } L_M^2(\Omega \times D). \quad (5.4.3)$$

Thus we have shown that the mapping

$$\mathcal{I} : v \in L^2(\Omega, L_M^2(D)) \mapsto V := \sum_{j=1}^{\infty} (v(\cdot), \varphi_j) \varphi_j(\cdot) \in L_M^2(\Omega \times D)$$

is correctly defined. Next, we prove that  $\mathcal{I}$  is a bijective isometry, and this will imply that the spaces  $L^2(\Omega; L_M^2(D))$  and  $L_M^2(\Omega \times D)$  are isometrically isomorphic.

We begin by showing that  $\mathcal{I}$  is injective. As  $\mathcal{I}$  is linear it suffices to prove that if  $\mathcal{I}(v) = 0$  then  $v = 0$ . Indeed, if  $\mathcal{I}(v) = 0$ , then

$$\sum_{j=1}^{\infty} (v(\underline{x}), \varphi_j) \varphi_j(\underline{q}) = 0 \quad \text{for a.e. } (\underline{x}, \underline{q}) \in \Omega \times D.$$

Since  $\{\varphi_j\}_{j=1}^\infty$  is an orthonormal system in  $L_M^2(D)$ , it follows that  $(v(\underline{x}), \varphi_j) = 0$  for a.e.  $\underline{x} \in \Omega$  and all  $j = 1, 2, \dots$ . The completeness of the orthonormal system  $\{\varphi_j\}_{j=1}^\infty$  in  $L_M^2(D)$  now implies that  $v(\underline{x}) = 0$  in  $L_M^2(D)$  for a.e.  $\underline{x} \in \Omega$ , i.e.  $v = 0$  in  $L^2(\Omega; L_M^2(D))$ .

Next we show that  $\mathcal{I}$  is surjective. Suppose that  $V \in L_M^2(\Omega \times D)$ . Then, by Fubini's theorem,  $V(\underline{x}, \cdot) \in L_M^2(D)$  for a.e.  $\underline{x} \in \Omega$ . Since  $\{\varphi_j\}_{j=1}^\infty$  is a complete orthonormal system in  $L_M^2(D)$ , it follows that

$$V(\underline{x}, \cdot) = \sum_{j=1}^{\infty} (V(\underline{x}, \cdot), \varphi_j) \varphi_j(\cdot).$$

On defining  $v(\underline{x}) := V(\underline{x}, \cdot) \in L_M^2(D)$ , we have that  $\mathcal{I}(v) = V$ . Hence  $\mathcal{I}$  is surjective.

Finally, we show that  $\mathcal{I}$  is an isometry. Clearly

$$\begin{aligned} \|V\|_{L_M^2(\Omega \times D)}^2 &\stackrel{(5.4.3)}{=} \lim_{N \rightarrow \infty} \|V_N\|_{L_M^2(\Omega \times D)}^2 \\ &\stackrel{(5.4.1)}{=} \lim_{N \rightarrow \infty} \sum_{j=1}^N \int_{\Omega} |(v(\underline{x}), \varphi_j)|^2 d\underline{x} \\ &\stackrel{(5.4.2)}{=} \int_{\Omega} \sum_{j=1}^{\infty} |(v(\underline{x}), \varphi_j)|^2 d\underline{x}. \end{aligned}$$

Applying Parseval's identity in  $L_M^2(D)$  to the infinite series under the last integral sign, we deduce that

$$\|V\|_{L_M^2(\Omega \times D)}^2 = \int_{\Omega} \|v(\underline{x})\|_{L_M^2(D)}^2 d\underline{x} = \|v\|_{L^2(\Omega; L_M^2(D))}^2.$$

Thus we have shown that  $\|\mathcal{I}v\|_{L_M^2(\Omega \times D)} = \|v\|_{L^2(\Omega; L_M^2(D))}$ , whereby  $\mathcal{I}$  is an isometry.

**Isometric isomorphism of  $H_M^{0,1}(\Omega \times D)$  and  $L^2(\Omega; H_M^1(D))$ .** Let us begin by observing that  $L_M^2(\Omega \times D) \subset L_{\text{loc}}^1(\Omega \times D)$ , and therefore any  $V$  in  $L_M^2(\Omega \times D)$  can be considered to be an element of  $\mathcal{D}'(\Omega \times D)$ , the space of  $\mathbb{R}$ -valued distributions on  $\Omega \times D$ . Let  $\nabla_q$  denote the distributional gradient with respect to  $q$ , defined on  $\mathcal{D}'(\Omega \times D)$ . We define

$$H_M^{0,1}(\Omega \times D) := \{V \in L_M^2(\Omega \times D) : \nabla_q V \in L_M^2(\Omega \times D)\}.$$

A completely identical argument to the one above shows that  $H_M^{0,1}(\Omega \times D)$  is isometrically isomorphic to  $L^2(\Omega; H_M^1(D))$ ; the only change that is required is to replace  $L_M^2(D)$  by  $H_M^1(D)$  throughout and to take  $\{\varphi_j\}_{j=1}^\infty$  to be a complete orthonormal system in the inner product  $(\cdot, \cdot)$  of the (separable) Hilbert space  $H_M^1(D)$ , instead of  $L_M^2(D)$ .

**Isometric isomorphism of  $H_M^{1,0}(\Omega \times D)$  and  $H^1(\Omega; L_M^2(D))$ .** Concerning the isometric isomorphism of  $H_M^{1,0}(\Omega \times D)$  and  $H^1(\Omega; L_M^2(D))$  we proceed as follows. Given  $v \in H^1(\Omega; L_M^2(D)) \subset L^2(\Omega; L_M^2(D))$ , we define, as in the proof of the isometric isomorphism of  $L^2(\Omega; L_M^2(D))$  and  $L_M^2(\Omega \times D)$  above, the function

$$V : (\underline{x}, \underline{q}) \in \Omega \times D \mapsto V(\underline{x}, \underline{q}) := \sum_{j=1}^{\infty} (v(\underline{x}), \varphi_j) \varphi_j(\underline{q}) \in \mathbb{R},$$

where  $\{\varphi_j\}_{j=1}^\infty$  is a complete orthonormal system in  $L_M^2(D)$ . We showed above that  $V \in L_M^2(\Omega \times D)$ , and  $\|V\|_{L_M^2(\Omega \times D)} = \|v\|_{L^2(\Omega; L_M^2(D))}$ .

Let  $\nabla_x$  denote the distributional gradient with respect to  $\underline{x}$ , defined on  $\mathcal{D}'(\Omega \times D)$ , and let  $\underline{D}_x$  denote the distributional gradient, defined on  $\mathcal{D}'(\Omega; L_M^2(D))$ , the space of  $L_M^2(D)$ -valued distributions on  $\Omega$ . Applying  $\nabla_x$  to

$$V = \sum_{j=1}^{\infty} (v, \varphi_j) \varphi_j \quad \text{in } \mathcal{D}'(\Omega \times D)$$

and noting that

$$\nabla_x V = \sum_{j=1}^{\infty} (\underline{D}_x v, \varphi_j) \varphi_j,$$

it follows from the isometric isomorphism of  $L_M^2(\Omega \times D)$  and  $L^2(\Omega; L_M^2(D))$  that

$$\begin{aligned} \|V\|_{H_M^{1,0}(\Omega \times D)}^2 &= \|V\|_{L_M^2(\Omega \times D)}^2 + \|\nabla_x V\|_{L_M^2(\Omega \times D)}^2 \\ &= \|v\|_{L^2(\Omega; L_M^2(D))}^2 + \|\underline{D}_x v\|_{L^2(\Omega; L_M^2(D))}^2 \\ &= \|v\|_{H^1(\Omega; L_M^2(D))}^2, \end{aligned}$$

which shows that  $H_M^{1,0}(\Omega \times D)$  and  $H^1(\Omega; L_M^2(D))$  are isometrically isomorphic.

### 5.4.3 Step 3: Compact embedding of $H_M^1(\Omega \times D)$ into $L_M^2(\Omega \times D)$

We use the results of Step 2 to identify the space  $L_M^2(\Omega \times D)$  with  $L^2(\Omega; L_M^2(D))$  and the space  $H_M^1(\Omega \times D) = H_M^{1,0}(\Omega \times D) \cap H_M^{0,1}(\Omega \times D)$  with  $H^1(\Omega; L_M^2(D)) \cap L^2(\Omega; H_M^1(D))$ .

Upon doing so, the compact embedding of  $H_M^1(\Omega \times D)$  into  $L_M^2(\Omega \times D)$  directly follows from the compact embedding of  $H^1(\Omega; L_M^2(D)) \cap L^2(\Omega; H_M^1(D))$  into  $L^2(\Omega; L_M^2(D))$ , implied by Theorem 2 on p.1499 in the paper of Shakhmurov [113].

## Chapter 6

# Finite element approximation of Navier–Stokes–Fokker–Planck systems

### 6.1 Introduction

This chapter is concerned with the construction and convergence analysis of a Galerkin finite element approximation to weak solutions of a system of nonlinear partial differential equations that arises from the kinetic theory of dilute polymer solutions. The solvent is an incompressible, viscous, isothermal Newtonian fluid confined to an open set  $\Omega \subset \mathbb{R}^d$ ,  $d = 2$  or  $3$ , with boundary  $\partial\Omega$ . For the sake of simplicity of presentation we shall suppose that  $\Omega$  has solid boundary  $\partial\Omega$ ; the velocity field  $\underline{u}$  will then satisfy the no-slip boundary condition  $\underline{u} = \mathbf{0}$  on  $\partial\Omega$ . The polymer chains, which are suspended in the solvent, are assumed not to interact with each other. The conservation of momentum and mass equations for the solvent then have the form of the incompressible Navier–Stokes equations in which the elastic *extra-stress* tensor  $\underline{\tau}$  (*i.e.*, the polymeric part of the Cauchy stress tensor,) appears as a source term:

Find  $\underline{u} : (\underline{x}, t) \in \bar{\Omega} \times [0, T] \mapsto \underline{u}(\underline{x}, t) \in \mathbb{R}^d$  and  $p : (\underline{x}, t) \in \Omega \times (0, T] \mapsto p(\underline{x}, t) \in \mathbb{R}$  such that

$$\frac{\partial \underline{u}}{\partial t} + (\underline{u} \cdot \nabla_x) \underline{u} - \nu \Delta_x \underline{u} + \nabla_x p = \underline{f} + \nabla_x \cdot \underline{\tau} \quad \text{in } \Omega \times (0, T], \quad (6.1.1a)$$

$$\nabla_x \cdot \underline{u} = 0 \quad \text{in } \Omega \times (0, T], \quad (6.1.1b)$$

$$\underline{u} = \mathbf{0} \quad \text{on } \partial\Omega \times (0, T], \quad (6.1.1c)$$

$$\underline{u}(\underline{x}, 0) = \underline{u}^0(\underline{x}) \quad \forall \underline{x} \in \Omega; \quad (6.1.1d)$$

where  $\underline{u}$  is the velocity field,  $p$  is the pressure,  $\nu \in \mathbb{R}_{>0}$  is the viscosity of the solvent, and  $\underline{f}$  is the density of body forces acting on the fluid.

The extra stress tensor  $\underline{\tau}$  is defined via a weighted average of  $\psi$ , the probability density function of the (random) conformation vector of the polymer molecules (cf. (6.1.3) below); the progressive Kolmogorov equation satisfied by  $\psi$  is a Fokker–Planck type second-order parabolic equation whose transport coefficients depend on the velocity field  $\underline{u}$ .

Kinetic theories of polymeric fluids ignore quantum mechanical and atomistic effects, and focus on ‘coarse-grained’ models of the polymeric *conformations*, *i.e.*, the orientation and the degree of stretching experienced by polymer molecules. The coarsest in the hierarchy of kinetic models of dilute polymers is the *dumbbell model*, which describes the polymer molecule by two beads connected by a massless elastic spring [21]; the elastic force  $\underline{F} : D \subseteq \mathbb{R}^d \rightarrow \mathbb{R}^d$  of the spring connecting the two beads is defined by a (sufficiently smooth) *spring potential*  $U : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$  through

$$\underline{F}(\underline{q}) = H U'(\tfrac{1}{2}|\underline{q}|^2) \underline{q}, \quad \underline{q} \in D, \quad (6.1.2)$$

where  $H \in \mathbb{R}_{>0}$  is a spring constant. The elongation (or conformation) vector  $\underline{q}$ , whose direction and length define the direction and length of the polymer chain represented by the dumbbell, is assumed to be confined to a balanced convex open set  $D \subset \mathbb{R}^d$ ; the term *balanced* means that  $\underline{q} \in D$ , and  $-\underline{q} \in D$  whenever  $\underline{q} \in D$ . Typically,  $D$  is an open  $d$ -dimensional ball of fixed radius  $r_D > 0$ , or an ellipse with fixed half-axes, or the whole of  $\mathbb{R}^d$ . Our analytical results in this chapter are concerned with the physically realistic case when  $D$  is bounded, although we shall also comment on the idealized situation when  $D = \mathbb{R}^d$ .

The governing equations of the dumbbell model considered here are (6.1.1a–d), where the *elastic extra-stress tensor*  $\underline{\tau}$  is defined by the *Kramers expression*:

$$\underline{\tau}(\underline{x}, t) = k_B \mathcal{T} \left( \int_D \underline{q} \underline{q}^T U' \left( \tfrac{1}{2}|\underline{q}|^2 \right) \psi(\underline{x}, \underline{q}, t) d\underline{q} - \rho(\underline{x}, t) \underline{I} \right); \quad (6.1.3)$$

here  $k_B$  is the *Boltzmann constant* and  $\mathcal{T}$  is the *absolute temperature*. Further,

$$\rho(\underline{x}, t) = \int_D \psi(\underline{x}, \underline{q}, t) d\underline{q} \quad (6.1.4)$$

signifies *density*, and the *probability density function*  $\psi(\underline{x}, \underline{q}, t)$  is a solution to the Fokker–Planck equation

$$\frac{\partial \psi}{\partial t} + (\underline{y} \cdot \nabla_{\underline{x}}) \psi + \nabla_{\underline{q}} \cdot ((\nabla_{\underline{x}} \underline{y}) \underline{q} \psi) = \varepsilon \Delta_{\underline{x}} \psi + \frac{1}{2\lambda} \nabla_{\underline{q}} \cdot (\nabla_{\underline{q}} \psi + U' \underline{q} \psi). \quad (6.1.5)$$

Here  $\lambda \in \mathbb{R}_{>0}$  and  $\varepsilon \in \mathbb{R}_{>0}$  are fixed positive real numbers, called the *relaxation time* and the *centre-of-mass diffusion coefficient*, respectively. We refer to [11] for the derivation of the model; see also the recent paper of Schieber [108] for a justification of the presence of the  $\underline{x}$ -dissipative centre-of-mass diffusion term  $\varepsilon \Delta_{\underline{x}} \psi$  on the right-hand side of (6.1.5).

When  $D$  is  $B(0, b^{\frac{1}{2}})$ , a ball of radius  $b^{\frac{1}{2}}$  in  $\mathbb{R}^d$  centred at the origin, a typical spring force  $\underline{F}(\underline{q})$  for a finitely-extensible model, such as the FENE (finitely-extensible nonlinear elastic) model for example in which

$$U(s) = -\frac{b}{2} \ln \left( 1 - \frac{2s}{b} \right), \quad s \in [0, \tfrac{b}{2}),$$

explodes as  $\underline{q}$  approaches  $\partial D$ ; see Section 6.2.2 below. Parabolic PDEs with unbounded coefficients are studied, for example, in the monographs of Cerrai [31] and Lorenzi and Bertoldi [89]; see also the article of Da Prato and Lunardi [105] and references therein. We note in passing that, on letting  $b \rightarrow +\infty$ , the FENE potential converges to the (linear) Hookean spring potential  $U(s) = s$  while  $D$  then becomes the whole of  $\mathbb{R}^d$ , — corresponding to a mathematically

simple(r) albeit physically unrealistic scenario in which a polymer chain can have arbitrarily large elongation.

We note in passing that in contrast with the case of Hookean dumbbells, the FENE model does not have an exact closure at the macroscopic level, though Du, Yu, and Liu [42] and Yu, Du, and Liu [125] have recently considered the analysis of approximate closures of the FENE model. Previously, El-Kareh and Leal [45] had proposed a macroscopic model, with added dissipation in the equation which governs the evolution of the conformation tensor  $\underline{A}(\underline{x}, t) := \int_D \underline{q} \underline{q}^T \psi(\underline{x}, \underline{q}, t) d\underline{q}$  in order to account for Brownian motion across streamlines; the model can be thought of as an approximate macroscopic closure of a FENE-type microscopic-macroscopic model with centre-of-mass diffusion.

An early effort to show the existence and uniqueness of local-in-time solutions to a family of bead-spring type polymeric flow models is due to Renardy [106]. While the class of potentials  $\underline{F}(\underline{q})$  considered by Renardy [106] (cf. hypotheses (F) and (F') on pp. 314–315) does include the case of Hookean dumbbells, it excludes the practically relevant case of the FENE model (see Section 6.2.2 below). More recently, E, Li, and Zhang [43] and Li, Zhang, and Zhang [81] have revisited the question of local existence of solutions for dumbbell models.

The existence of global weak solutions to the coupled Navier–Stokes–Fokker–Planck systems of the form (6.1.1a)–(6.1.5) with FENE type potentials, and related systems of partial differential equations, have been studied by Barrett, Schwab, and Süli [10], Constantin [37], Lions and Masmoudi [87], Barrett and Süli [11], [13], Otto and Tzavaras [102], and Masmoudi [96]. We refer to [13] for a detailed survey of the relevant literature.

For a survey of numerical algorithms for the approximation of kinetic models of dilute polymers see, for example, Section 4 of the survey article of Li and Zhang [82]; for recent progress on deterministic algorithms for the approximation of Fokker–Planck and coupled Navier–Stokes–Fokker–Planck systems, see, for example, Lozinski et al. [92, 93], and Knezevic and Süli [71, 72].

This chapter should be seen as a continuation of the discussion in the previous chapter, which was based on our recent work [12], and [13]; in [13], under very general assumptions on the finite-dimensional spaces used for the purpose of spatial discretization, including, in particular, classical conforming finite element spaces and spectral Galerkin subspaces, we showed the convergence of a (sub)sequence of numerical approximations to a weak solution of the coupled Navier–Stokes–Fokker–Planck system (6.1.1a)–(6.1.5), for a large class of unbounded spring potentials, including the FENE potential, in the case of the corotational model, where  $\underline{\nabla}_x \underline{u}$  in the Fokker–Planck equation is replaced by its skew-symmetric part  $\frac{1}{2}(\underline{\nabla}_x \underline{u} - (\underline{\nabla}_x \underline{u})^T)$ .

Here, we shall be concerned with the general noncorotational model (6.1.1a)–(6.1.5), but where a cut-off function  $\beta^L(\cdot) := \min(\cdot, L)$ , with  $L \gg 1$ , is introduced into the drag and convective terms of (6.1.5). In contrast with the previous chapter where we used the subscript  $_L$  to indicate the presence of the cut-off, here we shall use the superscript,  $^L$ . The chapter is organized as follows. Section 6.2 is devoted to the statement of the problem, including our structural assumptions on the admissible class of nonlinear spring potentials. In addition, we review the energy law satisfied by the system. In Section 6.3, we introduce the appropriate function spaces for the problem. Finally, in Section 6.4 we introduce our Galerkin finite element method for this coupled Navier–Stokes–Fokker–Planck system with microscopic cut-off, which involves an additional regularization parameter  $\delta > 0$ . We show the existence of this numerical approximation, and that it satisfies a discrete analogue of the

energy law for the continuous system. We then pass to the limit as the spatial discretization parameter  $h$  and the time step parameter  $\Delta t$ , as well as the regularization parameter  $\delta$ , tend to zero; using a weak-compactness argument in Maxwellian-weighted Sobolev spaces we show that a subsequence of the sequence  $\{\underline{u}_{\delta,h}^{\Delta t}, \hat{\psi}_{\delta,h}^{\Delta t}\}_{\delta>0, h>0, \Delta t>0}$  of numerical approximations to the velocity field  $\underline{u}$  and the scaled probability density function  $\hat{\psi} = \psi/M$ , where  $M$  is the normalized Maxwellian

$$M(\underline{q}) = Z^{-1} \exp(-U(\frac{1}{2}|\underline{q}|^2)), \quad (6.1.6)$$

where

$$Z := \int_D \exp(-U(\frac{1}{2}|\underline{q}|^2)) \, d\underline{q},$$

converges to a weak solution  $\{\underline{u}, \hat{\psi}\}$  of the coupled Navier–Stokes–Fokker–Planck system with microscopic cut-off. We close the chapter with an Appendix, where we use the Brascamp–Lieb inequality to construct a quasi-interpolation operator in Maxwellian-weighted Sobolev spaces. By applying an extension of the Bramble–Hilbert lemma due to Tartar, we prove sharp approximation error bounds; we also establish an, apparently new, elliptic regularity result in the Maxwellian-weighted  $H^2$  norm on  $D$ ; we then use these results to show that the orthogonal projection operator in the Maxwellian-weighted  $L^2$  inner product is stable in the Maxwellian-weighted  $H^1$  norm, — a result that plays a crucial role in our convergence proof of the numerical method.

The passage to the limit in this chapter is performed under minimal regularity assumptions on the data. The definition of the sequence of approximating solutions is completely constructive in the sense that it is based on a fully-discrete and practically implementable Galerkin finite element method. To the best of our knowledge this is the first rigorous result concerning the convergence of a sequence of numerical approximations to a global weak solution of the coupled Navier–Stokes–Fokker–Planck model in the case of a general, non-rotational, drag term.

## 6.2 Polymer models

We term polymer models under consideration here microscopic-macroscopic type models, since the continuum mechanical *macroscopic* equations of incompressible fluid flow are coupled to a *microscopic* model: the Fokker–Planck equation describing the statistical properties of particles in the continuum. We first present these equations and collect the assumptions on the parameters in the model.

### 6.2.1 Microscopic-macroscopic polymer models

Let  $\Omega \subset \mathbb{R}^d$  be a bounded open set with a Lipschitz-continuous boundary  $\partial\Omega$ , and suppose that the set  $D \subseteq \mathbb{R}^d$ ,  $d = 2$  or  $3$ , of admissible elongation vectors  $\underline{q}$  in (6.1.5) is a balanced convex open set. For the sake of simplicity of presentation, we shall suppose that  $D$  is a bounded open ball in  $\mathbb{R}^d$ . Gathering (6.1.1a–d), (6.1.3) and (6.1.5), we then consider the following initial-boundary-value problem:

**(P)** Find  $\underline{u} : (\underline{x}, t) \in \bar{\Omega} \times [0, T] \mapsto \underline{u}(\underline{x}, t) \in \mathbb{R}^d$  and  $p : (\underline{x}, t) \in \Omega \times (0, T) \mapsto p(\underline{x}, t) \in \mathbb{R}$  such

that

$$\frac{\partial \underline{u}}{\partial t} + (\underline{u} \cdot \underline{\nabla}_x) \underline{u} - \nu \Delta_x \underline{u} + \underline{\nabla}_x p = \underline{f} + \underline{\nabla}_x \cdot \underline{\tau}(\psi) \quad \text{in } \Omega \times (0, T], \quad (6.2.1a)$$

$$\underline{\nabla}_x \cdot \underline{u} = 0 \quad \text{in } \Omega \times (0, T], \quad (6.2.1b)$$

$$\underline{u} = 0 \quad \text{on } \partial\Omega \times (0, T], \quad (6.2.1c)$$

$$\underline{u}(\underline{x}, 0) = \underline{u}^0(\underline{x}) \quad \forall \underline{x} \in \Omega; \quad (6.2.1d)$$

where  $\nu \in \mathbb{R}_{>0}$  is the given viscosity,  $\underline{f}$  is the given density of the body forces acting on the fluid, and  $\underline{\tau}(\psi) : (\underline{x}, t) \in \Omega \times (0, T) \mapsto \underline{\tau}(\psi)(\underline{x}, t) \in \mathbb{R}^{d \times d}$  is the symmetric extra-stress tensor, dependent on a probability density function  $\psi : (\underline{x}, \underline{q}, t) \in \Omega \times D \times (0, T) \mapsto \psi(\underline{x}, \underline{q}, t) \in \mathbb{R}$ , defined as

$$\underline{\tau}(\psi) = k_B T (\underline{C}(\psi) - \rho(\psi) \underline{I}). \quad (6.2.2)$$

Here  $k_B, T \in \mathbb{R}_{>0}$  are, respectively, the Boltzmann constant and the absolute temperature,  $\underline{I}$  is the unit  $d \times d$  tensor,

$$\underline{C}(\psi)(\underline{x}, t) = \int_D \psi(\underline{x}, \underline{q}, t) U'(\frac{1}{2}|\underline{q}|^2) \underline{q} \underline{q}^T d\underline{q} \quad \text{and} \quad \rho(\psi)(\underline{x}, t) = \int_D \psi(\underline{x}, \underline{q}, t) d\underline{q}. \quad (6.2.3)$$

In addition, the real-valued, continuous, nonnegative and strictly monotonic increasing function  $U$ , defined on a relatively open subset of  $[0, \infty)$ , is an elastic potential which gives the elastic force  $\underline{F} : D \rightarrow \mathbb{R}^d$  on the springs *via* (6.1.2).

The probability density  $\psi(\underline{x}, \underline{q}, t)$  represents the probability at time  $t$  of finding the centre of mass of a dumbbell in the volume element  $\underline{x} + d\underline{x}$  and having the endpoint of its elongation vector within the volume element  $\underline{q} + d\underline{q}$ . Hence  $\rho(\psi)(\underline{x}, t)$  is the density of the polymer chains located at  $\underline{x}$  at time  $t$ . The function  $\psi$  satisfies the following Fokker–Planck equation, together with suitable boundary and initial conditions:

$$\frac{\partial \psi}{\partial t} + (\underline{u} \cdot \underline{\nabla}_x) \psi + \underline{\nabla}_q \cdot ((\underline{\nabla}_x \underline{u}) \underline{q} \psi) = \frac{1}{2\lambda} \underline{\nabla}_q \cdot (\underline{\nabla}_q \psi + U' \underline{q} \psi) + \varepsilon \Delta_x \psi \quad \text{in } \Omega \times D \times (0, T], \quad (6.2.4a)$$

$$\left[ \frac{1}{2\lambda} (\underline{\nabla}_q \psi + U' \underline{q} \psi) - (\underline{\nabla}_x \underline{u}) \underline{q} \psi \right] \cdot \underline{n}_{\partial D} = 0 \quad \text{on } \Omega \times \partial D \times (0, T], \quad (6.2.4b)$$

$$\varepsilon \underline{\nabla}_x \psi \cdot \underline{n}_{\partial\Omega} = 0 \quad \text{on } \partial\Omega \times D \times (0, T], \quad (6.2.4c)$$

$$\psi(\underline{x}, \underline{q}, 0) = \psi^0(\underline{x}, \underline{q}) \geq 0 \quad \forall (\underline{x}, \underline{q}) \in \Omega \times D; \quad (6.2.4d)$$

where  $\underline{n}_{\partial D}$  and  $\underline{n}_{\partial\Omega}$  are the unit outward normal vectors to  $\partial D$  and  $\partial\Omega$ , respectively, and  $U' := U'(\frac{1}{2}|\underline{q}|^2)$ . Here  $\int_D \psi^0(\underline{x}, \underline{q}) d\underline{q} = 1$  for a.e.  $\underline{x} \in \Omega$ . The boundary conditions for  $\psi$  on  $\Omega \times \partial D \times (0, T]$  and  $\partial\Omega \times D \times (0, T]$  have been chosen so as to ensure that  $\rho(\psi)(\underline{x}, t) = \int_D \psi(\underline{x}, \underline{q}, t) d\underline{q} = \int_D \psi^0(\underline{x}, \underline{q}) d\underline{q} = 1$  for a.e.  $(\underline{x}, t) \in \Omega_T$ . In (6.2.4a–c) the parameters  $\varepsilon, \lambda \in \mathbb{R}_{>0}$ , with  $\lambda$  characterizing the elastic relaxation property of the fluid, and  $(\underline{\nabla}_x \underline{u})(\underline{x}, t) \in \mathbb{R}^{d \times d}$  with  $\{\underline{\nabla}_x \underline{u}\}_{ij} = \frac{\partial u_i}{\partial x_j}$ .

On introducing the (normalized) Maxwellian (6.1.6), we have that

$$M \nabla_q M^{-1} = -M^{-1} \nabla_q M = \nabla_q U = U' q. \quad (6.2.5)$$

Thus, the Fokker–Planck system (6.2.4a–d) can be rewritten in terms of the scaled probability density function  $\hat{\psi} = \psi/M$  as

$$M \left[ \frac{\partial \hat{\psi}}{\partial t} + (\underline{u} \cdot \nabla_x) \hat{\psi} \right] + \nabla_q \cdot ((\nabla_x \underline{u}) \underline{q} M \hat{\psi}) = \frac{1}{2\lambda} \nabla_q \cdot (M \nabla_q \hat{\psi}) + \varepsilon M \Delta_x \hat{\psi} \quad \text{in } \Omega \times D \times (0, T], \quad (6.2.6a)$$

$$M \left[ \frac{1}{2\lambda} (\nabla_q \hat{\psi} - (\nabla_x \underline{u}) \underline{q} \hat{\psi}) \cdot \underline{n}_{\partial D} \right] = 0 \quad \text{on } \Omega \times \partial D \times (0, T], \quad (6.2.6b)$$

$$\varepsilon M \nabla_x \hat{\psi} \cdot \underline{n}_{\partial \Omega} = 0 \quad \text{on } \partial \Omega \times D \times (0, T], \quad (6.2.6c)$$

$$M \hat{\psi}(x, \underline{q}, 0) = M \hat{\psi}^0(x, \underline{q}) = \psi^0(x, \underline{q}) \geq 0 \quad \forall (x, \underline{q}) \in \Omega \times D. \quad (6.2.6d)$$

### 6.2.2 FENE model

We present an example of a spring potential: the FENE potential, where  $D$  is a bounded open ball in  $\mathbb{R}^d$ .

In this widely used model

$$D = B(\underline{0}, b^{\frac{1}{2}}) \quad \text{and} \quad U(s) = -\frac{b}{2} \ln \left( 1 - \frac{2s}{b} \right),$$

$$\text{and hence} \quad e^{-U(\frac{1}{2}|\underline{q}|^2)} = \left( 1 - \frac{|\underline{q}|^2}{b} \right)^{\frac{b}{2}}. \quad (6.2.7)$$

Here  $B(\underline{0}, s)$  is the bounded open ball of radius  $s > 0$  in  $\mathbb{R}^d$  centred at the origin, and  $b > 0$  is an input parameter. Hence the length  $|\underline{q}|$  of the elongation vector  $\underline{q}$  cannot exceed  $b^{\frac{1}{2}}$ .

Letting  $b \rightarrow \infty$  in (6.2.7) leads to the so-called Hookean dumbbell model where

$$D = \mathbb{R}^d \quad \text{and} \quad U(s) = s, \quad \text{and therefore} \quad e^{-U(\frac{1}{2}|\underline{q}|^2)} = e^{-\frac{1}{2}|\underline{q}|^2}. \quad (6.2.8)$$

This particular kinetic model, with  $\varepsilon \in \mathbb{R}_{>0}$ , corresponds formally to a dissipative Oldroyd-B type model; see [11] for details.

### 6.2.3 General structural assumptions on the potential

As has been noted above, the choice of  $D = \mathbb{R}^d$  (corresponding to the Hookean model) is physically unrealistic; thus, we shall henceforth suppose for simplicity that  $D = B(\underline{0}, r_D)$  is a bounded open ball in  $\mathbb{R}^d$  of radius  $r_D \in \mathbb{R}_{>0}$  centred at the origin. We assume that  $\underline{q} \mapsto U(\frac{1}{2}|\underline{q}|^2) \in C^\infty(D)$ ; that  $\underline{q} \mapsto U(\frac{1}{2}|\underline{q}|^2)$  is nonnegative, convex and has a positive definite Hessian at each  $\underline{q} \in D$ ; that  $\underline{q} \mapsto U'(\frac{1}{2}|\underline{q}|^2)$  is positive on  $D$ ; and that there exist constants  $c_i > 0$ ,  $i = 1 \rightarrow 5$ , such that the Maxwellian  $M$  and the associated elastic potential  $U$  satisfy

$$c_1 [\text{dist}(q, \partial D)]^\zeta \leq M(q) \leq c_2 [\text{dist}(q, \partial D)]^\zeta \quad \forall q \in D, \quad (6.2.9a)$$

$$c_3 \leq [\text{dist}(q, \partial D)] U'(\frac{1}{2}|q|^2) \leq c_4, \quad [U'(\frac{1}{2}|q|^2)]^2 \leq c_5 U''(\frac{1}{2}|q|^2) \quad \forall q \in D. \quad (6.2.9b)$$

It is an easy matter to show that the Maxwellian  $M$  and the elastic potential  $U$  of the FENE dumbbell model satisfy conditions (6.2.9a,b) with  $D = B(\underline{0}, b^{\frac{1}{2}})$  and  $\zeta = \frac{b}{2}$ . Since  $[U(q)]^2 = (-\ln M(q) + \text{Const.})^2$ , it follows from (6.2.9a,b) that if  $\zeta > 1$ , then

$$\int_D M [1 + U^2 + |U'|^2] \, dq < \infty. \quad (6.2.10)$$

We shall therefore suppose that  $\zeta > 1$ . For the FENE model (6.2.7),  $\zeta = \frac{b}{2}$ , and so the condition  $\zeta > 1$  translates into the requirement that  $b > 2$ . It is interesting to note that in the, equivalent, stochastic version of the FENE model, a solution to the system of stochastic differential equations associated with the Fokker–Planck equation exists and has trajectorial uniqueness if, and only if,  $b > 2$  (cf. [62] for details). Thus, the assumption  $\zeta > 1$  can be seen as the weakest reasonable requirement on the decay-rate of  $M$  as  $\text{dist}(q, \partial D) \rightarrow 0_+$ .

#### 6.2.4 Formal estimates

We end this section by identifying formally the energy structure for (P). Multiplying (6.2.1a) by  $\underline{u}$ , integrating over  $\Omega$ , and noting (6.2.1b,c) yields that

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \left[ \int_{\Omega} |u|^2 \, dx \right] + \nu \int_{\Omega} |\nabla_x u|^2 \, dx - \int_{\Omega} f \cdot u \, dx &= - \int_{\Omega} \tau(M \hat{\psi}) : \nabla_x u \, dx \\ &= -k_B T \int_{\Omega} C(M \hat{\psi}) : \nabla_x u \, dx. \end{aligned} \quad (6.2.11)$$

Let  $\mathcal{F}(s) := (\ln s - 1)s + 1$  for  $s > 0$ , with  $\mathcal{F}(0) := 1$ . Multiplying the Fokker–Planck equation (6.2.6a) with  $\mathcal{F}'(\hat{\psi}) \equiv \ln \hat{\psi}$ , on assuming that  $\hat{\psi} > 0$ , integrating over  $\Omega \times D$  yields that

$$\begin{aligned} \frac{d}{dt} \left[ \int_{\Omega \times D} M \mathcal{F}(\hat{\psi}) \, dq \, dx \right] + \int_{\Omega \times D} M \left[ \frac{1}{2\lambda} \nabla_q \hat{\psi} \cdot \nabla_q [\mathcal{F}'(\hat{\psi})] + \varepsilon \nabla_x \hat{\psi} \cdot \nabla_x [\mathcal{F}'(\hat{\psi})] \right] \, dq \, dx \\ = \int_{\Omega \times D} M \hat{\psi} [(\nabla_x u) q] \cdot \nabla_q [\mathcal{F}'(\hat{\psi})] \, dq \, dx. \end{aligned} \quad (6.2.12)$$

It follows, on noting that  $\mathcal{F}''(s) = s^{-1} > 0$  for  $s > 0$  and hence that  $\hat{\psi} \nabla_q [\mathcal{F}'(\hat{\psi})] = \nabla_q \hat{\psi}$ , (6.2.5), (6.2.1b) and  $M = 0$  on  $\partial D$  that

$$\begin{aligned} \int_{\Omega \times D} M \hat{\psi} [(\nabla_x u) q] \cdot \nabla_q [\mathcal{F}'(\hat{\psi})] \, dq \, dx &= \int_{\Omega \times D} M [(\nabla_x u) q] \cdot \nabla_q \hat{\psi} \, dq \, dx \\ &= \int_{\Omega \times D} M U' q \cdot [(\nabla_x u) q] \hat{\psi} \, dq \, dx \\ &= \int_{\Omega} C(M \hat{\psi}) : \nabla_x u \, dx, \end{aligned} \quad (6.2.13)$$

on recalling (6.2.3). Combining (6.2.11)–(6.2.13), we obtain the following energy law for (P):

$$\begin{aligned} \frac{d}{dt} \left[ \frac{1}{2} \int_{\Omega} |u|_{\sim}^2 dx + k_B \mathcal{T} \int_{\Omega \times D} M \mathcal{F}(\hat{\psi}) dq dx \right] + \nu \int_{\Omega} |\nabla_x u|_{\sim}^2 dx \\ + k_B \mathcal{T} \int_{\Omega \times D} M \left[ \frac{1}{2\lambda} \nabla_q \hat{\psi} \cdot \nabla_q [\mathcal{F}'(\hat{\psi})] + \varepsilon \nabla_x \hat{\psi} \cdot \nabla_x [\mathcal{F}'(\hat{\psi})] \right] dq dx = \int_{\Omega} f \cdot u dx. \end{aligned} \quad (6.2.14)$$

To make the above rigorous, and for computational purposes, we replace the convex function  $\mathcal{F} \in C(\mathbb{R}_{\geq 0}) \cap C^\infty(\mathbb{R}_{> 0})$  by the the convex regularization  $\mathcal{F}_\delta^L \in C^{2,1}(\mathbb{R})$  defined, for any  $\delta \in (0, 1)$  and  $L > 1$ , by

$$\mathcal{F}_\delta^L(s) := \begin{cases} \frac{s^2 - \delta^2}{2\delta} + (\ln \delta - 1)s + 1 & s \leq \delta, \\ \mathcal{F}(s) \equiv (\ln s - 1)s + 1 & \delta \leq s \leq L, \\ \frac{s^2 - L^2}{2L} + (\ln L - 1)s + 1 & L \leq s. \end{cases} \quad (6.2.15)$$

Hence, it follows that

$$[\mathcal{F}_\delta^L]'(s) = \begin{cases} \frac{s}{\delta} + \ln \delta - 1 & s \leq \delta, \\ \ln s & \delta \leq s \leq L, \\ \frac{s}{L} + \ln L - 1 & L \leq s, \end{cases} \quad \text{and} \quad [\mathcal{F}_\delta^L]''(s) = \begin{cases} \delta^{-1} & s \leq \delta, \\ s^{-1} & \delta \leq s \leq L, \\ L^{-1} & L \leq s. \end{cases} \quad (6.2.16)$$

In addition, we introduce

$$\beta_\delta^L(s) := [[\mathcal{F}_\delta^L]''(s)]^{-1} = \begin{cases} \delta & s \leq \delta, \\ s & \delta \leq s \leq L, \\ L & L \leq s. \end{cases} \quad (6.2.17)$$

It follows from (6.2.17) for any sufficiently smooth  $\hat{\varphi}$  that

$$\beta_\delta^L(\hat{\varphi}) \nabla_x ([\mathcal{F}_\delta^L]'(\hat{\varphi})) = \nabla_x \hat{\varphi} \quad \text{and} \quad \beta_\delta^L(\hat{\varphi}) \nabla_q ([\mathcal{F}_\delta^L]'(\hat{\varphi})) = \nabla_q \hat{\varphi}. \quad (6.2.18)$$

Let  $\{u_\delta^L, \hat{\psi}_\delta^L\}$  solve problem  $(P_\delta^L)$ , which is a regularization of the problem (P) where the drag term  $\nabla_q \cdot ((\nabla_x u) q \hat{\psi})$  in the Fokker–Planck equation (6.2.6a) is replaced by

$$\nabla_q \cdot ((\nabla_x u_\delta^L) q \beta_\delta^L(\hat{\psi}_\delta^L)). \quad (6.2.19)$$

Multiplying the Fokker–Planck equation in  $(P_\delta^L)$  with  $[\mathcal{F}_\delta^L]'(\hat{\psi}_\delta^L)$ , integrating over  $\Omega \times D$ , noting (6.2.18) yields, similarly to (6.2.12) and (6.2.13), that

$$\begin{aligned} \frac{d}{dt} \left[ \int_{\Omega \times D} M \mathcal{F}_\delta^L(\hat{\psi}_\delta^L) dq dx \right] + \frac{1}{2\lambda} \int_{\Omega \times D} M \nabla_q \hat{\psi}_\delta^L \cdot \nabla_q [\mathcal{F}_\delta^L]'(\hat{\psi}_\delta^L) dq dx \\ + \varepsilon \int_{\Omega \times D} M \nabla_x \hat{\psi}_\delta^L \cdot \nabla_x [\mathcal{F}_\delta^L]'(\hat{\psi}_\delta^L) dq dx = \int_{\Omega} C(M \hat{\psi}_\delta^L) : \nabla_x u_\delta^L dx. \end{aligned} \quad (6.2.20)$$

Combining (6.2.20) and the  $(P_\delta^L)$  version of (6.2.11), we obtain the following energy law for  $(P_\delta^L)$ , the regularized analogue of (6.2.14),

$$\begin{aligned} & \frac{d}{dt} \left[ \frac{1}{2} \int_{\Omega} |u_\delta^L|^2 dx + k_B \mathcal{T} \int_{\Omega \times D} M \mathcal{F}_\delta^L(\widehat{\psi}_\delta^L) dq dx \right] + \nu \int_{\Omega} |\nabla_x u_\delta^L|^2 dx \\ & + k_B \mathcal{T} \int_{\Omega \times D} M \left[ \frac{1}{2\lambda} \nabla_q \widehat{\psi}_\delta^L \cdot \nabla_q \left[ [\mathcal{F}_\delta^L]'(\widehat{\psi}_\delta^L) \right] + \varepsilon \nabla_x \widehat{\psi}_\delta^L \cdot \nabla_x \left[ [\mathcal{F}_\delta^L]'(\widehat{\psi}_\delta^L) \right] \right] dq dx \\ & = \int_{\Omega} f \cdot u_\delta^L dx. \end{aligned} \quad (6.2.21)$$

On noting that  $[\mathcal{F}_\delta^L]'' \geq L^{-1}$ , and

$$\min\{\mathcal{F}_\delta^L(s), s [\mathcal{F}_\delta^L]'(s)\} \geq \begin{cases} \frac{s^2}{2\delta} & \text{if } s \leq 0, \\ \frac{s^2}{4L} - C(L) & \text{if } s \geq 0, \end{cases} \quad (6.2.22)$$

one can establish from (6.2.21), on assuming that  $\widehat{\psi}^0 \leq L$ , that

$$\sup_{t \in (0, T)} \left[ \int_{\Omega} |u_\delta^L|^2 dx + \nu \int_{\Omega_T} |\nabla_x u_\delta^L|^2 dx dt + \delta^{-1} \sup_{t \in (0, T)} \left[ \int_{\Omega \times D} M |[\widehat{\psi}_\delta^L]_-|^2 dq dx \right] \right] \leq C. \quad (6.2.23)$$

In addition, one can establish that

$$\begin{aligned} & \sup_{t \in (0, T)} \left[ \int_{\Omega \times D} M |\widehat{\psi}_\delta^L|^2 dq dx \right] + \frac{1}{\lambda} \int_0^T \int_{\Omega \times D} M \left| \nabla_q \widehat{\psi}_\delta^L \right|^2 dq dx dt \\ & + \varepsilon \int_0^T \int_{\Omega \times D} M \left| \nabla_x \widehat{\psi}_\delta^L \right|^2 dq dx dt + \sup_{t \in (0, T)} \left[ \int_{\Omega} |C(M \widehat{\psi}_\delta^L)|^2 dx \right] \leq C(L, T). \end{aligned} \quad (6.2.24)$$

The above formal bounds have been made rigorous and the existence of a global-in-time weak solution  $\{u_\delta^L, \widehat{\psi}_\delta^L\}$  to  $(P_\delta^L)$  has been established in [12], see also the previous chapter. Moreover, one can take the limit  $\delta \rightarrow 0_+$  in problem  $(P_\delta^L)$  to establish the existence of a global-in-time weak solution  $\{u^L, \widehat{\psi}^L\}$  to problem  $(P^L)$ , which is a regularization of the problem (P) where the drag term  $\nabla_q \cdot ((\nabla_x u) q \widehat{\psi})$  in the Fokker–Planck equation (6.2.6a) is replaced by

$$\nabla_q \cdot ((\nabla_x u) q \beta^L(\widehat{\psi}_\delta^L)) \quad \text{with} \quad \beta^L(s) := \begin{cases} s & s \leq L, \\ L & L \leq s. \end{cases} \quad (6.2.25)$$

Once again, see [12] and the previous chapter.

The aim of this chapter is to construct a finite element approximation of problem  $(P_\delta^L)$ , which mimics the energy law (6.2.21) at a discrete level. Moreover, show that this approximation converges to a weak solution of  $(P^L)$ , as the spatial discretization parameter  $h$  and the time step parameter  $\Delta t$ , as well as the regularization parameter  $\delta$ , tend to zero.

### 6.3 Function spaces

Assuming that  $\partial\Omega \in C^{0,1}$ , let

$$\mathbb{H} := \{w \in \mathbb{L}^2(\Omega) : \nabla_x \cdot w = 0\} \quad \text{and} \quad \mathbb{Y} := \{w \in \mathbb{H}_0^1(\Omega) : \nabla_x \cdot w = 0\}, \quad (6.3.1)$$

where the divergence operator  $\nabla_x \cdot$  is to be understood in the sense of vector-valued distributions on  $\Omega$ . Here, and throughout, we adopt, for example, the notation  $\underline{L}^2(\Omega) \equiv [L^2(\Omega)]^d$  and  $\underline{H}_0^1(\Omega) \equiv [H_0^1(\Omega)]^d$ . Let  $\underline{V}'$  be the dual of  $\underline{V}$ . Let  $\underline{S} : \underline{V}' \rightarrow \underline{V}$  be such that  $\underline{S} \underline{v}$  is the unique solution to the Helmholtz–Stokes problem

$$\int_{\Omega} \underline{S} \underline{v} \cdot \underline{w} \, d\tilde{x} + \int_{\Omega} \nabla_x \underline{S} \underline{v} : \nabla_x \underline{w} \, d\tilde{x} = \langle \underline{v}, \underline{w} \rangle_V \quad \forall \underline{w} \in \underline{V}, \quad (6.3.2)$$

where  $\langle \cdot, \cdot \rangle_V$  denotes the duality pairing between  $\underline{V}'$  and  $\underline{V}$ . We note that

$$\langle \underline{v}, \underline{S} \underline{v} \rangle_V = \|\underline{S} \underline{v}\|_{\underline{H}^1(\Omega)}^2 \quad \forall \underline{v} \in \underline{V}' \supset (\underline{H}_0^1(\Omega))' \equiv \underline{H}^{-1}(\Omega), \quad (6.3.3)$$

and  $\|\underline{S} \cdot\|_{\underline{H}^1(\Omega)}$  is a norm on  $\underline{V}'$ . Here, and throughout, we adopt, for example, the notation  $\|\cdot\|_{\underline{H}^1(\Omega)}$  for the norm, and  $|\cdot|_{\underline{H}^1(\Omega)}$  for the semi-norm, on  $\underline{H}^1(\Omega)$  or  $\underline{H}^{-1}(\Omega)$ . We require also the duality pairing  $\langle \cdot, \cdot \rangle_{\underline{H}_0^1(\Omega)}$  between  $\underline{H}^{-1}(\Omega)$  and  $\underline{H}_0^1(\Omega)$ .

For later purposes, we recall the following well-known Gagliardo–Nirenberg inequality. Let  $r \in [2, \infty)$  if  $d = 2$ , and  $r \in [2, 6]$  if  $d = 3$  and  $\theta = d \left(\frac{1}{2} - \frac{1}{r}\right)$ . Then, there is a constant  $C$ , depending only on  $\Omega$ ,  $r$  and  $d$ , such that the following inequality holds for all  $\eta \in \underline{H}^1(\Omega)$ :

$$\|\eta\|_{L^r(\Omega)} \leq C \|\eta\|_{L^2(\Omega)}^{1-\theta} \|\eta\|_{\underline{H}^1(\Omega)}^{\theta}. \quad (6.3.4)$$

We make the following assumptions on the given initial data and the cut-off parameter  $L$  occurring in (6.2.15):

$$\underline{u}^0 \in \underline{H} \quad \text{and} \quad \hat{\psi}^0 := M^{-1} \psi^0 \in L^\infty(\Omega \times D) \quad \text{with} \quad 0 \leq \hat{\psi}^0 \leq L \quad \text{a.e. in } \Omega \times D; \quad (6.3.5a)$$

and the body force density

$$\underline{f} \in L^2(0, T; \underline{H}^{-1}(\Omega)). \quad (6.3.5b)$$

Let  $L_M^2(\Omega \times D)$  be the Maxwellian-weighted  $L^2$  space over  $\Omega \times D$  with norm

$$\|\hat{\varphi}\|_{L_M^2(\Omega \times D)} := \left\{ \int_{\Omega \times D} M |\hat{\varphi}|^2 \, d\tilde{q} \, d\tilde{x} \right\}^{\frac{1}{2}}.$$

Similarly, we consider  $L_M^2(D)$ , the Maxwellian-weighted  $L^2$  space over  $D$ . On introducing

$$\|\hat{\varphi}\|_{\underline{H}_M^1(\Omega \times D)} := \left\{ \int_{\Omega \times D} M \left[ |\hat{\varphi}|^2 + |\nabla_x \hat{\varphi}|^2 + |\nabla_q \hat{\varphi}|^2 \right] \, d\tilde{q} \, d\tilde{x} \right\}^{\frac{1}{2}}, \quad (6.3.6)$$

we then set

$$\hat{X} \equiv \underline{H}_M^1(\Omega \times D) := \left\{ \hat{\varphi} \in L_{\text{loc}}^1(\Omega \times D) : \|\hat{\varphi}\|_{\underline{H}_M^1(\Omega \times D)} < \infty \right\}. \quad (6.3.7)$$

It follows that

$$C^\infty(\overline{\Omega \times D}) \text{ is dense in } \hat{X}. \quad (6.3.8)$$

This can be shown, for example, by a simple adaptation of Lemma 3.1 in Barrett, Schwab, and Süli [10], which appeals to fundamental results on weighted Sobolev spaces in Triebel [120] and Kufner [77]. We have from Sobolev embedding that

$$L^s(\Omega; L_M^2(D)) \hookrightarrow H^1(\Omega; L_M^2(D)), \quad (6.3.9)$$

where  $s \in [1, \infty)$  if  $d = 2$  or  $s \in [1, 6]$  if  $d = 3$ . Similarly to (6.3.4) we have, with  $r$  and  $\theta$  as defined there, that there exists a constant  $C$ , depending only on  $\Omega$ ,  $r$  and  $d$ , such that

$$\|\hat{\varphi}\|_{L^r(\Omega; L_M^2(D))} \leq C \|\hat{\varphi}\|_{L^2(\Omega; L_M^2(D))}^{1-\theta} \|\hat{\varphi}\|_{H^1(\Omega; L_M^2(D))}^\theta \quad \forall \hat{\varphi} \in H^1(\Omega; L_M^2(D)). \quad (6.3.10)$$

In addition, we note that the embeddings

$$L_M^2(D) \hookrightarrow H_M^1(D), \quad (6.3.11a)$$

$$L_M^2(\Omega \times D) \equiv L^2(\Omega; L_M^2(D)) \hookrightarrow H_M^1(\Omega \times D) \equiv L^2(\Omega; H_M^1(D)) \cap H^1(\Omega; L_M^2(D)) \quad (6.3.11b)$$

are compact if  $\zeta \geq 1$  in (6.2.9a); see the Appendix to Chapter 5 or the Appendix in [12].

Let  $\hat{X}'$  be the dual space of  $\hat{X}$  with  $L_M^2(\Omega \times D)$  being the pivot space. Then, similarly to (6.3.2), let  $\mathcal{G} : \hat{X}' \rightarrow \hat{X}$  be such that  $\mathcal{G} \hat{\eta}$  is the unique solution of

$$\int_{\Omega \times D} M \left[ (\mathcal{G} \hat{\eta}) \hat{\varphi} + \nabla_q (\mathcal{G} \hat{\eta}) \cdot \nabla_q \hat{\varphi} + \nabla_x (\mathcal{G} \hat{\eta}) \cdot \nabla_x \hat{\varphi} \right] dq dx = \langle M \hat{\eta}, \hat{\varphi} \rangle_{\hat{X}} \quad \forall \hat{\varphi} \in \hat{X}, \quad (6.3.12)$$

where  $\langle M \cdot, \cdot \rangle_{\hat{X}}$  denotes the duality pairing between  $\hat{X}'$  and  $\hat{X}$ . Then, similarly to (6.3.3), we have that

$$\langle M \hat{\eta}, \mathcal{G} \hat{\eta} \rangle_{\hat{X}} = \|\mathcal{G} \hat{\eta}\|_{\hat{X}}^2 \quad \forall \hat{\eta} \in \hat{X}', \quad (6.3.13)$$

and  $\|\mathcal{G} \cdot\|_{\hat{X}}$  is a norm on  $\hat{X}'$ .

We recall the following compactness result, see, e.g., Temam [119] and Simon [115]. Let  $\mathcal{Y}_0$ ,  $\mathcal{Y}$  and  $\mathcal{Y}_1$  be Banach spaces,  $\mathcal{Y}_i$ ,  $i = 0, 1$ , reflexive, with a compact embedding  $\mathcal{Y}_0 \hookrightarrow \mathcal{Y}$  and a continuous embedding  $\mathcal{Y} \hookrightarrow \mathcal{Y}_1$ . Then, for  $\alpha_i > 1$ ,  $i = 0, 1$ , the embedding

$$\{ \eta \in L^{\alpha_0}(0, T; \mathcal{Y}_0) : \frac{\partial \eta}{\partial t} \in L^{\alpha_1}(0, T; \mathcal{Y}_1) \} \hookrightarrow L^{\alpha_0}(0, T; \mathcal{Y}) \quad (6.3.14)$$

is compact.

We note for future reference that (6.2.3) and (6.2.10) yield that, for  $\hat{\varphi} \in L_M^2(\Omega \times D)$ ,

$$\begin{aligned} \int_{\Omega} |C(M \hat{\varphi})|^2 dx &= \int_{\Omega} \sum_{i=1}^d \sum_{j=1}^d \left( \int_D M \hat{\varphi} U' q_i q_j dq \right)^2 dx \\ &\leq \left( \int_D M |U'|^2 |q|^4 dq \right) \left( \int_{\Omega \times D} M |\hat{\varphi}|^2 dq dx \right) \leq C \left( \int_{\Omega \times D} M |\hat{\varphi}|^2 dq dx \right). \end{aligned} \quad (6.3.15)$$

In [12] (see also Chapter 5), for any  $\varepsilon > 0$ ,  $L > 1$  and  $T > 0$  existence of a solution to the following weak formulation was established:

(P<sup>L</sup>) Find  $\underline{u}^L \in L^\infty(0, T; \underline{\mathbb{L}}^2(\Omega)) \cap L^2(0, T; \underline{\mathbb{V}}) \cap W^{1, \frac{4}{d}}(0, T; \underline{\mathbb{V}}')$  and  $\widehat{\psi}^L \in L^\infty(0, T; L_M^2(\Omega \times D)) \cap L^2(0, T; \widehat{\mathbb{X}}) \cap W^{1, \frac{4}{d}}(0, T; \widehat{\mathbb{X}}')$  with  $\underline{C}(M \widehat{\psi}^L) \in L^\infty(0, T; \underline{\mathbb{L}}^2(\Omega))$ , such that  $\underline{u}^L(\cdot, 0) = \underline{u}^0(\cdot)$ ,  $\widehat{\psi}^L(\cdot, \cdot, 0) = \widehat{\psi}^0(\cdot, \cdot)$  and

$$\begin{aligned} & \int_0^T \left\langle \frac{\partial \underline{u}^L}{\partial t}, \underline{w} \right\rangle_{\underline{\mathbb{V}}} dt + \int_{\Omega_T} \left[ \left[ (\underline{u}^L \cdot \nabla_x) \underline{u}^L \right] \cdot \underline{w} + \nu \nabla_x \underline{u}^L : \nabla_x \underline{w} \right] dx dt \\ & = \int_0^T \langle \underline{f}, \underline{w} \rangle_{\mathbb{H}_0^1(\Omega)} dt - k_B \mathcal{T} \int_{\Omega_T} C(M \widehat{\psi}^L) : \nabla_x \underline{w} dx dt \quad \forall \underline{w} \in L^{\frac{4}{4-d}}(0, T; \underline{\mathbb{V}}); \end{aligned} \quad (6.3.16a)$$

$$\begin{aligned} & \int_0^T \left\langle \frac{\partial \widehat{\psi}^L}{\partial t}, \widehat{\varphi} \right\rangle_{\widehat{\mathbb{X}}} dt + \int_0^T \int_{\Omega \times D} M \left[ \varepsilon \nabla_x \widehat{\psi}^L - \underline{u}^L \widehat{\psi}^L \right] \cdot \nabla_x \widehat{\varphi} dq dx dt \\ & + \int_0^T \int_{\Omega \times D} \left[ \frac{M}{2\lambda} \nabla_q \widehat{\psi}^L - (\nabla_x \underline{u}^L) q \beta^L(\widehat{\psi}^L) \right] \cdot \nabla_q \widehat{\varphi} dq dx dt = 0 \quad \forall \widehat{\varphi} \in L^{\frac{4}{4-d}}(0, T; \widehat{\mathbb{X}}). \end{aligned} \quad (6.3.16b)$$

**Remark 6.3.1** If  $d = 2$ , then  $\underline{u}^L \in C([0, T]; \underline{\mathbb{H}})$  (cf. Lemma 1.2 on p. 176 of Temam [119]), whereas if  $d = 3$ , then  $\underline{u}^L$  is weakly continuous only as a mapping from  $[0, T]$  into  $\underline{\mathbb{H}}$  (similarly as in Theorem 3.1 on p. 191 in Temam [119]). It is in the latter, weaker sense that the imposition of the initial condition to the  $\underline{u}^L$ -equation will be understood for  $d = 2, 3$ : that is,  $\lim_{t \rightarrow 0^+} \int_{\Omega} (\underline{u}^L(\underline{x}, t) - \underline{u}^0(\underline{x})) \cdot \underline{v}(\underline{x}) d\underline{x} = 0$  for all  $\underline{v} \in \underline{\mathbb{H}}$ . Similarly, for the initial conditions of the  $\widehat{\psi}^L$ -equation for  $d = 2, 3$ :  $\lim_{t \rightarrow 0^+} \int_{\Omega \times D} M (\widehat{\psi}^L(\underline{x}, \underline{q}, t) - \widehat{\psi}^0(\underline{x}, \underline{q})) \widehat{\varphi}(\underline{x}, \underline{q}) dq d\underline{x} = 0$  for all  $\widehat{\varphi} \in L_M^2(\Omega \times D)$ .  $\diamond$

**Remark 6.3.2** Since the test functions in  $\underline{\mathbb{V}}$  are divergence-free, the pressure has been eliminated in (6.3.16a,b); it can be recovered in a very weak sense following the same procedure as for the incompressible Navier–Stokes equations discussed on p.208 in Temam [119]; *i.e.*, one obtains that  $\int_0^t p^L(\cdot, s) ds \in C([0, T]; L^2(\Omega))$ .  $\diamond$

## 6.4 Finite element approximation

Let us denote the measure of a bounded open region  $\omega \subset \mathbb{R}^d$  by  $\underline{m}(\omega)$ . We make the following assumption on  $\Omega$  and the partitions of  $\Omega$  and  $D$ .

(A1) For ease of exposition, we shall assume that  $\Omega$  is a convex polytope. Let  $\{\mathcal{T}_h^x\}_{h>0}$  be a quasiuniform family of partitions of  $\Omega$  into disjoint open nonobtuse simplices  $\kappa_x$ , so that

$$\overline{\Omega} \equiv \bigcup_{\kappa_x \in \mathcal{T}_h^x} \overline{\kappa_x} \quad \text{with } h_{\kappa_x} := \text{diam}(\kappa_x), \quad h_x := \max_{\kappa_x \in \mathcal{T}_h^x} h_{\kappa_x} \leq \text{diam}(\Omega) h \quad \text{and} \quad \underline{m}(\kappa_x) \geq C h^d.$$

Let  $\{\mathcal{T}_h^q\}_{h>0}$  be a quasiuniform family of partitions of  $D \equiv B(\underline{0}, r_D)$ ,  $r_D \in \mathbb{R}_{>0}$ , into disjoint open nonobtuse simplices  $\kappa_q$ , with possibly one curved edge,  $d = 2$ , or face,  $d = 3$ , on  $\partial D$ ; so

that

$$\bar{D} \equiv \bigcup_{\kappa_q \in \mathcal{T}_h^q} \bar{\kappa}_q \quad \text{with} \quad h_{\kappa_q} := \text{diam}(\kappa_q), \quad h_q := \max_{\kappa_q \in \mathcal{T}_h^q} h_{\kappa_q} \leq \text{diam}(D)h \quad \text{and} \quad \underline{m}(\kappa_q) \geq Ch^d.$$

A “simplex”  $\kappa_q$  with a curved edge/face is nonobtuse if it is convex and the enclosed simplex with the same vertices is nonobtuse, in the sense that all of its dihedral angles are  $\leq \pi/2$ . It follows from the above that

$$\frac{h_x}{h_q} + \frac{h_q}{h_x} \leq C \quad \text{as} \quad h \rightarrow 0_+. \tag{6.4.1}$$

We note that such nonobtuse simplicial partitions of  $\Omega$  and  $D$  are easily constructed in the case  $d = 2$ . For the construction of nonobtuse three-dimensional simplicial partitions we refer to the papers of Korotov and Krížek, [74] and [75], for example; the reader should note, however, that in [74] the authors use the term *acute* when they mean *nonobtuse*. Elsewhere in the computational geometry literature the term *acute* is reserved for a simplicial partition where all dihedral angles of any simplex in the partition are  $< \pi/2$ , which is a more restrictive requirement (especially in the case of  $d = 3$ ) than what we assume here; see, for example, the articles of Brandts, Korotov, Krížek and Šolc [25], Eppstein, Sullivan and Üngör [47], and Itoh and Zamfirescu [59], and references therein. Nonobtuse simplicial partitions are sometimes also called *weakly acute* (cf. [110], p. 363).

We adopt the standard notation for  $L^2$  inner products, with  $\eta_i \in L^2(\Omega \times D)$ ,

$$(\eta_1, \eta_2)_\Omega := \int_\Omega \eta_1 \eta_2 \, dx \quad \forall \eta_i \in L^2(\Omega) \quad \text{and} \quad (\eta_1, \eta_2)_{\Omega \times D} := \int_{\Omega \times D} \eta_1 \eta_2 \, dq \, dx, \tag{6.4.2}$$

which are naturally extended to vector/matrix functions.

Let  $\mathbb{P}_k^x$  and  $\mathbb{P}_k^q$  denote polynomials of degree less than or equal to  $k$  in  $x$  and  $q$ , respectively. We approximate the pressure and velocity with the lowest order Taylor–Hood element; that is,

$$\mathbf{R}_h := \{ \eta_h \in C(\bar{\Omega}) : \eta_h|_{\kappa_x} \in \mathbb{P}_1^x \quad \forall \kappa_x \in \mathcal{T}_h^x \}, \tag{6.4.3a}$$

$$\mathbb{W}_h := \{ \underset{\sim}{w}_h \in [C(\bar{\Omega})]^d : \underset{\sim}{w}_h|_{\kappa_x} \in [\mathbb{P}_2^x]^d \quad \forall \kappa_x \in \mathcal{T}_h^x \text{ and } \underset{\sim}{w}_h = 0 \text{ on } \partial\Omega \} \subset [\mathbf{H}_0^1(\Omega)]^d, \tag{6.4.3b}$$

$$\mathbb{V}_h := \{ \underset{\sim}{v}_h \in \mathbb{W}_h : (\nabla_x \cdot \underset{\sim}{v}_h, \eta_h)_\Omega = 0 \quad \forall \eta_h \in \mathbf{R}_h \}. \tag{6.4.3c}$$

It is well-known that  $\mathbf{R}_h$  and  $\mathbb{W}_h$  satisfy the inf-sup condition

$$\sup_{\underset{\sim}{w}_h \in \mathbb{W}_h} \frac{(\nabla_x \cdot \underset{\sim}{w}_h, r_h)_\Omega}{\|\underset{\sim}{w}_h\|_{\mathbf{H}^1(\Omega)}} \geq C_0 \|r_h\|_{L^2(\Omega)} \quad \forall r_h \in \mathbf{R}_h, \tag{6.4.4}$$

see e.g. [27, §VI.6]. Hence for all  $v \in \mathbb{V}$ , there exists a sequence  $\{\underset{\sim}{v}_h\}_{h>0}$ , with  $\underset{\sim}{v}_h \in \mathbb{V}_h$ , such that

$$\lim_{h \rightarrow 0_+} \|\underset{\sim}{v} - \underset{\sim}{v}_h\|_{\mathbf{H}^1(\Omega)} = 0. \tag{6.4.5}$$

We require the  $L^2$  projector  $Q_h : \mathbb{V} \rightarrow \mathbb{V}_h$  defined by

$$(\underset{\sim}{v} - Q_h \underset{\sim}{v}, \underset{\sim}{w}_h)_\Omega = 0 \quad \forall \underset{\sim}{w}_h \in \mathbb{V}_h. \tag{6.4.6}$$

We note that the convexity of  $\Omega$  and the quasiuniformity of  $\{\mathcal{T}_h^x\}_{h>0}$  imply that  $\mathcal{Q}_h$  is uniformly  $H^1(\Omega)$  stable; that is,

$$\|\mathcal{Q}_h v\|_{H^1(\Omega)} \leq C \|v\|_{H^1(\Omega)} \quad \forall v \in \mathcal{V}, \quad (6.4.7)$$

see [56].

For the approximation of the advection term in the Navier–Stokes equation we note that, for all  $v \in \mathcal{V}$  and  $w, z \in H^1(\Omega)$ , we have that

$$((v \cdot \nabla_x)w, z)_\Omega \equiv \frac{1}{2} \left[ ((v \cdot \nabla_x)w, z)_\Omega - ((v \cdot \nabla_x)z, w)_\Omega \right]. \quad (6.4.8)$$

In addition, the choice  $w = z$  leads to both sides of (6.4.8) vanishing. Obviously as  $\mathcal{V}_h \not\subset \mathcal{V}$ , the discrete analogue of the above does not hold; that is, it is *not* generally true that, for all  $v_h \in \mathcal{V}_h$ ,  $w_h, z_h \in \mathcal{W}_h$ ,

$$((v_h \cdot \nabla_x)w_h, z_h)_\Omega \equiv \frac{1}{2} \left[ ((v_h \cdot \nabla_x)w_h, z_h)_\Omega - ((v_h \cdot \nabla_x)z_h, w_h)_\Omega \right]. \quad (6.4.9)$$

We note that the right-hand side of (6.4.9) vanishes if  $w_h = z_h$ , which is not necessarily true for the left-hand side. Hence, we use the right-hand side form of (6.4.9) for the approximation of the advection term in the Navier–Stokes equation.

To approximate  $\hat{X}$ , we first introduce

$$\hat{X}_h^x := \{\hat{\varphi}_h^x \in C(\overline{\Omega}) : \hat{\varphi}_h^x|_{\kappa_x} \in \mathbb{P}_1^x \quad \forall \kappa_x \in \mathcal{T}_h^x\} \subset W^{1,\infty}(\Omega), \quad (6.4.10a)$$

$$\hat{X}_h^q := \{\hat{\varphi}_h^q \in C(\overline{D}) : \hat{\varphi}_h^q|_{\kappa_q} \in \mathbb{P}_1^q \quad \forall \kappa_q \in \mathcal{T}_h^q\} \subset W^{1,\infty}(D). \quad (6.4.10b)$$

We then set

$$\hat{X}_h := \hat{X}_h^x \otimes \hat{X}_h^q \subset \hat{X}. \quad (6.4.11)$$

We note from (6.4.3a,d), (6.4.10a) and (6.4.11) that, for any  $v_h \in \mathcal{V}_h$  and any  $q \in \overline{D}$ ,

$$(\nabla_x \cdot v_h, \hat{\varphi}_h(\cdot, q))_\Omega = 0 \quad \forall \hat{\varphi}_h \in \hat{X}_h. \quad (6.4.12)$$

We note that for (6.4.12) to hold in general, we require that  $\hat{X}_h^x \subseteq R_h$ .

We introduce the interpolation operators  $\pi_h^x : C(\overline{\Omega}) \rightarrow \hat{X}_h^x$  and  $\pi_h^q : C(\overline{D}) \rightarrow \hat{X}_h^q$  such that

$$\pi_h^x \hat{\varphi}_h^x(P_i^x) = \hat{\varphi}_h^x(P_i^x), \quad i = 1 \rightarrow I^x, \quad \text{and} \quad \pi_h^q \hat{\varphi}_h^q(P_i^q) = \hat{\varphi}_h^q(P_i^q), \quad i = 1 \rightarrow I^q, \quad (6.4.13)$$

where  $\{P_i^x\}_{i=1}^{I^x}$  and  $\{P_i^q\}_{i=1}^{I^q}$  are the nodes (vertices) of  $\mathcal{T}_h^x$  and  $\mathcal{T}_h^q$ , respectively. The associated basis functions are

$$\chi_i^x \in \hat{X}_h^x \quad \text{such that} \quad \chi_i^x(P_j^x) = \delta_{ij} \quad \text{for } i, j = 1 \rightarrow I^x, \quad (6.4.14a)$$

$$\text{and} \quad \chi_i^q \in \hat{X}_h^q \quad \text{such that} \quad \chi_i^q(P_j^q) = \delta_{ij} \quad \text{for } i, j = 1 \rightarrow I^q. \quad (6.4.14b)$$

We introduce also  $\pi_h : C(\overline{\Omega \times D}) \rightarrow \hat{X}_h$  such that

$$(\pi_h \hat{\varphi})(P_i^x, P_j^q) = \hat{\varphi}(P_i^x, P_j^q) \quad \text{for } i = 1 \rightarrow I^x, \quad j = 1 \rightarrow I^q. \quad (6.4.15)$$

Of course, we have that  $\pi_h \equiv \pi_h^x \pi_h^q \equiv \pi_h^q \pi_h^x$ . The vector versions of the above interpolation operators are

$$\tilde{\pi}_h^x : [C(\overline{\Omega})]^d \rightarrow [\hat{X}_h^x]^d, \quad \tilde{\pi}_h^q : [C(\overline{D})]^d \rightarrow [\hat{X}_h^q]^d \quad \text{and} \quad \tilde{\pi}_h : [C(\overline{\Omega \times D})]^d \rightarrow [\hat{X}_h]^d. \quad (6.4.16)$$

We require also the local interpolation operators

$$\begin{aligned} \pi_{h,\kappa_x}^x &\equiv \pi_h^x |_{\kappa_x}, & \pi_{h,\kappa_q}^q &\equiv \pi_h^q |_{\kappa_q}, & \pi_{h,\kappa_x \times \kappa_q} &\equiv \pi_h |_{\kappa_x \times \kappa_q}, & \tilde{\pi}_{h,\kappa_x}^x &\equiv \tilde{\pi}_h^x |_{\kappa_x}, \\ \tilde{\pi}_{h,\kappa_q}^q &\equiv \tilde{\pi}_h^q |_{\kappa_q} & \text{and} & & \tilde{\pi}_{h,\kappa_x \times \kappa_q} &\equiv \tilde{\pi}_h |_{\kappa_x \times \kappa_q} & \forall \kappa_x \in \mathcal{T}_h^x, & \forall \kappa_q \in \mathcal{T}_h^q. \end{aligned} \quad (6.4.17)$$

For any  $\hat{\varphi}_h \in \hat{X}_h$ , there exist  $[\Xi_\delta^q(\hat{\varphi}_h)](\underline{x}, \underline{q})$ ,  $[\Xi_\delta^x(\hat{\varphi}_h)](\underline{x}, \underline{q}) \in \mathbb{R}^{d \times d}$  for a.e.  $(\underline{x}, \underline{q}) \in \Omega \times D$  such that on  $\kappa_x \times \kappa_q$ , for all  $\kappa_x \in \mathcal{T}_h^x$ ,  $\kappa_q \in \mathcal{T}_h^q$ ,

$$\Xi_\delta^x(\hat{\varphi}_h) \in [\mathbb{P}_1^q]^{d \times d} \quad \text{and} \quad \tilde{\pi}_{h,\kappa_x \times \kappa_q} \left[ \Xi_\delta^x(\hat{\varphi}_h) \tilde{\nabla}_x (\pi_h [ [\mathcal{F}_\delta^L]'(\hat{\varphi}_h) ]) \right] = \tilde{\nabla}_x \hat{\varphi}_h; \quad (6.4.18a)$$

$$\Xi_\delta^q(\hat{\varphi}_h) \in [\mathbb{P}_1^x]^{d \times d} \quad \text{and} \quad \tilde{\pi}_{h,\kappa_x \times \kappa_q} \left[ \Xi_\delta^q(\hat{\varphi}_h) \tilde{\nabla}_q (\pi_h [ [\mathcal{F}_\delta^L]'(\hat{\varphi}_h) ]) \right] = \tilde{\nabla}_q \hat{\varphi}_h. \quad (6.4.18b)$$

Hence (6.4.18a,b) are discrete analogues of the relations (6.2.18). We now give the construction of  $\Xi_\delta^x(\cdot)$  and  $\Xi_\delta^q(\cdot)$ . Let  $\{\underline{e}_i\}_{i=1}^d$  be the orthonormal vectors in  $\mathbb{R}^d$ , such that the  $j^{\text{th}}$  component of  $\underline{e}_i$  is  $\delta_{ij}$ ,  $i, j = 1 \rightarrow d$ . Let  $\tilde{\kappa}$  be the standard reference simplex in  $\mathbb{R}^d$  with vertices  $\{\tilde{P}_i\}_{i=0}^d$ , where  $\tilde{P}_0$  is the origin and  $\tilde{P}_i = \underline{e}_i$ ,  $i = 1 \rightarrow d$ . Given  $\hat{\varphi}_h \in \hat{X}_h$ ,  $\kappa_x \in \mathcal{T}_h^x$  with vertices  $\{\underline{P}_{i_j}^x\}_{j=0}^d$  and  $\kappa_q \in \mathcal{T}_h^q$  with vertices  $\{\underline{P}_{i_j}^q\}_{j=0}^d$ , then for a fixed vertex  $\underline{P}_{i_k}^q$  of  $\kappa_q$ , let  $\Lambda_\delta^x(\underline{P}_{i_k}^q) \in \mathbb{R}^{d \times d}$  be diagonal with entries

$$[\Lambda_\delta^x(\underline{P}_{i_k}^q)]_{jj} = \begin{cases} \frac{\hat{\varphi}_h(\underline{P}_{i_j}^x, \underline{P}_{i_k}^q) - \hat{\varphi}_h(\underline{P}_{i_0}^x, \underline{P}_{i_k}^q)}{[\mathcal{F}_\delta^L]'(\hat{\varphi}_h(\underline{P}_{i_j}^x, \underline{P}_{i_k}^q)) - [\mathcal{F}_\delta^L]'(\hat{\varphi}_h(\underline{P}_{i_0}^x, \underline{P}_{i_k}^q))} & \text{if } \hat{\varphi}_h(\underline{P}_{i_j}^x, \underline{P}_{i_k}^q) \neq \hat{\varphi}_h(\underline{P}_{i_0}^x, \underline{P}_{i_k}^q), \\ \frac{1}{[\mathcal{F}_\delta^L]''(\hat{\varphi}_h(\underline{P}_{i_j}^x, \underline{P}_{i_k}^q))} = \beta_\delta^L(\hat{\varphi}_h(\underline{P}_{i_j}^x, \underline{P}_{i_k}^q)) & \text{if } \hat{\varphi}_h(\underline{P}_{i_j}^x, \underline{P}_{i_k}^q) = \hat{\varphi}_h(\underline{P}_{i_0}^x, \underline{P}_{i_k}^q), \end{cases} \quad j = 1 \rightarrow d. \quad (6.4.19)$$

Let  $B_{\kappa_x} \in \mathbb{R}^{d \times d}$  be such that the linear mapping  $\mathcal{B}_{\kappa_x} : \underline{y} \in \mathbb{R}^d \mapsto \underline{P}_{i_0}^x + B_{\kappa_x} \underline{y}$  maps the vertex  $\tilde{P}_j$  to  $\underline{P}_{i_j}^x$ ,  $j = 0 \rightarrow d$ , and hence  $\tilde{\kappa}$  to  $\kappa_x$ . For any  $\hat{\varphi}_h \in \hat{X}_h$ , let  $\hat{\varphi}_{h,y}^x(\underline{x}) \equiv \hat{\varphi}_h(\mathcal{B}_{\kappa_x} \underline{y})$  for all  $\underline{y} \in \tilde{\kappa}$ . Hence it follows that

$$\tilde{\nabla}_x \hat{\varphi}_h^x = [B_{\kappa_x}^T]^{-1} \tilde{\nabla}_y \hat{\varphi}_{h,y}^x. \quad (6.4.20)$$

Therefore, for  $k = 0 \rightarrow d$ ,

$$\Xi_\delta^x(\underline{P}_{i_k}^q) = [B_{\kappa_x}^T]^{-1} \Lambda_\delta^x(\underline{P}_{i_k}^q) B_{\kappa_x}^T \quad (6.4.21)$$

is such that

$$\Xi_\delta^x(\underline{P}_{i_k}^q) \tilde{\nabla}_x \pi_h [ [\mathcal{F}_\delta^L]'(\hat{\varphi}_h) ](\underline{x}, \underline{P}_{i_k}^q) = \tilde{\nabla}_x \hat{\varphi}_h(\underline{x}, \underline{P}_{i_k}^q) \quad \forall \underline{x} \in \kappa_x. \quad (6.4.22)$$

Finally, on recalling (6.4.14b), we set

$$\Xi_{\delta}^x(x, q) = \sum_{k=0}^d \Xi_{\delta}^x(P_{i_k}^q) \chi_{i_k}^q(q) \quad \forall x \in \kappa_x, \quad \forall q \in \kappa_q. \quad (6.4.23)$$

Hence  $\Xi_{\delta}^x$  satisfies (6.4.18a). A similar construction yields  $\Xi_{\delta}^q$  satisfying (6.4.18b). The only difference is for those  $\kappa_q$  with a curved side or face, the corresponding linear mapping  $\mathcal{B}_{\kappa_q}$  maps  $\tilde{\kappa}$  to the enclosed simplex with the same vertices as  $\kappa_q$ .

As  $\mathcal{T}_x^h, \mathcal{T}_q^h$  are quasiuniform partitions, we have from (6.4.23), (6.4.21) and (6.4.19), and their  $\Xi_{\delta}^q$  counterparts that, for all  $\hat{\varphi}_h \in \hat{X}_h$ ,

$$\|\Xi_{\delta}^x(\hat{\varphi}_h)\|_{L^\infty(\Omega \times D)}^2 + \|\Xi_{\delta}^q(\hat{\varphi}_h)\|_{L^\infty(\Omega \times D)}^2 \leq C L^2. \quad (6.4.24)$$

We note that the construction of  $\Xi_{\delta}^x(\cdot)$  and  $\Xi_{\delta}^q(\cdot)$  satisfying (6.4.18a,b) is an extension of ideas used in e.g. [53] and [9] for the finite element approximation of fourth-order degenerate nonlinear parabolic equations, such as the thin film equation.

As the partitions  $\mathcal{T}_x^h$  and  $\mathcal{T}_q^h$  are nonobtuse, we have that

$$\nabla_x \chi_i^x \cdot \nabla_x \chi_j^x \leq 0 \quad \text{on } \kappa_x \quad i \neq j, \quad i, j = 1 \rightarrow I^x, \quad \forall \kappa_x \in \mathcal{T}_x^h; \quad (6.4.25a)$$

$$\text{and} \quad \nabla_q \chi_i^q \cdot \nabla_q \chi_j^q \leq 0 \quad \text{on } \kappa_q \quad i \neq j, \quad i, j = 1 \rightarrow I^q, \quad \forall \kappa_q \in \mathcal{T}_q^h. \quad (6.4.25b)$$

It follows from (6.4.25a,b) and the convexity of  $\mathcal{F}_d^L$  that, for all  $\kappa_x \in \mathcal{T}_x^h, \kappa_q \in \mathcal{T}_q^h$  and for all  $\hat{\varphi} \in \hat{X}_h$ ,

$$\begin{aligned} & \delta \int_{\kappa_x \times \kappa_q} M \pi_{h, \kappa_x \times \kappa_q} \left[ \left| \nabla_x (\pi_{h, \kappa_x \times \kappa_q} [\mathcal{F}_\delta^L]'(\hat{\varphi}_h)) \right|^2 \right] dq dx \\ & \leq \int_{\kappa_x \times \kappa_q} M \pi_{h, \kappa_x \times \kappa_q} \left[ \nabla_x \hat{\varphi}_h \cdot \nabla_x (\pi_{h, \kappa_x \times \kappa_q} [\mathcal{F}_\delta^L]'(\hat{\varphi}_h)) \right] dq dx; \end{aligned} \quad (6.4.26a)$$

$$\begin{aligned} \text{and} \quad & \delta \int_{\kappa_x \times \kappa_q} M \pi_{h, \kappa_x \times \kappa_q} \left[ \left| \nabla_q (\pi_{h, \kappa_x \times \kappa_q} [\mathcal{F}_\delta^L]'(\hat{\varphi}_h)) \right|^2 \right] dq dx \\ & \leq \int_{\kappa_x \times \kappa_q} M \pi_{h, \kappa_x \times \kappa_q} \left[ \nabla_q \hat{\varphi}_h \cdot \nabla_q (\pi_{h, \kappa_x \times \kappa_q} [\mathcal{F}_\delta^L]'(\hat{\varphi}_h)) \right] dq dx. \end{aligned} \quad (6.4.26b)$$

Let  $0 = t_0 < t_1 < \dots < t_{N-1} < t_N = T$  be a partition of the time interval  $[0, T]$  into time steps  $\Delta t^n = t^n - t^{n-1}$ ,  $n = 1 \rightarrow N$ . We set  $\Delta t = \max_{n=1 \rightarrow N} \Delta t^n$ . We make the following assumptions on the time steps  $\{\Delta t^n\}_{n=1}^N$  and the discrete initial data.

**(A2)** We assume that there exists  $C \in \mathbb{R}_{>0}$  such that

$$\Delta t^n \leq C \Delta t^{n-1}, \quad n = 2 \rightarrow N, \quad \text{as } \Delta t \rightarrow 0_+. \quad (6.4.27)$$

With  $\Delta t_1$  and  $C$  as above, let  $\Delta t_0 \in \mathbb{R}_{>0}$  be such that  $\Delta t_1 \leq C \Delta t_0$ . Given initial data satisfying (6.3.5a), we choose  $u_h^0 \in \mathcal{V}_h$  and  $\hat{\psi}_h^0 \in \hat{X}_h$  such that

$$(u_h^0, v_h)_\Omega + \Delta t_0 (\nabla_x u_h^0, \nabla_x v_h)_\Omega = (u_h^0, v_h)_\Omega \quad \forall v_h \in \mathcal{V}_h, \quad (6.4.28a)$$

$$(M, \pi_h[\hat{\psi}_h^0 \hat{\varphi}_h])_{\Omega \times D} = (M \hat{\psi}_h^0, \hat{\varphi}_h)_{\Omega \times D} \quad \forall \hat{\varphi}_h \in \hat{X}_h. \quad (6.4.28b)$$

It follows from (6.4.28a,b) and (6.3.5a) that

$$\int_{\Omega} \left[ |\tilde{u}_h^0|^2 + \Delta t_0 |\nabla_x \tilde{u}_h^0|^2 \right] dx \leq C \quad \text{and} \quad 0 \leq \hat{\psi}_h^0 \leq L. \quad (6.4.29)$$

We set

$$\tilde{f}^n(\cdot) := \frac{1}{\Delta t^n} \int_{t^{n-1}}^{t^n} f(\cdot, t) dt \in \tilde{\mathbf{H}}^{-1}(\Omega). \quad (6.4.30)$$

It is easily deduced from (6.3.5b) and (6.4.30) that

$$\sum_{n=1}^N \Delta t^n \|\tilde{f}^n\|_{\tilde{\mathbf{H}}^{-1}(\Omega)}^r \leq \int_0^T \|f\|_{\tilde{\mathbf{H}}^{-1}(\Omega)}^r dt \leq C \quad \text{for any } r \in [1, 2], \quad (6.4.31a)$$

$$\text{and} \quad \tilde{f}^{\Delta t, +} \rightarrow \tilde{f} \quad \text{strongly in } L^2(0, T; \tilde{\mathbf{V}}') \text{ as } \Delta t \rightarrow 0_+, \quad (6.4.31b)$$

where  $\tilde{f}^{\Delta t, +}(\cdot, t) := \tilde{f}^n(\cdot)$  for  $t \in (t^{n-1}, t^n]$ ,  $n = 1 \rightarrow N$ .

Our numerical approximation of  $(\mathbf{P}_\delta^L)$  is then defined as follows.

$(\mathbf{P}_\delta^{h, \Delta t})$  For  $n = 1 \rightarrow N$ , given  $\{\tilde{u}_{\delta, h}^{n-1}, \hat{\psi}_{\delta, h}^{n-1}\} \in \tilde{\mathcal{V}}_h \times \hat{\mathbf{X}}_h$ , find  $\{\tilde{u}_{\delta, h}^n, \hat{\psi}_{\delta, h}^n\} \in \tilde{\mathcal{V}}_h \times \hat{\mathbf{X}}_h$  such that

$$\begin{aligned} & \left( \frac{\tilde{u}_{\delta, h}^n - \tilde{u}_{\delta, h}^{n-1}}{\Delta t^n}, w_h \right)_{\Omega} + \nu (\nabla_x \tilde{u}_{\delta, h}^n, \nabla_x w_h)_{\Omega} \\ & + \frac{1}{2} \left[ ((\tilde{u}_{\delta, h}^{n-1} \cdot \nabla_x) \tilde{u}_{\delta, h}^n, w_h)_{\Omega} - ((\tilde{u}_{\delta, h}^{n-1} \cdot \nabla_x) w_h, \tilde{u}_{\delta, h}^n)_{\Omega} \right] \\ & = \langle \tilde{f}^n, w_h \rangle_{\tilde{\mathbf{H}}_0^1(\Omega)} - k_B \mathcal{T} (C(M \hat{\psi}_{\delta, h}^n), \nabla_x w_h)_{\Omega} \quad \forall w_h \in \tilde{\mathcal{V}}_h, \end{aligned} \quad (6.4.32a)$$

$$\begin{aligned} & \left( M, \pi_h \left[ \frac{\hat{\psi}_{\delta, h}^n - \hat{\psi}_{\delta, h}^{n-1}}{\Delta t^n} \hat{\varphi}_h + \varepsilon \nabla_x \hat{\psi}_{\delta, h}^n \cdot \nabla_x \hat{\varphi}_h + \frac{1}{2\lambda} \nabla_q \hat{\psi}_{\delta, h}^n \cdot \nabla_q \hat{\varphi}_h \right] \right)_{\Omega \times D} \\ & = \left( M (\nabla_x \tilde{u}_{\delta, h}^n) q, \pi_h \left[ \Xi_{\delta}^q(\hat{\psi}_{\delta, h}^n) \nabla_q \hat{\varphi}_h \right] \right)_{\Omega \times D} + \left( M \tilde{u}_{\delta, h}^n, \pi_h \left[ \Xi_{\delta}^x(\hat{\psi}_{\delta, h}^n) \nabla_x \hat{\varphi}_h \right] \right)_{\Omega \times D} \\ & \quad \forall \hat{\varphi}_h \in \hat{\mathbf{X}}_h; \end{aligned} \quad (6.4.32b)$$

where for ease of notation, we write  $\pi_h$  and  $\tilde{\pi}_h$  in (6.4.32b) whereas it should really be  $\pi_{h, \kappa_x \times \kappa_q}$  and  $\tilde{\pi}_{h, \kappa_x \times \kappa_q}$ , respectively, on each  $\kappa_x \times \kappa_q$  of  $\Omega \times D$ . In addition, we have suppressed the dependence of the solution  $\{\tilde{u}_{\delta, h}^n, \hat{\psi}_{\delta, h}^n\}$  on  $L$  through the dependence of  $\Xi_{\delta}^x$  and  $\Xi_{\delta}^q$  on  $\mathcal{F}_\delta^L$ . This is because we will not be passing to the limit  $L \rightarrow \infty$ , but only to the limit  $\delta \rightarrow 0_+$  in addition to letting the discretization parameters  $h, \Delta t \rightarrow 0_+$ .

We note that the approximations  $\tilde{u}_{\delta, h}^n$  and  $\hat{\psi}_{\delta, h}^n$  at time level  $t^n$  to the velocity field and the scaled probability distribution satisfy a coupled nonlinear system, (6.4.32a,b). We will show existence of a solution to (6.4.32a,b) below, see Theorem 6.4.2, via a Brouwer fixed point theorem. First, assuming existence, we show that  $(\mathbf{P}_\delta^{h, \Delta t})$  satisfies a discrete analogue of the energy equality (6.2.21). For all the following lemmas and theorems we assume the assumptions (A1) and (A2) hold.

**Lemma 6.4.1** For  $n = 1 \rightarrow N$ , a solution  $\{\underline{u}_{\delta,h}^n, \hat{\psi}_{\delta,h}^n\} \in \underline{\mathcal{Y}}_h \times \hat{\mathcal{X}}_h$  of (6.4.32a,b), if it exists, satisfies

$$\begin{aligned}
& \frac{1}{2} \left[ \|\underline{u}_{\delta,h}^n\|_{\mathbb{L}^2(\Omega)}^2 + \|\underline{u}_{\delta,h}^n - \underline{u}_{\delta,h}^{n-1}\|_{\mathbb{L}^2(\Omega)}^2 \right] + k_B \mathcal{T} (M, \pi_h[\mathcal{F}_\delta^L(\hat{\psi}_{\delta,h}^n)])_{\Omega \times D} + \Delta t^n \nu \|\underline{\nabla}_x \underline{u}_{\delta,h}^n\|_{\mathbb{L}^2(\Omega)}^2 \\
& + \Delta t^n k_B \mathcal{T} \left( M, \pi_h \left[ \varepsilon \underline{\nabla}_x \hat{\psi}_{\delta,h}^n \cdot \underline{\nabla}_x (\pi_h[[\mathcal{F}_\delta^L]'(\hat{\psi}_{\delta,h}^n)]) + \frac{1}{2\lambda} \underline{\nabla}_q \hat{\psi}_{\delta,h}^n \cdot \underline{\nabla}_q (\pi_h[[\mathcal{F}_\delta^L]'(\hat{\psi}_{\delta,h}^n)]) \right] \right)_{\Omega \times D} \\
& \leq \frac{1}{2} \|\underline{u}_{\delta,h}^{n-1}\|_{\mathbb{L}^2(\Omega)}^2 + k_B \mathcal{T} (M, \pi_h[\mathcal{F}_\delta^L(\hat{\psi}_{\delta,h}^{n-1})])_{\Omega \times D} + \Delta t^n \langle f^n, \underline{u}_{\delta,h}^n \rangle_{\mathbb{H}_0^1(\Omega)} \\
& \leq \frac{1}{2} \|\underline{u}_{\delta,h}^{n-1}\|_{\mathbb{L}^2(\Omega)}^2 + k_B \mathcal{T} (M, \pi_h[\mathcal{F}_\delta^L(\hat{\psi}_{\delta,h}^{n-1})])_{\Omega \times D} \\
& \quad + \Delta t^n \left[ \frac{\nu}{2} \|\underline{\nabla}_x \underline{u}_{\delta,h}^n\|_{\mathbb{L}^2(\Omega)}^2 + C \|f^n\|_{\mathbb{H}^{-1}(\Omega)}^2 \right]. \tag{6.4.33}
\end{aligned}$$

**Proof.** On choosing  $w_h = \underline{u}_{\delta,h}^n$  in (6.4.32a), it follows that

$$\begin{aligned}
& \frac{1}{2} \int_{\Omega} \left[ |\underline{u}_{\delta,h}^n|^2 + |\underline{u}_{\delta,h}^n - \underline{u}_{\delta,h}^{n-1}|^2 - |\underline{u}_{\delta,h}^{n-1}|^2 \right] d\underline{x} + \Delta t^n \nu \int_{\Omega} |\underline{\nabla}_x \underline{u}_{\delta,h}^n|^2 d\underline{x} \\
& = \Delta t^n \left[ \langle f^n, \underline{u}_{\delta,h}^n \rangle_{\mathbb{H}_0^1(\Omega)} - k_B \mathcal{T} (C(M \hat{\psi}_{\delta,h}^n), \underline{\nabla}_x \underline{u}_{\delta,h}^n)_{\Omega} \right], \tag{6.4.34}
\end{aligned}$$

where we have noted the simple identity

$$2(s_1 - s_2)s_1 = s_1^2 + (s_1 - s_2)^2 - s_2^2 \quad \forall s_1, s_2 \in \mathbb{R}. \tag{6.4.35}$$

On choosing  $\hat{\varphi}_h = \pi_h[[\mathcal{F}_\delta^L]'(\hat{\psi}_{\delta,h}^n)]$  in (6.4.32b), and noting the convexity of  $\mathcal{F}_\delta^L$ , (6.4.18a,b), (6.2.5), (6.4.12) and (6.2.3), we have that

$$\begin{aligned}
& (M, \pi_h[\mathcal{F}_\delta^L(\hat{\psi}_{\delta,h}^n) - \mathcal{F}_\delta^L(\hat{\psi}_{\delta,h}^{n-1})])_{\Omega \times D} \\
& + \Delta t^n \left( M, \pi_h \left[ \varepsilon \underline{\nabla}_x \hat{\psi}_{\delta,h}^n \cdot \underline{\nabla}_x (\pi_h[[\mathcal{F}_\delta^L]'(\hat{\psi}_{\delta,h}^n)]) + \frac{1}{2\lambda} \underline{\nabla}_q \hat{\psi}_{\delta,h}^n \cdot \underline{\nabla}_q (\pi_h[[\mathcal{F}_\delta^L]'(\hat{\psi}_{\delta,h}^n)]) \right] \right)_{\Omega \times D} \\
& \leq (M (\underline{\nabla}_x \underline{u}_{\delta,h}^n) q, \underline{\nabla}_q \hat{\psi}_{\delta,h}^n)_{\Omega \times D} + (M \underline{u}_{\delta,h}^n, \underline{\nabla}_x \hat{\psi}_{\delta,h}^n)_{\Omega \times D} \\
& = (M U' q \cdot [(\underline{\nabla}_x \underline{u}_{\delta,h}^n) q], \hat{\psi}_{\delta,h}^n)_{\Omega \times D} - 2 (M \underline{\nabla}_x \cdot \underline{u}_{\delta,h}^n, \hat{\psi}_{\delta,h}^n)_{\Omega \times D} \\
& = (C(M \hat{\psi}_{\delta,h}^n), \underline{\nabla}_x \underline{u}_{\delta,h}^n)_{\Omega}. \tag{6.4.36}
\end{aligned}$$

Combining (6.4.34) and (6.4.36) yields the first inequality (6.4.33). The second inequality follows from using a Young's inequality and a Poincaré inequality.  $\square$

We now show using a Brouwer fixed point theorem that there exists a solution  $\{\underline{u}_{\delta,h}^n, \hat{\psi}_{\delta,h}^n\}$  at time level  $t^n$  to (6.4.32a,b).

**Theorem 6.4.2** Given  $\{\underline{u}_{\delta,h}^{n-1}, \hat{\psi}_{\delta,h}^{n-1}\} \in \underline{\mathcal{Y}}_h \times \hat{\mathcal{X}}_h$  and for any time step  $\Delta t^n > 0$ , there exists at least one solution  $\{\underline{u}_{\delta,h}^n, \hat{\psi}_{\delta,h}^n\} \in \underline{\mathcal{Y}}_h \times \hat{\mathcal{X}}_h$  to (6.4.32a,b).

**Proof.** We define the inner product,  $((\cdot, \cdot))$ , on the Hilbert space  $\underline{\mathcal{Y}}_h \times \hat{\mathcal{X}}_h$  as follows:

$$((\{u_h, \hat{\psi}_h\}, \{w_h, \hat{\varphi}_h\})) := (u_h, w_h)_{\Omega} + (M, \pi_h[\hat{\psi}_h \hat{\varphi}_h])_{\Omega \times D} \quad \forall \{u_h, \hat{\psi}_h\}, \{w_h, \hat{\varphi}_h\} \in \underline{\mathcal{Y}}_h \times \hat{\mathcal{X}}_h.$$

Given  $\{u_{\delta,h}^{n-1}, \hat{\psi}_{\delta,h}^{n-1}\} \in \mathbb{V}_h \times \hat{\mathbb{X}}_h$ , let  $\mathcal{H} : \mathbb{V}_h \times \hat{\mathbb{X}}_h \mapsto \mathbb{V}_h \times \hat{\mathbb{X}}_h$  be such that, for any  $\{u_h, \hat{\psi}_h\} \in \mathbb{V}_h \times \hat{\mathbb{X}}_h$ ,

$$\begin{aligned}
& ((\mathcal{H}(u_h, \hat{\psi}_h), \{w_h, \hat{\varphi}_h\})) \\
& := \left( \frac{u_h - u_{\delta,h}^{n-1}}{\Delta t^n}, w_h \right)_{\Omega} + \nu (\nabla_x u_h, \nabla_x w_h)_{\Omega} - \langle f^n, w_h \rangle_{\mathbb{H}_0^1(\Omega)} + k_B \mathcal{T} (C(M \hat{\psi}_h), \nabla_x w_h)_{\Omega} \\
& \quad + \frac{1}{2} \left[ ((u_{\delta,h}^{n-1} \cdot \nabla_x) u_h, w_h)_{\Omega} - ((u_{\delta,h}^{n-1} \cdot \nabla_x) w_h, u_h)_{\Omega} \right] \\
& \quad + \left( M, \pi_h \left[ \frac{\hat{\psi}_h - \hat{\psi}_{\delta,h}^{n-1}}{\Delta t^n} \hat{\varphi}_h + \varepsilon \nabla_x \hat{\psi}_h \cdot \nabla_x \hat{\varphi}_h + \frac{1}{2\lambda} \nabla_q \hat{\psi}_h \cdot \nabla_q \hat{\varphi}_h \right] \right)_{\Omega \times D} \\
& \quad - \left( M (\nabla_x u_h) q, \pi_h \left[ \Xi_{\delta}^q(\hat{\psi}_h) \nabla_q \hat{\varphi}_h \right] \right)_{\Omega \times D} - \left( M u_h, \pi_h \left[ \Xi_{\delta}^x(\hat{\psi}_h) \nabla_x \hat{\varphi}_h \right] \right)_{\Omega \times D} \\
& \quad \forall \{w_h, \hat{\varphi}_h\} \in \mathbb{V}_h \times \hat{\mathbb{X}}_h. \tag{6.4.37}
\end{aligned}$$

We note that a solution  $\{u_{\delta,h}^n, \hat{\psi}_{\delta,h}^n\}$  to (6.4.32a,b), if it exists, corresponds to a zero of  $\mathcal{H}$ ; that is,

$$((\mathcal{H}(u_{\delta,h}^n, \hat{\psi}_{\delta,h}^n), \{w_h, \hat{\varphi}_h\})) = 0 \quad \forall \{w_h, \hat{\varphi}_h\} \in \mathbb{V}_h \times \hat{\mathbb{X}}_h. \tag{6.4.38}$$

On noting the construction of  $\Xi_{\delta}^x$  and  $\Xi_{\delta}^q$ , (6.4.19)–(6.4.23), it is easily deduced that the mapping  $\mathcal{H}$  is continuous.

For any  $\{u_h, \hat{\psi}_h\} \in \mathbb{V}_h \times \hat{\mathbb{X}}_h$ , on choosing  $\{w_h, \hat{\varphi}_h\} = \{u_h, \pi_h[[\mathcal{F}_{\delta}^L]'](\hat{\psi}_h)]\}$ , we obtain analogously to (6.4.33), on noting (6.4.26a,b) and neglecting some nonnegative terms, that

$$\begin{aligned}
& ((\mathcal{H}(u_h, \hat{\psi}_h), \{u_h, \pi_h[[\mathcal{F}_{\delta}^L]'](\hat{\psi}_h)]\})) \\
& \geq \frac{1}{\Delta t^n} \left[ \frac{1}{2} \left( \|u_h\|_{L^2(\Omega)}^2 - \|u_{\delta,h}^{n-1}\|_{L^2(\Omega)}^2 \right) + k_B \mathcal{T} (M, \pi_h[\mathcal{F}_{\delta}^L(\hat{\psi}_h) - \mathcal{F}_{\delta}^L(\hat{\psi}_{\delta,h}^{n-1})])_{\Omega \times D} \right] \\
& \quad + \frac{\nu}{2} \|\nabla_x u_h\|_{L^2(\Omega)}^2 - C \|f^n\|_{\mathbb{H}^{-1}(\Omega)}^2. \tag{6.4.39}
\end{aligned}$$

Let us now assume that, for any  $\gamma \in \mathbb{R}_{>0}$ , the continuous mapping  $\mathcal{H}$  has no zero  $\{u_{\delta,h}^n, \hat{\psi}_{\delta,h}^n\}$  satisfying (6.4.38), which lies in the ball

$$\mathcal{Z}_{\gamma} := \{\{w_h, \hat{\varphi}_h\} \in \mathbb{V}_h \times \hat{\mathbb{X}}_h : |||\{w_h, \hat{\varphi}_h\}||| \leq \gamma\};$$

where

$$|||\{w_h, \hat{\varphi}_h\}||| := [((\{w_h, \hat{\varphi}_h\}, \{w_h, \hat{\varphi}_h\}))]^{\frac{1}{2}} = \left[ \|w_h\|_{L^2(\Omega)}^2 + (M, \pi_h[(\hat{\varphi}_h)^2])_{\Omega \times D} \right]^{\frac{1}{2}}.$$

Then, for such  $\gamma$ , we can define the continuous mapping  $\mathcal{E}_{\gamma} : \mathcal{Z}_{\gamma} \mapsto \mathcal{Z}_{\gamma}$  such that, for all  $\{w_h, \hat{\varphi}_h\} \in \mathcal{Z}_{\gamma}$ ,

$$\mathcal{E}_{\gamma}(w_h, \hat{\varphi}_h) := -\gamma \frac{\mathcal{H}(w_h, \hat{\varphi}_h)}{|||\mathcal{H}(w_h, \hat{\varphi}_h)|||}.$$

By the Brouwer fixed point theorem,  $\mathcal{E}_{\gamma}$  has at least one fixed point  $\{u_h^{\gamma}, \hat{\psi}_h^{\gamma}\}$  in  $\mathcal{Z}_{\gamma}$ ; hence it satisfies

$$|||\{u_h^{\gamma}, \hat{\psi}_h^{\gamma}\}||| = |||\mathcal{E}_{\gamma}(u_h^{\gamma}, \hat{\psi}_h^{\gamma})||| = \gamma. \tag{6.4.40}$$

It follows from (6.2.22) and (6.4.40) that

$$\begin{aligned} \frac{1}{2} \|\underline{u}_h^\gamma\|_{L^2(\Omega)}^2 + k_B \mathcal{T} (M, \pi_h[\mathcal{F}_\delta^L(\hat{\psi}_h^\gamma)])_{\Omega \times D} &\geq \frac{1}{2} \|\underline{u}_h^\gamma\|_{L^2(\Omega)}^2 + \frac{k_B \mathcal{T}}{4L} (M, \pi_h[(\hat{\psi}_h^\gamma)^2])_{\Omega \times D} - C(L) \\ &\geq \min \left\{ \frac{1}{2}, \frac{k_B \mathcal{T}}{4L} \right\} \|\{\underline{u}_h^\gamma, \hat{\psi}_h^\gamma\}\|^2 - C(L) \\ &= \min \left\{ \frac{1}{2}, \frac{k_B \mathcal{T}}{4L} \right\} \gamma^2 - C(L). \end{aligned} \quad (6.4.41)$$

Hence for all  $\gamma$  sufficiently large, it follows from (6.4.39) and (6.4.41) that

$$((\mathcal{H}(\underline{u}_h^\gamma, \hat{\psi}_h^\gamma), \{\underline{u}_h^\gamma, \pi_h[[\mathcal{F}_\delta^L]'(\hat{\psi}_h^\gamma)]\})) > 0. \quad (6.4.42)$$

On the other hand as  $\{\underline{u}_h^\gamma, \hat{\psi}_h^\gamma\}$  is a fixed point of  $\mathcal{E}_\gamma$ , we have that

$$((\mathcal{H}(\underline{u}_h^\gamma, \hat{\psi}_h^\gamma), \{\underline{u}_h^\gamma, \pi_h[[\mathcal{F}_\delta^L]'(\hat{\psi}_h^\gamma)]\})) = -\frac{\|\mathcal{H}(\underline{u}_h^\gamma, \hat{\psi}_h^\gamma)\|}{\gamma} \left[ \|\underline{u}_h^\gamma\|_{L^2(\Omega)}^2 + (M, \pi_h[\hat{\psi}_h^\gamma [\mathcal{F}_\delta^L]'(\hat{\psi}_h^\gamma)])_{\Omega \times D} \right]. \quad (6.4.43)$$

Similarly to (6.4.41), we have from (6.2.22) and (6.4.40) that

$$\|\underline{u}_h^\gamma\|_{L^2(\Omega)}^2 + (M, \pi_h[\hat{\psi}_h^\gamma [\mathcal{F}_\delta^L]'(\hat{\psi}_h^\gamma)])_{\Omega \times D} \geq \frac{1}{4L} \gamma^2 - C(L). \quad (6.4.44)$$

Therefore on combining (6.4.43) and (6.4.44), we have for all  $\gamma$  sufficiently large that

$$((\mathcal{H}(\underline{u}_h^\gamma, \hat{\psi}_h^\gamma), \{\underline{u}_h^\gamma, \pi_h[[\mathcal{F}_\delta^L]'(\hat{\psi}_h^\gamma)]\})) < 0, \quad (6.4.45)$$

which obviously contradicts (6.4.42). Hence the mapping  $\mathcal{H}$  has a zero in  $\mathcal{Z}_\gamma$  for  $\gamma$  sufficiently large.  $\square$

In order to establish a stability result for our approximation  $(\mathbf{P}_\delta^{h,\Delta t})$ , we need first to prove a number of auxiliary results. We note that, for all  $\kappa_x \in \mathcal{T}_x^h$  with vertices  $\{\underline{P}_{i_j}^x\}_{j=0}^d$ ,

$$\begin{aligned} |[\pi_{h,\kappa_x}^x \hat{\varphi}^x](\underline{x})|^2 &= \left| \sum_{j=0}^d \hat{\varphi}^x(\underline{P}_{i_j}^x) \chi_{i_j}^x(\underline{x}) \right|^2 \leq \sum_{j=0}^d [\hat{\varphi}^x(\underline{P}_{i_j}^x)]^2 \chi_{i_j}^x(\underline{x}) = [\pi_{h,\kappa_x}^x [(\hat{\varphi}^x)^2]](\underline{x}) \\ &\quad \forall \underline{x} \in \kappa_x, \quad \forall \hat{\varphi}^x \in C(\overline{\kappa_x}), \end{aligned} \quad (6.4.46a)$$

where we have used (6.4.14a) and that  $\chi_{i_j}^x$  are nonnegative, and  $\sum_{j=0}^d \chi_{i_j}^x(\underline{x}) = 1$  for all  $\underline{x} \in \kappa_x$ . Similarly, we have for all  $\kappa_x \in \mathcal{T}_x^h$ ,  $\kappa_q \in \mathcal{T}_q^h$  that

$$|[\pi_{h,\kappa_q}^q \hat{\varphi}^q](q)|^2 \leq [\pi_{h,\kappa_q}^q [(\hat{\varphi}^q)^2]](q) \quad \forall q \in \kappa_q, \quad \forall \hat{\varphi}^q \in C(\overline{\kappa_q}), \quad (6.4.46b)$$

$$|[\pi_{h,\kappa_x \times \kappa_q} \hat{\varphi}](\underline{x}, q)|^2 \leq [\pi_{h,\kappa_x \times \kappa_q} [(\hat{\varphi})^2]](\underline{x}, q) \quad \forall (\underline{x}, q) \in \kappa_x \times \kappa_q, \quad \forall \hat{\varphi} \in C(\overline{\kappa_x \times \kappa_q}), \quad (6.4.46c)$$

$$|[\pi_{h,\kappa_x \times \kappa_q} \hat{\varphi}](\underline{x}, q)|^2 \leq [\pi_{h,\kappa_x \times \kappa_q} [|\hat{\varphi}|^2]](\underline{x}, q) \quad \forall (\underline{x}, q) \in \kappa_x \times \kappa_q, \quad \forall \hat{\varphi} \in [C(\overline{\kappa_x \times \kappa_q})]^d. \quad (6.4.46d)$$

In addition, for all  $\kappa_x \in \mathcal{T}_x^h$ ,  $\kappa_q \in \mathcal{T}_q^h$  and for all  $\hat{\varphi}, \hat{\psi} \in C(\overline{\kappa_x \times \kappa_q})$ ,  $\hat{\varphi}, \hat{\psi} \in [C(\overline{\kappa_x \times \kappa_q})]^d$  the following inequalities are easily deduced for any  $\eta \in \mathbb{R}_{>0}$

$$|[\pi_{h,\kappa_x \times \kappa_q} [\hat{\varphi} \cdot \hat{\psi}]](\underline{x}, \underline{q})| \leq \frac{1}{2} [\pi_{h,\kappa_x \times \kappa_q} [\eta \hat{\varphi}^2 + \eta^{-1} \hat{\psi}^2]](\underline{x}, \underline{q}) \quad \forall(\underline{x}, \underline{q}) \in \kappa_x \times \kappa_q, \quad (6.4.47a)$$

and

$$|[\pi_{h,\kappa_x \times \kappa_q} [\hat{\varphi} \cdot \hat{\psi}]](\underline{x}, \underline{q})| \leq \frac{1}{2} [\pi_{h,\kappa_x \times \kappa_q} [\eta |\hat{\varphi}|^2 + \eta^{-1} |\hat{\psi}|^2]](\underline{x}, \underline{q}) \quad \forall(\underline{x}, \underline{q}) \in \kappa_x \times \kappa_q. \quad (6.4.47b)$$

The following interpolation stability results are easily established for all  $\kappa_x \in \mathcal{T}_h^x$ ,  $\kappa_q \in \mathcal{T}_h^q$

$$\|\tilde{\nabla}_x \pi_{h,\kappa_x}^x \hat{\varphi}^x\|_{L^\infty(\kappa_x)} \leq C \|\tilde{\nabla}_x \hat{\varphi}^x\|_{L^\infty(\kappa_x)} \quad \forall \hat{\varphi}^x \in W^{1,\infty}(\kappa_x), \quad (6.4.48a)$$

$$\|\tilde{\nabla}_q \pi_{h,\kappa_q}^q \hat{\varphi}^q\|_{L^\infty(\kappa_q)} \leq C \|\tilde{\nabla}_q \hat{\varphi}^q\|_{L^\infty(\kappa_q)} \quad \forall \hat{\varphi}^q \in W^{1,\infty}(\kappa_q). \quad (6.4.48b)$$

It follows from (6.4.48a,b) that

$$\begin{aligned} \sum_{i=1}^d \sum_{j=1}^d \left\| \frac{\partial^2}{\partial x_i \partial q_j} \pi_{h,\kappa_x \times \kappa_q} \hat{\varphi} \right\|_{L^\infty(\kappa_x \times \kappa_q)} &= \sum_{i=1}^d \sum_{j=1}^d \left\| \frac{\partial}{\partial x_i} \pi_{h,\kappa_x}^x \left[ \frac{\partial}{\partial q_j} \pi_{h,\kappa_q}^q \hat{\varphi} \right] \right\|_{L^\infty(\kappa_x \times \kappa_q)} \\ &\leq C \sum_{i=1}^d \sum_{j=1}^d \left\| \frac{\partial^2}{\partial x_i \partial q_j} \hat{\varphi} \right\|_{L^\infty(\kappa_x \times \kappa_q)} \quad \forall \hat{\varphi} \in W^{2,\infty}(\kappa_x \times \kappa_q). \end{aligned} \quad (6.4.49)$$

We recall the well-known approximation results for all  $\kappa_x \in \mathcal{T}_h^x$  and  $\kappa_q \in \mathcal{T}_h^q$

$$\|(I - \pi_{h,\kappa_x}^x) \hat{\varphi}^x\|_{L^\infty(\kappa_x)} \leq C h_x^2 |\hat{\varphi}^x|_{W^{2,\infty}(\kappa_x)} \quad \forall \hat{\varphi}^x \in W^{2,\infty}(\kappa_x), \quad (6.4.50a)$$

$$\|(I - \pi_{h,\kappa_q}^q) \hat{\varphi}^q\|_{L^\infty(\kappa_q)} \leq C h_q^2 |\hat{\varphi}^q|_{W^{2,\infty}(\kappa_q)} \quad \forall \hat{\varphi}^q \in W^{2,\infty}(\kappa_q). \quad (6.4.50b)$$

We require the following inverse bounds for all  $\hat{\varphi}_h^x \in \mathbb{P}_1^x$ ,  $\hat{\varphi}_h^q \in \mathbb{P}_1^q$  and for all  $\kappa_x^* \subset \kappa_x \in \mathcal{T}_h^x$ ,  $\kappa_q^* \subset \kappa_q \in \mathcal{T}_h^q$  with  $\underline{m}(\kappa_x) \leq C \underline{m}(\kappa_x^*)$ ,  $\underline{m}(\kappa_q) \leq C \underline{m}(\kappa_q^*)$ :

$$\|\hat{\varphi}_h^x\|_{L^\infty(\kappa_x)}^2 \leq C [\underline{m}(\kappa_x^*)]^{-1} \int_{\kappa_x^*} |\hat{\varphi}_h^x|^2 dx, \quad (6.4.51a)$$

$$\|\hat{\varphi}_h^q\|_{L^\infty(\kappa_q)}^2 \leq C [\underline{m}(\kappa_q^*)]^{-1} \int_{\kappa_q^*} |\hat{\varphi}_h^q|^2 dq, \quad (6.4.51b)$$

$$\int_{\kappa_x^*} |\tilde{\nabla}_x \hat{\varphi}_h^x|^2 dx \leq C h_x^{-2} \int_{\kappa_x^*} |\hat{\varphi}_h^x|^2 dx \leq C h_x^{-2} \int_{\kappa_x^*} \pi_{h,\kappa_x}^x [|\hat{\varphi}_h^x|^2] dx, \quad (6.4.51c)$$

$$\int_{\kappa_q^*} |\tilde{\nabla}_q \hat{\varphi}_h^q|^2 dq \leq C h_q^{-2} \int_{\kappa_q^*} |\hat{\varphi}_h^q|^2 dq \leq C h_q^{-2} \int_{\kappa_q^*} \pi_{h,\kappa_q}^q [|\hat{\varphi}_h^q|^2] dq. \quad (6.4.51d)$$

The bounds (6.4.51a,b) are standard inverse bounds in the case  $\kappa_x^* \equiv \kappa_x$  and  $\kappa_q^* \equiv \kappa_q$ . However, these results are easily generalized to  $\kappa_x^* \subset \kappa_x$  and  $\kappa_q^* \subset \kappa_q$  under the stated conditions. The first inequalities in (6.4.51c,d) then follow immediately from (6.4.51a,b), respectively; whereas the second inequalities in (6.4.51c,d) follow from (6.4.46a,b), respectively. The following bounds follow immediately from (6.4.51a,b) under the same stated conditions:

$$\int_{\kappa_x^*} \pi_{h,\kappa_x}^x [|\hat{\varphi}_h^x|^2] dx \leq C \int_{\kappa_x^*} |\hat{\varphi}_h^x|^2 dx \quad \text{and} \quad \int_{\kappa_q^*} \pi_{h,\kappa_q}^q [|\hat{\varphi}_h^q|^2] dq \leq C \int_{\kappa_q^*} |\hat{\varphi}_h^q|^2 dq. \quad (6.4.52)$$

In addition, we require the following weighted bounds.

**Lemma 6.4.3** For all  $\kappa_q \in \mathcal{T}_q^h$  and for all  $\hat{\varphi}_h^q \in \mathbb{P}_1^q$  we have that

$$\int_{\kappa_q} M |\nabla_q \hat{\varphi}_h^q|^2 dq \leq C h_q^{-2} \int_{\kappa_q} M |\hat{\varphi}_h^q|^2 dq \leq C h_q^{-2} \int_{\kappa_q} M \pi_{h,\kappa_q}^q [|\hat{\varphi}_h^q|^2] dq, \quad (6.4.53a)$$

$$\int_{\kappa_q} M \pi_{h,\kappa_q}^q [|\hat{\varphi}_h^q|^2] dq \leq \left( \int_{\kappa_q} M dq \right) \|\hat{\varphi}_h^q\|_{L^\infty(\kappa_q)}^2 \leq C \int_{\kappa_q} M |\hat{\varphi}_h^q|^2 dq. \quad (6.4.53b)$$

**Proof.** If  $\kappa_q$  has no vertices on  $\partial D$ , let  $q_{\min}$  be the nearest point of  $\kappa_q$  to  $\partial D$ . It follows from the quasiuniformity of  $\mathcal{T}_q^h$  that  $\text{dist}(q_{\min}, \partial D) \geq C h_q$ , and hence, on noting (6.2.9a), it follows that

$$\frac{\max_{q \in \kappa_q} M(q)}{\min_{q \in \kappa_q} M(q)} \leq \frac{c_2 [\text{dist}(q_{\min}, \partial D) + h_q]^\zeta}{c_1 [\text{dist}(q_{\min}, \partial D)]^\zeta} \leq C. \quad (6.4.54)$$

The first inequality in (6.4.53a) then follows immediately from (6.4.51d) and (6.4.54). Similarly, (6.4.53b) follows immediately from (6.4.51b) and (6.4.54).

If  $\kappa_q$  has vertices on  $\partial D$ , we introduce, for appropriate  $C_i \in \mathbb{R}_{>0}$ ,

$$\kappa_q^* := \{q \in \kappa_q : \text{dist}(q, \partial D) \geq C_1 h_q\} \subset \kappa_q \quad \text{and} \quad \underline{m}(\kappa_q) \leq C_2 \underline{m}(\kappa_q^*). \quad (6.4.55)$$

Similarly to (6.4.54), we have from (6.4.54) and (6.2.9a) that

$$\frac{\max_{q \in \kappa_q} M(q)}{\min_{q \in \kappa_q^*} M(q)} \leq C. \quad (6.4.56)$$

It follows from (6.4.55), (6.2.9a), (6.4.56) and (6.4.51d) that

$$\begin{aligned} \int_{\kappa_q} M |\nabla_q \hat{\varphi}_h^q|^2 dq &\leq C_2 \int_{\kappa_q^*} M |\nabla_q \hat{\varphi}_h^q|^2 dq \\ &\leq C h_q^{-2} \int_{\kappa_q^*} M |\hat{\varphi}_h^q|^2 dq \\ &\leq C h_q^{-2} \int_{\kappa_q} M |\hat{\varphi}_h^q|^2 dq, \end{aligned} \quad (6.4.57)$$

and hence the first inequality in (6.4.53a). Similarly, the bound (6.4.53b) in this case follows immediately from (6.4.51b), (6.4.55) and (6.4.56).

Finally, the second inequality in (6.4.53a) follows in both cases from (6.4.46b).  $\square$

In addition, we require the following inverse inequalities.

**Lemma 6.4.4** For all  $\hat{\varphi}_h \in \mathbb{P}_1^x \otimes \mathbb{P}_1^q$  and for all  $\kappa_x \in \mathcal{T}_h^x$ ,  $\kappa_q \in \mathcal{T}_h^q$  we have that

$$\begin{aligned} \int_{\kappa_x \times \kappa_q} M \pi_{h, \kappa_x \times \kappa_q} [|\nabla_x \hat{\varphi}_h|^2] \, dq \, dx &\leq \int_{\kappa_x \times \kappa_q} M |\nabla_x \hat{\varphi}_h|^2 \, dq \, dx \\ &\leq C h_x^{-2} \int_{\kappa_x \times \kappa_q} M |\hat{\varphi}_h|^2 \, dq \, dx, \end{aligned} \quad (6.4.58a)$$

$$\begin{aligned} \int_{\kappa_x \times \kappa_q} M \pi_{h, \kappa_x \times \kappa_q} [|\nabla_q \hat{\varphi}_h|^2] \, dq \, dx &\leq \int_{\kappa_x \times \kappa_q} M |\nabla_q \hat{\varphi}_h|^2 \, dq \, dx \\ &\leq C h_q^{-2} \int_{\kappa_x \times \kappa_q} M |\hat{\varphi}_h|^2 \, dq \, dx. \end{aligned} \quad (6.4.58b)$$

**Proof.** The first inequalities in (6.4.58a,b) follow immediately from (6.4.53b) and (6.4.52), respectively. The second inequalities in (6.4.58a,b) follow immediately from the first inequalities in (6.4.51c) and (6.4.53a), respectively.  $\square$

We require the following results.

**Lemma 6.4.5** For all  $\kappa_x \in \mathcal{T}_h^x$ ,  $\kappa_q \in \mathcal{T}_h^q$  and for all  $\hat{\psi}_h, \hat{\varphi}_h \in \hat{X}_h$  we have that

$$\begin{aligned} &\left| \int_{\kappa_x \times \kappa_q} M (I - \pi_{h, \kappa_x \times \kappa_q}) [\nabla_q \hat{\psi}_h \cdot \nabla_q \hat{\varphi}_h] \, dq \, dx \right| \\ &\leq C h_x \left( \int_{\kappa_x \times \kappa_q} M |\nabla_q \hat{\psi}_h|^2 \, dq \, dx \right)^{\frac{1}{2}} \left( \sum_{i=1}^d \sum_{j=1}^d \int_{\kappa_x \times \kappa_q} M \left| \frac{\partial^2 \hat{\varphi}_h}{\partial x_i \partial q_j} \right|^2 \, dq \, dx \right)^{\frac{1}{2}}, \end{aligned} \quad (6.4.59a)$$

$$\begin{aligned} &\left| \int_{\kappa_x \times \kappa_q} M (I - \pi_{h, \kappa_x \times \kappa_q}) [\nabla_x \hat{\psi}_h \cdot \nabla_x \hat{\varphi}_h] \, dq \, dx \right| \\ &\leq C h_q \left( \int_{\kappa_x \times \kappa_q} M |\nabla_x \hat{\psi}_h|^2 \, dq \, dx \right)^{\frac{1}{2}} \left( \sum_{i=1}^d \sum_{j=1}^d \int_{\kappa_x \times \kappa_q} M \left| \frac{\partial^2 \hat{\varphi}_h}{\partial x_i \partial q_j} \right|^2 \, dq \, dx \right)^{\frac{1}{2}}, \end{aligned} \quad (6.4.59b)$$

and

$$\begin{aligned} &\left| \int_{\kappa_x \times \kappa_q} M (I - \pi_{h, \kappa_x \times \kappa_q}) [\hat{\psi}_h \hat{\varphi}_h] \, dq \, dx \right| \\ &\leq C h_x^2 \left( \int_{\kappa_x \times \kappa_q} M |\nabla_x \hat{\psi}_h|^2 \, dq \, dx \right)^{\frac{1}{2}} \left( \int_{\kappa_x \times \kappa_q} M |\nabla_x \hat{\varphi}_h|^2 \, dq \, dx \right)^{\frac{1}{2}} \\ &\quad + C h_q^2 \left( \int_{\kappa_x \times \kappa_q} M |\nabla_q \hat{\psi}_h|^2 \, dq \, dx \right)^{\frac{1}{2}} \left( \int_{\kappa_x \times \kappa_q} M |\nabla_q \hat{\varphi}_h|^2 \, dq \, dx \right)^{\frac{1}{2}}. \end{aligned} \quad (6.4.59c)$$

**Proof.** As  $\nabla_q \hat{\psi}_h, \nabla_q \hat{\varphi}_h \in [\mathbb{P}_1^x]^d$  on  $\kappa_x \times \kappa_q$ , it follows from (6.4.50a) that

$$\begin{aligned}
& \left| \int_{\kappa_x \times \kappa_q} M (I - \pi_{h, \kappa_x \times \kappa_q}) [\nabla_q \hat{\psi}_h \cdot \nabla_q \hat{\varphi}_h] \, dq \, dx \right| \\
& \leq \left( \int_{\kappa_x \times \kappa_q} M \, dq \, dx \right) \|(I - \pi_{h, \kappa_x}^x) [\nabla_q \hat{\psi}_h \cdot \nabla_q \hat{\varphi}_h]\|_{L^\infty(\kappa_x)} \\
& \leq C h_x^2 \left( \int_{\kappa_x \times \kappa_q} M \, dq \, dx \right) \|\nabla_q \hat{\psi}_h \cdot \nabla_q \hat{\varphi}_h\|_{W^{2, \infty}(\kappa_x)} \\
& \leq C h_x^2 \left( \sum_{i=1}^d \sum_{j=1}^d \int_{\kappa_x \times \kappa_q} M \left| \frac{\partial^2 \hat{\psi}_h}{\partial x_i \partial q_j} \right|^2 \, dq \, dx \right)^{\frac{1}{2}} \left( \sum_{i=1}^d \sum_{j=1}^d \int_{\kappa_x \times \kappa_q} M \left| \frac{\partial^2 \hat{\varphi}_h}{\partial x_i \partial q_j} \right|^2 \, dq \, dx \right)^{\frac{1}{2}}.
\end{aligned} \tag{6.4.60}$$

The desired result (6.4.59a) then follows from (6.4.60) on applying (6.4.58a) to the first integral.

Similarly, as  $\nabla_x \hat{\psi}_h, \nabla_x \hat{\varphi}_h \in [\mathbb{P}_1^q]^d$  on  $\kappa_x \times \kappa_q$ , it follows from (6.4.50b) that

$$\begin{aligned}
& \left| \int_{\kappa_x \times \kappa_q} M (I - \pi_{h, \kappa_x \times \kappa_q}) [\nabla_x \hat{\psi}_h \cdot \nabla_x \hat{\varphi}_h] \, dq \, dx \right| \\
& \leq C h_q^2 \left( \sum_{i=1}^d \sum_{j=1}^d \int_{\kappa_x \times \kappa_q} M \left| \frac{\partial^2 \hat{\psi}_h}{\partial x_i \partial q_j} \right|^2 \, dq \, dx \right)^{\frac{1}{2}} \left( \sum_{i=1}^d \sum_{j=1}^d \int_{\kappa_x \times \kappa_q} M \left| \frac{\partial^2 \hat{\varphi}_h}{\partial x_i \partial q_j} \right|^2 \, dq \, dx \right)^{\frac{1}{2}}.
\end{aligned} \tag{6.4.61}$$

The desired result (6.4.59b) then follows from (6.4.61) on applying (6.4.58b) to the first integral.

To prove (6.4.59c), we first note that

$$I - \pi_{h, \kappa_x \times \kappa_q} \equiv (I - \pi_{h, \kappa_x}^x) + (I - \pi_{h, \kappa_q}^q) \pi_{h, \kappa_x}^x.$$

It follows from (6.4.50a) that

$$\begin{aligned}
& \int_{\kappa_x \times \kappa_q} M \|(I - \pi_{h, \kappa_x}^x) [\hat{\psi}_h \hat{\varphi}_h]\|_{L^\infty(\kappa_x)} \, dq \, dx \leq C h_x^2 \int_{\kappa_x \times \kappa_q} M |\nabla_x \hat{\psi}_h| |\nabla_x \hat{\varphi}_h| \, dq \, dx \\
& \leq C h_x^2 \left( \int_{\kappa_x \times \kappa_q} M |\nabla_x \hat{\psi}_h|^2 \, dq \, dx \right)^{\frac{1}{2}} \left( \int_{\kappa_x \times \kappa_q} M |\nabla_x \hat{\varphi}_h|^2 \, dq \, dx \right)^{\frac{1}{2}}.
\end{aligned} \tag{6.4.62}$$

It follows from (6.4.50b) and (6.4.51a) that

$$\begin{aligned}
& \left( \int_{\kappa_x \times \kappa_q} M \, \underset{\sim}{dq} \, \underset{\sim}{dx} \right) \left\| (I - \pi_{h, \kappa_x}^q) \pi_{h, \kappa_x}^x [\hat{\psi}_h \hat{\varphi}_h] \right\|_{L^\infty(\kappa_x \times \kappa_q)} \\
& \leq C h_q^2 \left( \int_{\kappa_x \times \kappa_q} M \, \underset{\sim}{dq} \, \underset{\sim}{dx} \right) \sum_{i=1}^d \sum_{j=1}^d \left\| \pi_{h, \kappa_x}^x \left[ \frac{\partial \hat{\psi}_h}{\partial q_i} \frac{\partial \hat{\varphi}_h}{\partial q_j} \right] \right\|_{L^\infty(\kappa_x)} \\
& \leq C h_q^2 \left( \int_{\kappa_x \times \kappa_q} M \, \underset{\sim}{dq} \, \underset{\sim}{dx} \right) \|\nabla_q \hat{\psi}_h\|_{L^\infty(\kappa_x)} \|\nabla_q \hat{\varphi}_h\|_{L^\infty(\kappa_x)} \\
& \leq C h_q^2 \left( \int_{\kappa_x \times \kappa_q} M |\nabla_q \hat{\psi}_h|^2 \, \underset{\sim}{dq} \, \underset{\sim}{dx} \right)^{\frac{1}{2}} \left( \int_{\kappa_x \times \kappa_q} M |\nabla_q \hat{\varphi}_h|^2 \, \underset{\sim}{dq} \, \underset{\sim}{dx} \right)^{\frac{1}{2}}. \tag{6.4.63}
\end{aligned}$$

Hence combining (6.4.62) and (6.4.63) yields the desired result (6.4.59c).  $\square$

**Lemma 6.4.6** For all  $\kappa_x \in \mathcal{T}_h^x$ ,  $\kappa_q \in \mathcal{T}_h^q$  and for all  $\hat{\psi}_h, \hat{\varphi}_h \in \hat{X}_h$  we have that

$$\begin{aligned}
& \int_{\kappa_x \times \kappa_q} M \left[ \left| (I - \pi_{h, \kappa_x \times \kappa_q}) [\Xi_\delta^q(\hat{\psi}_h) \nabla_q \hat{\varphi}_h] \right|^2 + \left| (I - \pi_{h, \kappa_x \times \kappa_q}) [\Xi_\delta^x(\hat{\psi}_h) \nabla_x \hat{\varphi}_h] \right|^2 \right] \, \underset{\sim}{dq} \, \underset{\sim}{dx} \\
& \leq C(L) (h_x^2 + h_q^2) \left( \sum_{i=1}^d \sum_{j=1}^d \int_{\kappa_x \times \kappa_q} M \left| \frac{\partial^2 \hat{\varphi}_h}{\partial x_i \partial q_j} \right|^2 \, \underset{\sim}{dq} \, \underset{\sim}{dx} \right). \tag{6.4.64}
\end{aligned}$$

**Proof.** As  $\Xi_\delta^q(\hat{\psi}_h) \in [\mathbb{P}_1^x]^{d \times d}$  and  $\nabla_q \hat{\varphi}_h \in [\mathbb{P}_1^x]^d$  on  $\kappa_x \times \kappa_q$ , it follows from (6.4.50a), (6.4.51c) and (6.4.24) that

$$\begin{aligned}
& \int_{\kappa_x \times \kappa_q} M \left| (I - \pi_{h, \kappa_x \times \kappa_q}) [\Xi_\delta^q(\hat{\psi}_h) \nabla_q \hat{\varphi}_h] \right|^2 \, \underset{\sim}{dq} \, \underset{\sim}{dx} \\
& \leq C h_x^4 \left( \int_{\kappa_x \times \kappa_q} M \, \underset{\sim}{dq} \, \underset{\sim}{dx} \right) \left( \sum_{i=1}^d \sum_{j=1}^d \|\nabla_x [\Xi_\delta^q(\hat{\psi}_h)]_{ij}\|_{L^\infty(\kappa_x)}^2 \right) \left( \sum_{i=1}^d \sum_{j=1}^d \left\| \frac{\partial^2 \hat{\varphi}_h}{\partial x_i \partial q_j} \right\|_{L^\infty(\kappa_x)}^2 \right) \\
& \leq C(L) h_x^2 \left( \sum_{i=1}^d \sum_{j=1}^d \int_{\kappa_x \times \kappa_q} M \left| \frac{\partial^2 \hat{\varphi}_h}{\partial x_i \partial q_j} \right|^2 \, \underset{\sim}{dq} \, \underset{\sim}{dx} \right). \tag{6.4.65}
\end{aligned}$$

Similarly, as  $\Xi_\delta^x(\hat{\psi}_h) \in [\mathbb{P}_1^q]^{d \times d}$  and  $\nabla_x \hat{\varphi}_h \in [\mathbb{P}_1^q]^d$  on  $\kappa_x \times \kappa_q$ , it follows from (6.4.50b), (6.4.53a) and (6.4.24) that

$$\begin{aligned}
& \int_{\kappa_x \times \kappa_q} M \left| (I - \pi_{h, \kappa_x \times \kappa_q}) [\Xi_\delta^x(\hat{\psi}_h) \nabla_x \hat{\varphi}_h] \right|^2 \, \underset{\sim}{dq} \, \underset{\sim}{dx} \\
& \leq C(L) h_q^2 \left( \sum_{i=1}^d \sum_{j=1}^d \int_{\kappa_x \times \kappa_q} M \left| \frac{\partial^2 \hat{\varphi}_h}{\partial x_i \partial q_j} \right|^2 \, \underset{\sim}{dq} \, \underset{\sim}{dx} \right). \tag{6.4.66}
\end{aligned}$$

Combining (6.4.65) and (6.4.66) yields the desired result (6.4.64).  $\square$

In addition, we introduce  $Q_h^M : \hat{X} \rightarrow \hat{X}_h$  and  $\tilde{Q}_h^M : \hat{X} \rightarrow \hat{X}_h$  such that

$$(M Q_h^M \hat{\psi}, \hat{\varphi}_h)_{\Omega \times D} = (M \hat{\psi}, \hat{\varphi}_h)_{\Omega \times D} \quad \forall \hat{\varphi}_h \in \hat{X}_h, \quad (6.4.67a)$$

$$(M, \pi_h[(\tilde{Q}_h^M \hat{\psi}) \hat{\varphi}_h])_{\Omega \times D} = (M \hat{\psi}, \hat{\varphi}_h)_{\Omega \times D} \quad \forall \hat{\varphi}_h \in \hat{X}_h. \quad (6.4.67b)$$

In the Appendix, it is shown that

$$\|Q_h^M \hat{\psi}\|_{\hat{X}}^2 \leq C \|\hat{\psi}\|_{\hat{X}}^2 \quad \forall \hat{\psi} \in \hat{X}. \quad (6.4.68)$$

We require a related result for  $\tilde{Q}_h^M$ .

**Lemma 6.4.7** *The following bounds hold:*

$$\|\tilde{Q}_h^M \hat{\psi}\|_{\hat{X}}^2 \leq \left( M, \pi_h \left[ |\tilde{Q}_h^M \hat{\psi}|^2 + |\nabla_x (\tilde{Q}_h^M \hat{\psi})|^2 + |\nabla_q (\tilde{Q}_h^M \hat{\psi})|^2 \right] \right)_{\Omega \times D} \leq C \|\hat{\psi}\|_{\hat{X}}^2 \quad \forall \hat{\psi} \in \hat{X}. \quad (6.4.69)$$

**Proof.** Given  $\hat{\psi} \in \hat{X}$ , let  $E = (Q_h^M - \tilde{Q}_h^M) \hat{\psi}$ . It follows from (6.4.46c), (6.4.67a,b), (6.4.59c), (6.4.68), (6.4.58a,b) that

$$\begin{aligned} (M, E^2)_{\Omega \times D} &\leq (M, \pi_h[E^2])_{\Omega \times D} = (M, (\pi_h - I)[(Q_h^M \hat{\psi}) E])_{\Omega \times D} \\ &\leq C \|\hat{\psi}\|_{\hat{X}} \left[ h_x^2 \left( \int_{\Omega \times D} M |\nabla_x E|^2 dq dx \right)^{\frac{1}{2}} + h_q^2 \left( \int_{\Omega \times D} M |\nabla_q E|^2 dq dx \right)^{\frac{1}{2}} \right] \\ &\leq C (h_x + h_q) \|\hat{\psi}\|_{\hat{X}} [(M, E^2)_{\Omega \times D}]^{\frac{1}{2}} \leq C (h_x + h_q)^2 \|\hat{\psi}\|_{\hat{X}}^2. \end{aligned} \quad (6.4.70)$$

It follows from (6.4.58a,b), (6.4.70) and (6.4.1) that

$$\|(Q_h^M - \tilde{Q}_h^M) \hat{\psi}\|_{\hat{X}}^2 \leq C \|\hat{\psi}\|_{\hat{X}}^2. \quad (6.4.71)$$

The desired result (6.4.69) then follows from (6.4.71), (6.4.68), (6.4.58a,b), (6.4.59c) and (6.4.46c,d).  $\square$

We are now in a position to prove the following stability result for  $(P_\delta^{h, \Delta t})$ .

**Lemma 6.4.8** *A solution  $\{u_{\delta,h}^n, \hat{\psi}_{\delta,h}^n\}_{n=1}^N$  of  $(P_\delta^{h, \Delta t})$  satisfies the following stability bounds:*

$$\begin{aligned} &\max_{n=1 \rightarrow N} \|u_{\delta,h}^n\|_{\tilde{L}^2(\Omega)}^2 + \max_{n=1 \rightarrow N} (M, \pi_h[\mathcal{F}_\delta^L(\hat{\psi}_{\delta,h}^n)])_{\Omega \times D} \\ &+ \sum_{n=1}^N \Delta t^n \|\nabla_x u_{\delta,h}^n\|_{\tilde{L}^2(\Omega)}^2 + \sum_{n=1}^N \|u_{\delta,h}^n - u_{\delta,h}^{n-1}\|_{\tilde{L}^2(\Omega)}^2 \\ &+ \delta \sum_{n=1}^N \Delta t^n \left[ (M, \pi_h[|\nabla_x (\pi_h[\mathcal{F}_\delta^L]'(\hat{\psi}_{\delta,h}^n))|^2])_{\Omega \times D} + (M, \pi_h[|\nabla_q (\pi_h[\mathcal{F}_\delta^L]'(\hat{\psi}_{\delta,h}^n))|^2])_{\Omega \times D} \right] \\ &\leq C \left[ \|u_{\delta,h}^0\|_{\tilde{L}^2(\Omega)}^2 + (M, \pi_h[\mathcal{F}_\delta^L(\hat{\psi}_{\delta,h}^0)])_{\Omega \times D} + \sum_{n=1}^n \Delta t^n \|f^n\|_{\tilde{H}^{-1}(\Omega)}^2 \right] \leq C, \end{aligned} \quad (6.4.72a)$$

$$\begin{aligned} &\max_{n=1 \rightarrow N} (M, \pi_h[|\hat{\psi}_{\delta,h}^n|^2])_{\Omega \times D} \\ &+ \sum_{n=1}^N \Delta t^n (M, \pi_h[|\nabla_q \hat{\psi}_{\delta,h}^n|^2 + |\nabla_x \hat{\psi}_{\delta,h}^n|^2])_{\Omega \times D} + \sum_{n=1}^N (M, \pi_h[|\hat{\psi}_{\delta,h}^n - \hat{\psi}_{\delta,h}^{n-1}|^2])_{\Omega \times D} \\ &\leq C(L) + C(M, \pi_h[|\hat{\psi}_h^0|^2])_{\Omega \times D} \leq C(L), \end{aligned} \quad (6.4.72b)$$

and

$$\begin{aligned} \max_{n=1 \rightarrow N} \left[ \int_{\Omega} |C(M \hat{\psi}_{\delta,h}^n)|^2 dx \right] + \sum_{n=1}^N \Delta t^n \left\| S \left( \frac{u_{\delta,h}^n - u_{\delta,h}^{n-1}}{\Delta t^n} \right) \right\|_{\mathbb{H}^1(\Omega)}^{\frac{4}{\vartheta}} \\ + \sum_{n=1}^N \Delta t^n \left\| \mathcal{G} \left( \frac{\hat{\psi}_{\delta,h}^n - \hat{\psi}_{\delta,h}^{n-1}}{\Delta t^n} \right) \right\|_{\hat{X}}^2 \leq C(L, T), \end{aligned} \quad (6.4.72c)$$

where

$$\vartheta \in (2, 4) \quad \text{if } d = 2 \quad \text{and} \quad \vartheta = 3 \quad \text{if } d = 3. \quad (6.4.73)$$

**Proof.** Summing (6.4.33) from  $n = 1 \rightarrow m$ , for  $m = 1 \rightarrow N$ , yields the desired result (6.4.72a) on noting (6.4.26a,b), (6.4.29), (6.2.15) and (6.4.31a).

On choosing  $\hat{\varphi}_h = \hat{\psi}_{\delta,h}^n$  in (6.4.32b) and noting (6.4.35), we obtain

$$\begin{aligned} \mathbb{T}^n &:= \left( M, \pi_h \left[ |\hat{\psi}_{\delta,h}^n|^2 + |\hat{\psi}_{\delta,h}^n - \hat{\psi}_{\delta,h}^{n-1}|^2 \right] \right)_{\Omega \times D} \\ &\quad + \Delta t^n \left( M, \pi_h \left[ 2\varepsilon |\nabla_x \hat{\psi}_{\delta,h}^n|^2 + \frac{1}{\lambda} |\nabla_q \hat{\psi}_{\delta,h}^n|^2 \right] \right)_{\Omega \times D} \\ &= \left( M, \pi_h \left[ |\hat{\psi}_{\delta,h}^{n-1}|^2 \right] \right)_{\Omega \times D} + 2 \Delta t^n \left( M \left( \nabla_x u_{\delta,h}^n \right)_q, \pi_h \left[ \Xi_{\delta}^q(\hat{\psi}_{\delta,h}^n) \nabla_q \hat{\psi}_{\delta,h}^n \right] \right)_{\Omega \times D} \\ &\quad + 2 \Delta t^n \left( M u_{\delta,h}^n, \pi_h \left[ \Xi_{\delta}^x(\hat{\psi}_{\delta,h}^n) \nabla_x \hat{\psi}_{\delta,h}^n \right] \right)_{\Omega \times D}. \end{aligned}$$

Hence, recalling (6.4.46d) and (6.4.24), for any  $\eta \in \mathbb{R}_{>0}$ , we have that

$$\begin{aligned} \mathbb{T}^n &\leq \left( M, \pi_h \left[ |\hat{\psi}_{\delta,h}^{n-1}|^2 \right] \right)_{\Omega \times D} + \Delta t^n C(\eta^{-1}) \left[ \|u_{\delta,h}^n\|_{L^2(\Omega)}^2 + \|\nabla_x u_{\delta,h}^n\|_{L^2(\Omega)}^2 \right] \\ &\quad + \Delta t^n \eta \left( M, \left| \pi_h \left[ \Xi_{\delta}^x(\hat{\psi}_{\delta,h}^n) \nabla_x \hat{\psi}_{\delta,h}^n \right] \right|^2 + \left| \pi_h \left[ \Xi_{\delta}^q(\hat{\psi}_{\delta,h}^n) \nabla_q \hat{\psi}_{\delta,h}^n \right] \right|^2 \right)_{\Omega \times D} \\ &\leq \left( M, \pi_h \left[ |\hat{\psi}_{\delta,h}^{n-1}|^2 \right] \right)_{\Omega \times D} + \Delta t^n C(\eta^{-1}) \left[ \|u_{\delta,h}^n\|_{L^2(\Omega)}^2 + \|\nabla_x u_{\delta,h}^n\|_{L^2(\Omega)}^2 \right] \\ &\quad + \Delta t_n C L^2 \eta \left( M, \pi_h \left[ |\nabla_x \hat{\psi}_{\delta,h}^n|^2 + |\nabla_q \hat{\psi}_{\delta,h}^n|^2 \right] \right)_{\Omega \times D}. \end{aligned} \quad (6.4.74)$$

On noting the definition of  $\mathbb{T}^n$ , summing (6.4.74) from  $n = 1 \rightarrow m$ , for  $m = 1 \rightarrow N$ , with  $\eta$  chosen sufficiently small, and recalling inequalities (6.4.72a) and (6.4.29), yields the desired result (6.4.72b).

The first bound in (6.4.72c) follows immediately from the first bound in (6.4.72b), (6.3.15) and (6.4.46c).

On choosing

$$\mathfrak{w}_h = \mathcal{Q}_h \left[ S \left( \frac{u_{\delta,h}^n - u_{\delta,h}^{n-1}}{\Delta t^n} \right) \right] \in \mathfrak{Y}_h$$

in (6.4.32a) yields, on noting (6.4.6), (6.3.3), (6.4.7) and Sobolev embedding, that

$$\begin{aligned}
\left\| S \left( \frac{u_{\delta,h}^n - u_{\delta,h}^{n-1}}{\Delta t^n} \right) \right\|_{\mathbf{H}^1(\Omega)}^2 &= \left( \frac{u_{\delta,h}^n - u_{\delta,h}^{n-1}}{\Delta t^n}, Q_h \left[ S \left( \frac{u_{\delta,h}^n - u_{\delta,h}^{n-1}}{\Delta t^n} \right) \right] \right)_{\Omega} \\
&= -\nu \left( \nabla_x u_{\delta,h}^n, \nabla_x \left[ Q_h \left[ S \left( \frac{u_{\delta,h}^n - u_{\delta,h}^{n-1}}{\Delta t^n} \right) \right] \right] \right)_{\Omega} \\
&\quad - k_B \mathcal{T} \left( C(M \hat{\psi}_{\delta,h}^n), \nabla_x \left[ Q_h \left[ S \left( \frac{u_{\delta,h}^n - u_{\delta,h}^{n-1}}{\Delta t^n} \right) \right] \right] \right)_{\Omega} \\
&\quad - \frac{1}{2} \left( (u_{\delta,h}^{n-1} \cdot \nabla_x) u_{\delta,h}^n, Q_h \left[ S \left( \frac{u_{\delta,h}^n - u_{\delta,h}^{n-1}}{\Delta t^n} \right) \right] \right)_{\Omega} \\
&\quad + \frac{1}{2} \left( u_{\delta,h}^n, (u_{\delta,h}^{n-1} \cdot \nabla_x) \left[ Q_h \left[ S \left( \frac{u_{\delta,h}^n - u_{\delta,h}^{n-1}}{\Delta t^n} \right) \right] \right] \right)_{\Omega} \\
&\quad + \left\langle f^n, Q_h \left[ S \left( \frac{u_{\delta,h}^n - u_{\delta,h}^{n-1}}{\Delta t^n} \right) \right] \right\rangle_{\mathbf{H}_0^1(\Omega)} \\
&\leq C \left[ \|C(M \hat{\psi}_{\delta,h}^n)\|_{L^2(\Omega)}^2 + \|\nabla_x u_{\delta,h}^n\|_{L^2(\Omega)}^2 + \| |u_{\delta,h}^{n-1}| |u_{\delta,h}^n| \|_{L^2(\Omega)}^2 \right. \\
&\quad \left. + \| |u_{\delta,h}^{n-1}| |\nabla_x u_{\delta,h}^n| \|_{L^{1+\theta}(\Omega)}^2 + \|f^n\|_{\mathbf{H}^{-1}(\Omega)}^2 \right], \tag{6.4.75}
\end{aligned}$$

for any  $\theta > 0$  if  $d = 2$  and for  $\theta = \frac{1}{5}$  if  $d = 3$ . Applying the Cauchy–Schwarz and the algebraic-geometric mean inequalities, in conjunction with (6.3.4) and a Poincaré inequality yields that

$$\begin{aligned}
\| |u_{\delta,h}^{n-1}| |u_{\delta,h}^n| \|_{L^2(\Omega)}^2 &\leq \|u_{\delta,h}^{n-1}\|_{L^4(\Omega)}^2 \|u_{\delta,h}^n\|_{L^4(\Omega)}^2 \leq \frac{1}{2} \sum_{m=n-1}^n \|u_{\delta,h}^m\|_{L^4(\Omega)}^4 \\
&\leq C \sum_{m=n-1}^n \left[ \|u_{\delta,h}^m\|_{L^2(\Omega)}^{4-d} \|\nabla_x u_{\delta,h}^m\|_{L^2(\Omega)}^d \right]. \tag{6.4.76}
\end{aligned}$$

Similarly, we have for any  $\theta \in (0, 1)$ , if  $d = 2$ , that

$$\begin{aligned}
\| |u_{\delta,h}^{n-1}| |\nabla_x u_{\delta,h}^n| \|_{L^{1+\theta}(\Omega)}^2 &\leq \|u_{\delta,h}^{n-1}\|_{L^{\frac{2(1+\theta)}{1-\theta}}(\Omega)}^2 \|\nabla_x u_{\delta,h}^n\|_{L^2(\Omega)}^2 \\
&\leq C \|u_{\delta,h}^{n-1}\|_{L^2(\Omega)}^{\frac{2(1-\theta)}{1+\theta}} \sum_{m=n-1}^n \|\nabla_x u_{\delta,h}^m\|_{L^2(\Omega)}^{\frac{2(1+3\theta)}{1+\theta}}; \tag{6.4.77a}
\end{aligned}$$

and if  $d = 3$ , ( $\theta = \frac{1}{5}$ ), that

$$\begin{aligned}
\| |u_{\delta,h}^{n-1}| |\nabla_x u_{\delta,h}^n| \|_{L^{\frac{6}{5}}(\Omega)}^2 &\leq \|u_{\delta,h}^{n-1}\|_{L^3(\Omega)}^2 \|\nabla_x u_{\delta,h}^n\|_{L^2(\Omega)}^2 \\
&\leq C \|u_{\delta,h}^{n-1}\|_{L^2(\Omega)} \sum_{m=n-1}^n \|\nabla_x u_{\delta,h}^m\|_{L^2(\Omega)}^3. \tag{6.4.77b}
\end{aligned}$$

On taking the  $\frac{2}{\vartheta}$  power of both sides of (6.4.75), recall (6.4.73), multiplying by  $\Delta t^n$ , summing from  $n = 1 \rightarrow N$  and noting (6.4.76), (6.4.77a) with  $\theta = (\vartheta - 2)/(6 - \vartheta)$ , (6.4.77b), (6.4.27), (6.4.31a), (6.4.72a,b), (6.4.29) and the first bound in (6.4.72c) yields that

$$\begin{aligned}
& \sum_{n=1}^N \Delta t^n \left\| \tilde{S} \left( \frac{u_{\delta,h}^n - u_{\delta,h}^{n-1}}{\Delta t^n} \right) \right\|_{\mathbb{H}^1(\Omega)}^{\frac{4}{\vartheta}} \\
& \leq C \left[ \sum_{n=1}^N \Delta t^n \|\tilde{C}(M \hat{\psi}_{\delta,h}^n)\|_{L^2(\Omega)}^{\frac{4}{\vartheta}} \right] + C(T) \left[ \sum_{n=1}^N \Delta t^n \left[ \|\nabla_x u_{\delta,h}^n\|_{L^2(\Omega)}^2 + \|f^n\|_{\mathbb{H}^{-1}(\Omega)}^2 \right] \right]^{\frac{2}{\vartheta}} \\
& \quad + C \left[ \max_{n=0 \rightarrow N} \left( \|u_{\delta,h}^n\|_{L^2(\Omega)}^2 \right)^{\frac{4}{\vartheta}-1} \right] \left[ \sum_{n=0}^N \Delta t^n \|\nabla_x u_{\delta,h}^n\|_{L^2(\Omega)}^2 \right] \\
& \leq C(L, T); \tag{6.4.78}
\end{aligned}$$

and hence the second bound in (6.4.72c).

On choosing

$$\hat{\varphi}_h = \tilde{Q}_h^M \left[ \mathcal{G} \left( \frac{\hat{\psi}_{\delta,h}^n - \hat{\psi}_{\delta,h}^{n-1}}{\Delta t^n} \right) \right] \in \hat{X}_h$$

in (6.4.32b) yields, on noting (6.4.67b), (6.3.13), (6.4.47b), (6.4.46d), (6.4.24) and (6.4.69), that

$$\begin{aligned}
& \left\| \mathcal{G} \left( \frac{\hat{\psi}_{\delta,h}^n - \hat{\psi}_{\delta,h}^{n-1}}{\Delta t^n} \right) \right\|_{\hat{X}}^2 = \left( M, \pi_h \left[ \left( \frac{\hat{\psi}_{\delta,h}^n - \hat{\psi}_{\delta,h}^{n-1}}{\Delta t^n} \right) \tilde{Q}_h^M \left[ \mathcal{G} \left( \frac{\hat{\psi}_{\delta,h}^n - \hat{\psi}_{\delta,h}^{n-1}}{\Delta t^n} \right) \right] \right] \right)_{\Omega \times D} \\
& = -\frac{1}{2\lambda} \left( M, \pi_h \left[ \nabla_q \hat{\psi}_{\delta,h}^n \cdot \nabla_q \left[ \tilde{Q}_h^M \left[ \mathcal{G} \left( \frac{\hat{\psi}_{\delta,h}^n - \hat{\psi}_{\delta,h}^{n-1}}{\Delta t^n} \right) \right] \right] \right] \right)_{\Omega \times D} \\
& \quad - \varepsilon \left( M, \pi_h \left[ \nabla_x \hat{\psi}_{\delta,h}^n \cdot \nabla_x \left[ \tilde{Q}_h^M \left[ \mathcal{G} \left( \frac{\hat{\psi}_{\delta,h}^n - \hat{\psi}_{\delta,h}^{n-1}}{\Delta t^n} \right) \right] \right] \right] \right)_{\Omega \times D} \\
& \quad + \left( M \left( \nabla_x u_{\delta,h}^n \right)_q, \pi_h \left[ \tilde{\Xi}_{\delta}^q(\hat{\psi}_{\delta,h}^n) \nabla_q \left[ \tilde{Q}_h^M \left[ \mathcal{G} \left( \frac{\hat{\psi}_{\delta,h}^n - \hat{\psi}_{\delta,h}^{n-1}}{\Delta t^n} \right) \right] \right] \right] \right)_{\Omega \times D} \\
& \quad + \left( M u_{\delta,h}^n, \pi_h \left[ \tilde{\Xi}_{\delta}^x(\hat{\psi}_{\delta,h}^n) \nabla_x \left[ \tilde{Q}_h^M \left[ \mathcal{G} \left( \frac{\hat{\psi}_{\delta,h}^n - \hat{\psi}_{\delta,h}^{n-1}}{\Delta t^n} \right) \right] \right] \right] \right)_{\Omega \times D} \\
& \leq C(L) \left[ \|u_{\delta,h}^n\|_{L^2(\Omega)}^2 + \|\nabla_x u_{\delta,h}^n\|_{L^2(\Omega)}^2 \right] \\
& \quad + C \left( M, \pi_h \left[ |\nabla_q \hat{\psi}_{\delta,h}^n|^2 + |\nabla_x \hat{\psi}_{\delta,h}^n|^2 \right] \right)_{\Omega \times D}. \tag{6.4.79}
\end{aligned}$$

Multiplying (6.4.79) by  $\Delta t^n$ , summing from  $n = 1 \rightarrow N$  and noting (6.4.72a,b) yields the desired result (6.4.72c).  $\square$

Now we introduce some definitions prior to passing to the limit  $\delta, h, \Delta t \rightarrow 0_+$ . Let

$$\tilde{u}_{\delta,h}^{\Delta t}(\cdot, t) := \frac{t - t^{n-1}}{\Delta t^n} \tilde{u}_{\delta,h}^n(\cdot) + \frac{t^n - t}{\Delta t^n} \tilde{u}_{\delta,h}^{n-1}(\cdot), \quad t \in [t^{n-1}, t^n], \quad n \geq 1, \quad (6.4.80a)$$

$$\tilde{u}_{\delta,h}^{\Delta t,+}(\cdot, t) := \tilde{u}_{\delta,h}^n(\cdot), \quad \tilde{u}_{\delta,h}^{\Delta t,-}(\cdot, t) := \tilde{u}_{\delta,h}^{n-1}(\cdot), \quad t \in (t^{n-1}, t^n], \quad n \geq 1, \quad (6.4.80b)$$

$$\text{and} \quad \Delta(t) := \Delta t^n, \quad t \in (t^{n-1}, t^n], \quad n \geq 1. \quad (6.4.80c)$$

We note for future reference that

$$\tilde{u}_{\delta,h}^{\Delta t} - \tilde{u}_{\delta,h}^{\Delta t,\pm} = (t - t^{n,\pm}) \frac{\partial \tilde{u}_{\delta,h}^{\Delta t}}{\partial t}, \quad t \in (t^{n-1}, t^n), \quad n \geq 1, \quad (6.4.81)$$

where  $t^{n,+} := t^n$  and  $t^{n,-} := t^{n-1}$ . Using the above notation, and introducing analogous notation for  $\{\hat{\psi}_{\delta,h}^n\}_{n=0}^N$  and  $\{\tilde{f}_n\}_{n=1}^N$ , (6.4.32a) multiplied by  $\Delta t^n$  and summed for  $n = 1 \rightarrow N$  can be restated as:

$$\begin{aligned} & \int_0^T \int_{\Omega} \left[ \frac{\partial \tilde{u}_{\delta,h}^{\Delta t}}{\partial t} \cdot \tilde{w}_h + \nu \nabla_x \tilde{u}_{\delta,h}^{\Delta t,+} : \nabla_x \tilde{w}_h \right] \tilde{d}x \tilde{d}t \\ & + \frac{1}{2} \int_0^T \int_{\Omega} \left[ \left[ (\tilde{u}_{\delta,h}^{\Delta t,-} \cdot \nabla_x) \tilde{u}_{\delta,h}^{\Delta t,+} \right] \cdot \tilde{w}_h - \left[ (\tilde{u}_{\delta,h}^{\Delta t,-} \cdot \nabla_x) \tilde{w}_h \right] \cdot \tilde{u}_{\delta,h}^{\Delta t,+} \right] \tilde{d}x \tilde{d}t \\ & = \int_0^T \left[ \langle \tilde{f}^+, \tilde{w}_h \rangle_{\mathbb{H}_0^1(\Omega)} - k_B \mathcal{T} \int_{\Omega} C(M \hat{\psi}_{\delta,h}^{\Delta t,+}) : \nabla_x \tilde{w}_h \tilde{d}x \right] \tilde{d}t \\ & \quad \forall \tilde{w}_h \in \mathbb{L}^{\frac{4}{4-\vartheta}}(0, T; \mathbb{V}_h), \end{aligned} \quad (6.4.82)$$

where  $\vartheta$  is as defined in (6.4.73). Similarly, (6.4.32b) multiplied by  $\Delta t^n$  and summed for  $n = 1 \rightarrow N$  can be restated as:

$$\begin{aligned} & \int_0^T \int_{\Omega \times D} M \pi_h \left[ \frac{\partial \hat{\psi}_{\delta,h}^{\Delta t}}{\partial t} \hat{\varphi}_h \right] \tilde{d}q \tilde{d}x \tilde{d}t - \int_0^T \int_{\Omega \times D} M \tilde{u}_{\delta,h}^{\Delta t,+} \cdot \tilde{\pi}_h \left[ \Xi_{\delta}^x(\hat{\psi}_{\delta,h}^{\Delta t,+}) \nabla_x \hat{\varphi}_h \right] \tilde{d}q \tilde{d}x \tilde{d}t \\ & + \int_0^T \int_{\Omega \times D} M \pi_h \left[ \frac{1}{2\lambda} \nabla_q \hat{\psi}_{\delta,h}^{\Delta t,+} \cdot \nabla_q \hat{\varphi}_h + \varepsilon \nabla_x \hat{\psi}_{\delta,h}^{\Delta t,+} \cdot \nabla_x \hat{\varphi}_h \right] \tilde{d}q \tilde{d}x \tilde{d}t \\ & - \int_0^T \int_{\Omega \times D} M (\nabla_x \tilde{u}_{\delta,h}^{\Delta t,+} q) \cdot \tilde{\pi}_h \left[ \Xi_{\delta}^q(\hat{\psi}_{\delta,h}^{\Delta t,+}) \nabla_q \hat{\varphi}_h \right] \tilde{d}q \tilde{d}x \tilde{d}t = 0 \\ & \quad \forall \hat{\varphi}_h \in \mathbb{L}^2(0, T; \hat{\mathbb{X}}_h). \end{aligned} \quad (6.4.83)$$

It follows from (6.4.72a–c), (6.4.80a–c), (6.2.22) and (6.4.46c,d) that

$$\begin{aligned} & \sup_{t \in (0, T)} \left[ \|\tilde{u}_{\delta,h}^{\Delta t(\pm)}\|_{\mathbb{L}^2(\Omega)}^2 \right] + \frac{1}{\delta} \sup_{t \in (0, T)} \left[ (M, \pi_h [ [\hat{\psi}_{\delta,h}^{\Delta t(\pm)} ]_-^2 ] )_{\Omega \times D} \right] + \int_0^T \|\nabla_x \tilde{u}_{\delta,h}^{\Delta t(\pm)}\|_{\mathbb{L}^2(\Omega)}^2 \tilde{d}t \\ & + \delta \int_0^T \left[ (M, |\nabla_x (\pi_h [ [\mathcal{F}_{\delta}^L]'(\hat{\psi}_{\delta,h}^n)])|^2 + |\nabla_q (\pi_h [ [\mathcal{F}_{\delta}^L]'(\hat{\psi}_{\delta,h}^n)])|^2 )_{\Omega \times D} \right] \tilde{d}t \\ & + \int_0^T \int_{\Omega} \frac{|\tilde{u}_{\delta,h}^{\Delta t,+} - \tilde{u}_{\delta,h}^{\Delta t,-}|^2}{\Delta(t)} \tilde{d}x \tilde{d}t \leq C \end{aligned} \quad (6.4.84a)$$

and

$$\begin{aligned}
& \sup_{t \in (0, T)} \left[ (M, |\hat{\psi}_{\delta, h}^{\Delta t, (\pm)}|^2)_{\Omega \times D} \right] + \int_0^T (M, |\nabla_q \hat{\psi}_{\delta, h}^{\Delta t, +}|^2 + |\nabla_x \hat{\psi}_{\delta, h}^{\Delta t, +}|^2)_{\Omega \times D} dt \\
& + \int_0^T \left[ \int_{\Omega \times D} M \frac{|\hat{\psi}_{\delta, h}^{\Delta t, +} - \hat{\psi}_{\delta, h}^{\Delta t, -}|^2}{\Delta(t)} dq dx \right] dt + \sup_{t \in (0, T)} \left[ \|C(\hat{\psi}_{\delta, h}^{\Delta t, (\pm)})\|_{L^2(\Omega)}^2 \right] \\
& + \int_0^T \left\| \tilde{S} \frac{\partial u_{\delta, h}^{\Delta t}}{\partial t} \right\|_{\mathbb{H}^1(\Omega)}^{\frac{4}{\vartheta}} dt + \int_0^T \left\| \tilde{\mathcal{G}} \frac{\partial \hat{\psi}_{\delta, h}^{\Delta t}}{\partial t} \right\|_{\hat{X}}^2 dt \leq C(L, T), \tag{6.4.84b}
\end{aligned}$$

where  $\vartheta$  is as defined in (6.4.73). In the above and throughout, the notation  $u_{\delta, h}^{\Delta t, (\pm)}$  means  $u_{\delta, h}^{\Delta t}$  with or without the superscripts  $\pm$ , and similarly  $\hat{\psi}_{\delta, h}^{\Delta t, (\pm)}$ .

Before proving a convergence result for  $(P_\delta^{h, \Delta t})$ , we need the following result.

**Lemma 6.4.9** *For all  $\kappa_x \in \mathcal{T}_h^x$ ,  $\kappa_q \in \mathcal{T}_h^q$  and for all  $\hat{\varphi}_h \in \hat{X}_h$  we have that*

$$\begin{aligned}
& \int_{\kappa_x \times \kappa_q} M |\Xi_\delta^x(\hat{\varphi}_h) - \beta^L(\hat{\varphi}_h) I|^2 dq dx \\
& \leq C \left( \delta^2 + h_x^2 \int_{\kappa_x \times \kappa_q} M |\nabla_x \hat{\varphi}_h|^2 dq dx + \int_{\kappa_x \times \kappa_q} M \pi_{h, \kappa_x \times \kappa_q} [[\hat{\varphi}_h]_-]^2 dq dx \right), \tag{6.4.85a}
\end{aligned}$$

$$\begin{aligned}
& \int_{\kappa_x \times \kappa_q} M |\Xi_\delta^q(\hat{\varphi}_h) - \beta^L(\hat{\varphi}_h) I|^2 dq dx \\
& \leq C \left( \delta^2 + h_q^2 \int_{\kappa_x \times \kappa_q} M |\nabla_q \hat{\varphi}_h|^2 dq dx + \int_{\kappa_x \times \kappa_q} M \pi_{h, \kappa_x \times \kappa_q} [[\hat{\varphi}_h]_-]^2 dq dx \right). \tag{6.4.85b}
\end{aligned}$$

**Proof.** Firstly, we have from (6.4.21), (6.4.19), (6.2.17) and (6.4.53b) that

$$\begin{aligned}
& \int_{\kappa_x \times \kappa_q} M |\Xi_\delta^x(\hat{\varphi}_h) - \beta_\delta^L(\hat{\varphi}_h) I|^2 dq dx \leq \left( \int_{\kappa_x \times \kappa_q} M dq dx \right) \|\Lambda_\delta^x(\hat{\varphi}_h) - \beta_\delta^L(\hat{\varphi}_h) I\|_{L^\infty(\kappa_x \times \kappa_q)}^2 \\
& \leq C h_x^2 \left( \int_{\kappa_x \times \kappa_q} M dq dx \right) \|\nabla_x [\beta_\delta^L(\hat{\varphi}_h)]\|_{L^\infty(\kappa_x \times \kappa_q)}^2 \\
& \leq C h_x^2 \int_{\kappa_x \times \kappa_q} M |\nabla_x \hat{\varphi}_h|^2 dq dx. \tag{6.4.86}
\end{aligned}$$

Similarly, we have from (6.4.21), (6.4.19), (6.2.17) and (6.4.51a) that

$$\int_{\kappa_x \times \kappa_q} M |\Xi_\delta^q(\hat{\varphi}_h) - \beta_\delta^L(\hat{\varphi}_h) I|^2 dq dx \leq C h_q^2 \int_{\kappa_x \times \kappa_q} M |\nabla_q \hat{\varphi}_h|^2 dq dx. \tag{6.4.87}$$

Next we note from (6.2.17) and (6.2.25) that, for all  $s \in \mathbb{R}$ ,

$$|\beta_\delta^L(s) - \beta^L(s)| \leq \delta - [s]_-. \tag{6.4.88}$$

In addition, we note that

$$[\hat{\varphi}_h]_{-}(x, q) \geq \pi_{h, \kappa_x \times \kappa_q} [[\hat{\varphi}_h]_{-}](\tilde{x}, \tilde{q}) \quad \forall (x, q) \in \kappa_x \times \kappa_q. \quad (6.4.89)$$

Hence (6.4.88), (6.4.89) and (6.4.46c) yield that

$$\begin{aligned} \int_{\kappa_x \times \kappa_q} M |\beta_{\delta}^L(\hat{\varphi}_h) - \beta^L(\hat{\varphi}_h)|^2 dq dx &\leq \int_{\kappa_x \times \kappa_q} M |\delta - [\hat{\varphi}_h]_{-}|^2 dq dx \\ &\leq \int_{\kappa_x \times \kappa_q} M |\delta - \pi_{h, \kappa_x \times \kappa_q} [\hat{\varphi}_h]_{-}|^2 dq dx \\ &\leq C \left[ \delta^2 + \int_{\kappa_x \times \kappa_q} M |\pi_{h, \kappa_x \times \kappa_q} [\hat{\varphi}_h]_{-}|^2 dq dx \right]. \end{aligned} \quad (6.4.90)$$

Combining (6.4.86), (6.4.87) and (6.4.90) yields the desired results (6.4.85a,b).  $\square$

We are now in a position to prove the following convergence result for  $(P_{\delta}^{h, \Delta t})$ .

**Theorem 6.4.10** *There exists a subsequence of  $\{ \{ u_{\delta, h}^{\Delta t}, \hat{\psi}_{\delta, h}^{\Delta t} \} \}_{\delta > 0, h > 0, \Delta t > 0}$ , and functions  $u \in L^{\infty}(0, T; \underline{L}^2(\Omega)) \cap L^2(0, T; \underline{V}) \cap W^{1, \frac{4}{\vartheta}}(0, T; \underline{V}')$  and  $\hat{\psi} \in L^{\infty}(0, T; L_M^2(\Omega \times D)) \cap L^2(0, T; \hat{X}) \cap H^1(0, T; \hat{X}')$  with  $\hat{\psi} \geq 0$  a.e. in  $\Omega \times D \times (0, T)$  such that, as  $\delta, h, \Delta t \rightarrow 0_+$ ,*

$$u_{\delta, h}^{\Delta t(\pm)} \rightharpoonup u \quad \text{weak* in } L^{\infty}(0, T; \underline{L}^2(\Omega)), \quad (6.4.91a)$$

$$u_{\delta, h}^{\Delta t(\pm)} \rightharpoonup u \quad \text{weakly in } L^2(0, T; H_0^1(\Omega)), \quad (6.4.91b)$$

$$S \frac{\partial u_{\delta, h}}{\partial t} \rightharpoonup S \frac{\partial u}{\partial t} \quad \text{weakly in } L^{\frac{4}{\vartheta}}(0, T; \underline{V}), \quad (6.4.91c)$$

$$u_{\delta, h}^{\Delta t(\pm)} \rightarrow u \quad \text{strongly in } L^2(0, T; L^r(\Omega)), \quad (6.4.91d)$$

and

$$M^{\frac{1}{2}} \hat{\psi}_{\delta, h}^{\Delta t(\pm)} \rightharpoonup M^{\frac{1}{2}} \hat{\psi} \quad \text{weak* in } L^{\infty}(0, T; L^2(\Omega \times D)), \quad (6.4.92a)$$

$$M^{\frac{1}{2}} \nabla_q \hat{\psi}_{\delta, h}^{\Delta t, +} \rightharpoonup M^{\frac{1}{2}} \nabla_q \hat{\psi} \quad \text{weakly in } L^2(0, T; L^2(\Omega \times D)), \quad (6.4.92b)$$

$$M^{\frac{1}{2}} \nabla_x \hat{\psi}_{\delta, h}^{\Delta t, +} \rightharpoonup M^{\frac{1}{2}} \nabla_x \hat{\psi} \quad \text{weakly in } L^2(0, T; L^2(\Omega \times D)), \quad (6.4.92c)$$

$$\mathcal{G} \frac{\partial \hat{\psi}_{\delta, h}^{\Delta t}}{\partial t} \rightharpoonup \mathcal{G} \frac{\partial \hat{\psi}}{\partial t} \quad \text{weakly in } L^2(0, T; \hat{X}), \quad (6.4.92d)$$

$$M^{\frac{1}{2}} \hat{\psi}_{\delta, h}^{\Delta t(\pm)} \rightarrow M^{\frac{1}{2}} \hat{\psi} \quad \text{strongly in } L^2(0, T; L^2(\Omega \times D)), \quad (6.4.92e)$$

$$M^{\frac{1}{2}} \Xi_{\delta}^x(\hat{\psi}_{\delta, h}^{\Delta t(\pm)}) \rightarrow M^{\frac{1}{2}} \beta^L(\hat{\psi}) I \quad \text{strongly in } L^2(0, T; L^2(\Omega \times D)), \quad (6.4.92f)$$

$$M^{\frac{1}{2}} \Xi_{\delta}^q(\hat{\psi}_{\delta, h}^{\Delta t(\pm)}) \rightarrow M^{\frac{1}{2}} \beta^L(\hat{\psi}) I \quad \text{strongly in } L^2(0, T; L^2(\Omega \times D)), \quad (6.4.92g)$$

$$C(M \hat{\psi}_{\delta, h}^{\Delta t(\pm)}) \rightharpoonup C(M \hat{\psi}) \quad \text{strongly in } L^2(0, T; L^2(\Omega)); \quad (6.4.92h)$$

where  $\vartheta$  is defined by (6.4.73) and  $r \in [1, \infty)$  if  $d = 2$  and  $r \in [1, 6)$  if  $d = 3$ .

Furthermore,  $\{u, \hat{\psi}\}$  solves the following problem:

(P) Find functions

$$u \in L^\infty(0, T; \underline{L}^2(\Omega)) \cap L^2(0, T; \underline{Y}) \cap W^{1, \frac{4}{\vartheta}}(0, T; \underline{Y}')$$

and

$$\hat{\psi} \in L^\infty(0, T; L_M^2(\Omega \times D)) \cap L^2(0, T; \hat{X}) \cap H^1(0, T; \hat{X}'),$$

with  $\hat{\psi} \geq 0$  a.e. in  $\Omega \times D \times (0, T)$  and  $\underline{C}(M \hat{\psi}) \in L^\infty(0, T; \underline{L}^2(\Omega))$ , such that  $\underline{u}(\cdot, 0) = \underline{u}^0(\cdot)$ ,  $\hat{\psi}(\cdot, \cdot, 0) = \hat{\psi}^0(\cdot, \cdot)$  and

$$\begin{aligned} & \int_0^T \left\langle \frac{\partial \underline{u}}{\partial t}, \underline{w} \right\rangle_{\underline{V}} dt + \int_{\Omega_T} \left[ \left[ (\underline{u} \cdot \underline{\nabla}_x) \underline{u} \right] \cdot \underline{w} + \nu \underline{\nabla}_x \underline{u} : \underline{\nabla}_x \underline{w} \right] dx dt \\ & = \int_0^T \langle f, \underline{w} \rangle_{\underline{H}_0^1(\Omega)} dt - k_B \mathcal{T} \int_{\Omega_T} \underline{C}(M \hat{\psi}) : \underline{\nabla}_x \underline{w} dx dt \quad \forall \underline{w} \in L^{\frac{4}{4-\vartheta}}(0, T; \underline{V}); \end{aligned} \quad (6.4.93a)$$

$$\begin{aligned} & \int_0^T \left\langle \frac{\partial \hat{\psi}}{\partial t}, \hat{\varphi} \right\rangle_{\hat{X}} dt + \int_0^T \int_{\Omega \times D} M \left[ \varepsilon \underline{\nabla}_x \hat{\psi} - \underline{u} \beta^L(\hat{\psi}) \right] \cdot \underline{\nabla}_x \hat{\varphi} dq dx dt \\ & + \int_0^T \int_{\Omega \times D} \left[ \frac{M}{2\lambda} \underline{\nabla}_q \hat{\psi} - (\underline{\nabla}_x \underline{u}) q \beta^L(\hat{\psi}) \right] \cdot \underline{\nabla}_q \hat{\varphi} dq dx dt = 0 \quad \forall \hat{\varphi} \in L^2(0, T; \hat{X}). \end{aligned} \quad (6.4.93b)$$

**Proof.** The results (6.4.91a-c) follow immediately from the bounds (6.4.84a,b) on noting the notation (6.4.80a-c). The denseness of  $\bigcup_{h>0} R_h$  in  $L^2(\Omega)$  and (6.4.3c) yield that  $\underline{u} \in L^2(0, T; \underline{V})$ . The strong convergence result (6.4.91d) for  $\underline{u}_{\delta, h}^{\Delta t}$  follows immediately from (6.4.91a-c), (6.3.3) and (6.3.14), on noting that  $\underline{V} \subset \underline{H}_0^1(\Omega)$  is compactly embedded in  $\underline{L}^r(\Omega)$  for the stated values of  $r$ . We now prove (6.4.91d) for  $\underline{u}_{\delta, h}^{\Delta t, \pm}$ . First we obtain from the bound on the last term on the left-hand side of (6.4.84a) and (6.4.81) that

$$\|\underline{u}_{\delta, h}^{\Delta t} - \underline{u}_{\delta, h}^{\Delta t, \pm}\|_{L^2(0, T; L^2(\Omega))}^2 \leq C \Delta t. \quad (6.4.94)$$

Second, we note from Sobolev embedding that, for all  $\eta \in L^2(0, T; H^1(\Omega))$ ,

$$\|\eta\|_{L^2(0, T; L^r(\Omega))} \leq \|\eta\|_{L^2(0, T; L^2(\Omega))}^\theta \|\eta\|_{L^2(0, T; L^s(\Omega))}^{1-\theta} \leq C \|\eta\|_{L^2(0, T; L^2(\Omega))}^\theta \|\eta\|_{L^2(0, T; H^1(\Omega))}^{1-\theta} \quad (6.4.95)$$

for all  $r \in [2, s)$ , with any  $s \in (2, \infty)$  if  $d = 2$  or any  $s \in (2, 6]$  if  $d = 3$ , and

$$\theta = [2(s-r)]/[r(s-2)] \in (0, 1].$$

Hence, combining (6.4.94), (6.4.95) and (6.4.91d) for  $\underline{u}_{\delta, h}^{\Delta t}$  yields (6.4.91d) for  $\underline{u}_{\delta, h}^{\Delta t, \pm}$ .

The result (6.4.92a) follows immediately from the bounds on the first and third terms on the left-hand side of (6.4.84b). It follows immediately from the bound on the second term on the left-hand side of (6.4.84b) that (6.4.92b) holds for some limit  $\underline{g} \in L^2(0, T; \underline{L}^2(\Omega \times D))$ , which we need to identify. However for any  $\underline{\eta} \in L^2(0, T; \underline{C}_0^\infty(\Omega \times D))$ , it follows from (6.2.5)

and the compact support of  $\eta$  on  $D$  that  $[\nabla_q \cdot (M^{\frac{1}{2}} \eta)]/M^{\frac{1}{2}} \in L^2(0, T; L^2(\Omega \times D))$  and hence the above convergence implies, on noting (6.4.92a), that

$$\begin{aligned} \int_0^T \int_{\Omega \times D} \underline{g} \cdot \underline{\eta} \, d\underline{q} \, d\underline{x} \, dt &\leftarrow - \int_0^T \int_{\Omega \times D} M^{\frac{1}{2}} \hat{\psi}_{\delta, h}^{\Delta t, +} \frac{\nabla_q \cdot (M^{\frac{1}{2}} \eta)}{M^{\frac{1}{2}}} \, d\underline{q} \, d\underline{x} \, dt \\ &\rightarrow - \int_0^T \int_{\Omega \times D} M^{\frac{1}{2}} \hat{\psi} \frac{\nabla_q \cdot (M^{\frac{1}{2}} \eta)}{M^{\frac{1}{2}}} \, d\underline{q} \, d\underline{x} \, dt \end{aligned} \quad (6.4.96)$$

as  $\delta, h, \Delta t \rightarrow 0_+$ . Hence the desired result (6.4.92b) follows from (6.4.96) on noting the denseness of  $C_0^\infty(\Omega \times D)$  in  $L^2(\Omega \times D)$ . Similar arguments also prove (6.4.92c,d) on noting (6.4.92a) and the second and sixth bounds in (6.4.84b). The strong convergence result (6.4.92e) for  $\hat{\psi}_{\delta, h}^{\Delta t}$  follows immediately from (6.4.92a–c), (6.3.13), (6.3.14) and (6.3.11b). Similarly to (6.4.94), the third bound in (6.4.84b) then yields that (6.4.92e) holds for  $\hat{\psi}_{\delta, h}^{\Delta t(\pm)}$ . The desired results (6.4.92f,g) follow immediately from (6.4.85a,b) the second bounds in (6.4.84a,b), (6.2.25) and (6.4.92e). The desired result (6.4.92h) follows immediately from (6.4.92a), (6.2.3) and (6.3.15). Finally, the nonnegativity of  $\hat{\psi}$  follows from (6.4.92e) and the second bound in (6.4.84a).

It remains to prove that  $\{u, \hat{\psi}\}$  solve (P). It follows from (6.4.5), (6.4.84a,b), (6.4.91a–d), (6.4.92h), (6.4.31b), (6.3.2) and (6.4.8) that we may pass to the limit,  $\delta, h, \Delta t \rightarrow 0_+$ , in (6.4.82) to obtain that  $u \in L^\infty(0, T; \underline{L}^2(\Omega)) \cap L^2(0, T; \underline{Y}) \cap W^{1, \frac{4}{\vartheta}}(0, T; \underline{Y}')$  and  $\underline{C}(M \hat{\psi}) \in L^\infty(0, T; \underline{L}^2(\Omega))$  satisfy (6.4.93a). It also follows from (6.4.28a), (6.4.5), (6.4.84a) and (6.4.91d) that  $u(\cdot, 0) = u^0(\cdot)$  in the required sense; recall Remark 6.3.1.

It follows from (6.4.92a–g), (6.4.91b,d), (6.3.12), (6.4.59a–c), (6.4.64), (6.4.84a,b), (6.4.48a,b), (6.4.49) and (6.4.50a,b) that we may pass to the limit  $\delta, h, \Delta t \rightarrow 0_+$  in (6.4.83) with  $\hat{\varphi}_h = \pi_h \hat{\varphi}$  to obtain (6.4.93b) for any  $\hat{\varphi} \in C_0^\infty(0, T; C(\overline{\Omega \times D}))$ . In order to pass to the limit on the first term in (6.4.83), we note that

$$\begin{aligned} &\int_0^T \int_{\Omega \times D} M \pi_h \left[ \frac{\partial \hat{\psi}_{\delta, h}^{\Delta t}}{\partial t} [\pi_h \hat{\varphi}] \right] \, d\underline{q} \, d\underline{x} \, dt \\ &= \int_0^T \int_{\Omega \times D} M \frac{\partial \hat{\psi}_{\delta, h}^{\Delta t}}{\partial t} [\pi_h \hat{\varphi}] \, d\underline{q} \, d\underline{x} \, dt + \int_0^T \int_{\Omega \times D} M (I - \pi_h) \left[ \hat{\psi}_{\delta, h}^{\Delta t} \frac{\partial [\pi_h \hat{\varphi}]}{\partial t} \right] \, d\underline{q} \, d\underline{x} \, dt. \end{aligned} \quad (6.4.97)$$

The desired result (6.4.93b) then follows from noting that  $C_0^\infty(0, T; C(\overline{\Omega \times D}))$  is dense in  $L^2(0, T; \hat{X})$ , on recalling (6.3.8). Finally, it follows from (6.4.28b), (6.4.59c), (6.4.51c), (6.4.53a), (6.4.50a,b), (6.3.8), (6.4.84b) and (6.4.92e) that  $\hat{\psi}(\cdot, \cdot, 0) = \hat{\psi}^0(\cdot, \cdot)$  in the required sense; recall Remark 6.3.1.  $\square$

**Remark 6.4.11** We note that (P), (6.4.93a,b), differs slightly from (P<sup>L</sup>), (6.3.16a,b), in that  $u \in W^{1, \frac{4}{\vartheta}}(0, T, \underline{Y}')$  for the stated value of  $\vartheta$ , recall (6.4.73), is slightly weaker than  $u^L \in W^{1, \frac{4}{d}}(0, T, \underline{Y}')$  in the case  $d = 2$  with the subsequent slight strengthening of the regularity of the test functions in (6.4.93a). In addition,  $\hat{\psi}^L$  in the convective term in (6.3.16b) is replaced by  $\beta^L(\hat{\psi})$  in (6.4.93b). It does not appear possible to construct a variation of the finite element approximation (P <sub>$\delta$</sub>  <sup>$h, \Delta t$</sup> ) that converges to the former version of the convective

term, and at the same time converges to the other terms in (6.4.93b). The presence of the cut-off  $\beta^L(\cdot)$  in this convective term improves the regularity in time of  $\hat{\psi}$  in (6.4.93a,b), to that in (6.3.16a,b), and hence the weakening of the regularity in time of the test functions in (6.4.93b).  $\diamond$

**Remark 6.4.12** Finally, it follows from (6.4.84a) and (6.4.91a,b) that

$$\sup_{t \in (0, T)} \left[ \|u\|_{L^2(\Omega)}^2 \right] + \int_0^T \|\nabla_x u\|_{L^2(\Omega)}^2 dt \leq C. \tag{6.4.98}$$

Hence, although we have introduced a cut-off  $L \gg 1$  to  $\hat{\psi}$  in the drag and convective terms, and added diffusion in the  $\underline{x}$  direction with a positive coefficient  $\varepsilon \ll 1$  in the Fokker–Planck equation compared to the standard polymer model; the bound (6.4.98) on  $u$ , the variable of real physical interest, is independent of the parameters  $L$  and  $\varepsilon$ .  $\diamond$

## 6.5 Appendix: Maxwellian Sobolev norm quasi-interpolation

The aim of this Appendix is to prove the stability result (6.4.68). To do so, we first need to show certain quasi-interpolation results in Maxwellian-weighted Sobolev spaces. The starting point for the construction of the relevant quasi-interpolation operators is the Brascamp–Lieb inequality stated below.

Suppose that  $D$  is a convex open set,  $D \subset \mathbb{R}^d$  (e.g., a bounded open ball in  $\mathbb{R}^d$  centred at the origin; or, more specifically, in the case of the FENE model,  $D = B(\underline{0}; b^{\frac{1}{2}})$ ,  $b > 2$ ). Consider a probability measure  $\mu$  supported on  $D$  with density  $e^{-V(\underline{q})}$ ,  $\underline{q} \in D$ , with respect to the Lebesgue measure  $d\underline{q}$  on  $\mathbb{R}^d$ , where  $V$  is a convex function on  $D$ ;  $\mu$  is usually referred to as a *Gibbs measure*. In particular,

$$\mu(B) = \int_B d\mu = \int_B e^{-V(\underline{q})} d\underline{q},$$

for any  $\mu$ -measurable set  $B \subset D$ , with  $\mu(D) = 1$ . The following geometric functional inequality comes from the paper of Bobkov & Ledoux [24].

**Theorem 6.5.1 (Brascamp–Lieb inequality)** *Assume that  $V$  is a twice continuously differentiable and convex function on a convex open set  $D \subset \mathbb{R}^d$ , such that, for each  $\underline{q} \in D$ , the Hessian*

$$\underline{H}(\underline{q}) := \left( \frac{\partial^2 V(\underline{q})}{\partial q_i \partial q_j} \right)$$

*is positive definite. Then, for any sufficiently smooth function  $f$ ,*

$$\text{Var}_\mu[f] := \mathbb{E}_\mu[(f - \mathbb{E}_\mu[f])^2] \leq \int_D [\underline{H}^{-1}(\underline{q}) \nabla_{\underline{q}} f] \cdot \nabla_{\underline{q}} f d\mu, \quad \text{where } \mathbb{E}_\mu[f] = \int_D f d\mu.$$

In terms of simpler notation, the Brascamp–Lieb inequality can be restated as follows:

$$\int_D \left[ f(\underline{q}) - \int_D f(\underline{p}) e^{-V(\underline{p})} d\underline{p} \right]^2 e^{-V(\underline{q})} d\underline{q} \leq \int_D [\underline{H}^{-1}(\underline{q}) \nabla_{\underline{q}} f] \cdot \nabla_{\underline{q}} f e^{-V(\underline{q})} d\underline{q},$$

for any sufficiently smooth function  $f$ .

### 6.5.1 The univariate case

Suppose that  $d = 1$ ,  $D := (0, q_1) \subset \mathbb{R}$ , and  $V(q) := \ln \left[ \frac{q_1}{\alpha+1} \left( \frac{q_1}{q} \right)^\alpha \right]$  with  $\alpha > 0$ . Clearly,  $\int_D e^{-V(q)} dq = 1$ . By the Brascamp–Lieb inequality,

$$\int_0^{q_1} \left[ f(q) - \frac{\alpha+1}{q_1^{\alpha+1}} \int_0^{q_1} f(p) p^\alpha dp \right]^2 q^\alpha dq \leq \frac{1}{\alpha} \int_0^{q_1} |f'(q)|^2 q^{\alpha+2} dq \leq \frac{q_1^2}{\alpha} \int_0^{q_1} |f'(q)|^2 q^\alpha dq. \quad (6.5.1)$$

Let us consider the nonuniform partition  $0 = q_0 < q_1 < \dots < q_N = 1$  of the interval  $[0, 1]$ , with  $h_q := \max_{k=1 \rightarrow N} (q_k - q_{k-1})$ , and let  $\hat{X}_h^q$  denote the set of all continuous piecewise linear functions defined on this partition. For  $m \in \mathbb{Z}_{\geq 0}$  and a nonempty open interval  $(a, b) \subset \mathbb{R}_{>0}$ , let

$$\mathbf{H}^m((a, b); q^\alpha) := \left\{ \hat{\varphi} \in \mathbf{H}_{\text{loc}}^m(a, b) : \|\hat{\varphi}\|_{\mathbf{H}^m((a, b); q^\alpha)}^2 := \sum_{k=0}^m \int_a^b |\hat{\varphi}^{(k)}(q)|^2 q^\alpha dq < \infty \right\}.$$

When  $m = 0$ , we write  $\mathbf{L}^2((a, b); q^\alpha)$  instead of  $\mathbf{H}^0((a, b); q^\alpha)$ .

For  $\hat{\psi} \in \mathbf{H}^1((0, 1); q^\alpha)$ , let  $I_h^q \hat{\psi} \in \hat{X}_h^q$  denote the continuous piecewise linear (quasi-)interpolant of  $\hat{\psi}$ , defined by

$$(I_h^q \hat{\psi})(q) := \begin{cases} \hat{\psi}(q_1) + (q - q_1) \frac{\alpha+1}{q_1^{\alpha+1}} \int_0^{q_1} \hat{\psi}'(p) p^\alpha dp, & q \in [0, q_1], \\ \frac{\hat{\psi}(q_k) - \hat{\psi}(q_{k-1})}{q_k - q_{k-1}} (q - q_{k-1}) + \hat{\psi}(q_{k-1}), & q \in [q_{k-1}, q_k], \quad k = 2 \rightarrow N. \end{cases}$$

We note that since  $\mathbf{H}^1((0, 1); q^\alpha) \subset \mathbf{C}(0, 1]$ , the definition is meaningful. Observe, further, that  $(I_h^q \hat{\psi})(q_k) = \hat{\psi}(q_k)$ ,  $k = 1 \rightarrow N$ ; i.e. the function  $I_h^q \hat{\psi}$  interpolates  $\hat{\psi}$  at  $q = q_k$ ,  $k = 1 \rightarrow N$ , but not at  $q = q_0 = 0$ . In the interval  $[0, q_1]$  the function  $I_h^q \hat{\psi}$  has been chosen so as to ensure that  $(I_h^q \hat{\psi})'(q) = \frac{\alpha+1}{q_1^{\alpha+1}} \int_0^{q_1} \hat{\psi}'(p) p^\alpha dp$  and  $(I_h^q \hat{\psi})(q_1) = \hat{\psi}(q_1)$ . Hence, on applying the inequality (6.5.1),

$$\int_0^{q_1} \left[ \hat{\psi}'(q) - (I_h^q \hat{\psi})'(q) \right]^2 q^\alpha dq \leq \frac{q_1^2}{\alpha} \int_0^{q_1} |\hat{\psi}''(q)|^2 q^\alpha dq.$$

On the remaining subintervals in the partition, using  $q_{k-1}^\alpha \leq q^\alpha \leq q_k^\alpha$  and a standard error bound for the linear interpolant of  $\hat{\psi} \in \mathbf{H}^2(q_{k-1}, q_k)$ ,  $k = 2 \rightarrow N$ , we have that

$$\int_{q_{k-1}}^{q_k} \left[ \hat{\psi}'(q) - (I_h^q \hat{\psi})'(q) \right]^2 q^\alpha dq \leq \left( \frac{q_k}{q_{k-1}} \right)^\alpha \frac{(q_k - q_{k-1})^2}{\pi^2} \int_{q_{k-1}}^{q_k} |\hat{\psi}''(q)|^2 q^\alpha dq, \quad k = 2 \rightarrow N.$$

On summing our bounds through  $k = 1 \rightarrow N$  and noting that  $q_1 \leq h_q$  and  $q_k - q_{k-1} \leq h_q$  for  $k = 1 \rightarrow N$ , we obtain

$$\int_0^1 \left[ \hat{\psi}'(q) - (I_h^q \hat{\psi})'(q) \right]^2 q^\alpha dq \leq \max \left( \frac{h_q^2}{\alpha}, \max_{k=2 \rightarrow N} \left( \frac{q_k}{q_{k-1}} \right)^\alpha \frac{h_q^2}{\pi^2} \right) \int_0^1 |\hat{\psi}''(q)|^2 q^\alpha dq.$$

We shall henceforth assume that the partition  $0 = q_0 < q_1 < \dots < q_N = 1$  is such that there exists a fixed constant  $C_0 > 1$  such that

$$\max_{k=2 \rightarrow N} \frac{q_k}{q_{k-1}} \leq C_0. \quad (6.5.2)$$

Now, letting  $C_\alpha := \max(\frac{1}{\alpha}, \frac{1}{\pi^2} C_0^\alpha)$ , we get

$$\int_0^1 \left[ \hat{\psi}'(q) - (I_h^q \hat{\psi})'(q) \right]^2 q^\alpha dq \leq C_\alpha h_q^2 \int_0^1 |\hat{\psi}''(q)|^2 q^\alpha dq. \quad (6.5.3)$$

We note the weighted Poincaré inequality for all  $\hat{v} \in H^1((0, 1); q^\alpha)$  with  $\hat{v}(1) = 0$

$$\begin{aligned} \int_0^1 |\hat{v}(q)|^2 q^\alpha dq &= \int_0^1 \left( \int_q^1 \hat{v}'(t) t^{\frac{\alpha}{2}} t^{-\frac{\alpha}{2}} dt \right)^2 q^\alpha dq \\ &\leq \left( \int_0^1 q^\alpha \left( \int_q^1 t^{-\alpha} dt \right) dq \right) \int_0^1 |\hat{v}'(q)|^2 q^\alpha dq = \frac{1}{2(\alpha+1)} \int_0^1 |\hat{v}'(q)|^2 q^\alpha dq, \end{aligned} \quad (6.5.4)$$

which, in fact, holds for any  $\alpha > -1$ . Applying (6.5.4) with  $\hat{v} = \hat{\psi} - I_h^q \hat{\psi}$ , and noting (6.5.3) we deduce that

$$\int_0^1 \left[ \hat{\psi}(q) - (I_h^q \hat{\psi})(q) \right]^2 q^\alpha dq \leq \frac{C_\alpha}{2(\alpha+1)} h_q^2 \int_0^1 |\hat{\psi}''(q)|^2 q^\alpha dq,$$

and therefore

$$\|\hat{\psi} - I_h^q \hat{\psi}\|_{H^1((0,1);q^\alpha)}^2 \leq C_\alpha \left( 1 + \frac{1}{2(\alpha+1)} \right) h_q^2 \|\hat{\psi}''\|_{L^2((0,1);q^\alpha)}^2. \quad (6.5.5)$$

Let  $P_h^q$  denote the orthogonal projector from  $H^1((0, 1); q^\alpha)$  onto  $\hat{X}_h^q$  with respect to the  $q^\alpha$ -weighted  $H^1(0, 1)$  inner product

$$a(\hat{\psi}, \hat{\varphi}) := \int_0^1 \hat{\psi}'(q) \hat{\varphi}'(q) q^\alpha dq + \int_0^1 \hat{\psi}(q) \hat{\varphi}(q) q^\alpha dq,$$

where  $\alpha > 0$ . That is,

$$a(\hat{\psi} - P_h^q \hat{\psi}, \hat{\varphi}_h) = 0 \quad \forall \hat{\varphi}_h \in \hat{X}_h^q. \quad (6.5.6)$$

Now, consider the following boundary-value problem: Find  $\hat{z} \in H^1((0, 1); q^\alpha)$  such that

$$a(\hat{\varphi}, \hat{z}) = \ell(\hat{\varphi}) \quad \forall \hat{\varphi} \in H^1((0, 1); q^\alpha), \quad (6.5.7)$$

where

$$\ell(\hat{\varphi}) := \int_0^1 \hat{g}(q) \hat{\varphi}(q) q^\alpha dq, \quad \text{with} \quad \hat{g} := \hat{\psi} - P_h^q \hat{\psi}.$$

The existence of a unique weak solution  $\hat{z} \in H^1((0, 1); q^\alpha)$  to (6.5.7) follows from the Lax–Milgram theorem. Hence, on taking  $\hat{\varphi} = \hat{z}$  in (6.5.7), we obtain

$$\begin{aligned} \|\hat{z}\|_{H^1((0,1);q^\alpha)}^2 &= a(\hat{z}, \hat{z}) \leq \|\hat{z}\|_{L^2((0,1);q^\alpha)} \|\hat{\psi} - P_h^q \hat{\psi}\|_{L^2((0,1);q^\alpha)} \\ &\leq \|\hat{z}\|_{H^1((0,1);q^\alpha)} \|\hat{\psi} - P_h^q \hat{\psi}\|_{L^2((0,1);q^\alpha)}, \end{aligned}$$

and therefore

$$\|\hat{z}\|_{H^1((0,1);q^\alpha)} \leq \|\hat{\psi} - P_h^q \hat{\psi}\|_{L^2((0,1);q^\alpha)}.$$

Problem (6.5.7) is the weak form of the following boundary value problem:

$$-\hat{z}'' - \frac{\alpha}{q} \hat{z}' + \hat{z} = \hat{\psi} - P_h^q \hat{\psi}, \quad q \in (0, 1), \quad \lim_{q \rightarrow 0^+} q^\alpha \hat{z}'(q) = 0, \quad \hat{z}'(1) = 0.$$

Formally differentiating this equation, multiplying the resulting equation by  $\hat{z}' q^\alpha$ , integrating over  $q \in (0, 1)$  and integrating by parts in the first term on the left-hand side and on the right-hand side yields

$$\begin{aligned} & \int_0^1 |\hat{z}''|^2 q^\alpha dq + \alpha \int_0^1 |\hat{z}'|^2 q^{\alpha-2} dq + \int_0^1 |\hat{z}'|^2 q^\alpha dq \\ &= - \int_0^1 (\hat{\psi} - P_h^q \hat{\psi}) \hat{z}'' q^\alpha dq - \alpha \int_0^1 (\hat{\psi} - P_h^q \hat{\psi}) \hat{z}' q^{\alpha-1} dq. \end{aligned}$$

This formal argument can be made rigorous by replacing  $q^\alpha$  with  $(q + \delta)^\alpha$ ,  $\delta > 0$ , in the definitions of  $a(\cdot, \cdot)$  and  $\ell(\cdot)$  above, and passing to the limit  $\delta \rightarrow 0_+$ ; we refer to Section 6.5.6 for the details of an analogous, but rigorous, multidimensional argument. Hence,

$$\|\hat{z}''\|_{L^2((0,1);q^\alpha)}^2 + \alpha \|\hat{z}'\|_{L^2((0,1);q^{\alpha-2})}^2 + 2 \|\hat{z}'\|_{L^2((0,1);q^\alpha)}^2 \leq (1 + \alpha) \|\hat{\psi} - P_h^q \hat{\psi}\|_{L^2((0,1);q^\alpha)}^2. \quad (6.5.8)$$

Now, by (6.5.7) with  $\hat{\varphi} = \hat{\psi} - P_h^q \hat{\psi}$ , the definition (6.5.6) of the projector  $P_h^q$ , and the bound (6.5.5),

$$\begin{aligned} \|\hat{\psi} - P_h^q \hat{\psi}\|_{L^2((0,1);q^\alpha)}^2 &= a(\hat{\psi} - P_h^q \hat{\psi}, \hat{z}) = a(\hat{\psi} - P_h^q \hat{\psi}, \hat{z} - P_h^q \hat{z}) \\ &\leq \|\hat{\psi} - P_h^q \hat{\psi}\|_{H^1((0,1);q^\alpha)} \|\hat{z} - P_h^q \hat{z}\|_{H^1((0,1);q^\alpha)} \\ &\leq \|\hat{\psi} - P_h^q \hat{\psi}\|_{H^1((0,1);q^\alpha)} \|\hat{z} - I_h^q \hat{z}\|_{H^1((0,1);q^\alpha)} \\ &\leq \|\hat{\psi} - P_h^q \hat{\psi}\|_{H^1((0,1);q^\alpha)} \left[ C_\alpha \left( 1 + \frac{1}{2(\alpha + 1)} \right) \right]^{1/2} h_q \|\hat{z}''\|_{L^2((0,1);q^\alpha)}. \end{aligned}$$

Thus, by (6.5.8),

$$\|\hat{\psi} - P_h^q \hat{\psi}\|_{L^2((0,1);q^\alpha)} \leq [C_\alpha \left( \frac{3}{2} + \alpha \right)]^{1/2} h_q \|\hat{\psi} - P_h^q \hat{\psi}\|_{H^1((0,1);q^\alpha)}; \quad (6.5.9)$$

and, denoting by  $Q_h^q$  the orthogonal projection in the inner product of  $L^2((0, 1); q^\alpha)$  onto  $\hat{X}_h^q$ , trivially

$$\|\hat{\psi} - Q_h^q \hat{\psi}\|_{L^2((0,1);q^\alpha)} \leq [C_\alpha \left( \frac{3}{2} + \alpha \right)]^{1/2} h_q \|\hat{\psi} - P_h^q \hat{\psi}\|_{H^1((0,1);q^\alpha)}. \quad (6.5.10)$$

Now,

$$\|\hat{\psi}' - (Q_h^q \hat{\psi})'\|_{L^2((0,1);q^\alpha)} \leq \|\hat{\psi}' - (P_h^q \hat{\psi})'\|_{L^2((0,1);q^\alpha)} + \|(P_h^q \hat{\psi})' - (Q_h^q \hat{\psi})'\|_{L^2((0,1);q^\alpha)}.$$

Let us, at this point, strengthen the mesh-regularity hypothesis (6.5.2) by assuming that the partition  $0 = q_0 < q_1 < \dots < q_N = 1$  is quasiuniform. Then, by the inverse inequality

$$\int_{q_{k-1}}^{q_k} |(\hat{\varphi}_h)'|^2 q^\alpha dq \leq C_{\text{inv}}^2 h_q^{-2} \int_{q_{k-1}}^{q_k} |\hat{\varphi}_h|^2 q^\alpha dq \quad \forall \hat{\varphi}_h \in \hat{X}_h^q,$$

whose proof is identical to that of the first inequality stated in (6.4.53a), we have that

$$\begin{aligned} \|\hat{\psi}' - (Q_h^q \hat{\psi})'\|_{L^2((0,1);q^\alpha)} &\leq \|\hat{\psi}' - (P_h^q \hat{\psi})'\|_{L^2((0,1);q^\alpha)} + C_{\text{inv}} h_q^{-1} \|P_h^q \hat{\psi} - Q_h^q \hat{\psi}\|_{L^2((0,1);q^\alpha)} \\ &\leq \|\hat{\psi}' - (P_h^q \hat{\psi})'\|_{L^2((0,1);q^\alpha)} + 2C_{\text{inv}} h_q^{-1} \|\hat{\psi} - P_h^q \hat{\psi}\|_{L^2((0,1);q^\alpha)}. \end{aligned}$$

This, together with (6.5.9) and (6.5.10) yields

$$\|\hat{\psi} - Q_h^q \hat{\psi}\|_{\mathbb{H}^1((0,1);q^\alpha)}^2 \leq [2 + (h_q^2 + 8C_{\text{inv}}^2) C_\alpha (\frac{3}{2} + \alpha)] \|\hat{\psi} - P_h^q \hat{\psi}\|_{\mathbb{H}^1((0,1);q^\alpha)}^2, \quad (6.5.11)$$

which in turn implies, by the triangle inequality and the fact that

$$\|P_h^q \hat{\psi}\|_{\mathbb{H}^1((0,1);q^\alpha)} \leq \|\hat{\psi}\|_{\mathbb{H}^1((0,1);q^\alpha)},$$

the existence of a positive constant  $C$ , independent of  $h$ , such that

$$\|Q_h^q \hat{\psi}\|_{\mathbb{H}^1((0,1);q^\alpha)} \leq C \|\hat{\psi}\|_{\mathbb{H}^1((0,1);q^\alpha)} \quad \forall \hat{\psi} \in \hat{\mathbb{H}}^1((0,1);q^\alpha).$$

This is the univariate counterpart of the desired stability result (6.4.68).

**Remark 6.5.2** Supposing that  $\hat{\psi} \in \mathbb{H}^2((0,1);q^\alpha)$ , we have that

$$\|\hat{\psi} - P_h^q \hat{\psi}\|_{\mathbb{H}^1((0,1);q^\alpha)}^2 \leq \|\hat{\psi} - I_h^q \hat{\psi}\|_{\mathbb{H}^1((0,1);q^\alpha)}^2 \leq C_\alpha \left(1 + \frac{1}{2(\alpha + 1)}\right) h_q^2 \int_0^1 |\hat{\psi}''(q)|^2 q^\alpha dq.$$

Thus, (6.5.11) implies that an analogous bound holds for  $\hat{\psi} - Q_h^q \hat{\psi}$  in the  $\|\cdot\|_{\mathbb{H}^1((0,1);q^\alpha)}$  norm.  $\diamond$

## 6.5.2 Multiple dimensions

In multiple space dimensions the proof of the stability result (6.4.68) proceeds in a similar manner as in the univariate case discussed above, except for two technical complications. The first is that  $D$  is ball, and therefore  $D$  has a curved boundary  $\partial D$ ; the second is that an open (possibly, curved) simplex  $\kappa_q$  in the partition of  $D$ , whose closure has nonempty intersection with  $\partial D$ , may intersect  $\partial D$  in  $d$  different configurations: with exactly one curved  $(d - k)$ -dimensional face contained in  $\partial D$ ,  $k = 1 \rightarrow d - 1$ , accounting for  $d - 1$  different configurations, and with exactly one vertex contained in  $\partial D$ , accounting for the  $d^{\text{th}}$  configuration. Each of the  $d$  possible configurations necessitates a different local definition of the quasi-interpolation operator  $I_h^q$ , which we use in the proof of the stability result (6.4.68). Since the two-dimensional case is sufficiently representative of the general argument, we shall restrict ourselves to showing (6.4.68) in the bivariate case. The proof in the case of  $d = 3$  is identical; in Section 6.5.5 we shall indicate the essential alterations that have to be made to the arguments presented herein to obtain the corresponding bounds in the case of  $d = 3$ .

## 6.5.3 Two dimensions: flat boundary

We begin by assuming that the boundary of  $D \subset \mathbb{R}^2$  is flat, e.g. that it is the straight line  $q_1 = 0$  in the  $\tilde{q} = (q_1, q_2)$ -plane. For ease of exposition we shall, intermittently, write  $x$  and  $y$  instead of  $q_1$  and  $q_2$ , i.e.  $x := q_1$  and  $y := q_2$ .

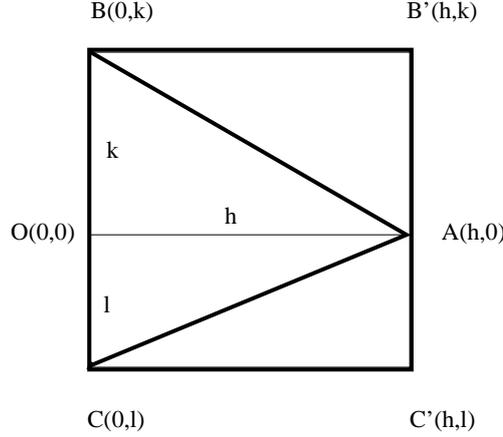


Figure 6.1: The nonobtuse open triangle  $\kappa = \triangle ABC$  in the  $(x, y) := (q_1, q_2)$ -plane, with  $A=(h, 0)$ ,  $B=(0, k)$ ,  $C=(0, l)$ , in configuration 1-flat, that is with two points, B and C, on the line  $x = 0$  along which the weight function  $(x, y) \mapsto x^\alpha$  vanishes.

**Two dimensions: configuration 1-flat.** Consider a nonobtuse open triangle  $\kappa = \triangle ABC$ , as in Figure 6.1, with  $A=(h, 0)$ ,  $B=(0, k)$ ,  $C=(0, l)$ , contained in the rectangle  $R(\kappa) := (0, h) \times (l, k) = \square B'BCC'$ , with  $B' = (h, k)$  and  $C' = (h, l)$ , where  $l \leq 0 \leq k$ ,  $k - l > 0$  and  $h > 0$ . Here, B and C belong to the line  $x = 0$  along which the weight-function  $(x, y) \mapsto x^\alpha$  vanishes;  $\alpha > 0$ . We define,

$$\hat{\Phi}(0, k) := \hat{\varphi}(h, k) - h \frac{\alpha + 1}{h^{\alpha+1}} \int_0^h \hat{\varphi}_x(x, k) x^\alpha dx$$

and

$$\hat{\Phi}(0, l) := \hat{\varphi}(h, l) - h \frac{\alpha + 1}{h^{\alpha+1}} \int_0^h \hat{\varphi}_x(x, l) x^\alpha dx.$$

We then define  $p_{\hat{\varphi}}$  as the affine function whose values at the points A, B and C are, respectively,  $\hat{\varphi}(h, 0)$ ,  $\hat{\Phi}(0, k)$  and  $\hat{\Phi}(0, l)$ . Thus,  $p_{\hat{\varphi}}$  interpolates  $\hat{\varphi}$  at A, while at the points B and C the values of  $p_{\hat{\varphi}}$  are based on extrapolating from the points B' and C', respectively, by means of the univariate quasi-interpolant  $I_h^q$ . Thus,

$$p_{\hat{\varphi}}(x, y) := \hat{\varphi}(h, 0) \frac{x}{h} + \hat{\Phi}(0, k) \left(1 - \frac{x}{h} - \frac{y}{l}\right) \frac{l}{l-k} + \hat{\Phi}(0, l) \left(1 - \frac{x}{h} - \frac{y}{k}\right) \frac{k}{k-l},$$

which implies that the partial derivatives of  $p_{\hat{\varphi}}$  with respect to  $x$  and  $y$  are:

$$(p_{\hat{\varphi}})_x(x, y) = \hat{\varphi}(h, 0) \frac{1}{h} + \hat{\Phi}(0, k) \left(-\frac{1}{h}\right) \frac{l}{l-k} + \hat{\Phi}(0, l) \left(-\frac{1}{h}\right) \frac{k}{k-l},$$

and

$$(p_{\hat{\varphi}})_y(x, y) = \hat{\Phi}(0, k) \left(-\frac{1}{l}\right) \frac{l}{l-k} + \hat{\Phi}(0, l) \left(-\frac{1}{k}\right) \frac{k}{k-l}.$$

We define the linear functionals

$$L_1(\hat{\varphi}) := \hat{\varphi}_x - (p_{\hat{\varphi}})_x \quad \text{and} \quad L_2(\hat{\varphi}) := \hat{\varphi}_y - (p_{\hat{\varphi}})_y.$$

By direct computation,  $\hat{\Phi}(0, k) = \hat{\varphi}(0, k)$  and  $\hat{\Phi}(0, l) = \hat{\varphi}(0, l)$  all  $\hat{\varphi} \in \mathbb{P}_1$ , and hence  $p_{\hat{\varphi}} \equiv \hat{\varphi}$  and  $L_i(\hat{\varphi}) \equiv 0$  for all  $\hat{\varphi} \in \mathbb{P}_1$ ,  $i = 1, 2$ . Further,

$$|\hat{\Phi}(0, k)| \leq \frac{\alpha + 1}{h^{\alpha+1}} \int_0^h |\hat{\varphi}(h, k) - h \hat{\varphi}_x(x, k)| x^\alpha dx. \quad (6.5.12)$$

Now,

$$\begin{aligned} \hat{\varphi}(h, k) h^\alpha &= \hat{\varphi}(x, k) x^\alpha + \int_x^h \frac{d}{dt} (\hat{\varphi}(t, k) t^\alpha) dt \\ &= \hat{\varphi}(x, k) x^\alpha + \int_x^h \hat{\varphi}_x(t, k) t^\alpha dt + \alpha \int_x^h \hat{\varphi}(t, k) t^{\alpha-1} dt. \end{aligned}$$

Therefore, by integration over the interval  $x \in [0, h]$ , integration by parts in the third integral on the right-hand side, and applying the Cauchy–Schwarz inequality,

$$\begin{aligned} |\hat{\varphi}(h, k)| h^{\alpha+1} &\leq \int_0^h |\hat{\varphi}(x, k)| x^\alpha dx + \int_0^h \int_x^h |\hat{\varphi}_x(t, k)| t^\alpha dt dx + \alpha \int_0^h \int_x^h |\hat{\varphi}(t, k)| t^{\alpha-1} dt dx \\ &= (\alpha + 1) \int_0^h |\hat{\varphi}(x, k)| x^\alpha dx + \int_0^h \int_x^h |\hat{\varphi}_x(t, k)| t^\alpha dt dx \\ &\leq (\alpha + 1) \left( \frac{h^{\alpha+1}}{\alpha + 1} \right)^{1/2} \left( \int_0^h |\hat{\varphi}(x, k)|^2 x^\alpha dx \right)^{1/2} \\ &\quad + h \left( \frac{h^{\alpha+1}}{\alpha + 1} \right)^{1/2} \left( \int_0^h |\hat{\varphi}_x(x, k)|^2 x^\alpha dx \right)^{1/2}. \end{aligned}$$

Thus,

$$\begin{aligned} |\hat{\varphi}(h, k)| &\leq \left( \frac{\alpha + 1}{h^{\alpha+1}} \right)^{1/2} \left( \int_0^h |\hat{\varphi}(x, k)|^2 x^\alpha dx \right)^{1/2} \\ &\quad + h^{-\alpha} \left( \frac{h^{\alpha+1}}{\alpha + 1} \right)^{1/2} \left( \int_0^h |\hat{\varphi}_x(x, k)|^2 x^\alpha dx \right)^{1/2}. \end{aligned} \quad (6.5.13)$$

To bound the first term on the right-hand side, note that, for any  $y \in [l, k]$ ,

$$|\hat{\varphi}(x, k)|^2 = |\hat{\varphi}(x, y)|^2 + 2 \int_y^k \hat{\varphi}(x, s) \hat{\varphi}_y(x, s) ds,$$

and hence

$$\begin{aligned} \int_0^h |\hat{\varphi}(x, k)|^2 x^\alpha dx &= \int_0^h |\hat{\varphi}(x, y)|^2 x^\alpha dx + 2 \int_0^h \left( \int_y^k \hat{\varphi}(x, s) x^{\frac{\alpha}{2}} \hat{\varphi}_y(x, s) x^{\frac{\alpha}{2}} ds \right) dx \\ &\leq \int_0^h |\hat{\varphi}(x, y)|^2 x^\alpha dx + 2 \int_0^h \int_l^k |\hat{\varphi}(x, y)| x^{\frac{\alpha}{2}} |\hat{\varphi}_y(x, y)| x^{\frac{\alpha}{2}} dx dy. \end{aligned}$$

Thus, on integrating over all  $y \in [l, k]$  (recall that  $l \leq 0 \leq k$ ,  $k - l > 0$  and  $h > 0$ ),

$$\begin{aligned} (k - l) \int_0^h |\hat{\varphi}(x, k)|^2 x^\alpha dx &\leq \int_0^h \int_l^k |\hat{\varphi}(x, y)|^2 x^\alpha dx dy \\ &\quad + 2(k - l) \left( \int_0^h \int_l^k |\hat{\varphi}(x, y)|^2 x^\alpha dx dy \right)^{1/2} \left( \int_0^h \int_l^k |\hat{\varphi}_y(x, y)|^2 x^\alpha dx dy \right)^{1/2}, \end{aligned}$$

which then implies that

$$\begin{aligned} \int_0^h |\hat{\varphi}(x, k)|^2 x^\alpha dx &\leq \frac{1}{k-l} \int_0^h \int_l^k |\hat{\varphi}(x, y)|^2 x^\alpha dx dy \\ &+ 2 \left( \int_0^h \int_l^k |\hat{\varphi}(x, y)|^2 x^\alpha dx dy \right)^{1/2} \left( \int_0^h \int_l^k |\hat{\varphi}_y(x, y)|^2 x^\alpha dx dy \right)^{1/2}. \end{aligned}$$

Analogously,

$$\begin{aligned} \int_0^h |\hat{\varphi}_x(x, k)|^2 x^\alpha dx &\leq \frac{1}{k-l} \int_0^h \int_l^k |\hat{\varphi}_x(x, y)|^2 x^\alpha dx dy \\ &+ 2 \left( \int_0^h \int_l^k |\hat{\varphi}_x(x, y)|^2 x^\alpha dx dy \right)^{1/2} \left( \int_0^h \int_l^k |\hat{\varphi}_{xy}(x, y)|^2 x^\alpha dx dy \right)^{1/2}. \end{aligned} \quad (6.5.14)$$

Substituting the last two bounds into (6.5.13) it follows that

$$|\hat{\varphi}(h, k)| \leq C(h, k-l) \|\hat{\varphi}\|_{\mathbf{H}^2((0,h) \times (l,k); x^\alpha)}.$$

Further, (6.5.14) implies that

$$\int_0^h |\hat{\varphi}_x(x, k)| x^\alpha dx \leq \left( \frac{h^{\alpha+1}}{\alpha+1} \right)^{1/2} \left( \int_0^h |\hat{\varphi}_x(x, k)|^2 x^\alpha dx \right)^{1/2} \leq C(h, k-l) \|\hat{\varphi}\|_{\mathbf{H}^2((0,h) \times (l,k); x^\alpha)}.$$

Substituting the last two bounds into (6.5.12), we deduce that

$$|\hat{\Phi}(0, k)| \leq C(h, k-l) \|\hat{\varphi}\|_{\mathbf{H}^2((0,h) \times (l,k); x^\alpha)}.$$

Analogously,

$$|\hat{\varphi}(h, l)| \leq C(h, k-l) \|\hat{\varphi}\|_{\mathbf{H}^2((0,h) \times (l,k); x^\alpha)} \quad \text{and} \quad |\hat{\Phi}(0, l)| \leq C(h, k-l) \|\hat{\varphi}\|_{\mathbf{H}^2((0,h) \times (l,k); x^\alpha)},$$

as well as

$$|\hat{\varphi}(h, 0)| \leq C(h, k-l) \|\hat{\varphi}\|_{\mathbf{H}^2((0,h) \times (l,k); x^\alpha)}.$$

These inequalities imply that, for  $i = 1, 2$ ,

$$\begin{aligned} \|L_i(\hat{\varphi})\|_{\mathbf{L}^2(\kappa; x^\alpha)} &\leq \|L_i(\hat{\varphi})\|_{\mathbf{L}^2((0,h) \times (l,k); x^\alpha)} \\ &\leq \left[ 1 + \max \left( \frac{3}{h}, \frac{2}{k-l} \right) \left( \frac{h^{\alpha+1}}{\alpha+1} (k-l) \right)^{\frac{1}{2}} C(h, k-l) \right] \|\hat{\varphi}\|_{\mathbf{H}^2((0,h) \times (l,k); x^\alpha)}. \end{aligned}$$

Recall that  $L_i(\hat{\varphi}) \equiv 0$  for all  $\hat{\varphi} \in \mathbb{P}_1$ ,  $i = 1, 2$ .

Let  $\tilde{\mathbf{A}} = (1, 0)$ ,  $\tilde{\mathbf{B}} = (0, b)$ ,  $\tilde{\mathbf{C}} = (0, c)$  denote the counterparts of  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{C}$ , respectively, with  $c < 0 < b$ , in the open reference triangle  $\tilde{\kappa}$ , obtained by rescaling the open triangle  $\kappa = \triangle ABC$  by  $h$ , i.e.  $b = k/h$  and  $c = l/h$ , and let  $\rho := (k-l)/h = b-c (> 0)$ . We define  $\tilde{x} = x/h$  and  $\tilde{y} = y/h$ ,  $\tilde{\varphi}(\tilde{x}, \tilde{y}) := \hat{\varphi}(x, y)$ ,  $\tilde{p}_{\tilde{\varphi}}(\tilde{x}, \tilde{y}) := p_{\hat{\varphi}}(x, y)$ . Finally, we define  $\tilde{L}_i$  by

$$\tilde{L}_i(\tilde{\varphi})(\tilde{x}, \tilde{y}) := h L_i(\hat{\varphi})(x, y), \quad i = 1, 2.$$

Thus,

$$\tilde{L}_1(\tilde{\varphi})(\tilde{x}, \tilde{y}) = \tilde{\varphi}_{\tilde{x}}(\tilde{x}, \tilde{y}) - (\tilde{p}_{\tilde{\varphi}})_{\tilde{x}}(\tilde{x}, \tilde{y}), \quad \tilde{L}_2(\tilde{\varphi})(\tilde{x}, \tilde{y}) = \tilde{\varphi}_{\tilde{y}}(\tilde{x}, \tilde{y}) - (\tilde{p}_{\tilde{\varphi}})_{\tilde{y}}(\tilde{x}, \tilde{y}).$$

Then,  $\tilde{L}_i(\tilde{\varphi}) \equiv 0$  for all  $\tilde{\varphi} \in \mathbb{P}_1$ . In addition, repeating the bounds above with  $h$ ,  $k$  and  $l$  replaced by 1,  $b$  and  $c$ , noting that all constants in the bounds depend continuously on  $\rho = b - c$ , we deduce the existence of a positive constant  $C(\rho)$ , which depends continuously on  $\rho$ , such that

$$\|\tilde{L}_i(\tilde{\varphi})\|_{L^2(\tilde{\kappa}; \tilde{x}^\alpha)} \leq \|\tilde{L}_i(\tilde{\varphi})\|_{L^2((0,1) \times (c,b); \tilde{x}^\alpha)} \leq C(\rho) \|\tilde{\varphi}\|_{H^2((0,1) \times (c,b); \tilde{x}^\alpha)}, \quad i = 1, 2.$$

Note that  $\rho$  depends only on the shape of  $\kappa$ ; in particular, it is independent of the size of  $\kappa$ .

Let us recall the following generalization of the Bramble–Hilbert Lemma, due to Tartar (cf. Ciarlet [34], Section 3.1, Exercise 3.1.1).

**Lemma 6.5.3 (L. Tartar)** *Let  $V$  be a Banach space, and let  $V_1$ ,  $V_2$  and  $W$  be three normed linear spaces. Suppose that  $A_i \in \mathcal{L}(V; V_i)$ ,  $i = 1, 2$ , and that  $A_1$  is compact. Suppose, further, that there exists a positive constant  $c_0$  such that*

$$\|v\|_V \leq c_0 (\|A_1 v\|_{V_1} + \|A_2 v\|_{V_2}) \quad \forall v \in V.$$

Finally, suppose that  $L \in \mathcal{L}(V; W)$  is such that

$$v \in \ker A_2 \implies Lv = 0.$$

Then, the following statements hold.

(i)  $\mathbb{P} := \ker A_2$  is a finite-dimensional linear space.

(ii) There exists a positive constant  $c_1$  such that

$$\inf_{p \in \mathbb{P}} \|v - p\|_V \leq c_1 \|A_2 v\|_{V_2} \quad \forall v \in V.$$

(iii) There exists a positive constant  $C$  such that

$$\|Lv\|_W \leq C \|A_2 v\|_{V_2} \quad \forall v \in V.$$

We shall apply this result with  $\alpha \geq 1$ ,  $V := H^2((0, 1) \times (c, b); \tilde{x}^\alpha)$ ,  $V_1 := H^1((0, 1) \times (c, b); \tilde{x}^\alpha)$ ,  $V_2 := [L^2((0, 1) \times (c, b); \tilde{x}^\alpha)]^4$ ,  $W := L^2((0, 1) \times (c, b); \tilde{x}^\alpha)$ ,  $A_2 : \tilde{v} \in H^2((0, 1) \times (c, b); \tilde{x}^\alpha) \mapsto (\tilde{v}_{\tilde{x}\tilde{x}}, \tilde{v}_{\tilde{x}\tilde{y}}, \tilde{v}_{\tilde{y}\tilde{x}}, \tilde{v}_{\tilde{y}\tilde{y}})$ ,  $A_1 := \text{Id}$ , and  $L = \tilde{L}_i$ ,  $i = 1, 2$ , together with the compact embedding

$$H^2((0, 1) \times (c, b); \tilde{x}^\alpha) \hookrightarrow H^1((0, 1) \times (c, b); \tilde{x}^\alpha),$$

which requires the restriction  $\alpha \geq 1$  (cf. Lemma 5.2 in Antoci [5])

Thus, we deduce that

$$\|\tilde{\varphi}_{\tilde{x}} - (\tilde{p}_{\tilde{\varphi}})_{\tilde{x}}\|_{L^2((0,1) \times (c,b); \tilde{x}^\alpha)} \leq C(\rho) |\tilde{\varphi}|_{H^2((0,1) \times (c,b); \tilde{x}^\alpha)}$$

and

$$\|\tilde{\varphi}_{\tilde{y}} - (\tilde{p}_{\tilde{\varphi}})_{\tilde{y}}\|_{L^2((0,1) \times (c,b); \tilde{x}^\alpha)} \leq C(\rho) |\tilde{\varphi}|_{H^2((0,1) \times (c,b); \tilde{x}^\alpha)},$$

where  $|\cdot|_{\mathbb{H}^2((0,1)\times(c,b);\tilde{x}^\alpha)}$  is the semi-norm on  $\mathbb{H}^2((0,1)\times(c,b);\tilde{x}^\alpha)$ .

After returning from the scaled variables  $\tilde{x}$  and  $\tilde{y}$  to the original variables  $x = h\tilde{x}$  and  $y = h\tilde{y}$  and combining the resulting inequalities into a single inequality, we obtain

$$\|\nabla(\hat{\varphi} - p_{\hat{\varphi}})\|_{\mathbb{L}^2((0,h)\times(l,k);x^\alpha)} \leq C(\rho) h |\hat{\varphi}|_{\mathbb{H}^2((0,h)\times(l,k);x^\alpha)}.$$

In other words,

$$\|\nabla(\hat{\varphi} - p_{\hat{\varphi}})\|_{\mathbb{L}^2(R(\kappa);x^\alpha)} \leq C(\rho) h |\hat{\varphi}|_{\mathbb{H}^2(R(\kappa);x^\alpha)}, \quad (6.5.15)$$

whereupon

$$\|\nabla(\hat{\varphi} - p_{\hat{\varphi}})\|_{\mathbb{L}^2(\kappa;x^\alpha)} \leq C(\rho) h |\hat{\varphi}|_{\mathbb{H}^2(R(\kappa);x^\alpha)}, \quad (6.5.16)$$

where  $R(\kappa) := (0, h) \times (l, k)$ ,  $\rho := (k - l)/h$  and  $\alpha \geq 1$ .

Using that, for  $(x, y) \in \kappa$ ,  $0 \leq x/h \leq 1$  and  $|y|/(k - l) \leq 1$ , one can obtain a similar bound on  $\hat{\varphi} - p_{\hat{\varphi}}$  in the  $x^\alpha$ -weighted  $\mathbb{L}^2$  norm on  $\kappa$ . The only difference is that then

$$L(\hat{\varphi}) := \hat{\varphi} - p_{\hat{\varphi}} \quad \text{and} \quad \tilde{L}(\tilde{\varphi})(\tilde{x}, \tilde{y}) := L(\hat{\varphi})(x, y),$$

with the same definitions of  $p_{\hat{\varphi}}$ ,  $\tilde{\varphi}$ ,  $\tilde{p}_{\tilde{\varphi}}$ ,  $\tilde{x}$  and  $\tilde{y}$  as before. We recall that  $p_{\hat{\varphi}} \equiv \hat{\varphi}$  for all  $\hat{\varphi} \in \mathbb{P}_1$ , and hence  $L(\hat{\varphi}) \equiv 0$  for all  $\hat{\varphi} \in \mathbb{P}_1$  and therefore  $\tilde{L}(\tilde{\varphi}) \equiv 0$  for all  $\tilde{\varphi} \in \mathbb{P}_1$ . We still have that

$$\|\tilde{L}(\tilde{\varphi})\|_{\mathbb{L}^2(\tilde{\kappa};\tilde{x}^\alpha)} \leq \|\tilde{L}(\tilde{\varphi})\|_{\mathbb{L}^2((0,1)\times(c,b);\tilde{x}^\alpha)} \leq C(\rho) \|\tilde{\varphi}\|_{\mathbb{H}^2((0,1)\times(c,b);\tilde{x}^\alpha)}.$$

Hence, Lemma 6.5.3, with the same choice of  $V$ ,  $V_1$ ,  $V_2$ ,  $W$ ,  $A_1$  and  $A_2$  as before, and  $\alpha \geq 1$ , implies that

$$\|\tilde{\varphi} - \tilde{p}_{\tilde{\varphi}}\|_{\mathbb{L}^2((0,1)\times(c,b);\tilde{x}^\alpha)} \leq C(\rho) |\tilde{\varphi}|_{\mathbb{H}^2((0,1)\times(c,b);\tilde{x}^\alpha)}.$$

After returning from the scaled variables  $\tilde{x} = x/h$  and  $\tilde{y} = y/h$  to the original variables  $x$  and  $y$ , we obtain that

$$\|\hat{\varphi} - p_{\hat{\varphi}}\|_{\mathbb{L}^2((0,h)\times(l,k);x^\alpha)} \leq C(\rho) h^2 |\hat{\varphi}|_{\mathbb{H}^2((0,h)\times(l,k);x^\alpha)}.$$

In other words,

$$\|\hat{\varphi} - p_{\hat{\varphi}}\|_{\mathbb{L}^2(R(\kappa);x^\alpha)} \leq C(\rho) h^2 |\hat{\varphi}|_{\mathbb{H}^2(R(\kappa);x^\alpha)},$$

whereupon

$$\|\hat{\varphi} - p_{\hat{\varphi}}\|_{\mathbb{L}^2(\kappa;x^\alpha)} \leq C(\rho) h^2 |\hat{\varphi}|_{\mathbb{H}^2(R(\kappa);x^\alpha)}, \quad (6.5.17)$$

with  $R(\kappa) := (0, h) \times (l, k)$ ,  $\rho := (k - l)/h$  and  $\alpha \geq 1$ . The constant  $C(\rho)$  is a continuous function of  $\rho$  in each of these bounds.

**Two dimensions: configuration 2-flat.** The alternative configuration of the triangle  $\kappa = \triangle ABC$  is:  $A=(0,0)$ ,  $B=(h,k)$  and  $C=(h,l)$ , with only one point,  $A$ , on the line  $x = 0$  along which the weight-function  $(x, y) \mapsto x^\alpha$  vanishes. In this case, we define  $p_{\hat{\varphi}}$  as the affine function that interpolates  $\hat{\varphi}$  at  $B$  and  $C$ , and has the value

$$\hat{\Phi}(0,0) = \hat{\varphi}(h,0) - h \frac{\alpha + 1}{h^{\alpha+1}} \int_0^h \hat{\varphi}_x(x,0) x^\alpha dx$$

at  $A=(0,0)$ , extrapolated from  $(h,0)$  using the univariate quasi-interpolation operator. Thus,

$$\begin{aligned} p_{\hat{\varphi}}(x,y) &= \hat{\Phi}(0,0) \left(1 - \frac{x}{h}\right) + \hat{\varphi}(h,k) \left(y - \frac{l}{h}x\right) \frac{1}{k-l} + \hat{\varphi}(h,l) \left(y - \frac{k}{h}x\right) \frac{1}{l-k}, \\ (p_{\hat{\varphi}})_x(x,y) &= -\hat{\Phi}(0,0) \frac{1}{h} + \hat{\varphi}(h,k) \left(-\frac{l}{h}\right) \frac{1}{k-l} + \hat{\varphi}(h,l) \left(-\frac{k}{h}\right) \frac{1}{l-k}, \\ (p_{\hat{\varphi}})_y(x,y) &= \hat{\varphi}(h,k) \frac{1}{k-l} + \hat{\varphi}(h,l) \frac{1}{l-k}. \end{aligned}$$

Again, we define

$$L_1(\hat{\varphi}) := \hat{\varphi}_x - (p_{\hat{\varphi}})_x \quad \text{and} \quad L_2(\hat{\varphi}) := \hat{\varphi}_y - (p_{\hat{\varphi}})_y,$$

and we observe that  $\hat{\Phi}(0,0) = \hat{\varphi}(0,0)$  for all  $\hat{\varphi} \in \mathbb{P}_1$ , and hence  $p_{\hat{\varphi}} \equiv \hat{\varphi}$  and  $L_i(\hat{\varphi}) \equiv 0$  for all  $\hat{\varphi} \in \mathbb{P}_1$ ,  $i = 1, 2$ .

The rest of the argument is the same as in the case of configuration 1-flat, and leads to the same final bound:

$$\|\nabla(\hat{\varphi} - p_{\hat{\varphi}})\|_{L^2(\kappa; x^\alpha)} \leq C(\rho) h |\hat{\varphi}|_{H^2(R(\kappa); x^\alpha)}, \quad (6.5.18)$$

where again  $R(\kappa) := (0, h) \times (l, k)$ ,  $\rho := (k-l)/h$  and  $\alpha \geq 1$ . Also, as in the case of configuration 1-flat,

$$\|\hat{\varphi} - p_{\hat{\varphi}}\|_{L^2(\kappa; x^\alpha)} \leq C(\rho) h^2 |\hat{\varphi}|_{H^2(R(\kappa); x^\alpha)}, \quad (6.5.19)$$

with  $R(\kappa) := (0, h) \times (l, k)$ ,  $\rho := (k-l)/h$  and  $\alpha \geq 1$ . The constant  $C(\rho)$  is a continuous function of  $\rho$  in each of these bounds.

#### 6.5.4 Two dimensions: curved boundary

Now suppose that  $D$  is an open disc in  $\mathbb{R}^2$  of radius  $r_D \in \mathbb{R}_{>0}$ , centred at the origin. Suppose, further, that  $\{\mathcal{T}_h^q\}_{h>0}$  is a quasiuniform family of partitions of  $D$  (in the sense of Hypothesis (A1) from Section 6.4, with  $d = 2$ ,) into disjoint open nonobtuse triangles  $\kappa_q$ , with possibly one curved edge on  $\partial D$ . We focus our attention on elements  $\kappa_q$  that are in contact with  $\partial D$ . There are again two possible configurations, which will be considered separately. We shall assume throughout the section that the potential  $U$  and the associated Maxwellian  $M$  satisfy on  $D$  the assumptions stated at the start of Section 6.2.3, including (6.2.9a), with  $\zeta \geq 1$ , and (6.2.9b).

**Two dimensions: configuration 1-curved.** We consider a circle  $\mathcal{C} \subset D$ , concentric with  $\partial D$ , which is a distance  $h$  away from  $\partial D$ ; cf. Figure 6.2. The analogue of configuration 1-flat is an open curved nonobtuse triangle  $\kappa_q := \triangle ABC$ , with one curved edge  $BC \subset \partial D$  and with  $A \in \mathcal{C}$ . Let  $B'$  and  $C'$  be points on  $\mathcal{C}$  such that  $BB'$  and  $CC'$  are aligned with the directions of the normal vectors to  $\partial D$  at  $B$  and  $C$ , respectively. We mimic the construction of the quasi-interpolant  $p_{\hat{\varphi}}$  of  $\hat{\varphi}$  described in the previous section.

Note that, for  $\hat{\varphi} \in H_M^2(D)$  and any pair of points  $Q_1$  and  $Q_2$  in  $D$ ,

$$\hat{\varphi}(Q_2) = \hat{\varphi}(Q_1) - \int_0^1 \frac{d}{d\tau} \hat{\varphi}((1-\tau)Q_2 + \tau Q_1) d\tau.$$

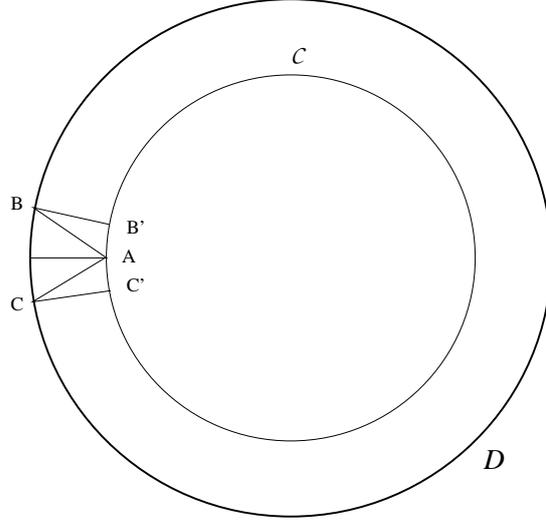


Figure 6.2: The domain  $D$ , the circle  $C \subset D$ , with  $\text{dist}(\partial D, C) = h$ , and  $C', A, B' \in C$  and the open curved nonobtuse triangle  $\kappa_q = \triangle ABC$  in configuration 1-curved.

Motivated by this identity, for  $Q_1 \in D$  and  $Q_2 \in \bar{D}$ , we define

$$\hat{\Phi}(Q_2) := \hat{\varphi}(Q_1) - \frac{\int_0^1 M((1-\tau)Q_2 + \tau Q_1) \frac{d}{d\tau} \hat{\varphi}((1-\tau)Q_2 + \tau Q_1) d\tau}{\int_0^1 M((1-\tau)Q_2 + \tau Q_1) d\tau}. \quad (6.5.20)$$

**Remark 6.5.4** In one space dimension, with  $M(q) = q^\alpha$ ,  $q \in [0, h]$ ,  $Q_2 = 0$ ,  $Q_1 = h$ , and performing the change of variable  $q = \tau h$ , (6.5.20) yields our univariate extrapolation operator:

$$\hat{\Phi}(0) := \hat{\varphi}(h) - h \frac{\int_0^h q^\alpha \hat{\varphi}'(q) dq}{\int_0^h q^\alpha dq} = \hat{\varphi}(h) - h \frac{\alpha + 1}{h^{\alpha+1}} \int_0^h q^\alpha \hat{\varphi}'(q) dq. \quad \diamond$$

In multiple space dimensions the formula (6.5.20), after performing the  $\tau$ -differentiation under the integral sign, becomes

$$\hat{\Phi}(Q_2) := \hat{\varphi}(Q_1) - (Q_1 - Q_2) \cdot \frac{\int_0^1 M((1-\tau)Q_2 + \tau Q_1) (\nabla_q \hat{\varphi})((1-\tau)Q_2 + \tau Q_1) d\tau}{\int_0^1 M((1-\tau)Q_2 + \tau Q_1) d\tau}.$$

In particular, in the two-dimensional setting considered here, and with reference to Figure 6.2,

$$\hat{\Phi}(B) := \hat{\varphi}(B') - (B' - B) \cdot \frac{\int_0^1 M((1-\tau)B + \tau B') (\nabla_q \hat{\varphi})((1-\tau)B + \tau B') d\tau}{\int_0^1 M((1-\tau)B + \tau B') d\tau}$$

and

$$\hat{\Phi}(C) := \hat{\varphi}(C') - (C' - C) \cdot \frac{\int_0^1 M((1-\tau)C + \tau C') (\nabla_q \hat{\varphi})((1-\tau)C + \tau C') d\tau}{\int_0^1 M((1-\tau)C + \tau C') d\tau}.$$

We then define the affine function  $p_{\hat{\varphi}}$  on  $\kappa_q = \triangle ABC$  by

$$p_{\hat{\varphi}}(\underline{q}) := \hat{\varphi}(\mathbf{A}) \psi_{\mathbf{A}}(\underline{q}) + \hat{\Phi}(\mathbf{B}) \psi_{\mathbf{B}}(\underline{q}) + \hat{\Phi}(\mathbf{C}) \psi_{\mathbf{C}}(\underline{q}), \quad \underline{q} = (q_1, q_2) \in \kappa_q, \quad (6.5.21)$$

where  $\{\psi_{\mathbf{A}}, \psi_{\mathbf{B}}, \psi_{\mathbf{C}}\}$  is the  $\mathbb{P}_1^q$  local (nodal/Lagrange) basis associated with the triangle  $\triangle ABC$ .

Let  $R(\kappa_q)$  denote the curvilinear rectangle  $\mathbf{B}'\mathbf{BCC}'$  depicted in Figure 6.2. Our aim is to show that, in analogy with (6.5.15),

$$\|\nabla_q(\hat{\varphi} - p_{\hat{\varphi}})\|_{L_M^2(R(\kappa_q))} \leq C(\rho) h |\hat{\varphi}|_{H_M^2(R(\kappa_q))},$$

where  $\rho$  is a positive constant dependent only on the shape of  $\kappa_q$ ; this will in turn imply that

$$\|\nabla_q(\hat{\varphi} - p_{\hat{\varphi}})\|_{L_M^2(\kappa_q)} \leq C(\rho) h |\hat{\varphi}|_{H_M^2(R(\kappa_q))}.$$

Using polar co-ordinates, the curvilinear rectangle  $R(\kappa_q)$  in the  $\underline{q} := (q_1, q_2)$  domain can be mapped into the rectangular domain

$$R_{\text{polar}}(\kappa_q) := \{(r, \theta) : -r_D < r < -r_D + h, \theta_C < \theta < \theta_B\}.$$

Let us therefore perform the following change of independent variables:

$$q_1 = r \cos \theta, \quad q_2 = r \sin \theta, \quad r \in (-r_D, -r_D + h), \quad \theta \in (\theta_C, \theta_B); \quad (6.5.22)$$

thus,  $r = -|q|$ . Naturally,  $0 < h \ll 1 < r_D$ , and we can therefore assume without loss of generality that  $-r_D + h \leq -\frac{1}{2}$ ; therefore,  $r = 0$  is, uniformly in  $h$ , separated from the range  $(-r_D, -r_D + h)$  of  $r$ , whereby the change of variables (6.5.22) is a smooth bijective diffeomorphism from  $R(\kappa_q)$  to  $R_{\text{polar}}(\kappa_q)$ .

By virtue of (6.2.9a) we may assume without loss of generality that  $M(\underline{q}) = (r_D - |q|)^\alpha$ , with  $\alpha = \zeta \geq 1$  and  $\zeta$  as in (6.2.9a), and  $|q| \in (r_D - h, r_D)$ . In polar co-ordinates, with  $|q| = -r$ , we therefore define  $N(r) := (r_D + r)^\alpha$  for  $r \in (-r_D, -r_D + h)$ , where  $\alpha = \zeta \geq 1$ .

Now, on noting that  $M(\underline{q}) = N(r)$  with  $r = -|q| \in (-r_D, -r_D + h)$ , we have that

$$\hat{\Phi}(\mathbf{B}) = \hat{\Phi}(-r_D, \theta_B) = \hat{\varphi}(-r_D + h, \theta_B) - h \frac{\int_0^1 N(-r_D + \tau h) \hat{\varphi}_r(-r_D + \tau h, \theta_B) d\tau}{\int_0^1 N(-r_D + \tau h) d\tau}.$$

Hence,

$$\hat{\Phi}(\mathbf{B}) = \hat{\Phi}(-r_D, \theta_B) = \hat{\varphi}(-r_D + h, \theta_B) - h \frac{\alpha + 1}{h^{\alpha+1}} \int_0^h t^\alpha \hat{\varphi}_r(-r_D + t, \theta_B) dt. \quad (6.5.23)$$

Analogously,

$$\hat{\Phi}(\mathbf{C}) = \hat{\Phi}(-r_D, \theta_C) = \hat{\varphi}(-r_D + h, \theta_C) - h \frac{\alpha + 1}{h^{\alpha+1}} \int_0^h t^\alpha \hat{\varphi}_r(-r_D + t, \theta_C) dt, \quad (6.5.24)$$

while

$$\hat{\varphi}(\mathbf{A}) = \hat{\varphi}(-r_D + h, \theta_A). \quad (6.5.25)$$

It is clear from (6.5.23) that if the restriction of  $\hat{\varphi}$  to the closed line segment connecting  $B'$  to  $B$  is a linear function, and therefore  $\hat{\varphi}_r$  is constant along this line segment, then  $\hat{\Phi}(B) = \hat{\varphi}(-r_D, \theta_B) = \hat{\varphi}(B)$ . Analogously, (6.5.24) implies that if the restriction of  $\hat{\varphi}$  to the closed line segment connecting  $C'$  to  $C$  is a linear function, then  $\hat{\Phi}(C) = \hat{\varphi}(-r_D, \theta_C) = \hat{\varphi}(C)$ .

Hence, if  $\hat{\varphi} \in \mathbb{P}_1^q$ , then (6.5.21) implies that  $p_{\hat{\varphi}}(q) = \hat{\varphi}(A) \psi_A(q) + \hat{\varphi}(B) \psi_B(q) + \hat{\varphi}(C) \psi_C(q)$ , the standard linear nodal/Lagrange interpolant of  $\hat{\varphi}$ , whereby  $\nabla_q(\hat{\varphi} - p_{\hat{\varphi}}) \equiv 0$ . Equivalently, letting

$$L_1(\hat{\varphi}) = (\hat{\varphi})_{q_1} - (p_{\hat{\varphi}})_{q_1}, \quad L_2(\hat{\varphi}) = (\hat{\varphi})_{q_2} - (p_{\hat{\varphi}})_{q_2},$$

we have that  $L_i(\hat{\varphi}) \equiv 0$  for all  $\hat{\varphi} \in \mathbb{P}_1^q$ ,  $i = 1, 2$ .

Since the formulae (6.5.23), (6.5.24), (6.5.25) are essentially the same as those corresponding to  $\hat{\Phi}(B) = \hat{\Phi}(0, k)$ ,  $\hat{\Phi}(C) = \hat{\Phi}(0, l)$  and  $\hat{\varphi}(A) = \hat{\varphi}(h, 0)$  in the case of configuration 1-flat in the previous section, defining  $\rho := (\theta_B - \theta_C)/h$ , changing variables to the rectangular region  $R_{\text{polar}}(\kappa_q)$ , rescaling this by  $1/h$  as in the previous section, applying Lemma 6.5.3, and then rescaling by  $h$  to return from  $R_{\text{polar}}(\kappa_q)$  to  $R(\kappa_q)$  yields

$$\|\nabla_q(\hat{\varphi} - p_{\hat{\varphi}})\|_{L_M^2(R(\kappa_q))} \leq C(\rho) h |\hat{\varphi}|_{H_M^2(R(\kappa_q))}.$$

Hence,

$$\|\nabla_q(\hat{\varphi} - p_{\hat{\varphi}})\|_{L_M^2(R(\kappa_q))} \leq C(\rho) h |\hat{\varphi}|_{H_M^2(R(\kappa_q))}, \quad (6.5.26)$$

with  $\rho := (\theta_B - \theta_C)/h$ .

Next, we prove that

$$\|\hat{\varphi} - p_{\hat{\varphi}}\|_{L_M^2(R(\kappa_q))} \leq C(\rho) h^2 |\hat{\varphi}|_{H_M^2(R(\kappa_q))}. \quad (6.5.27)$$

This time, we define  $L(\hat{\varphi}) := \hat{\varphi} - p_{\hat{\varphi}}$  where, again,  $p_{\hat{\varphi}}(q) = \hat{\varphi}(A) \psi_A(q) + \hat{\Phi}(B) \psi_B(q) + \hat{\Phi}(C) \psi_C(q)$ . Once again, if  $\hat{\varphi} \in \mathbb{P}_1^q$ , then  $p_{\hat{\varphi}}$  is just the standard linear nodal/Lagrange interpolant of  $\hat{\varphi}$  and therefore  $L(\hat{\varphi}) \equiv 0$ . The rest of the argument is the same as in the case of the error estimate in the  $M$ -weighted  $H^1$  seminorm above. Thus, on applying Lemma 6.5.3 and a scaling argument in the same way as before,

$$\|\hat{\varphi} - p_{\hat{\varphi}}\|_{L_M^2(R(\kappa_q))} \leq C(\rho) h^2 |\hat{\varphi}|_{H_M^2(R(\kappa_q))}, \quad (6.5.28)$$

where, again,  $\rho := (\theta_B - \theta_C)/h$ . The constant  $C(\rho)$  is a continuous function of  $\rho$  in each of these bounds.

**Two dimensions: configuration 2-curved.** The alternative configuration of the triangle  $\kappa_q = \triangle ABC$  is that  $A \in \partial D$  while  $B, C \in \mathcal{C}$ . In this case, we define  $p_{\hat{\varphi}}$  as the affine function that interpolates  $\hat{\varphi}$  at  $B$  and  $C$ , and has the value

$$\hat{\Phi}(A) := \hat{\varphi}(A') - (A' - A) \cdot \frac{\int_0^1 M((1-\tau)A + \tau A') (\nabla \hat{\varphi})((1-\tau)A + \tau A') d\tau}{\int_0^1 M((1-\tau)A + \tau A') d\tau}$$

at  $A$ . Here  $A'$  is the point on  $\mathcal{C}$  where the line segment, normal to  $\partial D$ , connecting  $A$  to the centre of the disc  $D$  intersects  $\mathcal{C}$ ; thus the segment  $AA'$  is orthogonal to  $\partial D$ . The value  $\hat{\Phi}(A)$  is therefore obtained by extrapolating  $\hat{\varphi}$  from  $A'$ . Thus,

$$p_{\hat{\varphi}}(q) = \hat{\Phi}(A) \psi_A(q) + \hat{\varphi}(B) \psi_B(q) + \hat{\varphi}(C) \psi_C(q).$$

Again, we define,

$$L_1(\hat{\varphi}) := \hat{\varphi}_{q_1} - (p_{\hat{\varphi}})_{q_1} \quad \text{and} \quad L_2(\hat{\varphi}) := \hat{\varphi}_{q_2} - (p_{\hat{\varphi}})_{q_2},$$

and we observe that  $L_i(\hat{\varphi}) \equiv 0$ ,  $i = 1, 2$ , for all  $\hat{\varphi} \in \mathbb{P}_1^q$ . The rest of the argument is the same as in the case of configuration 1-curved, and leads to the same final bound:

$$\|\nabla_q(\hat{\varphi} - p_{\hat{\varphi}})\|_{L_M^2(\kappa_q)} \leq C(\rho) h |\hat{\varphi}|_{\mathbb{H}_M^2(R(\kappa_q))}, \quad (6.5.29)$$

where now  $R(\kappa_q)$  is the curvilinear rectangle  $BB'C'C$ , whose curved edges  $B'C' \subset \partial D$ ,  $BC \subset \mathcal{C}$ ; here  $B'$  and  $C'$  are the points on  $\partial D$  where the line segments passing through the centre of the disc  $D$  and the points  $B$  and  $C$ , respectively, extended beyond  $B$  and  $C$ , respectively, intersect  $\partial D$ . Clearly, each of the line segments  $BB'$  and  $CC'$  is orthogonal to  $\partial D$  as in the case of configuration 1-curved. The definition of  $\rho$  is the same as in the case of configuration 1-curved, i.e.  $\rho := (\theta_B - \theta_C)/h$ .

Arguing in the same way as in the case of the  $M$ -weighted  $L^2$  norm bound derived above in the case of configuration 1-curved, we also have that, with  $\rho := (\theta_B - \theta_C)/h$ ,

$$\|\hat{\varphi} - p_{\hat{\varphi}}\|_{L_M^2(\kappa_q)} \leq C(\rho) h^2 |\hat{\varphi}|_{\mathbb{H}_M^2(R(\kappa_q))}. \quad (6.5.30)$$

The constant  $C(\rho)$  is a continuous function of  $\rho$  in each of these bounds.

**Two dimensions: global interpolation bound.** Let  $h_q$  denote the maximum diameter of any triangle  $\kappa_q$  in the quasiuniform and nonobtuse family of partitions  $\{\mathcal{T}_h^q\}_{h>0}$  of  $D$ . Each triangle  $\kappa_q \in \mathcal{T}_h^q$  whose closure intersects  $\partial D$  is either in configuration 1-curved or in configuration 2-curved; on such triangles we define  $p_{\hat{\varphi}}$  as above. Any triangle  $\kappa_q \in \mathcal{T}_h^q$  that is neither in configuration 1-curved or configuration 2-curved is such that the closure of  $\kappa_q$  is contained in the open disc  $D$ ; on such triangles, referred to as being in configuration 0, we define  $p_{\hat{\varphi}}$  as the standard nodal interpolant of  $\hat{\varphi}$ . For  $\hat{\varphi} \in \mathbb{H}_M^2(D)$ , we then define the global quasi-interpolant  $I_h^q \hat{\varphi} := p_{\hat{\varphi}}$ . Note, in particular, that  $I_h^q \hat{\varphi}$  is a continuous piecewise linear function on  $\bar{D}$  with the following properties: suppose that  $P$  is a vertex of a triangle  $\kappa_q \in \mathcal{T}_h^q$ ; if  $P \in D$ , then  $(I_h^q \hat{\varphi})(P) = \hat{\varphi}(P)$ ; if, on the other hand,  $P \in \partial D$ , then  $(I_h^q \hat{\varphi})(P) = \hat{\Phi}(P)$ , the value extrapolated from  $P' \in D$  using the formula

$$\hat{\Phi}(P) := \hat{\varphi}(P') - (P' - P) \cdot \frac{\int_0^1 M((1-\tau)P + \tau P') (\nabla_q \hat{\varphi})(1-\tau)P + \tau P') d\tau}{\int_0^1 M((1-\tau)P + \tau P') d\tau},$$

where  $P'$  is the unique point of intersection of the line segment that connects  $P \in \partial D$  to the centre of  $D$  with the circle  $\mathcal{C} \subset D$  concentric with  $\partial D$  and such that  $\text{dist}(\partial D, \mathcal{C}) = h$  and  $0 < h \ll r_D$ .

By virtue of (6.5.26) (on triangles  $\kappa_q \subset D$  in configuration 1-curved), (6.5.29) (on triangles  $\kappa_q \in D$  in configuration 2-curved), and classical interpolation results on the remaining triangles  $\kappa_q \in \mathcal{T}_h$  (in configuration 0) whose closure does not intersect  $\partial D$ , together with upper and lower bounds on  $M$  on triangles in configuration 0 and recalling (6.4.54), to relate the  $M$ -weighted  $L^2$ ,  $H^1$  and  $H^2$  norms to standard (nonweighted)  $L^2$ ,  $H^1$  and  $H^2$  norms, we deduce that

$$\|\nabla_q(\hat{\psi} - I_h^q \hat{\psi})\|_{L_M^2(D)} \leq C h_q |\hat{\psi}|_{\mathbb{H}_M^2(D)} \quad \text{and} \quad \|\hat{\psi} - I_h^q \hat{\psi}\|_{L_M^2(D)} \leq C h_q^2 |\hat{\psi}|_{\mathbb{H}_M^2(D)},$$

whereby

$$\|\hat{\psi} - I_h^q \hat{\psi}\|_{H_M^1(D)} \leq C h_q |\hat{\psi}|_{H_M^2(D)}. \quad (6.5.31)$$

Here we made use of the fact that the parameter  $\rho$  appearing in the bounds on the triangles  $\kappa_q \in \mathcal{T}_h^q$  in configuration 1-curved and configuration 2-curved belongs to a compact subinterval of  $\mathbb{R}_{>0}$ , independent of  $h_q$ , due to our assumption that  $\{\mathcal{T}_h^q\}_{h>0}$  is a quasiuniform family of nonobtuse partitions; since the constants  $C(\rho)$  featuring in those bounds are continuous functions of  $\rho$ , it follows that the constant  $C$  in (6.5.31) depends only on the shape-regularity parameters of  $\{\mathcal{T}_h^q\}_{h>0}$ , which, in particular, fix the range of  $\rho$ .

### 6.5.5 Three dimensions

We briefly comment on the modifications that need to be made to our arguments above when  $d = 3$ . Consider a family of quasiuniform nonobtuse partitions  $\{\mathcal{T}_h^q\}_{h>0}$ , in the sense of (A1) in Section 6.4, of the ball  $D = B(Q, r_D) \subset \mathbb{R}^3$ . Excluding the case of configuration 0, when the closure of a simplex  $\kappa_q \in \mathcal{T}_h^q$  has empty intersection with  $\partial D$ , there are now three different configurations to consider, corresponding to the cases when the closure of  $\kappa_q$  has one, two or three vertices on  $\partial D$ .

Let us suppose, for example, that the open nonobtuse simplex  $\kappa_q \in \mathcal{T}_h^q$  has three vertices A, B and C on the sphere  $\partial D$ , while the fourth vertex D is in the interior of the domain  $D$ , on a sphere  $\mathcal{C}$  concentric with  $\partial D$ , that is a distance  $h$  away from  $\partial D$ . We raise the inward normals from A, B, C to  $\partial D$ , and consider the points  $A'$ ,  $B'$ ,  $C'$  in the interior of the ball  $D$  that are on the respective normals to  $\partial D$  at A, B and C, and a distance  $h$  away from A, B and C, respectively; i.e.  $A'$ ,  $B'$ ,  $C'$  are on the sphere  $\mathcal{C}$ . The tetrahedron  $\kappa_q = ABCD$  is then contained in the curved triangular prismoid  $R(\kappa_q) := ABCA'B'C'$ , with curved faces ABC and  $A'B'C'$ .

Given a function  $\hat{\varphi} \in H_M^2(D)$ , we then extrapolate  $\hat{\varphi}$  from  $A'$ ,  $B'$  and  $C'$  using (6.5.20) to define  $\hat{\Phi}(A)$ ,  $\hat{\Phi}(B)$  and  $\hat{\Phi}(C)$ , and define  $p_{\hat{\varphi}}$  as the affine function of  $q$  on the simplex ABCD whose nodal values are  $\hat{\Phi}(A)$ ,  $\hat{\Phi}(B)$ ,  $\hat{\Phi}(C)$  and  $\hat{\varphi}(D)$ . We note in particular that if  $\hat{\varphi} \in \mathbb{P}_1^q$ , then  $p_{\hat{\varphi}} = \hat{\varphi}$ . Using spherical polar co-ordinates we map the curved triangular prismoid  $R(\kappa_q)$  containing the simplex  $\kappa_q = ABCD$  into a right triangular prism  $R_{\text{polar}}(\kappa_q)$ , and then argue as in the case of  $d = 2$  above, using Lemma 6.5.3, to deduce the analogue of (6.5.31) in the case of  $d = 3$ .

### 6.5.6 Stability of the Maxwellian-weighted $L^2$ projector in the Maxwellian-weighted $H^1$ norm

Now we are ready to discuss the question of stability, in the  $M$ -weighted  $H^1$  norm, of the orthogonal projector in the  $M$ -weighted  $L^2$  inner product on  $D \subset \mathbb{R}^d$ ,  $d = 2, 3$ . We begin by considering the following auxiliary problem: Let  $\hat{g} \in L_M^2(D)$ ; find  $\hat{z} \in H_M^1(D)$  such that

$$a(\hat{z}, \hat{\varphi}) = \ell(\hat{\varphi}) \quad \forall \hat{\varphi} \in H_M^1(D), \quad (6.5.32)$$

where, for  $\hat{\zeta}, \hat{\varphi} \in H_M^1(D)$ ,

$$a(\hat{\zeta}, \hat{\varphi}) := \int_D M \left( \nabla_q \hat{\zeta} \cdot \nabla_q \hat{\varphi} + \hat{\zeta} \hat{\varphi} \right) dq \quad \text{and} \quad \ell(\hat{\varphi}) := \int_D M \hat{g} \hat{\varphi} dq.$$

The existence of a unique solution  $\hat{z} \in \mathbf{H}_M^1(D)$  to (6.5.32) follows by the Lax–Milgram theorem. Note that

$$\|\hat{z}\|_{\mathbf{L}_M^2(D)} \leq \|\hat{z}\|_{\mathbf{H}_M^1(D)} \leq \|\hat{g}\|_{\mathbf{L}_M^2(D)}.$$

We begin by showing the following elliptic regularity result for (6.5.32):  $\hat{z} \in \mathbf{H}_D^2(M)$ , and the bound stated in (6.5.41) below holds. To this end, for  $\delta > 0$  we define

$$U_\delta(s) := U\left(\left(\frac{r_D}{r_D + \delta}\right)^2 s\right), \quad s \in [0, \tfrac{1}{2}r_D^2] \quad \text{and} \quad M_\delta(\underline{q}) := Z^{-1} \exp(-U_\delta(\tfrac{1}{2}|\underline{q}|^2)), \quad \underline{q} \in D,$$

where, as in (6.1.6) (i.e. with no  $\delta$ -dependence in the definition of  $Z$ ),

$$Z := \int_D \exp(-U(\tfrac{1}{2}|q|^2)) \, dq.$$

Note that since  $U'(s) > 0$  for  $s \in [0, \frac{1}{2}r_D^2]$ , we have  $0 \leq U_\delta(s) \leq U(s)$  for all  $s \in [0, \frac{1}{2}r_D^2]$ , with strict inequalities for  $s \neq 0$ , and  $M(\underline{q}) \leq M_\delta(\underline{q})$  for  $\underline{q} \in D$ , with strict inequality for  $\underline{q} \neq \mathbf{0}$ . The fact that, thereby,  $\int_D M_\delta(\underline{q}) \, d\underline{q}$  is strictly greater than 1 rather than equal to 1 is of no significance. For  $\hat{g} \in \mathbf{L}_M^2(D)$  and  $\delta > 0$ , we define

$$\hat{g}_\delta(\underline{q}) := \left(\frac{M(\underline{q})}{M_\delta(\underline{q})}\right)^{\frac{1}{2}} \hat{g}(\underline{q}), \quad \underline{q} \in D,$$

and note that  $\hat{g}_\delta \in \mathbf{L}_{M_\delta}^2(D)$  with  $\|\hat{g}_\delta\|_{\mathbf{L}_{M_\delta}^2(D)} = \|\hat{g}\|_{\mathbf{L}_M^2(D)}$ .

We consider the following problem: For  $\hat{g} \in \mathbf{L}_M^2(D)$  and  $\delta > 0$ , and with  $M_\delta$  and  $\hat{g}_\delta$  as defined above, find  $\hat{z}_\delta \in \mathbf{H}_{M_\delta}^1(D)$  such that

$$a_\delta(\hat{z}_\delta, \hat{\varphi}) = \ell_\delta(\hat{\varphi}) \quad \forall \hat{\varphi} \in \mathbf{H}_{M_\delta}^1(D), \quad (6.5.33)$$

where, for  $\hat{\zeta}, \hat{\varphi} \in \mathbf{H}_{M_\delta}^1(D)$ ,

$$a_\delta(\hat{\zeta}, \hat{\varphi}) := \int_D M_\delta \left( \nabla_q \hat{\zeta} \cdot \nabla_q \hat{\varphi} + \hat{\zeta} \hat{\varphi} \right) \, d\underline{q} \quad \text{and} \quad \ell_\delta(\hat{\varphi}) := \int_D M_\delta \hat{g}_\delta \hat{\varphi} \, d\underline{q}.$$

We note that, for  $\delta > 0$  and  $\underline{q} \in D$ ,  $0 < Z^{-1} \exp(-U_\delta(\frac{1}{2}r_D^2)) \leq M_\delta(\underline{q}) \leq Z^{-1}$ , and therefore  $\mathbf{L}_{M_\delta}^2(D)$  and  $\mathbf{H}_{M_\delta}^1(D)$  are homeomorphic to  $\mathbf{L}^2(D)$  and  $\mathbf{H}^1(D)$ , respectively, with equivalent respective norms, so they can be identified with  $\mathbf{L}^2(D)$  and  $\mathbf{H}^1(D)$ , respectively.

As in the case of (6.5.32), the existence of a unique solution  $\hat{z}_\delta \in \mathbf{H}_{M_\delta}^1(D)$  to (6.5.33) follows by the Lax–Milgram theorem, and

$$\|\hat{z}_\delta\|_{\mathbf{L}_{M_\delta}^2(D)} \leq \|\hat{z}_\delta\|_{\mathbf{H}_{M_\delta}^1(D)} \leq \|\hat{g}_\delta\|_{\mathbf{L}_{M_\delta}^2(D)} = \|\hat{g}\|_{\mathbf{L}_M^2(D)}. \quad (6.5.34)$$

Also, by (standard) elliptic regularity theory,  $\hat{z}_\delta \in \mathbf{H}_{M_\delta}^1(D) = \mathbf{H}^1(D)$  belongs to  $\mathbf{H}^2(D) = \mathbf{H}_{M_\delta}^2(D)$  for all  $\delta > 0$ .

Since  $\mathbf{C}_0^\infty(D) \subset \mathbf{H}_{M_\delta}^1(D)$  for any  $\delta > 0$ , on choosing  $\hat{\varphi} \in \mathbf{C}_0^\infty(D)$  in (6.5.33), it follows that

$$-\nabla_q \cdot (M_\delta \nabla_q \hat{z}_\delta) + M_\delta \hat{z}_\delta = M_\delta \hat{g}_\delta \quad \text{in } \mathcal{D}'(D), \quad (6.5.35)$$

i.e. in the sense of distributions on  $D$ . As  $M_\delta \in C^\infty(D)$ , multiplication by  $M_\delta$  of elements of  $\mathcal{D}'(D)$  is correctly defined; thus, by the Leibniz rule for differentiation of the product of a  $C^\infty(D)$  function and an element of  $\mathcal{D}'(D)$ , (6.5.35) yields

$$-M_\delta \Delta_q \hat{z}_\delta - \nabla_q M_\delta \cdot \nabla_q \hat{z}_\delta + M_\delta \hat{z}_\delta = M_\delta \hat{g} \quad \text{in } \mathcal{D}'(D). \quad (6.5.36)$$

Noting that  $M_\delta$  and  $U'_\delta$  satisfy an identity analogous to (6.2.5), and that since  $M_\delta^{-1} \in C^\infty(D)$  multiplication by  $M_\delta^{-1}$  in  $\mathcal{D}'(D)$  is meaningful, multiplying (6.5.36) by  $M_\delta^{-1}$  we deduce that

$$-\Delta_q \hat{z}_\delta + U'_\delta q \cdot \nabla_q \hat{z}_\delta + \hat{z}_\delta = \hat{g}_\delta \quad \text{in } \mathcal{D}'(D). \quad (6.5.37)$$

As  $q \mapsto U'_\delta(\frac{1}{2}|q|^2)q$  belongs to  $[C^\infty(D)]^d$ , the dot-product in the second term of (6.5.37) is meaningful as an operation in  $[\mathcal{D}'(D)]^d$ . Taking the partial derivative in  $\mathcal{D}'(D)$  of (6.5.37) with respect to  $q_i$ , the  $i$ th component of  $q$ , gives

$$-\Delta_q \frac{\partial \hat{z}_\delta}{\partial q_i} + q_i U''_\delta q \cdot \nabla_q \hat{z}_\delta + U'_\delta \frac{\partial \hat{z}_\delta}{\partial q_i} + U'_\delta q \cdot \nabla_q \frac{\partial \hat{z}_\delta}{\partial q_i} + \frac{\partial \hat{z}_\delta}{\partial q_i} = \frac{\partial \hat{g}_\delta}{\partial q_i} \quad \text{in } \mathcal{D}'(D), \quad i \in \{1, \dots, d\}. \quad (6.5.38)$$

For  $\hat{\varphi} \in C_0^\infty(D)$ , we have  $M_\delta \frac{\partial \hat{\varphi}}{\partial q_i} \in C_0^\infty(D)$ , and therefore (6.5.38) implies that

$$\begin{aligned} & \left\langle -\Delta_q \frac{\partial \hat{z}_\delta}{\partial q_i}, M_\delta \frac{\partial \hat{\varphi}}{\partial q_i} \right\rangle + \left\langle q_i U''_\delta q \cdot \nabla_q \hat{z}_\delta, M_\delta \frac{\partial \hat{\varphi}}{\partial q_i} \right\rangle + \left\langle U'_\delta \frac{\partial \hat{z}_\delta}{\partial q_i}, M_\delta \frac{\partial \hat{\varphi}}{\partial q_i} \right\rangle \\ & + \left\langle U'_\delta q \cdot \nabla_q \frac{\partial \hat{z}_\delta}{\partial q_i}, M_\delta \frac{\partial \hat{\varphi}}{\partial q_i} \right\rangle + \left\langle \frac{\partial \hat{z}_\delta}{\partial q_i}, M_\delta \frac{\partial \hat{\varphi}}{\partial q_i} \right\rangle = \left\langle \frac{\partial \hat{g}_\delta}{\partial q_i}, M_\delta \frac{\partial \hat{\varphi}}{\partial q_i} \right\rangle \\ & \quad \forall \hat{\varphi} \in C_0^\infty(D), \quad i \in \{1, \dots, d\}; \quad (6.5.39) \end{aligned}$$

where  $\langle \cdot, \cdot \rangle$  denotes the duality pairing on  $\mathcal{D}'(D) \times C_0^\infty(D)$ . Writing  $\Delta_q = \nabla_q \cdot \nabla_q$  in the first term on the left-hand side of (6.5.39), passing  $\nabla_q$  to the test function in this term, using the Leibniz rule in  $C^\infty(D)$ , noting (6.2.5) and that  $U'_\delta \in C^\infty(D)$ , whereby multiplication in  $\mathcal{D}'(D)$  by  $U'_\delta$  is legitimate, and observing that one of the two terms that result upon the use of the Leibniz rule from the first term on the left-hand side of (6.5.39) cancels with the fourth term on the left-hand side of (6.5.39), gives

$$\begin{aligned} & \left\langle \nabla_q \frac{\partial \hat{z}_\delta}{\partial q_i}, M_\delta \nabla_q \frac{\partial \hat{\varphi}}{\partial q_i} \right\rangle + \left\langle q_i U''_\delta q \cdot \nabla_q \hat{z}_\delta, M_\delta \frac{\partial \hat{\varphi}}{\partial q_i} \right\rangle \\ & + \left\langle U'_\delta \frac{\partial \hat{z}_\delta}{\partial q_i}, M_\delta \frac{\partial \hat{\varphi}}{\partial q_i} \right\rangle + \left\langle \frac{\partial \hat{z}_\delta}{\partial q_i}, M_\delta \frac{\partial \hat{\varphi}}{\partial q_i} \right\rangle = \left\langle \frac{\partial \hat{g}_\delta}{\partial q_i}, M_\delta \frac{\partial \hat{\varphi}}{\partial q_i} \right\rangle \end{aligned}$$

for all  $\hat{\varphi} \in C_0^\infty(D)$ ,  $i \in \{1, \dots, d\}$ . Summing over  $i = 1 \rightarrow d$ , we deduce the identity

$$\begin{aligned} \mathcal{A}_\delta(\hat{z}_\delta, \hat{\varphi}) & := \int_D M_\delta \nabla_q \nabla_q \hat{z}_\delta : \nabla_q \nabla_q \hat{\varphi} \, dq + \int_D M_\delta U''_\delta (q \cdot \nabla_q \hat{z}_\delta) (q \cdot \nabla_q \hat{\varphi}) \, dq \\ & + \int_D M_\delta (U'_\delta + 1) \nabla_q \hat{z}_\delta \cdot \nabla_q \hat{\varphi} \, dq \\ & = - \int_D \hat{g}_\delta \nabla_q \cdot (M_\delta \nabla_q \hat{\varphi}) \, dq = - \int_D M_\delta \hat{g}_\delta \Delta_q \hat{\varphi} \, dq + \int_D M_\delta \hat{g}_\delta U'_\delta q \cdot \nabla_q \hat{\varphi} \, dq =: \mathcal{L}_\delta(\hat{\varphi}). \end{aligned}$$

for all  $\varphi \in C_0^\infty(D)$ . Consider the norm  $\|\cdot\|_{\mathcal{H}_{M_\delta}^2(D)}$  defined by

$$\|\hat{\zeta}\|_{\mathcal{H}_{M_\delta}^2(D)}^2 := \int_D M_\delta \left[ |\nabla_q \nabla_q \hat{\zeta}|^2 + U_\delta'' |q \cdot \nabla_q \hat{\zeta}|^2 + (U_\delta' + 1) |\nabla_q \hat{\zeta}|^2 + |\hat{\zeta}|^2 \right] dq.$$

We observe that  $\|\cdot\|_{\mathcal{H}_{M_\delta}^2(D)}$  is an equivalent norm on  $H_{M_\delta}^2(D) = H^2(D)$  and, in particular,  $\|\hat{z}_\delta\|_{\mathcal{H}_{M_\delta}^2(D)} < \infty$ . Next, we show that  $\|\hat{z}_\delta\|_{\mathcal{H}_{M_\delta}^2(D)}$  is, in fact, bounded, independent of  $\delta > 0$ . Recalling (6.5.34) we have that

$$\begin{aligned} \|\hat{z}_\delta\|_{\mathcal{H}_{M_\delta}^2(D)}^2 &= \mathcal{A}(\hat{z}_\delta, \hat{z}_\delta) + (M_\delta \hat{z}_\delta, \hat{z}_\delta)_D = \mathcal{L}_\delta(\hat{z}_\delta) + (M_\delta \hat{z}_\delta, \hat{z}_\delta)_D = \mathcal{L}_\delta(\hat{z}_\delta) + \|\hat{z}_\delta\|_{L_{M_\delta}^2(D)}^2 \\ &\leq \|\hat{g}_\delta\|_{L_{M_\delta}^2(D)} \|\Delta_q \hat{z}_\delta\|_{L_{M_\delta}^2(D)} + \|\hat{g}_\delta\|_{L_{M_\delta}^2(D)} \|U_\delta' q \cdot \nabla_q \hat{z}_\delta\|_{L_{M_\delta}^2(D)} + \|\hat{g}_\delta\|_{L_{M_\delta}^2(D)} \|\hat{z}_\delta\|_{L_{M_\delta}^2(D)}. \end{aligned}$$

Since  $\|\Delta_q \hat{z}_\delta\|_{L_{M_\delta}^2(D)} \leq d^{\frac{1}{2}} \|\nabla_q \nabla_q \hat{z}_\delta\|_{L_{M_\delta}^2(D)}$  and, thanks to (6.2.9b),  $[U_\delta'(s)]^2 \leq c_5 U_\delta''(s)$ ,  $s \in [0, \frac{1}{2}r_D^2)$ , we thus have that

$$\|\hat{z}_\delta\|_{\mathcal{H}_{M_\delta}^2(D)}^2 \leq (d + c_5 + 1)^{\frac{1}{2}} \|\hat{g}_\delta\|_{L_{M_\delta}^2(D)} \|\hat{z}_\delta\|_{\mathcal{H}_{M_\delta}^2(D)},$$

which implies that

$$\|\hat{z}_\delta\|_{\mathcal{H}_{M_\delta}^2(D)}^2 \leq \|\hat{z}_\delta\|_{\mathcal{H}_{M_\delta}^2(D)}^2 \leq (d + c_5 + 1) \|\hat{g}_\delta\|_{L_{M_\delta}^2(D)}^2 = (d + c_5 + 1) \|\hat{g}\|_{L_M^2(D)}^2.$$

Since  $M(q) \leq M_\delta(q)$  for all  $q \in D$  and  $\delta > 0$ , we deduce that

$$\|\hat{z}_\delta\|_{H_M^2(D)}^2 \leq (d + c_5 + 1) \|\hat{g}\|_{L_M^2(D)}^2.$$

Since  $\{\hat{z}_\delta\}_{\delta>0}$  is bounded in  $H_M^2(D)$ , there exists  $\hat{z}_0 \in H_M^2(D)$  and a subsequence, still denoted  $\{\hat{z}_\delta\}_{\delta>0}$ , such that  $\hat{z}_\delta \rightharpoonup \hat{z}_0$  weakly in  $H_M^2(D)$  as  $\delta \rightarrow 0_+$ . By the weak lower semicontinuity of the norm function  $\hat{\zeta} \mapsto \|\hat{\zeta}\|_{H_M^2(D)}$ ,

$$\|\hat{z}_0\|_{H_M^2(D)}^2 \leq \|\hat{z}_0\|_{H_M^2(D)}^2 \leq (d + c_5 + 1) \|\hat{g}\|_{L_M^2(D)}^2. \quad (6.5.40)$$

Since for  $\zeta \geq 1$  (cf. (6.2.9a)) the space  $H_M^2(D)$  is compactly embedded into  $H_M^1(D)$  (see Lemma 5.2 in Antoci [5]),  $\{\hat{z}_\delta\}_{\delta>0}$  is strongly convergent to  $\hat{z}_0$  in  $H_M^1(D)$  as  $\delta \rightarrow 0_+$ . Noting that  $\{M_\delta\}_{\delta>0}$  converges to  $M$  uniformly on  $\bar{D}$  as  $\delta \rightarrow 0_+$  it follows that, as  $\delta \rightarrow 0_+$ ,

$$\ell_\delta(\hat{\varphi}) = \int_D M_\delta \hat{g}_\delta \hat{\varphi} dq = \int_D (M_\delta)^{\frac{1}{2}} \hat{g} M^{\frac{1}{2}} \hat{\varphi} dq \rightarrow \int_D M^{\frac{1}{2}} \hat{g} M^{\frac{1}{2}} \hat{\varphi} dq = \int_D M \hat{g} \hat{\varphi} dq = \ell(\hat{\varphi})$$

for all  $\hat{\varphi} \in C^\infty(\bar{D})$ , and  $a_\delta(\hat{z}_\delta, \hat{\varphi}) \rightarrow a(\hat{z}_0, \hat{\varphi})$  for all  $\hat{\varphi} \in C^\infty(\bar{D})$ . Hence, passage to the limit  $\delta \rightarrow 0_+$  in (6.5.33) yields  $a(\hat{z}_0, \hat{\varphi}) = \ell(\hat{\varphi})$  for all  $\hat{\varphi} \in C^\infty(\bar{D})$ . Since  $C^\infty(\bar{D})$  is dense in  $H_M^1(D)$ , also  $a(\hat{z}_0, \hat{\varphi}) = \ell(\hat{\varphi})$  for all  $\hat{\varphi} \in H_M^1(D)$ . However,  $\hat{z} \in H_M^1(D)$  is the unique solution to (6.5.32), and therefore  $\hat{z} = \hat{z}_0 \in H_M^2(D)$ , and then by (6.5.40),

$$\|\hat{z}\|_{H_M^2(D)}^2 \leq \|\hat{z}\|_{H_M^2(D)}^2 \leq (d + c_5 + 1) \|\hat{g}\|_{L_M^2(D)}^2. \quad (6.5.41)$$

That completes the proof of the elliptic regularity result that we need in order to proceed with the proof of stability, in the  $M$ -weighted  $H^1$  norm, of the orthogonal projector in the  $M$ -weighted  $L^2$  inner product on  $D$ .

Taking  $g = \hat{\psi} - P_h^q \hat{\psi}$  in (6.5.32), where  $P_h^q$  denotes the orthogonal projector in the  $M$ -weighted  $H^1$  inner product on  $D$ , we have from the symmetry of the bilinear form  $a(\cdot, \cdot)$ , the definitions of  $\hat{z}$  and  $P_h^q$ , the Cauchy–Schwarz inequality and (6.5.31) that

$$\begin{aligned} \|\hat{\psi} - P_h^q \hat{\psi}\|_{L_M^2(D)}^2 &\leq a(\hat{\psi} - P_h^q \hat{\psi}, \hat{z}) = a(\hat{\psi} - P_h^q \hat{\psi}, \hat{z} - P_h^q \hat{z}) \\ &\leq \|\hat{\psi} - P_h^q \hat{\psi}\|_{H_M^1(D)} \|\hat{z} - P_h^q \hat{z}\|_{H_M^1(D)} \\ &\leq C h_q \|\hat{\psi} - P_h^q \hat{\psi}\|_{H_M^1(D)} |\hat{z}|_{H_M^2(D)}. \end{aligned}$$

The elliptic regularity result (6.5.41) with  $\hat{g} = \hat{\psi} - P_h^q \hat{\psi}$  gives

$$|\hat{z}|_{H_M^2(D)} \leq (d + c_5 + 1)^{\frac{1}{2}} \|\hat{\psi} - P_h^q \hat{\psi}\|_{L_M^2(D)}.$$

We thus have that

$$\|\hat{\psi} - P_h^q \hat{\psi}\|_{L_M^2(D)} \leq C h_q \|\hat{\psi} - P_h^q \hat{\psi}\|_{H_M^1(D)}. \quad (6.5.42)$$

Now, by the first inverse inequality in the  $M$ -weighted  $H^1$  norm on  $D$  stated in (6.4.53a), and (6.5.42),

$$\begin{aligned} \|\hat{\psi} - Q_h^q \hat{\psi}\|_{H_M^1(D)} &\leq \|\hat{\psi} - P_h^q \hat{\psi}\|_{H_M^1(D)} + \|P_h^q \hat{\psi} - Q_h^q \hat{\psi}\|_{H_M^1(D)} \\ &\leq \|\hat{\psi} - P_h^q \hat{\psi}\|_{H_M^1(D)} + C_{\text{inv}} h_q^{-1} \|P_h^q \hat{\psi} - Q_h^q \hat{\psi}\|_{L_M^2(D)} \\ &\leq \|\hat{\psi} - P_h^q \hat{\psi}\|_{H_M^1(D)} + C_{\text{inv}} h_q^{-1} \|\hat{\psi} - P_h^q \hat{\psi}\|_{L_M^2(D)} + C_{\text{inv}} h_q^{-1} \|\hat{\psi} - Q_h^q \hat{\psi}\|_{L_M^2(D)} \\ &\leq \|\hat{\psi} - P_h^q \hat{\psi}\|_{H_M^1(D)} + 2 C_{\text{inv}} h_q^{-1} \|\hat{\psi} - P_h^q \hat{\psi}\|_{L_M^2(D)} \leq (1 + C) \|\hat{\psi} - P_h^q \hat{\psi}\|_{H_M^1(D)}. \end{aligned}$$

In particular the last inequality implies that

$$\|\hat{\psi} - Q_h^q \hat{\psi}\|_{H_M^1(D)} \leq 2(1 + C) \|\hat{\psi}\|_{H_M^1(D)} \quad \forall \hat{\psi} \in H_M^1(D)$$

and therefore also,

$$\|Q_h^q \hat{\psi}\|_{H_M^1(D)} \leq (3 + 2C) \|\hat{\psi}\|_{H_M^1(D)} \quad \forall \hat{\psi} \in H_M^1(D). \quad (6.5.43)$$

It remains to prove that the projector  $Q_h^M = Q_h^x Q_h^q = Q_h^q Q_h^x$ , where  $Q_h^x$  is the orthogonal projector in  $L^2(\Omega)$  onto  $X_h^x$  and  $Q_h^q$  is the orthogonal projector in  $L_M^2(D)$  onto  $X_h^q$ , is stable in the norm of  $\hat{X} := H^1(\Omega \times D; M)$ . Indeed,

$$\begin{aligned} \|Q_h^M \hat{\psi}\|_{\hat{X}}^2 &= \|Q_h^x Q_h^q \hat{\psi}\|_{\hat{X}}^2 = \int_{\Omega \times D} M \left[ |Q_h^x Q_h^q \hat{\psi}|^2 + |\nabla_x (Q_h^x Q_h^q \hat{\psi})|^2 + |\nabla_q (Q_h^x Q_h^q \hat{\psi})|^2 \right] dq dx \\ &\leq \int_D M \|Q_h^x (Q_h^q \hat{\psi})(\cdot, q)\|_{H^1(\Omega)}^2 dq + \int_{\Omega} \|Q_h^q (Q_h^x \hat{\psi})(x, \cdot)\|_{H_M^1(D)}^2 dx \\ &\leq C \left[ \int_D M \|Q_h^q \hat{\psi}(\cdot, q)\|_{H^1(\Omega)}^2 dq + \int_{\Omega} \|Q_h^x \hat{\psi}(x, \cdot)\|_{H_M^1(D)}^2 dx \right] \\ &\leq C \left[ \int_D M \|\hat{\psi}(\cdot, q)\|_{H^1(\Omega)}^2 dq + \int_{\Omega} \|\hat{\psi}(x, \cdot)\|_{H_M^1(D)}^2 dx \right] \\ &= C \|\hat{\psi}\|_{\hat{X}}^2, \end{aligned}$$

where in the transition to the third line we used the stability of  $Q_h^x$  in the  $H^1(\Omega)$  norm, and the stability of  $Q_h^q$  in the  $H_M^1(D)$  norm stated in (6.5.43). In the transition to the penultimate line we used Fubini's theorem to exchange the order of integration, together with the fact that  $Q_h^q$  is a contraction in the norm of  $L_M^2(D)$  and  $Q_h^x$  is a contraction in the norm of  $L^2(\Omega)$ .

# References

- [1] D. J. Acheson. *Elementary Fluid Dynamics*. Oxford University Press, 1990.
- [2] L. Ambrosio. Transport equation and Cauchy problem for  $BV$  vector fields. *Invent. Math.*, 158:227–260, 2004.
- [3] A. Ammar, B. Mokdad, F. Chinesta, and R. Keunings. A new family of solvers for some classes of multidimensional partial differential equations encountered in kinetic theory modeling of complex fluids. *J. Non-Newtonian Fluid Mech.*, 139:153–176, 2006.
- [4] A. Ammar, B. Mokdad, F. Chinesta, and R. Keunings. A new family of solvers for some classes of multidimensional partial differential equations encountered in kinetic theory modelling of complex fluids. part ii: Transient simulation using space-time separated representations. *J. Non-Newtonian Fluid Mech.*, 144:98–121, 2007.
- [5] F. Antoci. Some necessary and some sufficient conditions for the compactness of the embedding of weighted Sobolev spaces. *Ricerche Mat.*, 52(1):55–71, 2003.
- [6] A. Arnold, P. Markowich, G. Toscani, and A. Unterreiter. On convex Sobolev inequalities and the rate of convergence to equilibrium for Fokker-Planck type equations. *Comm. Partial Differential Equations*, 26(1-2):43–100, 2001.
- [7] F. G. Avkhadiev and K.-J. Wirths. Unified Poincaré and Hardy inequalities with sharp constants for convex domains. *ZAMM Z. Angew. Math. Mech.*, 87(8-9):632–642, 2007.
- [8] S. Balay, K. Buschelman, V. Eijkhout, W. D. Gropp, D. Kaushik, M. G. Knepley, L. C. McInnes, B. F. Smith, and H. Zhang. PETSc users manual. Technical Report ANL-95/11 - Revision 2.1.5, Argonne National Laboratory, 2004.
- [9] J. W. Barrett and R. Nürnberg. Convergence of a finite-element approximation of surfactant spreading on a thin film in the presence of van der Waals forces. *IMA J. Numer. Anal.*, 24:323–363, 2004.
- [10] J. W. Barrett, Ch. Schwab, and E. Süli. Existence of global weak solutions for some polymeric flow models. *Math. Models and Methods in Applied Sciences*, 15(3):939–983, 2005.
- [11] J. W. Barrett and E. Süli. Existence of global weak solutions to some regularized kinetic models for dilute polymers. *Multiscale Model. Simul.*, 6(2):506–546 (electronic), 2007.
- [12] J. W. Barrett and E. Süli. Existence of global weak solutions to dumbbell models for dilute polymers with microscopic cut-off. *M3AS: Mathematical Models and Methods in Applied Sciences*, 18(6):935–971, 2008.

- [13] J. W. Barrett and E. Süli. Numerical approximation of corotational dumbbell models for dilute polymers. *IMA J. Numer. Anal.*, Accepted for publication.
- [14] J. W. Barrett and E. Süli. Finite element approximation of kinetic dilute polymer models with microscopic cut-off. *M2AN Math. Model. Numer. Anal.*, Submitted for publication.
- [15] G. K. Batchelor. *An Introduction to Fluid Dynamics*. Cambridge University Press, 1967.
- [16] C. Bernardi and Y. Maday. Spectral methods. In P.G. Ciarlet and J.L. Lions, editors, *Handbook of Numerical Analysis*, volume V. Elsevier, 1997.
- [17] O. V. Besov, Ja. Kadlec, and A. Kufner. Certain properties of weight classes. *Dokl. Akad. Nauk SSSR*, 171:514–516, 1966.
- [18] O. V. Besov and A. Kufner. The density of smooth functions in weight spaces. *Czechoslovak Math. J.*, 18 (93):178–188, 1968.
- [19] A. V. Bhave, R. C. Armstrong, and R. A. Brown. Kinetic theory and rheology of dilute, nonhomogeneous polymer solutions. *J. Chem. Phys.*, 95:2988–3000, 1991.
- [20] B. Bialecki and R. Fernandes. An orthogonal spline collocation alternating direction implicit Crank–Nicolson method for linear parabolic problems on rectangles. *SIAM J. Numer. Anal.*, 36(5):1414–1434, 1999.
- [21] R. Bird, C. Curtiss, R. Armstrong, and O. Hassager. *Dynamics of Polymeric Liquids, Vol 2: Kinetic Theory*. John Wiley and Sons, 1987.
- [22] R. B. Bird, C. F. Curtiss, R. C. Armstrong, and O. Hassager. *Dynamics of Polymeric Liquids, Volume 1, Fluid Mechanics*. John Wiley and Sons, second edition, 1987.
- [23] R. B. Bird, C. F. Curtiss, R. C. Armstrong, and O. Hassager. *Dynamics of Polymeric Liquids, Volume 2, Kinetic Theory*. John Wiley and Sons, second edition, 1987.
- [24] S. Bobkov and M. Ledoux. From Brunn–Minkowski to Brascamp–Lieb and to logarithmic Sobolev inequalities. *Geom. Funct. Anal.*, 10:1028–1052, 2000.
- [25] J. Brandts, S. Korotov, M. Křížek, and J. Šolc. On acute and nonobtuse simplicial partitions. *Helsinki University of Technology, Institute of Mathematics, Research Reports*, A503, 2006.
- [26] S. C. Brenner and L. R. Scott. *The Mathematical Theory of Finite Element Methods*. Springer, second edition, 2002.
- [27] F. Brezzi and M. Fortin. *Mixed and Hybrid Finite Element Methods*. Springer-Verlag, 1991.
- [28] C. Canuto, A. Quarteroni, M. Y. Hussaini, and T. A. Zang. *Spectral Methods: Fundamentals in Single Domains*. Springer, 2006.

- [29] M. Celia and G. Pinder. An analysis of alternating-direction methods for parabolic equations. *Numerical Methods for Partial Differential Equations*, 1(1):57–70, 1985.
- [30] M. Celia and G. Pinder. Generalized alternating-direction collocation methods for parabolic equations. i. spatially varying coefficients. *Numerical Methods for Partial Differential Equations*, 6(3):193–214, 1990.
- [31] S. Cerrai. *Second-order PDEs in Finite and Infinite Dimension*, volume 1762 of *Lecture Notes in Mathematics*. Springer-Verlag, 2001.
- [32] C. Chauvière and A. Lozinski. Simulation of complex viscoelastic flows using Fokker–Planck equation: 3D FENE model. *J. Non-Newtonian Fluid Mech.*, 122:201–214, 2004.
- [33] C. Chauvière and A. Lozinski. Simulation of dilute polymer solutions using a Fokker–Planck equation. *Computers and Fluids*, 33:687–696, 2004.
- [34] P. G. Ciarlet. *The finite element method for elliptic problems*, volume 40 of *Classics in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2002. Reprint of the 1978 original [North-Holland, Amsterdam; MR0520174 (58 #25001)].
- [35] P. Clément. Approximation by finite element functions using local regularization. *Rev. Française Automat. Informat. Recherche Opérationnelle Sér. RAIRO Analyse Numérique*, 9(R-2):77–84, 1975.
- [36] W. T. Coffey, Y. P. Kalmykov, and J. T. Waldron. *The Langevin Equation: With Applications in Physics, Chemistry and Electrical Engineering*. World Scientific, 1996.
- [37] P. Constantin. Nonlinear Fokker–Planck Navier–Stokes systems. *Commun. Math. Sci.*, 3(4):531–544, 2005.
- [38] P. Delaunay, A. Lozinski, and R. G. Owens. Sparse tensor-product Fokker-Planck-based methods for nonlinear bead-spring chain models of dilute polymer solutions. *CRM Proceedings and Lecture Notes*, 41:73 – 89, 2007.
- [39] L. Desvillettes and C. Villani. On the trend to global equilibrium for spatially inhomogeneous kinetic systems: the Boltzmann equation. *Invent. Math.*, 159(2):245–316, 2005.
- [40] R. J. DiPerna and P.-L. Lions. Ordinary differential equations, transport theory and Sobolev spaces. *Invent. Math.*, 98:511–547, 1989.
- [41] J. Douglas and T. DuPont. Alternating-direction Galerkin methods on rectangles. *Numerical Solution of Partial Differential Equations, II (SYNSPADE 1970)*, pages 133–214, 1971.
- [42] Q. Du, C. Liu, and P. Yu. FENE dumbbell model and its several linear and nonlinear closure approximations. *Multiscale Model. Simul.*, 4(3):709–731, 2005.
- [43] W. E, T. Li, and P. Zhang. Well-posedness for the dumbbell model of polymeric fluids. *Comm. Math. Phys.*, 248(2):409–427, 2004.

- [44] H. Eisen, W. Heinrichs, and K. Witsch. Spectral collocation methods and polar coordinate singularities. *J. Comput. Phys.*, 96(2):241–257, 1991.
- [45] A. E. El Kareh and L. G. Leal. Existence of solutions for all Deborah numbers for a non-newtonian model modified to include diffusion. *J. Non-Newtonian Fluid Mech.*, 33:257–287, 1989.
- [46] H. Elman, D. Silvester, and A. Wathen. *Finite elements and fast iterative solvers*. Oxford Science Publications, 2005.
- [47] D. Eppstein, J. M. Sullivan, and A. Üngör. Tiling space and slabs with acute tetrahedra. *Comput. Geom.*, 27(3):237–255, 2004.
- [48] X. J. Fan. Molecular models and flow calculations: II. simulation of steady planar flow. *Acta Mechanica Sinica*, 5:216–226, 1989.
- [49] P. J. Flory. *Statistical Mechanics of Chain Molecules*. Wiley-Interscience, 1969.
- [50] C. Foias, D. D. Holm, and E. S. Titi. The Navier–Stokes-alpha model of fluid turbulence. *Phys. D*, 152/153:505–519, 2001. Advances in nonlinear mathematics and science.
- [51] V. Girault and P.-A. Raviart. *Finite Element Methods for Navier–Stokes Equations: Theory and Algorithms*. Springer, 1986.
- [52] M. Grosso, P. L. Maffettone, P. Halin, R. Keunings, and V. Legat. Flow of nematic polymers in eccentric cylinder geometry: influence of closure approximations. *J. Non-Newtonian Fluid Mech.*, 94:119–134, 2000.
- [53] G. Grün and M. Rumpf. Nonnegativity preserving numerical schemes for the thin film equation. *Numer. Math.*, 87:113–152, 2000.
- [54] P. Halin, G. Lielens, R. Keunings, and V. Legat. The Lagrangian particle method for macroscopic and micro-macro viscoelastic flow computations. *J. Non-Newtonian Fluid Mech.*, 79:387–403, 1998.
- [55] C. Helzel and F. Otto. Multiscale simulations of suspensions of rod-like molecules. *J. Comp. Phys.*, 216:52–75, 2006.
- [56] J. G. Heywood and R. Rannacher. Finite element approximation of the nonstationary Navier–Stokes problem I: Regularity of solutions and second-order error estimates for spatial discretization. *SIAM J. Numer. Anal.*, 19:275–311, 1982.
- [57] W. Huang and B. Guo. Fully discrete Jacobi-spherical harmonic spectral method for Navier-Stokes equations. *Appl. Math. Mech. (English Ed.)*, 29(4):453–476, 2008.
- [58] M. A. Hulsen, A. P. G. van Heel, and B. H. A. A. van den Brule. Simulation of viscoelastic flows using Brownian configuration fields. *J. Non-Newtonian Fluid Mech.*, 70:79–101, 1997.
- [59] J.-I. Itoh and T. Zamfirescu. Acute triangulations of the regular dodecahedral surface. *European J. Combin.*, 28(4):1072–1086, 2007.

- [60] B. Jourdain and T. Lelièvre. Mathematical analysis of a stochastic differential equation arising in the micro-macro modelling of polymeric fluids. *Probabilistic Methods in Fluids*, pages 205–223, 2003.
- [61] B. Jourdain, T. Lelièvre, and C. Le Bris. Numerical analysis of micro-macro simulations of polymeric fluid flows: A simple case. *Math. Models Methods Appl. Sci.*, 12:1205–1243, 2002.
- [62] B. Jourdain, T. Lelièvre, and C. Le Bris. Existence of solution for a micro-macro model of polymeric fluid: the FENE model. *J. Funct. Anal.*, 209(1):162–193, 2004.
- [63] B. Jourdain, T. Lelièvre, C. Le Bris, and F. Otto. Long-time asymptotics of a multiscale model for polymeric fluid flows. *Arch. Rat. Mech. Anal.*, 181:97–148, 2006.
- [64] P. Keast. Moderate-degree tetrahedral quadrature formulas. *Comput. Methods Appl. Mech. Engrg.*, 55(3):339–348, 1986.
- [65] R. Keunings. On the Peterlin approximation for finitely extensible dumbbells. *J. Non-Newtonian Fluid Mech.*, 68:85–100, 1997.
- [66] R. Keunings. A survey of computational rheology. In *XIIIth International Congress on Rheology*, Cambridge, UK, August 2000. Available at <http://www.mate.tue.nl/~hulsen>.
- [67] R. Keunings. Micro–macro methods for the multiscale simulation of viscoelastic flow using molecular models of kinetic theory. *Rheology Review*, pages 67–98, 2004.
- [68] B. S. Kirk, J. W. Peterson, R. M. Stogner, and G. F. Carey. libmesh: A C++ library for parallel adaptive mesh refinement/coarsening simulations. *Engineering with Computers*, 23(3–4):237–254, 2006.
- [69] J. G. Kirkwood. *Macromolecules*. Gordon and Breach, 1967.
- [70] D. Knezevic. *Analysis and Implementation of Numerical Methods for Simulating Dilute Polymeric Fluids*. PhD thesis, University of Oxford, 2008.
- [71] D. Knezevic and E. Süli. Spectral Galerkin approximation of Fokker–Planck equations with unbounded drift. *M2AN: Mathematical Modeling and Numerical Analysis*, Accepted for publication.
- [72] D. Knezevic and E. Süli. A heterogeneous alternating-direction method for a micro-macro dilute polymeric fluid model. *M2AN: Mathematical Modeling and Numerical Analysis*, Submitted for publication.
- [73] A. N. Kolmogorov. Über die analytischen Methoden in der Wahrscheinlichkeitsrechnung. *Math. Ann.*, 104, 1931.
- [74] S. Korotov and M. Křížek. Acute type refinements of tetrahedral partitions of polyhedral domains. *SIAM J. Numer. Anal.*, 39(2):724–733 (electronic), 2001.
- [75] S. Korotov and M. Křížek. Global and local refinement techniques yielding nonobtuse tetrahedral partitions. *Comput. Math. Appl.*, 50:1105–1113, 2005.

- [76] H. A. Kramers. The viscosity of macromolecules in a streaming fluid. *Physica*, 11(1), 1944.
- [77] A. Kufner. *Weighted Sobolev Spaces*. Teubner-Texte zur Mathematik. Teubner, 1980.
- [78] B. Lapeyre, É. Pardoux, and R. Sentis. *Introduction to Monte-Carlo Methods for Transport and Diffusion Equations*. Oxford University Press, 2003.
- [79] M. Laso and H. C. Öttinger. Calculation of viscoelastic flow using molecular models: the CONNFESSIT approach. *J. Non-Newtonian Fluid Mech.*, 47:1–20, 1993.
- [80] C. Le Bris and T. Lelièvre. Multiscale modelling of complex fluids: A mathematical initiation. In *Discrete and computational geometry (Tokyo, 2000)*, volume 66 of *Multiscale Modeling and Simulation in Science Series*, pages 49–138. Springer, Berlin, 2009.
- [81] T. Li, H. Zhang, and P. Zhang. Local existence for the dumbbell model of polymeric fluids. *Comm. Partial Differential Equations*, 29:903–923, 2004.
- [82] T. Li and P. Zhang. Mathematical analysis of multi-scale models of complex fluids. *Commun. Math. Sci.*, 5(1):1–51, 2007.
- [83] G. Lielens, P. Halin, I. Jaumain, R. Keunings, and V. Legat. New closure approximations for the kinetic theory of finitely extensible dumbbells. *J. Non-Newtonian Fluid Mech.*, 76:249–279, 1998.
- [84] F.-H. Lin, C. Liu, and P. Zhang. On a micro-macro model for polymeric fluids near equilibrium. *Comm. Pure Appl. Math.*, 60(6):838–866, 2007.
- [85] F.-H. Lin, P. Zhang, and Z. Zhang. On the global existence of smooth solution to the 2-D FENE dumbbell model. *Comm. Math. Phys.*, 277:531–553, 2008.
- [86] P.-L. Lions and N. Masmoudi. Global solutions for some Oldroyd models of non-Newtonian flows. *Chinese Ann. Math. Ser. B*, 21(2):131–146, 2000.
- [87] P.-L. Lions and N. Masmoudi. Global existence of weak solutions to some micro-macro models. *C. R. Math. Acad. Sci. Paris*, 345:15–20, 2007.
- [88] C. Liu and H. Liu. Boundary conditions for the microscopic FENE models. *SIAM J. Appl. Math.*, 68(5):1304–1315, 2008.
- [89] L. Lorenzi and M. Bertoldi. *Analytical methods for Markov semigroups*, volume 283 of *Pure and Applied Mathematics (Boca Raton)*. Chapman & Hall/CRC, Boca Raton, FL, 2007.
- [90] A. Lozinski. *Spectral methods for kinetic theory models of viscoelastic fluids*. PhD thesis, École Polytechnique Fédérale de Lausanne, 2003.
- [91] A. Lozinski and C. Chauvière. A fast solver for Fokker–Planck equation applied to viscoelastic flows calculation: 2D FENE model. *Journal of Computational Physics*, 189:607–625, 2003.

- [92] A. Lozinski, C. Chauvière, J. Fang, and R. G. Owens. Fokker–Planck simulations of fast flows of melts and concentrated polymer solutions in complex geometries. *J. Rheology*, 47:535–561, 2003.
- [93] A. Lozinski, R. G. Owens, and J. Fang. A Fokker–Planck-based numerical method for modelling non-homogeneous flows of dilute polymeric solutions. *J. Non-Newtonian Fluid Mech.*, 122:273–286, 2004.
- [94] J. N. Lyness and D. Jespersen. Moderate degree symmetric quadrature rules for the triangle. *J. Inst. Math. Appl.*, 15:19–32, 1975.
- [95] M. Marcus, V. J. Mizel, and Y. Pinchover. On the best constant for Hardy’s inequality in  $\mathbf{R}^n$ . *Trans. Amer. Math. Soc.*, 350(8):3237–3255, 1998.
- [96] N. Masmoudi. Well-posedness for the FENE dumbbell model of polymeric flows. *Comm. Pure Appl. Math.*, 61(12):1685–1714, 2008.
- [97] T. Matsushima and P. S. Marcus. A spectral method for polar coordinates. *J. Comput. Phys.*, 120:365–374, 1995.
- [98] D. A. McQuarrie. *Statistical Mechanics*. University Science Books, Second Edition, 2000.
- [99] R. Nayak. *Molecular simulation of liquid crystal polymer flow: a wavelet-finite element analysis*. PhD thesis, MIT, 1998.
- [100] J. G. Oldroyd. On the formulation of rheological equations of state. *Proc. Roy. Soc. London*, A200:523–541, 1950.
- [101] H. C. Öttinger. *Stochastic Processes in Polymeric Fluids*. Springer, 1996.
- [102] F. Otto and A. E. Tzavaras. Continuity of velocity gradients in suspensions of rod-like molecules. *Comm. Math. Phys.*, 277(3):729–758, 2008.
- [103] R. G. Owens and T. N. Phillips. *Computational Rheology*. Imperial College Press, 2002.
- [104] A. Papoulis. *Probability, Random Variables, and Stochastic Processes*. McGraw-Hill, 2 edition, 1984.
- [105] G. Da Prato and A. Lunardi. Elliptic operators with unbounded drift coefficients and Neumann boundary condition. *J. Differential Equations*, 198:35–52, 2004.
- [106] M. Renardy. An existence theorem for model equations resulting from kinetic theories of polymer solutions. *SIAM J. Math. Anal.*, 22(2):313–327, 1991.
- [107] P. E. Rouse. A theory of the linear viscoelastic properties of dilute solutions of coiling polymers. *J. Chem. Phys.*, 21:1272–1280, 1953.
- [108] J. D. Schieber. Generalized Brownian configuration field for Fokker–Planck equations including center-of-mass diffusion. *J. Non-Newtonian Fluid Mech.*, 135:179–181, 2006.
- [109] J. D. Schieber and H. C. Öttinger. The effects of bead inertia on the Rouse model. *J. Chem. Phys.*, 89(11), 1988.

- [110] W. H. A. Schilders and E. J. W. ter Maten, editors. *Numerical Methods in Electromagnetics*. Handbook of Numerical Analysis Vol. XIII. North-Holland, Amsterdam, 2005.
- [111] Ch. Schwab, E. Süli, and R.-A. Todor. Sparse finite element approximation of high-dimensional transport-dominated diffusion problems. *Math. Models Methods Appl. Sci.*, 42:777–820, 2008.
- [112] L. R. Scott and S. Zhang. Finite element interpolation of nonsmooth functions satisfying boundary conditions. *Math. Comp.*, 54(190):483–493, 1990.
- [113] V. B. Shakhmurov. Embedding and maximal regular differential operators in Sobolev-Lions spaces. *Acta Math. Sin. (Engl. Ser.)*, 22(5):1493–1508, 2006.
- [114] J. Shen. Efficient spectral Galerkin methods III: Polar and cylindrical geometries. *SIAM J. Sci. Comput.*, 18(6):1583–1604, 1997.
- [115] J. Simon. Compact sets in the space  $L^p(0, T; B)$ . *Ann. Mat. Pura Appl. (4)*, 146:65–96, 1987.
- [116] W. E. Stewart and J. P. Sørensen. Hydrodynamic interaction effects in rigid dumbbell suspensions. II. computations for steady shear flow. *Journal of Rheology*, 16(1):1–13, 1972.
- [117] E. Süli and D. F. Mayers. *An Introduction to Numerical Analysis*. Cambridge University Press, 2003.
- [118] R. Sureshkumar and A.N. Beris. Effect of artificial stress diffusivity on the stability of numerical calculations and the flow dynamics of time-dependent viscoelastic flows. *J. Non-Newtonian Fluid Mech.*, 60:53–80, 1995.
- [119] R. Temam. *Navier–Stokes Equations: Theory and Numerical Analysis*. North-Holland, Amsterdam, 3rd edition, 1984.
- [120] H. Triebel. *Interpolation Theory, Function Spaces, Differential Operators*. Second edition. Joh. Ambrosius Barth Publ., 1995.
- [121] W. T. M. Verkley. A spectral model for two-dimensional incompressible fluid flow in a circular basin I. Mathematical formulation. *J. Comput. Phys.*, 136(1):100–114, 1997.
- [122] T. von Petersdorff and C. Schwab. Numerical solution of parabolic equations in high dimensions. *M2AN, Mathematical Modelling and Numerical Analysis*, 38(1):93–127, 2004.
- [123] N. J. Walkington. Quadrature on simplices of arbitrary dimension. <http://www.math.cmu.edu/~nw0z/publications/00-CNA-023/023abs/>.
- [124] H. R. Warner. Kinetic theory and rheology of dilute suspensions of finitely extendible dumbbells. *Ind. Eng. Chem. Fundamentals*, pages 379–387, 1972.

- [125] P. Yu, Q. Du, and C. Liu. From micro to macro dynamics via a new closure approximation to the fene model of polymeric fluids. *Multiscale Model. Simul.*, 3:895–917, 2005.
- [126] H. Zhang and P. Zhang. Local existence for the FENE-dumbbell model of polymeric fluids. *Arch. Ration. Mech. Anal.*, 181(2):373–400, 2006.
- [127] H. Zhang and P. Zhang. Local existence for the FENE-dumbbell model of polymeric fluids. *Arch. Ration. Mech. Anal.*, 181(2):373–400, 2006.
- [128] Q. Zhou and A. Akhavan. A comparison of FENE and FENE-P dumbbell and chain models in turbulent flow. *J. Non-Newtonian Fluid Mech.*, 109:115–155, 2003.
- [129] O. C. Zienkiewicz, R. L. Taylor, and Zhu J. Z. *The Finite Element Method: Its basis and fundamentals*. Butterworth-Heinemann, 2005.
- [130] B. H. Zimm. Dynamics of polymer molecules in dilute solution: viscoelasticity, flow birefringence and dielectric loss. *J. Chem. Phys.*, 24:269–278, 1956.