



**HAL**  
open science

# Rigorous derivation of the macroscopic equations for the lattice Boltzmann method via the corresponding Finite Difference scheme

Thomas Bellotti

► **To cite this version:**

Thomas Bellotti. Rigorous derivation of the macroscopic equations for the lattice Boltzmann method via the corresponding Finite Difference scheme. 2022. hal-03659078v1

**HAL Id: hal-03659078**

**<https://hal.science/hal-03659078v1>**

Preprint submitted on 4 May 2022 (v1), last revised 14 Dec 2022 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Rigorous derivation of the macroscopic equations for the lattice Boltzmann method *via* the corresponding Finite Difference scheme

Thomas Bellotti

(thomas.bellotti@polytechnique.edu)

CMAP, CNRS, École polytechnique, Institut Polytechnique de Paris  
91128 Palaiseau Cedex, France.

May 4, 2022

## Abstract

Lattice Boltzmann schemes are efficient numerical methods to solve a broad range of problems under the form of conservation laws. However, they suffer from a chronic lack of clear theoretical foundations. In particular, the consistency analysis is still an open issue. We propose a rigorous derivation of the macroscopic equations for any lattice Boltzmann scheme under acoustic scaling. This is done by passing from a kinetic (lattice Boltzmann) to a macroscopic (Finite Difference) point of view at a fully discrete level in order to eliminate the non-conserved moments relaxing away from the equilibrium. We rewrite the lattice Boltzmann scheme as a multi-step Finite Difference scheme on the conserved variables, as introduced in our previous contribution. We then perform the usual consistency analysis for Finite Difference by exploiting its precise characterization using matrices of Finite Difference operators. Though we present the derivation until second-order under acoustic scaling, we provide all the elements to extend it to higher orders and to other scalings, since the kinetic-macroscopic connection is conducted at the fully discrete level. Finally, we show that our strategy yields, in a mathematically rigorous setting, the same results as previous works in the literature.

**Keywords:** Lattice Boltzmann, Finite Difference, macroscopic equations, consistency, Taylor expansions  
**2000 MSC:** 65M75, 65M06

## 1 Introduction

Lattice Boltzmann methods form a vast category of numerical schemes to address the approximation of the solution of Partial Differential Equations (PDEs) under the form of conservation laws, called macroscopic equations. These numerical schemes act in a kinetic fashion by employing a certain number  $q \in \mathbb{N}^*$  of discrete velocities, larger than the number  $N \in \mathbb{N}^*$  of macroscopic equations to be solved. The scheme proceeds *via* a kinetic-like algorithm made up of two distinct steps. The first one is a local non-linear collision phase on each site of the mesh, followed by a lattice-constrained transport which is inherently linear. The local nature of the collision phase allows for massive parallelization of the method and the fact that the “particles” are constrained to dwell on the lattice allows to implement the stream phase as a pointer shift in memory. This results in a very efficient numerical method capable of reaching problems of important size in terms of computational and memory cost. The historical seminal papers from the end of the eighties are [33] and [21], while for a general modern presentation of the lattice Boltzmann schemes and their extremely broad fields of application, including hyperbolic systems of conservation laws, the quasi-incompressible Navier-Stokes equations, multi-phase systems and porous media, the interested reader can consult [42], [19] and [30]. The presentation of this plethora of interesting applications is however beyond the scope of our contribution.

To our understanding, the highest price to pay for this highly efficient implementation of the method is the lack of pure theoretical understanding on why the overall procedure works well at approximating the solution of the target macroscopic equations. This is essentially due to the fact that – the standpoint of the lattice Boltzmann schemes being kinetic – the number of discrete velocities is larger than the number of target equations.

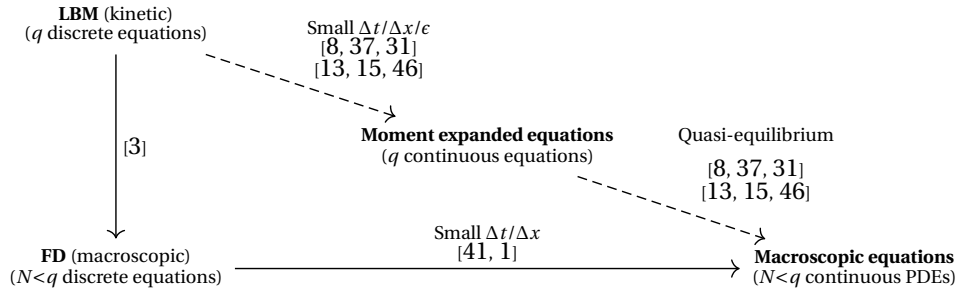


Figure 1: Different paths to recover the macroscopic equations. The formal approaches available in the literature [8, 37, 31, 13, 15, 46] rely on the path marked with dashed arrows. They perform Taylor expansions for small discretization parameters and then utilize the quasi-equilibrium of the non-conserved moments to get rid of them. Our way of proceeding is marked with full arrows: we eliminate exactly the non-conserved moments at the discrete level as in [3] and we perform the usual consistency analysis for Finite Difference schemes as in [41, 1].

Therefore, the formal analyses for lattice Boltzmann schemes available in the literature try to bridge the gap between a kinetic and a macroscopic point of view relying essentially on the quasi-equilibrium of the non-conserved variables. In particular, as far as the consistency with the macroscopic equations is concerned in the limit of small discretization parameters, two main approaches are at our disposal. The first one is based on the Chapman-Enskog expansion [7, 23] from statistical mechanics, shaped to the context of lattice Boltzmann schemes, see for example [8, 37, 31]. The second approach features the so-called equivalent equations introduced by Dubois [13, 15], consisting in performing a Taylor expansion of the scheme both for the conserved and non-conserved moments and progressively re-inject the developments order-by-order. This approach has proved to yield information in accordance with the numerical simulations, see [16, 17, 2]. Despite their proved empirical reliability and the fact that they yield the same results at the dominant orders (see [14] for instance) these two strategies are both formal: the Chapman-Enskog expansion relies on a multi-scale development with no clear mathematical foundation whereas the method of the equivalent equations writes the expansions also on the non-conserved variables, for which no equation under the form of a PDE is known. Other approaches known in the literature are the asymptotic analysis under parabolic scaling deployed in [27, 25, 26] as well as the Maxwell iteration method [46, 47], which shares strong bonds with the equivalent equations method presented before. The previous list of formal analysis techniques does not aim at being exhaustive (the interested reader can refer to [30]) and one should be aware that, despite efforts in this direction [6], there is no consensus on which is the right method to use [30].

A staple of all the previously mentioned approaches is that the expansion for the discretization parameters (time and space steps) tending to zero is performed on the kinetic numerical scheme, where both conserved and non-conserved variables are present. Eventually, the non-conserved variables are formally eliminated from the continuous formulation by scaling arguments, so to speak, using quasi-equilibrium. This corresponds to follow the diagonal path on Figure 1. In this contribution, we develop the other path, namely the top-down movement followed by the left-right one on Figure 1. In particular, in order to fill the hollow between lattice Boltzmann schemes and traditional approaches known to numerical analysts, such as Finite Difference schemes, we recently introduced [3] a formalism to recast any lattice Boltzmann scheme, regardless of its linearity, as a multi-step Finite Difference scheme solely on the conserved moments. It should be stressed that our standpoint, where lattice Boltzmann schemes are studied in terms of their Finite Difference counterpart, must not be seen as the right way of implementing them, because one would lose most of the previously mentioned computational efficiency coming from the kinetic vision. Conversely, our way of writing the scheme should be seen as a sort of one-way mathematical transform to pass from a kinetic standpoint to a macroscopic one in a purely discrete setting. The elimination of the non-conserved moments is carried exactly on the discrete formulation by algebraic devices, thus independently from the time-space scaling. The price to pay for the non-conserved moments relaxing away from the equilibria is the multi-step nature of the Finite Difference scheme. In our previous proposal

[3], it has been crucial to be able to provide, thanks to a systematic mathematical approach, a precise description of the main ingredient needed to reduce the lattice Boltzmann scheme to a Finite Difference scheme, namely the characteristic polynomial of matrices of Finite Difference operators. We are therefore allowed to utilize this characteristic polynomial as a tool satisfying certain properties alone from the particular underlying lattice Boltzmann scheme. Quite the opposite, using the algorithm proposed by [18], one is compelled to explicitly write down the corresponding Finite Difference scheme in order to perform the Taylor expansions to recover the continuous macroscopic equations. In our case, the mathematical understanding that we *a priori* have on the corresponding (macroscopic) Finite Difference schemes, regardless of the (kinetic) lattice Boltzmann scheme it stands for, allows the following theoretical discussion.

The main findings presented in the present paper are the following.

- We propose a procedure to rigorously analyze the consistency of any lattice Boltzmann scheme *via* its corresponding Finite Difference scheme.
- In the case of acoustic scaling between space and time discretization, we rigorously find the expression of the macroscopic PDEs approximated by any scheme until second order.
- Under acoustic scaling, these macroscopic PDEs are the same than the ones obtained by [15] until second order and to the ones from [47] at any order.

The paper is structured as follows. In Section 2, we set notations and assumptions concerning the lattice Boltzmann schemes we shall work with. Section 3 is devoted to recall the main results from our previous work [3] concerning the recast of any lattice Boltzmann scheme as a Finite Difference scheme. These results are then stated in a slightly different manner, facilitating the following analysis. The main result of the work is stated in Section 4 and comes under the form of two theorems, which are eventually proved in Section 5 under the assumption of dealing with one conservation law, for the sake of keeping the presentation and the notations as simple as possible. In Section 6, we indicate how the previous proof is easily extended to several conservation laws, whereas Section 7 is devoted to hint the links with some available approaches to find the macroscopic equations available in the literature. The conclusions and perspectives of this work are drawn in Section 8.

## 2 Lattice Boltzmann schemes

To start our contribution, we present the classical framework of the multiple-relaxation-times (known as MRT) lattice Boltzmann schemes, as introduced by [11]. For the sake of simplicity, we do not consider source terms which can be effortlessly introduced in the analysis. This fixes the perimeter of the schemes we shall be allowed to treat and study in the sequel.

### 2.1 Spatial and temporal discretization

We set the problem in spatial dimension  $d = 1, 2, 3$  considering the whole space  $\mathbb{R}^d$  since this work is not focused on the enforcement of boundary conditions. The space is discretized by a  $d$ -dimensional lattice denoted  $\mathcal{L} := \Delta x \mathbb{Z}^d$  with constant step  $\Delta x > 0$ . The time is uniformly discretized with step  $\Delta t > 0$ , rendering a time lattice  $\mathcal{T} := \Delta t \mathbb{N}$ . The role of the initial conditions is not investigated and is a subject on its own, see [44, 38]. We introduce the so-called “lattice velocity”  $\lambda > 0$  defined by  $\lambda := \Delta x / \Delta t$ . Observe that in the sequel, namely in Section 4, we shall introduce a particular relation between space step  $\Delta x$  and  $\Delta t$  when  $\Delta x \rightarrow 0$ , in order to provide the main results of the work, as done in [13, 15]. However, until the end of Section 3, the discussion remains valid for any choice of these parameters.

### 2.2 Discrete velocities

The discrete velocities are an essential ingredient of any lattice Boltzmann scheme. One has to choose  $(\mathbf{e}_j)_{j=1}^{j=q} \subset \mathbb{R}^d$  with  $q \in \mathbb{N}^*$ , discrete velocities, which are multiple of the lattice velocity  $\lambda$ , namely  $\mathbf{e}_j = \lambda \mathbf{c}_j$  for any  $j \in [1..q]^1$  with  $(\mathbf{c}_j)_{j=1}^{j=q} \subset \mathbb{Z}^d$ . Thus, the virtual particles are stuck to the lattice  $\mathcal{L}$  at each time step of the method. We denote the

<sup>1</sup>This shall be a notation to indicate closed intervals of integers, namely for  $a, b \in \mathbb{Z}$  with  $a \leq b$ , then  $[a..b] := \{a, a+1, \dots, b\}$ .

distribution density of the virtual particles moving with velocity  $\mathbf{e}_j$  by  $f_j = f_j(t, \mathbf{x})$  for every  $j \in [1..q]$ , depending on the space and time variable.

### 2.3 Lattice Boltzmann algorithm: collide and stream

As mentioned in the Introduction, any lattice Boltzmann scheme consists in a kinetic algorithm made up of two phases: a local collision phase performed on each site of the lattice  $\mathcal{L}$  and a stream phase where particles are exchanged between different sites of the lattice. Let us independently introduce them.

- *Collision phase.* We adopt the general point of view of the multiple-relaxation-times schemes, where the collision phase is written as a diagonal relaxation towards some equilibria in the moments basis, see [11]. We introduce a change of basis called moment matrix  $\mathbf{M} \in \text{GL}_q(\mathbb{R})$ . Gathering the distributions into  $\mathbf{f} = (f_1, \dots, f_q)^\top$ , the moments are recovered by  $\mathbf{m} = \mathbf{M}\mathbf{f}$  and *vice versa*. We also introduce
  - the matrix  $\mathbf{I} \in \text{GL}_q(\mathbb{R})$ , the identity matrix of size  $q$ ;
  - the matrix  $\mathbf{S} \in \mathcal{M}_q(\mathbb{R})$ , called relaxation matrix. This matrix is diagonal with  $N \in [1..q-1]$  being the number of conserved moments  $\mathbf{S} = \text{diag}(s_1, \dots, s_N, s_{N+1}, \dots, s_q)$ , where  $s_i \in \mathbb{R}$  for  $i \in [1..N]$  for the conserved moments and  $s_i \in ]0, 2]$  for  $i \in [N+1..q]$ , see [13], for the non-conserved ones. Observe that the relaxation parameters corresponding to the conserved moments do not play any role in the lattice Boltzmann algorithm, therefore the matrix  $\mathbf{S}$  can be singular, without any specific issue. In particular, we shall prove in Section 3.4 that the choice of relaxation parameter for the conserved variables does not have any influence on the outcomes presented in this work. For the sake of presentation, we start numbering the moments by the conserved ones;
  - we employ the notation  $\mathbf{m}^{\text{eq}}(t, \mathbf{x}) = \mathbf{m}^{\text{eq}}(m_1(t, \mathbf{x}), \dots, m_N(t, \mathbf{x}))$  for  $t \in \mathfrak{J}$  and  $\mathbf{x} \in \mathcal{L}$ , where  $\mathbf{m}^{\text{eq}} : \mathbb{R}^N \rightarrow \mathbb{R}^q$  are possibly non-linear functions of the conserved moments. In order to guarantee that the first  $N$  moments are conserved through the collision process, the constraints

$$m_i^{\text{eq}}(m_1, \dots, m_N) = m_i, \quad \forall i \in [1..N], \quad (1)$$

must hold [4].

Let  $t \in \mathfrak{J}$  and  $\mathbf{x} \in \mathcal{L}$ , the collision phase reads, denoting by  $\star$  any post-collision state

$$\mathbf{m}^\star(t, \mathbf{x}) = (\mathbf{I} - \mathbf{S})\mathbf{m}(t, \mathbf{x}) + \mathbf{S}\mathbf{m}^{\text{eq}}(t, \mathbf{x}). \quad (2)$$

- *Stream phase.* The stream phase is diagonal in the space of the distributions and consist in an exact upwind advection of the particle distribution densities. It can be written, for  $t \in \mathfrak{J}$  and  $\mathbf{x} \in \mathcal{L}$ , as

$$f_j(t + \Delta t, \mathbf{x}) = f_j^\star(t, \mathbf{x} - \mathbf{c}_j \Delta x), \quad (3)$$

for any  $j \in [1..q]$ .

## 3 Finite Difference formulation of a lattice Boltzmann scheme

Having defined the lattice Boltzmann schemes, we briefly introduce the setting allowing us to rewrite any lattice Boltzmann scheme (kinetic) as a multi-step Finite Difference scheme (macroscopic) on the  $N$  conserved moments of interest. The interested reader can refer to our previous contribution [3] for more details. Then, the formulation of the multi-step Finite Difference scheme is given using a more compact notation which is more suitable to the following discussion. We start with the assumptions needed in the sequel.

**Assumptions 3.1** (Finite Difference assumptions). The entries of  $\mathbf{M}$  and  $\mathbf{S}$  can depend on  $\Delta x$  and/or on  $\Delta t$  but cannot be a function of the space and time variables.

### 3.1 Algebraic setting

Let us first introduce the necessary algebraic setting. In particular, we define the shift operators associated with each discrete velocity as well as the derived Finite Difference operators in space. In the following Definition, the time variable does not play any role since kept frozen, thus it is not listed for the sake of readability.

**Definition 3.2** (Shift and Finite Difference operators in space). Let  $\mathbf{z} \in \mathbb{Z}^d$ , then the associated shift operator on the lattice  $\mathcal{L}$ , denoted  $\mathfrak{t}_{\mathbf{z}}$ , is defined in the following way. Take  $f : \mathcal{L} \rightarrow \mathbb{R}$  be any function defined on the lattice, then the action of  $\mathfrak{t}_{\mathbf{z}}$  is

$$(\mathfrak{t}_{\mathbf{z}}f)(\mathbf{x}) = f(\mathbf{x} - \mathbf{z}\Delta x), \quad \forall \mathbf{x} \in \mathcal{L}.$$

We also introduce  $\mathbb{T} := \{\mathfrak{t}_{\mathbf{z}} \mid \mathbf{z} \in \mathbb{Z}^d\} \cong \mathbb{Z}^d$ . The product  $\circ : \mathbb{T} \times \mathbb{T} \rightarrow \mathbb{T}$  of two shift operators is defined by

$$\mathfrak{t}_{\mathbf{z}} \circ \mathfrak{t}_{\mathbf{w}} := \mathfrak{t}_{\mathbf{z}+\mathbf{w}}, \quad \forall \mathbf{z}, \mathbf{w} \in \mathbb{Z}^d.$$

The set of Finite Difference operators on the lattice  $\mathcal{L}$  is defined as

$$\mathbb{D} := \mathbb{R}\mathbb{T} = \left\{ \sum_{\mathfrak{t} \in \mathbb{T}} \alpha_{\mathfrak{t}} \mathfrak{t}, \text{ where } \alpha_{\mathfrak{t}} \in \mathbb{R} \text{ and } \alpha_{\mathfrak{t}} = 0 \text{ almost everywhere} \right\}, \quad (4)$$

the group ring (or group algebra) of  $\mathbb{T}$  over  $\mathbb{R}$ . The sum  $+$  :  $\mathbb{D} \times \mathbb{D} \rightarrow \mathbb{D}$  and the product  $\circ$  :  $\mathbb{D} \times \mathbb{D} \rightarrow \mathbb{D}$  of two elements are defined by

$$\left( \sum_{\mathfrak{t} \in \mathbb{T}} \alpha_{\mathfrak{t}} \mathfrak{t} \right) + \left( \sum_{\mathfrak{t} \in \mathbb{T}} \beta_{\mathfrak{t}} \mathfrak{t} \right) = \sum_{\mathfrak{t} \in \mathbb{T}} (\alpha_{\mathfrak{t}} + \beta_{\mathfrak{t}}) \mathfrak{t}, \quad \left( \sum_{\mathfrak{t} \in \mathbb{T}} \alpha_{\mathfrak{t}} \mathfrak{t} \right) \circ \left( \sum_{\mathfrak{h} \in \mathbb{T}} \beta_{\mathfrak{h}} \mathfrak{h} \right) = \sum_{\mathfrak{t}, \mathfrak{h} \in \mathbb{T}} (\alpha_{\mathfrak{t}} \beta_{\mathfrak{h}}) (\mathfrak{t} \circ \mathfrak{h}).$$

Furthermore, the product of  $\sigma \in \mathbb{R}$  with elements of  $\mathbb{D}$  is given by

$$\sigma \left( \sum_{\mathfrak{t} \in \mathbb{T}} \alpha_{\mathfrak{t}} \mathfrak{t} \right) = \sum_{\mathfrak{t} \in \mathbb{T}} (\sigma \alpha_{\mathfrak{t}}) \mathfrak{t}.$$

In the sequel, the products  $\circ$  are generally understood.

**Remark 3.3.** We could achieve exactly the same construction, following Chapter 2 in [9], by considering functions on the lattice  $\mathcal{L}$  as sequences and the Finite Difference operators as sequences with compact support (whence the almost everywhere requirement in Equation (4)). Then, the product  $\circ$  can be seen as a convolution (Cauchy) product between compactly supported sequences and the action of a Finite Difference operator on a function as the convolution of a finitely supported sequence with a generic sequence.

Upon introducing the generating displacements along each axis  $x_k := \mathfrak{t}_{\mathbf{e}_k}$  where  $\mathbf{e}_k$  is the  $k$ -th vector of the canonical basis, for any  $k \in [1..d]$ , we can isomorphically identify  $\mathbb{D} \cong \mathbb{R}[x_1, x_1^{-1}, \dots, x_d, x_d^{-1}]$ , the ring of multivariate Laurent polynomials. On the other hand, the real numbers  $\mathbb{R}$  can be viewed as sub-ring of  $\mathbb{D}$ , being the constant polynomials. This identification can be somehow interpreted as the historical starting point of umbral calculus [39], also known as calculus of Finite Differences [34]: allow to interchange indices in sequences (operators or functions) with exponents (in polynomials). The stream phase Equation (3) can be recast under its non-diagonal form in the space of moments [46, 15] by introducing what we call the moments-stream matrix  $\mathbf{T} := \mathbf{M} \text{diag}(\mathfrak{t}_{c_1}, \dots, \mathfrak{t}_{c_q}) \mathbf{M}^{-1} \in \mathcal{M}_q(\mathbb{D})$  and merged with the collision phase Equation (2) to obtain the scheme, for any  $t \in \mathbb{J}$  and for any  $\mathbf{x} \in \mathcal{L}$

$$\mathbf{m}(t + \Delta t, \mathbf{x}) = \mathbf{A} \mathbf{m}(t, \mathbf{x}) + \mathbf{B} \mathbf{m}^{\text{eq}}(t, \mathbf{x}), \quad (5)$$

where  $\mathbf{A} := \mathbf{T}(\mathbf{I} - \mathbf{S}) \in \mathcal{M}_q(\mathbb{D})$  and  $\mathbf{B} := \mathbf{T}\mathbf{S} \in \mathcal{M}_q(\mathbb{D})$ .

### 3.2 Corresponding Finite Difference scheme

With this new compact algebraic form of any lattice Boltzmann scheme, namely Equation (5), we are able to recall the main results proved in [3]. These results encompass the findings from [43], [10] and [18]. The version for one conserved moment can be formulated as follows.

**Proposition 3.4** (Corresponding Finite Difference scheme for  $N = 1$ , [3]). *Consider  $N = 1$ . Then the lattice Boltzmann scheme given by Equation (5) corresponds to a multi-step explicit macroscopic Finite Difference scheme on the conserved moment  $m_1$  under the form*

$$m_1(t + \Delta t, \mathbf{x}) = - \sum_{k=0}^{q-1} c_k m_1(t + (1 - q + k)\Delta t, \mathbf{x}) + \left( \sum_{k=0}^{q-1} \left( \sum_{\ell=0}^k c_{q+\ell-k} \mathbf{A}^\ell \right) \mathbf{B} \mathbf{m}^{eq}(t - k\Delta t, \mathbf{x}) \right)_1, \quad (6)$$

for all  $t \in \mathfrak{Z}$  and for all  $\mathbf{x} \in \mathcal{L}$ , where  $(c_k)_{k=0}^{k=q} \subset \mathbb{D}$  are the coefficients of  $\chi_{\mathbf{A}} := \det(X\mathbf{I} - \mathbf{A}) = \sum_{k=0}^{k=q} c_k X^k$ , the characteristic polynomial of  $\mathbf{A}$ , with  $\det(\cdot)$  indicating the determinant of a matrix.

The proof – given in [3] – relies on the fact that  $\mathbb{D}$  is a commutative ring and that therefore the Cayley-Hamilton theorem [5], stipulating that any square matrix with entries in a commutative ring annihilates its characteristic polynomial, holds.

This result is easily generalized to the case of multiple conservation laws, namely  $N > 1$ . For this, let us introduce a new notation. For any square matrix  $\mathbf{C} \in \mathcal{M}_q(\mathfrak{R})$  on a commutative ring  $\mathfrak{R}$ , consider  $\mathbf{C}_I := (\sum_{i \in I} \mathbf{e}_i \otimes \mathbf{e}_i) \mathbf{C} (\sum_{i \in I} \mathbf{e}_i \otimes \mathbf{e}_i) \in \mathcal{M}_q(\mathfrak{R})$  for any  $I \subset [1..q]$ , corresponding to the matrix where only the entries with row and column indices in  $I$  are kept and the remaining ones are set to zero. Then we have the following statement.

**Proposition 3.5** (Corresponding Finite Difference scheme for  $N \geq 1$ , [3]). *Consider  $N \geq 1$ . Then the lattice Boltzmann scheme given by Equation (5) corresponds to a family of multi-step explicit macroscopic Finite Difference schemes on the conserved moments  $m_1, \dots, m_N$ . This is, for any  $i \in [1..N]$*

$$m_i(t + \Delta t, \mathbf{x}) = - \sum_{k=0}^{q-N} c_{i,k} m_i(t + (k - q + N)\Delta t, \mathbf{x}) + \left( \sum_{k=0}^{q-N} \left( \sum_{\ell=0}^k c_{i,q+1-N+\ell-k} \mathbf{A}_i^\ell \right) \mathbf{A}_i^\diamond \mathbf{m}(t - k\Delta t, \mathbf{x}) \right)_i \quad (7)$$

$$+ \left( \sum_{k=0}^{q-N} \left( \sum_{\ell=0}^k c_{i,q+1-N+\ell-k} \mathbf{A}_i^\ell \right) \mathbf{B} \mathbf{m}^{eq}(t - k\Delta t, \mathbf{x}) \right)_i,$$

for all  $t \in \mathfrak{Z}$  and  $\mathbf{x} \in \mathcal{L}$ , where  $\mathbf{A}_i := \mathbf{A}_{\{i\} \cup [N+1..q]}$  and  $\mathbf{A}_i^\diamond := \mathbf{A} - \mathbf{A}_i$  with  $(c_{i,k})_{k=0}^{k=q+1-N} \subset \mathbb{D}$  which are the coefficients of  $\chi_{\mathbf{A}_i} := \det(X\mathbf{I} - \mathbf{A}_i) = X^{N-1} \sum_{k=0}^{k=q+1-N} c_{i,k} X^k$ , the characteristic polynomial of  $\mathbf{A}_i$ .

This result is the natural generalization of Proposition 3.4 to the case  $N > 1$ , in the sense that each sub-problem for any  $i \in [1..N]$  deals with one conserved moment (the  $i$ -th) at each time, only trying to eliminate the non-conserved moments while keeping the conserved ones other than the  $i$ -th. This is achieved by using a tailored characteristic polynomial for each conserved moment in the problem.

Further comments on Proposition 3.4 and Proposition 3.5 are postponed to the following Section.

### 3.3 A more compact form of corresponding Finite Difference scheme

Although the asymptotic analysis we shall develop in Section 5 can be carried on the formulations from Proposition 3.4 and Proposition 3.5 previously introduced in [3], we propose a different formalism based on shift operators in time. Having utilized both approaches, the advantage of this new standpoint – which shall be adopted in this paper – is to easily deal with the asymptotic analysis of the coefficients of the characteristic polynomial and of the powers of the matrix  $\mathbf{A}$  on the right hand side of Equation (6). In particular, this allows for the straightforward generalization of the procedure above second-order. Furthermore, the links with other asymptotic analysis of lattice Boltzmann schemes from the literature – which we shall develop in Section 7 – become noticeably more transparent. To this end, we introduce the following Definition.

**Definition 3.6** (Shift operator in time). Let  $f : \mathfrak{Z} \rightarrow \mathbb{R}$  be any function defined on the time lattice, then the time shift operator  $z$  acts as

$$(zf)(t) = f(t + \Delta t), \quad \forall t \in \mathfrak{Z}.$$

With this, the scheme Equation (5) can be recast under the fully-operatorial form. For any  $t \in \mathfrak{Z}$  and for any  $\mathbf{x} \in \mathcal{L}$

$$(z\mathbf{I} - \mathbf{A})\mathbf{m}(t, \mathbf{x}) = \mathbf{B} \mathbf{m}^{eq}(t, \mathbf{x}), \quad (8)$$

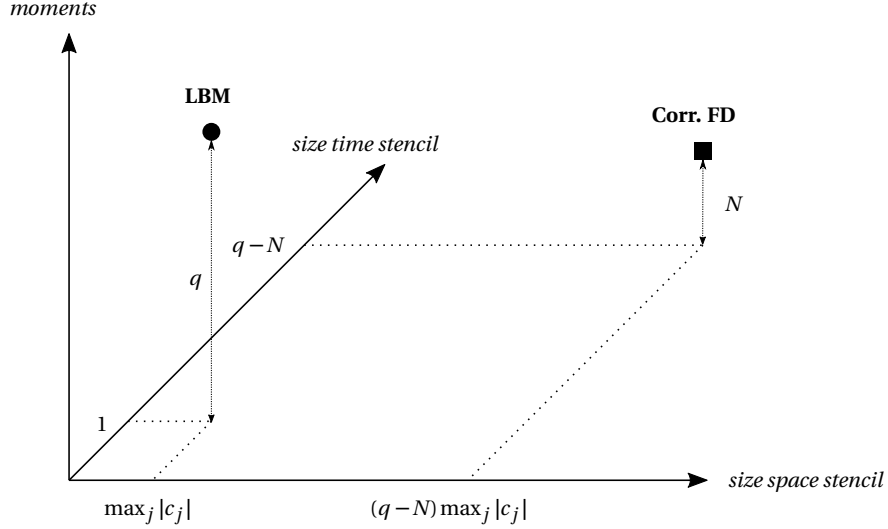


Figure 2: Comparison between lattice Boltzmann scheme (circle) and corresponding Finite Difference schemes (square) in terms of involved moments (respectively  $q$  and  $N$ ), number of time steps (respectively 1 and  $q - N$ ) and size of the maximal spatial stencil (respectively  $\max_j |c_j|$  and  $(q - N) \max_j |c_j|$ ).

which corresponds to taking the  $Z$ -transform [28] of the scheme in the variable  $z$ . Here, the inverse of the resolvent associated with  $\mathbf{A}$ , namely  $z\mathbf{I} - \mathbf{A} \in \mathcal{M}_q(\mathbb{R}[z] \otimes_{\mathbb{R}} \mathbb{D})$ , where  $\mathbb{R}[z] \otimes_{\mathbb{R}} \mathbb{D} \cong \mathbb{R}[z, x_1, x_1^{-1}, \dots, x_d, x_d^{-1}]$ , with  $\otimes_{\mathbb{R}}$  that indicates the tensor product of  $\mathbb{R}$ -algebras (see Chapter 16 in [32] or Chapter 2 in [29]), forms a commutative ring. In the sequel, we shall drop the time and the space variables when not strictly needed for the sake of readability, because the system given by Equation (8) is intrinsically time and space invariant thanks to Assumptions 3.1 and since we work on an unbounded domain, without considering the initial conditions.

**Proposition 3.4 bis** (Corresponding Finite Difference scheme for  $N = 1$ ). *Consider  $N = 1$ . Then the lattice Boltzmann scheme given by Equation (5) or Equation (8) corresponds to a multi-step explicit macroscopic Finite Difference scheme on the conserved moment  $m_1$  under the form*

$$\det(z\mathbf{I} - \mathbf{A})m_1 = (\text{adj}(z\mathbf{I} - \mathbf{A})\mathbf{B}\mathbf{m}^{eq})_1, \quad (9)$$

where  $\text{adj}(\cdot)$  indicates the adjugate matrix,<sup>2</sup> also known as classical adjoint, which is the transpose of the cofactor matrix [22].

Up to a temporal shift of the whole scheme, the corresponding multi-step explicit Finite Difference scheme by Equation (9) equals the one from Equation (6).

*Proof.* The proof can be done starting from Proposition 3.4. Alternatively using the fundamental relation between adjugate and determinant, see Chapter 0 in [22], which is a consequence of the Laplace formula, we have that for any  $\mathbf{C} \in \mathcal{M}_q(\mathfrak{R})$  where  $\mathfrak{R}$  is any commutative ring

$$\mathbf{C}\text{adj}(\mathbf{C}) = \text{adj}(\mathbf{C})\mathbf{C} = \det(\mathbf{C})\mathbf{I}. \quad (10)$$

Hence, multiplying Equation (8) by  $\text{adj}(z\mathbf{I} - \mathbf{A})$  yields  $\det(z\mathbf{I} - \mathbf{A})\mathbf{m} = \text{adj}(z\mathbf{I} - \mathbf{A})\mathbf{B}\mathbf{m}^{eq}$ . Selecting the first row gives Equation (9).  $\square$

**Remark 3.7** (From kinetic to macroscopic). We observe the following facts:

- The procedure can be reversed – when keeping all the lines in  $\det(z\mathbf{I} - \mathbf{A})\mathbf{m} = \text{adj}(z\mathbf{I} - \mathbf{A})\mathbf{B}\mathbf{m}^{eq}$  – using a multiplication by  $z\mathbf{I} - \mathbf{A}$  and then dividing by the polynomial  $\chi_{\mathbf{A}}(z) = \det(z\mathbf{I} - \mathbf{A})$ . In this way, one comes back

<sup>2</sup>It is worthwhile observing that the determinant and the adjugate matrix are defined for any square matrix with elements in a commutative ring.



to the lattice Boltzmann scheme by Equation (8). This can be done as long as one does not select and store only the first row as in Equation (9). Contrarily, if this selection is performed, the irreversible passage from the kinetic to the macroscopic formulation is accomplished. The non-conserved moments  $m_2, \dots, m_q$  are no longer defined and they cannot be recovered from Equation (9). This fact has been observed by [10]: the same macroscopic Finite Difference scheme can correspond to distinct lattice Boltzmann schemes which can have different evolution equations for the non-conserved moments  $m_2, \dots, m_q$ . This is not surprising, since for a given monic polynomial, one can find an infinite number of matrices of which is it the characteristic polynomial.

- Though – as previously emphasized – the non-conserved moments are no longer present in the macroscopic Finite Difference scheme by Equation (9), there is a residual shadow of their presence, namely the multi-step nature of the Finite Difference scheme, see Figure 2. In particular, each non-conserved moment  $m_i$  relaxing away from the equilibrium, namely with  $s_i \neq 1$ , for  $i \in [2..q]$ , adds a time step to the corresponding Finite Difference scheme solely acting on the conserved moment  $m_1$ .

**Remark 3.8** (Adjugate and characteristic polynomial). A time shift and a change of variable in Equation (6) allows to express  $\text{adj}(z\mathbf{I} - \mathbf{A})$  as a polynomial in  $z$  of degree  $q - 1$  computed from the characteristic polynomial. This relation is indeed classical and reads

$$\text{adj}(z\mathbf{I} - \mathbf{A}) = \sum_{k=0}^{q-1} \left( \sum_{\ell=0}^{q-1-k} c_{k+\ell+1} \mathbf{A}^\ell \right) z^k, \quad \text{where} \quad \det(z\mathbf{I} - \mathbf{A}) = \sum_{k=0}^q c_k z^k.$$

In the same way, we can restate Proposition 3.5 using the new formalism.

**Proposition 3.5 bis** (Corresponding Finite Difference scheme for  $N \geq 1$ ). *Consider  $N \geq 1$ . Then the lattice Boltzmann scheme given by Equation (5) or Equation (8) corresponds to a family of multi-step explicit macroscopic Finite Difference schemes on the conserved moments  $m_1, \dots, m_N$ . This is, for any  $i \in [1..N]$*

$$\det(z\mathbf{I} - \mathbf{A}_i) m_i = (\text{adj}(z\mathbf{I} - \mathbf{A}_i) \mathbf{A}_i^\diamond \mathbf{m})_i + (\text{adj}(z\mathbf{I} - \mathbf{A}_i) \mathbf{B} \mathbf{m}^{eq})_i. \quad (11)$$

*Up to a temporal shift of the whole scheme, the corresponding multi-step explicit Finite Difference scheme by Equation (11) equals the one from Equation (7).*

We could call the form of Finite Difference scheme from Proposition 3.5 and Proposition 3.5 bis “canonical” since we shall prove in Section 3.4 that it guarantees that the Finite Difference does not depend on the choice of relaxation parameters for the conserved variables, which do not play any role in the original lattice Boltzmann scheme either, as previously discussed.

**Remark 3.9** (Lack of scaling assumption). The results in Proposition 3.4, Proposition 3.5, Proposition 3.4 bis, Proposition 3.5 bis are fully discrete and do not make any assumption on the particular scaling between the time-step  $\Delta t$  and the space-step  $\Delta x$ . In particular, they can be employed both in the case of acoustic scaling, namely  $\Delta t \sim \Delta x$ , and in the case of diffusive scaling  $\Delta t \sim \Delta x^2$ .

The previous Remark signifies that the corresponding Finite Difference schemes can be utilized to assess the consistency of the underlying lattice Boltzmann scheme with respect to the macroscopic equations regardless of the particular scaling between time and space discretizations.

### 3.4 On the choice of relaxation parameters for the conserved moments

In Section 2, we have observed that the choice of relaxation parameters for the conserved moments, namely  $s_1, \dots, s_N$ , does not change the lattice Boltzmann scheme Equation (2). On the other hand, it could be argued that different choices for  $s_1, \dots, s_N$  can affect the formulations of the corresponding Finite Difference schemes resulting from Proposition 3.4 bis and Proposition 3.5 bis. We now show that, as one could hope, this is not the case for the Finite Difference schemes given by Proposition 3.5 bis.

**Proposition 3.10.** *The multi-step explicit macroscopic Finite Difference schemes given by Equation (11) in Proposition 3.5 bis do not depend on the choice of  $s_1, \dots, s_N$ , the relaxation parameters of the conserved moments.*

*Proof.* Fix the indices of the conserved moment  $i \in [1..N]$ . Let us decompose  $\mathbf{B}$ , the part of the lattice Boltzmann scheme dealing with the equilibria, as follows:  $\mathbf{B} = \mathbf{b}_i \otimes \mathbf{e}_i + \mathbf{B}|_{s_i=0}$  where  $\mathbf{b}_i = \mathbf{B}_{\cdot,i}$  is the  $i$ -th column of  $\mathbf{B}$ . The dependency of  $\mathbf{B}$  on the choice of  $s_i$  is now fully contained in  $\mathbf{b}_i$ . On the other hand  $\mathbf{B}|_{s_i=0}$  does not depend on it. The Finite Difference scheme from Proposition 3.5 bis can be therefore recast, upon rearranging and using well-known properties of the external product  $\otimes$ , as

$$(\det(\mathbf{zI} - \mathbf{A}_i) - \mathbf{e}_i^\top \text{adj}(\mathbf{zI} - \mathbf{A}_i) \mathbf{b}_i) m_i = (\text{adj}(\mathbf{zI} - \mathbf{A}_i) \mathbf{A}_i^\diamond \mathbf{m})_i + (\text{adj}(\mathbf{zI} - \mathbf{A}_i) \mathbf{B}|_{s_i=0} \mathbf{m}^{\text{eq}})_i. \quad (12)$$

The left hand side does not depend on  $s_j$  for  $j \in [1..N] \setminus \{i\}$  by construction of  $\mathbf{A}_i$  and  $\mathbf{b}_i$ . On the other hand, the right hand side does not depend on  $s_j$  for  $j \in [1..N] \setminus \{i\}$ , because  $(\mathbf{A}_i^\diamond)_{\cdot,j} + (\mathbf{B}|_{s_i=0})_{\cdot,j} = (\mathbf{A}_i^\diamond|_{s_j=0})_{\cdot,j}$ , where we have used Equation (2) and Equation (1). We are left to discuss the possible dependency of Equation (12) on  $s_i$ . For the left hand side, we need the following result concerning the determinant of matrices under rank-one updates, whose proof is analogous to that in [12].

**Lemma 3.11** (Matrix determinant). *Let  $\mathfrak{R}$  be a commutative ring,  $\mathbf{C} \in \mathcal{M}_q(\mathfrak{R})$  and  $\mathbf{u}, \mathbf{v} \in \mathfrak{R}^q$ , then  $\det(\mathbf{C} + \mathbf{u} \otimes \mathbf{v}) = \det(\mathbf{C}) + \mathbf{v}^\top \text{adj}(\mathbf{C}) \mathbf{u}$ .*

By this Lemma, we deduce that Equation (12) now reads

$$\det(\mathbf{zI} - (\mathbf{A}_i + \mathbf{b}_i \otimes \mathbf{e}_i)) m_i = (\text{adj}(\mathbf{zI} - \mathbf{A}_i) \mathbf{A}_i^\diamond \mathbf{m})_i + (\text{adj}(\mathbf{zI} - \mathbf{A}_i) \mathbf{B}|_{s_i=0} \mathbf{m}^{\text{eq}})_i. \quad (13)$$

Observe that  $\mathbf{A}_i + \mathbf{b}_i \otimes \mathbf{e}_i = \mathbf{A}_i|_{s_i=0}$ , thus the left hand side of Equation (13) does not depend on  $s_i$ . The right hand side of Equation (13) is independent of  $s_i$  because  $\mathbf{A}_i^\diamond$  does not depend on it and since the  $i$ -th row of  $\text{adj}(\mathbf{zI} - \mathbf{A}_i)$  – the transpose of the cofactor matrix of  $\mathbf{zI} - \mathbf{A}_i$  – cannot depend on  $s_i$ , because only the  $i$ -th column of  $\mathbf{zI} - \mathbf{A}_i$  depends on  $s_i$ . This concludes the proof.  $\square$

We have thus shown that the Finite Difference schemes from Proposition 3.5 bis do not depend on the choice of relaxation parameters for the conserved moments and so that we are allowed to take them equal to zero or any other value of specific convenience without loss of generality. In particular, the choice of taking  $s_i = 0$  for  $i \in [1..N]$  offers interesting simplifications in the computations to come in Section 5, in a way that shall be clearer by looking at the details. Moreover, this choice has the advantage of showing which moments are conserved at a glance.

## 4 Main results

Everything is in place to start the standard consistency analysis of Finite Difference schemes [45, 41, 1]. We start from the assumptions allowing us to identify each term once developing in formal power series of  $\Delta x$ , *i.e.* performing Taylor expansions.

**Assumptions 4.1** (Scaling assumptions). Assume that:

1. We utilize the acoustic scaling<sup>3</sup>  $\lambda > 0$  is a fixed real number as  $\Delta x \rightarrow 0$ , like in [13, 15, 46]. Thus, the only discretization parameter we shall consider in the sequel is  $\Delta x$ .
2. The change of basis  $\mathbf{M}$ , the relaxation matrix  $\mathbf{S}$  and the moments at equilibrium  $\mathbf{m}^{\text{eq}}$  are fixed as  $\Delta x \rightarrow 0$ .

We also introduce the spaces of differential operators which shall be obtained by taking the limit  $\Delta x \rightarrow 0$  as well as other tightly associated concepts.

**Definition 4.2** (Time-space differential operators). We define.

- The commutative ring of time-space differential operators  $\mathcal{D}$  by

$$\mathcal{D} := \mathbb{R}[\partial_t] \otimes_{\mathbb{R}} \mathbb{R}[\partial_{x_1}, \dots, \partial_{x_d}] \cong \mathbb{R}[\partial_t, \partial_{x_1}, \dots, \partial_{x_d}].$$

- Under the acoustic Assumptions 4.1, we consider the commutative ring of formal power series [36, 35]

$$\mathcal{S} := \mathcal{D}[[\Delta x]].$$

<sup>3</sup>Frequently,  $\lambda = 1$  is considered in the literature.

- For any  $\delta = \sum_{r=0}^{+\infty} \Delta x^r \delta^{(r)} \in \mathcal{S}$ , we indicate  $\delta = O(\Delta x^{r_o})$  for some  $r_o \in \mathbb{N}$  if  $\delta^{(r)} = 0$  for  $r \in [0..r_o - 1]$  and  $\delta^{(r_o)} \neq 0$ . The integer  $r_o$  is called “order” of the formal power series  $\delta$ , see Chapter 1 in [39].
- Finally, let  $d \in \mathbb{R}[z] \otimes_{\mathbb{R}} \mathbb{D}$  and  $\delta \in \mathcal{S}$ , then we indicate  $d \asymp \delta$ , called “asymptotic equivalence” of  $d$  and  $\delta$ , if for any smooth function of the time and space variables  $f : \mathbb{R} \times \mathbb{R}^d \rightarrow \mathbb{R}$ , we have

$$(df)(t, \mathbf{x}) = \sum_{r=0}^{+\infty} \Delta x^r (\delta^{(r)} f)(t, \mathbf{x}), \quad \forall (t, \mathbf{x}) \in \mathbb{R} \times \mathbb{R}^d, \quad \text{as } \Delta x \rightarrow 0.$$

The previous  $O(\cdot)$  notation and the notion of asymptotic equivalence are effortlessly extended to vectors and matrices in an entry-wise fashion. It shall be common and harmless not to distinguish between  $\mathcal{M}_q(\mathcal{S})$  and  $(\mathcal{M}_q(\mathcal{D}))[[\Delta x]]$ .

The momentum-velocity operator matrix  $\mathcal{G} \in \mathcal{M}_q(\mathcal{D})$ , introduced by [15] with slightly different notations, is defined as follows. It is indeed closely linked to the moment-stream matrix  $\mathbf{T} \in \mathcal{M}_q(\mathbb{D})$  that we have previously introduced.

**Definition 4.3** (Momentum-velocity operator matrix). The momentum-velocity operator matrix made up of first-order differential operators in space is given by

$$\mathcal{G} := \mathbf{M} \left( \sum_{|\mathbf{v}|=1} \text{diag}(\mathbf{e}_1^{\mathbf{v}}, \dots, \mathbf{e}_q^{\mathbf{v}}) \partial^{\mathbf{v}} \right) \mathbf{M}^{-1} \in \mathcal{M}_q(\mathcal{D}),$$

where the multi-index notation is employed.

This momentum-velocity operator matrix can be partitioned in four blocks with different meanings according to the different nature (conserved or not) of the corresponding moments, as for Equation (8) in [15]. We are now ready to state and then prove the main results of the present contribution. The Taylor expansions are applied to the solution of the corresponding Finite Difference schemes given by Proposition 3.5 or Proposition 3.5 bis, where non-conserved moments have been removed yielding purely macroscopic discrete equations.

**Theorem 4.4** (First order expansion). *Under Assumptions 4.1 and in the limit  $\Delta x \rightarrow 0$ , the conserved moments  $m_1, \dots, m_N$ , solution of the corresponding macroscopic Finite Difference schemes given by Proposition 3.5 or Proposition 3.5 bis, asymptotically satisfy the following system of macroscopic PDEs*

$$\partial_t m_i + \gamma_{1,i} = O(\Delta x), \quad \text{with} \quad \gamma_{1,i} = \gamma_{1,i}(m_1, \dots, m_N) := \sum_{j=1}^N \mathcal{G}_{ij} m_j + \sum_{j=N+1}^q \mathcal{G}_{ij} m_j^{eq}(m_1, \dots, m_N),$$

for  $i \in [1..N]$  and for  $(t, \mathbf{x}) \in \mathbb{R}^+ \times \mathbb{R}$ , where  $m_1, \dots, m_N : \mathbb{R}^+ \times \mathbb{R} \rightarrow \mathbb{R}$  are smooth functions.

The first term in  $\gamma_{1,i}$  represents the derivatives of fluxes of the conserved variables, which are necessarily linear, while the second one represents the derivatives of the fluxes given by the equilibria of the non-conserved moments, which can be non-linear. These limit equations can be refined with a second-order expansion yielding diffusion terms, where the so-called Hénon’s parameters [20] of type  $1/s_j - 1/2$  appear.

**Theorem 4.5** (Second order expansion). *Under Assumptions 4.1 and in the limit  $\Delta x \rightarrow 0$ , the conserved moments  $m_1, \dots, m_N$ , solution of the corresponding macroscopic Finite Difference schemes given by Proposition 3.5 or Proposition 3.5 bis, asymptotically satisfy the following system of macroscopic PDEs*

$$\partial_t m_i + \gamma_{1,i} + \frac{\Delta x}{\lambda} \sum_{j=N+1}^q \left( \frac{1}{s_j} - \frac{1}{2} \right) \mathcal{G}_{ij} \left( \sum_{\ell=1}^N \frac{dm_j^{eq}}{dm_\ell} \gamma_{1,\ell} - \sum_{\ell=1}^N \mathcal{G}_{j\ell} m_\ell - \sum_{\ell=N+1}^q \mathcal{G}_{j\ell} m_\ell^{eq} \right) = O(\Delta x^2),$$

for  $i \in [1..N]$  and for  $(t, \mathbf{x}) \in \mathbb{R}^+ \times \mathbb{R}$ , where  $m_1, \dots, m_N : \mathbb{R}^+ \times \mathbb{R} \rightarrow \mathbb{R}$  are smooth functions.

One can notice that the diffusive terms appear to be proportional to  $\Delta x$ . This is not surprising, since the only way of having a stable explicit Finite Difference scheme to simulate the heat equation under the acoustic scaling is to consider a diffusion coefficient proportional to  $\Delta x$ , in order to constrain the speed of propagation of information to remain finite in the limit  $\Delta x \rightarrow 0$ , see for instance Theorem 6.3.1 in [41].

In this contribution, we have deliberately neglected the behavior of the schemes close to the initial time  $t = 0$ . It is dictated by the choice of initial datum for the non conserved moments, which is not unique for lattice Boltzmann schemes since  $q > N$  but one only knows the  $N$  conserved moments at  $t = 0$ , being the initial datum of the macroscopic PDEs to be solved. The interested reader can consult [44, 38] for more information on this topic.

Let us sketch the main ideas of the proofs of Theorem 4.4 and Theorem 4.5:

- The result of Proposition 3.5 has allowed to eliminate the non-conserved moments from the discrete scheme, thus has completed the step represented by a vertical arrow in Figure 1. Contrarily to the existing approaches, we do not need (and we cannot, see Remark 3.7) to estimate the Taylor expansions of the non-conserved moments.
- We benefit from the clever formulation from Proposition 3.5 bis instead of that of Proposition 3.5. Indeed, considering  $\zeta \mathbf{I} - \mathcal{A}_i \asymp \mathbf{z} \mathbf{I} - \mathbf{A}_i$ , we are allowed to write, for every  $i \in [1..N]$

$$\det(\zeta \mathbf{I} - \mathcal{A}_i) m_i = (\text{adj}(\zeta \mathbf{I} - \mathcal{A}_i) \mathcal{A}_i^\circ m)_i + (\text{adj}(\zeta \mathbf{I} - \mathcal{A}_i) \mathcal{B} m^{\text{eq}})_i,$$

obtained by replacing matrices with entries in the ring  $\mathbb{R}[z] \otimes_{\mathbb{R}} \mathbb{D}$  of discrete operators by their asymptotic equivalents in the ring  $\mathcal{S}$ . Here, for example,  $\det(\zeta \mathbf{I} - \mathcal{A}_i) \in \mathcal{S}$ , and the expression perfectly makes sense because the determinant and the adjugate are well-defined polynomial functions of any square matrix on a commutative ring, like  $\mathcal{S}$ . Since the determinant and the adjugate are non-linear functions and thus mix different orders in the expansion  $\zeta \mathbf{I} - \mathcal{A}_i$ , if we want to recover a closed-form result at a given order of accuracy, we are compelled to utilize the Taylor expansions of the determinant and the adjugate. However, these expansions are well-known and can be easily computed at any order of accuracy.

Notice that, on the other hand, if we want to exploit the formulation of [3] stated in Proposition 3.5, we should characterize the asymptotic equivalents of any coefficient of the characteristic polynomial of  $\mathbf{A}_i$  and then combine them with the asymptotic equivalents of the time shifts  $z$  alone and the terms on the right hand side of Equation (7). Though this is actually feasible and we firstly did it, the computations are extremely involved<sup>4</sup> and very hard to generalize above second-order.

This justifies the use of the formulation from Proposition 3.5 bis to achieve the step denoted by an horizontal arrow in Figure 1.

## 5 Detailed proofs

The vast majority of rest of this work is devoted to the detailed proof of Theorem 4.4 and Theorem 4.5 for the scalar case  $N = 1$ . This choice has been adopted to keep the presentation and the involved notations as simple as possible. The idea behind the generalization to  $N > 1$  is eventually given in Section 6 and is straightforward except for the more involved notations.

Let us start by finding, for each shift operator from Definition 3.2, its asymptotically equivalent formal power series in  $\Delta x$ , see for instance [46, 15]. This is formalized by the following Lemma.

**Lemma 5.1** (Series expansion of a shift operator in space). *Let  $\mathbf{z} \in \mathbb{Z}^d$ , then the associated shift operator in space  $t_{\mathbf{z}} \in \mathbb{R}[z] \otimes_{\mathbb{R}} \mathbb{D}$  is asymptotically equivalent, in the limit of  $\Delta x \rightarrow 0$ , to the formal power series of differential operators of the form*

$$t_{\mathbf{z}} \asymp \sum_{|\mathbf{v}| \geq 0} \frac{(-\Delta x)^{|\mathbf{v}|} \mathbf{z}^{\mathbf{v}}}{\mathbf{v}!} \partial^{\mathbf{v}} \in \mathcal{S}.$$

*Proof.* Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be a smooth function of the spatial variable. Then performing a Taylor expansion for  $\Delta x \rightarrow 0$  yields

$$(t_{\mathbf{z}} f)(\mathbf{x}) = f(\mathbf{x} - \mathbf{z} \Delta x) = \sum_{|\mathbf{v}| \geq 0} \frac{(-\Delta x)^{|\mathbf{v}|} \mathbf{z}^{\mathbf{v}}}{\mathbf{v}!} \partial^{\mathbf{v}} f(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d.$$

□

<sup>4</sup>Probably, a deeper mastery of the elementary symmetric polynomials, the Newton's identities, the Bell polynomials and the Feddeev-Leverrier algorithm could simplify many reasonings.

The extension of Lemma 5.1 to any Finite Difference operator in  $\mathbb{D}$  according to Definition 3.2 is done by linearity. With this in mind, recalling the definition of  $\mathbf{T} \in \mathcal{M}_q(\mathbb{D})$ , the moments-stream matrix and using Assumptions 4.1, we have that

$$\mathbf{T} := \mathbf{M} \text{diag}(t_{c_1}, \dots, t_{c_q}) \mathbf{M}^{-1} \asymp \mathbf{M} \left( \sum_{|\mathbf{v}| \geq 0} \frac{(-\Delta x)^{|\mathbf{v}|}}{\mathbf{v}!} \text{diag}(\mathbf{c}_1^{\mathbf{v}}, \dots, \mathbf{c}_q^{\mathbf{v}}) \partial^{\mathbf{v}} \right) \mathbf{M}^{-1} =: \mathcal{T} \in \mathcal{M}_q(\mathcal{S}). \quad (14)$$

Accordingly, we introduce  $\mathcal{A} := \mathcal{T}(\mathbf{I} - \mathbf{S}) \in \mathcal{M}_q(\mathcal{S})$  and  $\mathcal{B} := \mathcal{T}\mathbf{S} \in \mathcal{M}_q(\mathcal{S})$  such that  $\mathcal{A} \asymp \mathbf{A}$  and  $\mathcal{B} \asymp \mathbf{B}$ . The tight bond between the momentum-velocity operator matrix  $\mathcal{G} \in \mathcal{M}_q(\mathcal{D})$  from [15] and our moments-stream matrix  $\mathbf{T} \in \mathcal{M}_q(\mathbb{D})$  and its asymptotic equivalent matrix  $\mathcal{T} \in \mathcal{M}_q(\mathcal{S})$  is given by the following Lemma.

**Lemma 5.2** (Link between  $\mathcal{G}$  and  $\mathcal{T}^{(r)}$ ). *For any order  $r \in \mathbb{N}$ , the matrix  $\mathcal{T}^{(r)} \in \mathcal{M}_q(\mathcal{D})$  is linked to  $\mathcal{G} \in \mathcal{M}_q(\mathcal{D})$  by*

$$\mathcal{T}^{(r)} = \frac{(-1)^r}{\lambda^r r!} \mathcal{G}^r.$$

Moreover, using the Assumptions 4.1, we also have

$$\mathcal{A}^{(r)} = \frac{(-1)^r}{\lambda^r r!} \mathcal{G}^r (\mathbf{I} - \mathbf{S}), \quad \mathcal{B}^{(r)} = \frac{(-1)^r}{\lambda^r r!} \mathcal{G}^r \mathbf{S}.$$

*Proof.* By Equation (21) in [15], we have that  $\mathbf{T} \asymp \mathcal{T} = \exp\left(-\frac{\Delta x}{\lambda} \mathcal{G}\right)$ . A Taylor expansion of the exponential function yields the result. Using Assumptions 4.1, one obtains that  $(\mathbf{I} - \mathbf{S})$  and  $\mathbf{S}$  do not perturb the orders of the expansion.  $\square$

As far as the time variable is concerned, we can complete by the development of the time shift operator  $\mathbf{z}$  in order to provide the overall expansion of the inverse of the resolvent  $\mathbf{z}\mathbf{I} - \mathbf{A} \in \mathcal{M}_q(\mathbb{R}[\mathbf{z}] \otimes_{\mathbb{R}} \mathbb{D})$ .

**Lemma 5.3** (Expansion of the inverse of the resolvent). *Under Assumptions 4.1 and in the limit of  $\Delta x \rightarrow 0$ , the inverse of the resolvent  $\mathbf{z}\mathbf{I} - \mathbf{A} \in \mathcal{M}_q(\mathbb{R}[\mathbf{z}] \otimes_{\mathbb{R}} \mathbb{D})$  is asymptotically equivalent to  $\zeta\mathbf{I} - \mathcal{A} \in \mathcal{M}_q(\mathcal{S})$ , where*

$$\zeta\mathbf{I} - \mathcal{A} = \sum_{r=0}^{+\infty} \frac{\Delta x^r}{\lambda^r r!} (\partial_t^r \mathbf{I} - (-1)^r \mathcal{G}^r (\mathbf{I} - \mathbf{S})) = \mathbf{S} + \frac{\Delta x}{\lambda} (\partial_t \mathbf{I} + \mathcal{G}(\mathbf{I} - \mathbf{S})) + \frac{\Delta x^2}{2\lambda^2} (\partial_{tt} \mathbf{I} - \mathcal{G}^2 (\mathbf{I} - \mathbf{S})) + O(\Delta x^3). \quad (15)$$

*Proof.* The standard Taylor expansion of  $\mathbf{z}$ , using the assumption on the acoustic scaling, gives the claim.  $\square$

The consistency analysis of the Finite Difference schemes from Proposition 3.5 bis could be carried on infinite formal power series of differential operators  $\mathcal{S}$  on the formulation

$$\det(\zeta\mathbf{I} - \mathcal{A}_i) m_i = (\text{adj}(\zeta\mathbf{I} - \mathcal{A}_i) \mathcal{A}_i^{\circ} \mathbf{m})_i + (\text{adj}(\zeta\mathbf{I} - \mathcal{A}_i) \mathcal{B} \mathbf{m}^{\text{eq}})_i, \quad (16)$$

for each  $i \in [1..N]$ , because the determinant and the adjugate perfectly make sense for any square matrix on a commutative ring, like  $\mathcal{S}$ . However, in order to prove Theorem 4.4 and Theorem 4.5, where formal power series are truncated at a certain order, we shall need Equation (15) from Lemma 5.3 as well as the Taylor expansions of the determinant and the adjugate matrix around a given matrix. Indeed, these are non-linear functions and thus mix different orders in the expansions  $\zeta\mathbf{I} - \mathcal{A} \in \mathcal{M}_q(\mathcal{S})$ . Since the product of the relaxation parameters for the non-conserved moments is a quantity which shall frequently appear in the computations to come, we fix a special notation for it, namely setting  $\Pi := \prod_{i=2}^{i=q} s_i \neq 0$ .

## 5.1 Determinant

We start by studying the expansion of the determinant up to second-order in the perturbation. For this, we need to characterize its derivatives. The expansion can be carried at higher order by employing the very same strategy.

**Lemma 5.4** (Derivatives and expansion of the determinant function). *Let  $\mathbf{C} \in GL_q(\mathfrak{R})$  and  $\mathbf{D}, \mathbf{E} \in \mathcal{M}_q(\mathfrak{R})$ , where  $\mathfrak{R}$  is a commutative ring. Then the determinant function*

$$\begin{aligned} \det: \mathcal{M}_q(\mathfrak{R}) &\rightarrow \mathfrak{R} \\ \mathbf{C} &\mapsto \det(\mathbf{C}), \end{aligned}$$

has the following derivatives.

$$D_C(\det(\mathbf{C}))(\mathbf{D}) = \det(\mathbf{C}) \operatorname{tr}(\mathbf{C}^{-1} \mathbf{D}), \quad (17)$$

$$D_{CC}(\det(\mathbf{C}))(\mathbf{D})(\mathbf{E}) = \det(\mathbf{C}) \left( \operatorname{tr}(\mathbf{C}^{-1} \mathbf{E}) \operatorname{tr}(\mathbf{C}^{-1} \mathbf{D}) - \operatorname{tr}(\mathbf{C}^{-1} \mathbf{E} \mathbf{C}^{-1} \mathbf{D}) \right), \quad (18)$$

where  $\operatorname{tr}(\cdot)$  indicates the trace, i.e. the sum of the diagonal entries. Equation (17) is known as Jacobi formula. Moreover, the second-order Taylor expansion of the determinant function reads

$$\det(\mathbf{C} + \mathbf{D}) = \det(\mathbf{C}) + D_C(\det(\mathbf{C}))(\mathbf{D}) + \frac{1}{2} D_{CC}(\det(\mathbf{C}))(\mathbf{D})(\mathbf{D}) + O(\|\mathbf{D}\|^3),$$

where the derivatives are given by Equation (17) and Equation (18).

*Proof.* The Jacobi formula Equation (17) is a standard result, see Chapter 0 in [22] or Chapter 5 in [48]. Let us prove Equation (18).

$$\begin{aligned} D_{CC}(\det(\mathbf{C}))(\mathbf{D})(\mathbf{E}) &:= D_C(D_C(\det(\mathbf{C}))(\mathbf{D}))(\mathbf{E}) = D_C(\det(\mathbf{C}) \operatorname{tr}(\mathbf{C}^{-1} \mathbf{D}))(\mathbf{E}), \\ &= D_C(\det(\mathbf{C}))(\mathbf{E}) \operatorname{tr}(\mathbf{C}^{-1} \mathbf{D}) + \det(\mathbf{C}) D_C(\operatorname{tr}(\mathbf{C}^{-1} \mathbf{D}))(\mathbf{E}), \\ &= \det(\mathbf{C}) \operatorname{tr}(\mathbf{C}^{-1} \mathbf{E}) \operatorname{tr}(\mathbf{C}^{-1} \mathbf{D}) + \det(\mathbf{C}) \operatorname{tr}(D_C(\mathbf{C}^{-1} \mathbf{D}))(\mathbf{E}), \\ &= \det(\mathbf{C}) \operatorname{tr}(\mathbf{C}^{-1} \mathbf{E}) \operatorname{tr}(\mathbf{C}^{-1} \mathbf{D}) - \det(\mathbf{C}) \operatorname{tr}(\mathbf{C}^{-1} \mathbf{E} \mathbf{C}^{-1} \mathbf{D}), \end{aligned}$$

where we have used, in this order, the product rule for derivatives, the Jacobi formula Equation (17), the linearity of the trace and the fact that  $D_C(\mathbf{C}^{-1})(\mathbf{D}) = -\mathbf{C}^{-1} \mathbf{D} \mathbf{C}^{-1}$ , see Chapter 5 in [48].  $\square$

**Remark 5.5** (On the invertibility assumption). There exists a form of the Jacobi formula Equation (17) for general  $\mathbf{C} \in \mathcal{M}_q(\mathfrak{R})$  without assuming invertibility, under the form  $D_C(\det(\mathbf{C}))(\mathbf{D}) = \operatorname{tr}(\operatorname{adj}(\mathbf{C}) \mathbf{D})$ , which is equivalent to Equation (17) since Equation (10) holds. Nevertheless, we decided to state Lemma 5.4 using the invertibility assumption. This is done, as we shall see, without loss of generality by taking advantage of some invertible approximation of real matrices and allows to easily find the formulae for higher order derivatives and expansions *via* basic differential calculus as illustrated in the previous proof.

In the sequel, we shall take  $\mathfrak{R} = \mathcal{S}$  and  $\mathbf{C} = \mathbf{S} \in \operatorname{GL}_q(\mathbb{R}) \subset \operatorname{GL}_q(\mathcal{S})$  and  $\mathbf{D} = O(\Delta x) \in \mathcal{M}_q(\mathcal{S})$ . To simplify the computations and relying on the findings of Section 3.4, we can consider  $\mathbf{S}$  singular by having  $s_1 = 0$ . To avoid the difficulties linked with singular matrices, in the spirit of Remark 5.5, we take advantage of the fact that the derivatives of the determinant (and the determinant itself) around  $\mathbf{C}$  are smooth (indeed, polynomial) functions of  $\mathbf{C}$ . Thus, we introduce the non-singular approximation  $\mathbf{S}$  where  $s_1 \neq 0$ , which is such that  $\mathbf{S} \rightarrow \mathbf{S}|_{s_1=0}$  as  $s_1 \rightarrow 0$  for any matricial topology.

We are now ready to use the expansion given by Lemma 5.3 into the terms stemming from Lemma 5.4 to find the leading order terms of the left hand side of Equation (9), namely of  $\det(\zeta \mathbf{I} - \mathcal{A}) \in \mathcal{S}$ . This is nothing but computing the Taylor series of composite functions (see the Faà di Bruno's formulae [24]) or the composition of formal series

$$\begin{aligned} \det(\zeta \mathbf{I} - \mathcal{A}) &= \det(\mathbf{S}) + \Delta x D_S(\det(\mathbf{S}))((\zeta \mathbf{I} - \mathcal{A})^{(1)}) \\ &\quad + \Delta x^2 (D_S(\det(\mathbf{S}))((\zeta \mathbf{I} - \mathcal{A})^{(2)}) + D_{SS}(\det(\mathbf{S}))((\zeta \mathbf{I} - \mathcal{A})^{(1)})(\zeta \mathbf{I} - \mathcal{A})^{(1)}) + O(\Delta x^3). \end{aligned}$$

- One clearly has  $\det(\mathbf{S}) = s_1 \Pi$ , because the matrix  $\mathbf{S}$  is diagonal. Thus, the Taylor expansion of  $\det(\zeta \mathbf{I} - \mathcal{A})$  does not contain zero-order terms if  $s_1 = 0$ .
- Let  $\mathbf{C} = \mathbf{S} \in \operatorname{GL}_q(\mathbb{R}) \subset \operatorname{GL}_q(\mathcal{S})$  and  $\mathbf{D} = \frac{\Delta x}{\lambda} (\partial_t \mathbf{I} + \mathcal{G}(\mathbf{I} - \mathbf{S})) + \frac{\Delta x^2}{2\lambda^2} (\partial_{tt} \mathbf{I} - \mathcal{G}^2(\mathbf{I} - \mathbf{S})) + O(\Delta x^3) \in \mathcal{M}_q(\mathcal{S})$  from Lemma 5.3. Using Equation (17) from Lemma 5.4 and performing elementary computations, we have

$$\begin{aligned} D_C(\det(\mathbf{C}))(\mathbf{D}) &= s_1 \Pi \left( \frac{\Delta x}{\lambda} \left( \frac{1}{s_1} (\partial_t + \mathcal{G}_{11}) + \sum_{i=2}^q \frac{1}{s_i} (\partial_t + (1 - s_i) \mathcal{G}_{ii}) \right) \right. \\ &\quad \left. + \frac{\Delta x^2}{2\lambda^2} \left( \frac{1}{s_1} \left( \partial_{tt} - \mathcal{G}_{11} \mathcal{G}_{11} - \sum_{\ell=2}^q \mathcal{G}_{1\ell} \mathcal{G}_{\ell 1} \right) + \sum_{i=2}^q \frac{1}{s_i} \left( \partial_{tt} - (1 - s_i) \sum_{\ell=1}^q \mathcal{G}_{i\ell} \mathcal{G}_{\ell i} \right) \right) \right) + O(\Delta x^3). \end{aligned} \quad (19)$$

We keep this expression without taking the limit in  $s_1$  for future use. On the other hand, taking the limit for  $s_1 \rightarrow 0$  yields the derivative around the singular matrix  $\mathbf{S}|_{s_1=0}$  instead of  $\mathbf{S} \in \text{GL}_q(\mathbb{R})$  for  $s_1 \neq 0$ .

$$\lim_{s_1 \rightarrow 0} D_{\mathbf{C}}(\det(\mathbf{C}))(\mathbf{D}) = \Pi \frac{\Delta x}{\lambda} \left( \partial_t + \mathcal{G}_{11} + \frac{\Delta x}{2\lambda} \left( \partial_{tt} - \mathcal{G}_{11} \mathcal{G}_{11} - \sum_{\ell=2}^q \mathcal{G}_{1\ell} \mathcal{G}_{\ell 1} \right) \right) + O(\Delta x^3). \quad (20)$$

This gives all the first-order term and part of the second-order term in the series  $\det(\zeta \mathbf{I} - \mathcal{A})$ .

- Let  $\mathbf{C} = \mathbf{S} \in \text{GL}_q(\mathbb{R}) \subset \text{GL}_q(\mathcal{S})$  and  $\mathbf{D} = \frac{\Delta x}{\lambda} (\partial_t \mathbf{I} + \mathcal{G}(\mathbf{I} - \mathbf{S})) + O(\Delta x^2) \in \mathcal{M}_q(\mathcal{S})$  from Lemma 5.3. Using Equation (18) from Lemma 5.4, we have, after some algebra

$$\begin{aligned} D_{\mathbf{C}\mathbf{C}}(\det(\mathbf{C}))(\mathbf{D})(\mathbf{D}) &= s_1 \Pi \frac{\Delta x^2}{\lambda^2} \left( \frac{2}{s_1} (\partial_t + \mathcal{G}_{11}) \sum_{i=2}^q \frac{1}{s_i} (\partial_t + (1 - s_i) \mathcal{G}_{ii}) + \left( \sum_{i=2}^q \frac{1}{s_i} (\partial_t + (1 - s_i) \mathcal{G}_{ii}) \right)^2 \right. \\ &\quad \left. - \frac{2}{s_1} \sum_{\ell=2}^q \left( \frac{1}{s_\ell} - 1 \right) \mathcal{G}_{1\ell} \mathcal{G}_{\ell 1} - \sum_{i=2}^q \frac{1}{s_i^2} (\partial_t + (1 - s_i) \mathcal{G}_{ii})^2 \right. \\ &\quad \left. - \sum_{i=2}^q \sum_{\substack{\ell=2 \\ \ell \neq i}}^q \left( \frac{1}{s_i} - 1 \right) \left( \frac{1}{s_\ell} - 1 \right) \mathcal{G}_{i\ell} \mathcal{G}_{\ell i} \right) + O(\Delta x^3). \end{aligned} \quad (21)$$

Once more, we take the limit for  $s_1 \rightarrow 0$  in order to find the desired result on the remaining second-order terms in the development  $\det(\zeta \mathbf{I} - \mathcal{A})$

$$\begin{aligned} \lim_{s_1 \rightarrow 0} D_{\mathbf{C}\mathbf{C}}(\det(\mathbf{C}))(\mathbf{D})(\mathbf{D}) &= 2\Pi \frac{\Delta x^2}{\lambda^2} \left( \left( \sum_{\ell=2}^q \frac{1}{s_\ell} \right) \partial_{tt} + \left( \sum_{\ell=2}^q \frac{1}{s_\ell} \right) \mathcal{G}_{11} \partial_t + \left( \sum_{i=2}^q \left( \frac{1}{s_i} - 1 \right) \mathcal{G}_{ii} \right) \partial_t + \mathcal{G}_{11} \left( \sum_{i=2}^q \left( \frac{1}{s_i} - 1 \right) \mathcal{G}_{ii} \right) \right. \\ &\quad \left. - \sum_{\ell=2}^q \left( \frac{1}{s_\ell} - 1 \right) \mathcal{G}_{1\ell} \mathcal{G}_{\ell 1} \right) + O(\Delta x^3). \end{aligned} \quad (22)$$

Putting Equation (20) and Equation (22) together in Lemma 5.4, with expansion around  $\mathbf{S}$ , allows to write  $\det(\zeta \mathbf{I} - \mathcal{A})$  up to third order. This is

$$\begin{aligned} \lim_{s_1 \rightarrow 0} \det(\zeta \mathbf{I} - \mathcal{A}) &= \Pi \frac{\Delta x}{\lambda} \left( \partial_t + \mathcal{G}_{11} + \frac{\Delta x}{\lambda} \left( \left( \frac{1}{2} + \sum_{\ell=2}^q \frac{1}{s_\ell} \right) \partial_{tt} + \left( \sum_{\ell=2}^q \frac{1}{s_\ell} \right) \mathcal{G}_{11} \partial_t + \left( \sum_{i=2}^q \left( \frac{1}{s_i} - 1 \right) \mathcal{G}_{ii} \right) \partial_t \right. \right. \\ &\quad \left. \left. - \frac{1}{2} \mathcal{G}_{11} \mathcal{G}_{11} - \sum_{\ell=2}^q \left( \frac{1}{s_\ell} - \frac{1}{2} \right) \mathcal{G}_{1\ell} \mathcal{G}_{\ell 1} + \mathcal{G}_{11} \left( \sum_{i=2}^q \left( \frac{1}{s_i} - 1 \right) \mathcal{G}_{ii} \right) \right) \right) + O(\Delta x^3). \end{aligned} \quad (23)$$

## 5.2 Adjugate

We now switch to the formal power series of the adjugate function of the inverse of the resolvent, in order to deal with the right hand side of the corresponding Finite Difference scheme given by Equation (9). Let us start by characterizing its derivatives.

**Lemma 5.6** (Derivatives and expansion of the adjugate function). *Let  $\mathbf{C} \in \text{GL}_q(\mathfrak{R})$  and  $\mathbf{D}, \mathbf{E} \in \mathcal{M}_q(\mathfrak{R})$ , where  $\mathfrak{R}$  is a commutative ring. Then the adjugate function*

$$\begin{aligned} \text{adj}: \mathcal{M}_q(\mathfrak{R}) &\rightarrow \mathcal{M}_q(\mathfrak{R}) \\ \mathbf{C} &\mapsto \text{adj}(\mathbf{C}), \end{aligned}$$

has the following derivatives.

$$D_{\mathbf{C}}(\text{adj}(\mathbf{C}))(\mathbf{D}) = \det(\mathbf{C}) \left( \text{tr}(\mathbf{C}^{-1} \mathbf{D}) \mathbf{I} - \mathbf{C}^{-1} \mathbf{D} \right) \mathbf{C}^{-1}, \quad (24)$$

$$\begin{aligned} D_{\mathbf{C}\mathbf{C}}(\text{adj}(\mathbf{C}))(\mathbf{D})(\mathbf{E}) &= \det(\mathbf{C}) \left( \left( \text{tr}(\mathbf{C}^{-1} \mathbf{E}) \text{tr}(\mathbf{C}^{-1} \mathbf{D}) - \text{tr}(\mathbf{C}^{-1} \mathbf{E} \mathbf{C}^{-1} \mathbf{D}) \right) \mathbf{C}^{-1} \right. \\ &\quad \left. + \mathbf{C}^{-1} \left( \mathbf{E} \mathbf{C}^{-1} \mathbf{D} + \mathbf{D} \mathbf{C}^{-1} \mathbf{E} - \text{tr}(\mathbf{C}^{-1} \mathbf{E}) \mathbf{D} - \text{tr}(\mathbf{C}^{-1} \mathbf{D}) \mathbf{E} \right) \mathbf{C}^{-1} \right). \end{aligned} \quad (25)$$

Moreover, the second-order Taylor expansion of the adjugate function reads

$$\text{adj}(\mathbf{C} + \mathbf{D}) = \text{adj}(\mathbf{C}) + D_{\mathbf{C}}(\text{adj}(\mathbf{C}))(\mathbf{D}) + \frac{1}{2}D_{\mathbf{C}\mathbf{C}}(\text{adj}(\mathbf{C}))(\mathbf{D})(\mathbf{D}) + O(\|\mathbf{D}\|^3),$$

where the derivatives are given by Equation (24) and Equation (25).

*Proof.* Since Equation (10) holds and  $\mathbf{C}$  is invertible, we have that  $\text{adj}(\mathbf{C}) = \det(\mathbf{C})\mathbf{C}^{-1}$ . Therefore

$$\begin{aligned} D_{\mathbf{C}}(\text{adj}(\mathbf{C}))(\mathbf{D}) &= D_{\mathbf{C}}(\det(\mathbf{C})\mathbf{C}^{-1})(\mathbf{D}) = D_{\mathbf{C}}(\det(\mathbf{C}))(\mathbf{D})\mathbf{C}^{-1} + \det(\mathbf{C})D_{\mathbf{C}}(\mathbf{C}^{-1})(\mathbf{D}), \\ &= \det(\mathbf{C})\text{tr}(\mathbf{C}^{-1}\mathbf{D})\mathbf{C}^{-1} - \det(\mathbf{C})\mathbf{C}^{-1}\mathbf{D}\mathbf{C}^{-1}, \end{aligned}$$

where we have used the rule for the derivative of a product, the Jacobi formula Equation (17) and the identity  $D_{\mathbf{C}}(\mathbf{C}^{-1})(\mathbf{D}) = -\mathbf{C}^{-1}\mathbf{D}\mathbf{C}^{-1}$ . For the second derivative, we have

$$\begin{aligned} D_{\mathbf{C}\mathbf{C}}(\text{adj}(\mathbf{C}))(\mathbf{D})(\mathbf{E}) &:= D_{\mathbf{C}}(D_{\mathbf{C}}(\text{adj}(\mathbf{C}))(\mathbf{D}))(\mathbf{E}) = D_{\mathbf{C}}(\det(\mathbf{C})(\text{tr}(\mathbf{C}^{-1}\mathbf{D})\mathbf{I} - \mathbf{C}^{-1}\mathbf{D})\mathbf{C}^{-1})(\mathbf{E}), \\ &= D_{\mathbf{C}}(\det(\mathbf{C}))(\mathbf{E})(\text{tr}(\mathbf{C}^{-1}\mathbf{D})\mathbf{I} - \mathbf{C}^{-1}\mathbf{D})\mathbf{C}^{-1} + \det(\mathbf{C})D_{\mathbf{C}}((\text{tr}(\mathbf{C}^{-1}\mathbf{D})\mathbf{I} - \mathbf{C}^{-1}\mathbf{D})\mathbf{C}^{-1})(\mathbf{E}), \\ &= \det(\mathbf{C})\text{tr}(\mathbf{C}^{-1}\mathbf{E})(\text{tr}(\mathbf{C}^{-1}\mathbf{D})\mathbf{I} - \mathbf{C}^{-1}\mathbf{D})\mathbf{C}^{-1} + \det(\mathbf{C})D_{\mathbf{C}}(\text{tr}(\mathbf{C}^{-1}\mathbf{D})\mathbf{I} - \mathbf{C}^{-1}\mathbf{D})(\mathbf{E})\mathbf{C}^{-1} \\ &\quad + \det(\mathbf{C})(\text{tr}(\mathbf{C}^{-1}\mathbf{D})\mathbf{I} - \mathbf{C}^{-1}\mathbf{D})D_{\mathbf{C}}(\mathbf{C}^{-1})(\mathbf{E}), \\ &= \det(\mathbf{C})\text{tr}(\mathbf{C}^{-1}\mathbf{E})(\text{tr}(\mathbf{C}^{-1}\mathbf{D})\mathbf{I} - \mathbf{C}^{-1}\mathbf{D})\mathbf{C}^{-1} \\ &\quad + \det(\mathbf{C})(\text{tr}(D_{\mathbf{C}}(\mathbf{C}^{-1})(\mathbf{E})\mathbf{D})\mathbf{I} - D_{\mathbf{C}}(\mathbf{C}^{-1})(\mathbf{E})\mathbf{D})\mathbf{C}^{-1} \\ &\quad - \det(\mathbf{C})(\text{tr}(\mathbf{C}^{-1}\mathbf{D})\mathbf{I} - \mathbf{C}^{-1}\mathbf{D})\mathbf{C}^{-1}\mathbf{E}\mathbf{C}^{-1}, \\ &= \det(\mathbf{C})\text{tr}(\mathbf{C}^{-1}\mathbf{E})(\text{tr}(\mathbf{C}^{-1}\mathbf{D})\mathbf{I} - \mathbf{C}^{-1}\mathbf{D})\mathbf{C}^{-1} - \det(\mathbf{C})(\text{tr}(\mathbf{C}^{-1}\mathbf{E}\mathbf{C}^{-1}\mathbf{D})\mathbf{I} - \mathbf{C}^{-1}\mathbf{E}\mathbf{C}^{-1}\mathbf{D})\mathbf{C}^{-1} \\ &\quad - \det(\mathbf{C})(\text{tr}(\mathbf{C}^{-1}\mathbf{D})\mathbf{I} - \mathbf{C}^{-1}\mathbf{D})\mathbf{C}^{-1}\mathbf{E}\mathbf{C}^{-1}, \end{aligned}$$

where we have used the rule for the derivative of a product, the Jacobi formula Equation (17), the linearity of the derivative and the trace and the identity  $D_{\mathbf{C}}(\mathbf{C}^{-1})(\mathbf{D}) = -\mathbf{C}^{-1}\mathbf{D}\mathbf{C}^{-1}$ . Upon rearrangement, this yields the result.  $\square$

**Remark 5.7.** We observe that, looking at Equation (24) and Equation (25) compared to Equation (17) and Equation (18), we have that

$$\begin{aligned} D_{\mathbf{C}}(\text{adj}(\mathbf{C}))(\mathbf{D}) &= D_{\mathbf{C}}(\det(\mathbf{C}))(\mathbf{D})\mathbf{C}^{-1} - \det(\mathbf{C})\mathbf{C}^{-1}\mathbf{D}\mathbf{C}^{-1}, \\ D_{\mathbf{C}\mathbf{C}}(\text{adj}(\mathbf{C}))(\mathbf{D})(\mathbf{E}) &= D_{\mathbf{C}\mathbf{C}}(\det(\mathbf{C}))(\mathbf{D})(\mathbf{E})\mathbf{C}^{-1} + \det(\mathbf{C})\mathbf{C}^{-1}(\mathbf{E}\mathbf{C}^{-1}\mathbf{D} + \mathbf{D}\mathbf{C}^{-1}\mathbf{E} - \text{tr}(\mathbf{C}^{-1}\mathbf{E})\mathbf{D} - \text{tr}(\mathbf{C}^{-1}\mathbf{D})\mathbf{E})\mathbf{C}^{-1}. \end{aligned}$$

This implies that we can reuse the computations we did for the determinant in the current treatment of the adjugate, as far as the first terms on the right hand sides are concerned. However, one must be careful that now they are multiplied by  $\mathbf{C}^{-1}$ .

If we had stopped the developments at first order, we could have used the first-order perturbation theory of the adjugate matrix as provided by Theorem 2.1 from [40]. However, to the best of our knowledge, no second-order perturbation theory for this matrix is available in the literature, thus we have been compelled to independently develop it using differential calculus. Lemma 5.6 is thus a generalization of the results from [40] and can therefore be used – beyond the application presented in this contribution – by researchers needing a second-order perturbation theory for the adjugate matrix.

Since we are ultimately interested, as one can notice from Equation (9), in multiplying the formal power series  $\text{adj}(\zeta\mathbf{I} - \mathbf{A}) \in \mathcal{M}_q(\mathcal{S})$  by  $\mathbf{B} \in \mathcal{M}_q(\mathcal{S})$  in a Cauchy-like fashion (the standard product of formal power series) and select the first row, see Proposition 3.4 bis, we perform the computations only for the first row of  $\text{adj}(\zeta\mathbf{I} - \mathbf{A})$ .

- Using the definition of the adjugate matrix in combination with the Laplace formula or using the explicit formula for the adjugate of an upper triangular matrix, see [22], we have

$$\text{adj}(\mathbf{S}) = \Pi \text{diag}\left(1, \frac{s_1}{s_2}, \dots, \frac{s_1}{s_q}\right), \quad \text{thus} \quad \lim_{s_1 \rightarrow 0} \text{adj}(\mathbf{S}) = \Pi \mathbf{e}_1 \otimes \mathbf{e}_1.$$

Hence, contrarily to the determinant, the zero-order term in  $\text{adj}(\zeta\mathbf{I} - \mathbf{A})$  is not zero for  $s_1 = 0$  but is still a singular one-rank diagonal matrix.



- Let  $\mathbf{C} = \mathbf{S} \in \text{GL}_q(\mathbb{R}) \subset \text{GL}_q(\mathcal{S})$  and  $\mathbf{D} = \frac{\Delta x}{\lambda} (\partial_t \mathbf{I} + \mathcal{G}(\mathbf{I} - \mathbf{S})) + \frac{\Delta x^2}{2\lambda^2} (\partial_{tt} \mathbf{I} - \mathcal{G}^2(\mathbf{I} - \mathbf{S})) + O(\Delta x^3) \in \mathcal{M}_q(\mathcal{S})$  from Lemma 5.3. We utilize the previous computations from Equation (19) as suggested in Remark 5.7, plus Equation (24).

$$\begin{aligned} \text{D}_{\mathbf{C}}(\text{adj}(\mathbf{C}))(\mathbf{D}) &= s_1 \Pi \left( \frac{\Delta x}{\lambda} \left( \frac{1}{s_1} (\partial_t + \mathcal{G}_{11}) \right) + \sum_{i=2}^q \frac{1}{s_i} (\partial_t + (1-s_i) \mathcal{G}_{ii}) \right) + \frac{\Delta x^2}{2\lambda^2} \left( \frac{1}{s_1} (\partial_{tt} - \mathcal{G}_{11} \mathcal{G}_{11} - \sum_{\ell=2}^q \mathcal{G}_{1\ell} \mathcal{G}_{\ell 1}) \right) \\ &\quad + \sum_{i=2}^q \frac{1}{s_i} \left( \partial_{tt} - (1-s_i) \sum_{\ell=1}^q \mathcal{G}_{i\ell} \mathcal{G}_{\ell i} \right) \text{diag} \left( \frac{1}{s_1}, \frac{1}{s_2}, \dots, \frac{1}{s_q} \right) \\ &\quad - s_1 \Pi \text{diag} \left( \frac{1}{s_1}, \frac{1}{s_2}, \dots, \frac{1}{s_q} \right) \mathbf{D} \text{diag} \left( \frac{1}{s_1}, \frac{1}{s_2}, \dots, \frac{1}{s_q} \right) + O(\Delta x^3). \end{aligned}$$

In this case, we do not even have to take the limit for  $s_1 \rightarrow 0$ , since all the terms in  $s_1$  cancel. Therefore, for the very first component, we get

$$\begin{aligned} (\text{D}_{\mathbf{C}}(\text{adj}(\mathbf{C}))(\mathbf{D}))_{11} &= \Pi \frac{\Delta x}{\lambda} \left( \left( \sum_{\ell=2}^q \frac{1}{s_\ell} \right) \partial_t + \sum_{i=2}^q \left( \frac{1}{s_i} - 1 \right) \mathcal{G}_{ii} + \frac{\Delta x}{2\lambda} \left( \left( \sum_{\ell=2}^q \frac{1}{s_\ell} \right) \partial_{tt} - \sum_{i=2}^q \left( \frac{1}{s_i} - 1 \right) \sum_{\ell=1}^q \mathcal{G}_{i\ell} \mathcal{G}_{\ell i} \right) \right) \\ &\quad + O(\Delta x^3). \end{aligned} \tag{26}$$

Now consider  $j \in [2..q]$ , then

$$(\text{D}_{\mathbf{C}}(\text{adj}(\mathbf{C}))(\mathbf{D}))_{1j} = -\Pi \frac{\Delta x}{\lambda} \left( \frac{1}{s_j} - 1 \right) \left( \mathcal{G}_{1j} - \frac{\Delta x}{2\lambda} \left( \mathcal{G}_{11} \mathcal{G}_{1j} + \sum_{\ell=2}^q \mathcal{G}_{1\ell} \mathcal{G}_{\ell j} \right) \right) + O(\Delta x^3), \tag{27}$$

This gives all the first-order terms on the first row of  $\text{adj}(\zeta \mathbf{I} - \mathcal{A})$  and part of the second-order terms.

- Let  $\mathbf{C} = \mathbf{S} \in \text{GL}_q(\mathbb{R}) \subset \text{GL}_q(\mathcal{S})$  and  $\mathbf{D} = \frac{\Delta x}{\lambda} (\partial_t \mathbf{I} + \mathcal{G}(\mathbf{I} - \mathbf{S})) + O(\Delta x^2) \in \mathcal{M}_q(\mathcal{S})$  from Lemma 5.3. We reuse computations from Equation (21) as well as Equation (25).

$$\begin{aligned} \text{D}_{\mathbf{C}\mathbf{C}}(\text{adj}(\mathbf{C}))(\mathbf{D})(\mathbf{D}) &= s_1 \Pi \frac{\Delta x^2}{\lambda^2} \left( \frac{2}{s_1} (\partial_t + \mathcal{G}_{11}) \sum_{i=2}^q \frac{1}{s_i} (\partial_t + (1-s_i) \mathcal{G}_{ii}) + \left( \sum_{i=2}^q \frac{1}{s_i} (\partial_t + (1-s_i) \mathcal{G}_{ii}) \right)^2 \right. \\ &\quad - \frac{2}{s_1} \sum_{\ell=2}^q \left( \frac{1}{s_\ell} - 1 \right) \mathcal{G}_{1\ell} \mathcal{G}_{\ell 1} - \sum_{i=2}^q \frac{1}{s_i^2} (\partial_t + (1-s_i) \mathcal{G}_{ii})^2 \\ &\quad \left. - \sum_{i=2}^q \sum_{\substack{\ell=2 \\ \ell \neq i}}^q \left( \frac{1}{s_i} - 1 \right) \left( \frac{1}{s_\ell} - 1 \right) \mathcal{G}_{i\ell} \mathcal{G}_{\ell i} \right) \text{diag} \left( \frac{1}{s_1}, \frac{1}{s_2}, \dots, \frac{1}{s_q} \right) \\ &\quad + 2s_1 \Pi \text{diag} \left( \frac{1}{s_1}, \frac{1}{s_2}, \dots, \frac{1}{s_q} \right) \left( \mathbf{D} \mathbf{S}^{-1} \mathbf{D} - \text{tr}(\mathbf{S}^{-1} \mathbf{D}) \mathbf{D} \right) \text{diag} \left( \frac{1}{s_1}, \frac{1}{s_2}, \dots, \frac{1}{s_q} \right) + O(\Delta x^3). \end{aligned}$$

Then we have, for the first matrix entry

$$\begin{aligned} (\text{D}_{\mathbf{C}\mathbf{C}}(\text{adj}(\mathbf{C}))(\mathbf{D})(\mathbf{D}))_{11} &= \Pi \frac{\Delta x^2}{\lambda^2} \left( \left( \sum_{i=2}^q \frac{1}{s_i} (\partial_t + (1-s_i) \mathcal{G}_{ii}) \right)^2 - \sum_{i=2}^q \frac{1}{s_i^2} (\partial_t + (1-s_i) \mathcal{G}_{ii})^2 \right. \\ &\quad \left. - \sum_{i=2}^q \sum_{\substack{\ell=2 \\ \ell \neq i}}^q \left( \frac{1}{s_i} - 1 \right) \left( \frac{1}{s_\ell} - 1 \right) \mathcal{G}_{i\ell} \mathcal{G}_{\ell i} \right) + O(\Delta x^3), \end{aligned} \tag{28}$$

independent from  $s_1$ . On the other hand, for  $j \in [2..q]$

$$\begin{aligned} (\text{D}_{\mathbf{C}\mathbf{C}}(\text{adj}(\mathbf{C}))(\mathbf{D})(\mathbf{D}))_{1j} &= 2\Pi \frac{\Delta x^2}{\lambda^2} \left( \frac{1}{s_j} - 1 \right) \left( \frac{1}{s_j} \mathcal{G}_{1j} (\partial_t + (1-s_j) \mathcal{G}_{jj}) + \sum_{\substack{\ell=2 \\ \ell \neq j}}^q \left( \frac{1}{s_\ell} - 1 \right) \mathcal{G}_{1\ell} \mathcal{G}_{\ell j} \right. \\ &\quad \left. - \mathcal{G}_{1j} \sum_{i=2}^q \frac{1}{s_i} (\partial_t + (1-s_i) \mathcal{G}_{ii}) \right) + O(\Delta x^3). \end{aligned} \tag{29}$$

Using Equation (26) and Equation (28), we have that the first entry on the first row of  $\text{adj}(\zeta \mathbf{I} - \mathcal{A})$  is

$$\begin{aligned} \lim_{s_1 \rightarrow 0} (\text{adj}(\zeta \mathbf{I} - \mathcal{A}))_{11} &= \Pi + \Pi \frac{\Delta x}{\lambda} \left( \left( \sum_{\ell=2}^q \frac{1}{s_\ell} \right) \partial_t + \sum_{i=2}^q \left( \frac{1}{s_i} - 1 \right) \mathcal{G}_{ii} \right) + \Pi \frac{\Delta x^2}{2\lambda^2} \left( \left( \sum_{\ell=2}^q \frac{1}{s_\ell} \right) \partial_{tt} - \sum_{i=2}^q \left( \frac{1}{s_i} - 1 \right) \sum_{\ell=1}^q \mathcal{G}_{i\ell} \mathcal{G}_{\ell i} \right. \\ &+ 2 \left( \sum_{i=2}^q \frac{1}{s_i} (\partial_t + (1 - s_i) \mathcal{G}_{ii}) \right)^2 - 2 \sum_{i=2}^q \frac{1}{s_i^2} (\partial_t + (1 - s_i) \mathcal{G}_{ii})^2 - 2 \sum_{i=2}^q \sum_{\substack{\ell=2 \\ \ell \neq i}}^q \left( \frac{1}{s_i} - 1 \right) \left( \frac{1}{s_\ell} - 1 \right) \mathcal{G}_{i\ell} \mathcal{G}_{\ell i} \Big) + O(\Delta x^3). \end{aligned} \quad (30)$$

On the other hand, using Equation (27) and Equation (29), for any  $j \in [2..q]$ , we write

$$\begin{aligned} \lim_{s_1 \rightarrow 0} (\text{adj}(\zeta \mathbf{I} - \mathcal{A}))_{1j} &= -\Pi \frac{\Delta x}{\lambda} \left( \frac{1}{s_j} - 1 \right) \mathcal{G}_{1j} + \Pi \frac{\Delta x^2}{2\lambda^2} \left( \frac{1}{s_j} - 1 \right) \left( \mathcal{G}_{11} \mathcal{G}_{1j} + \sum_{\ell=2}^q \mathcal{G}_{1\ell} \mathcal{G}_{\ell j} + \frac{2}{s_j} \mathcal{G}_{1j} (\partial_t + (1 - s_j) \mathcal{G}_{jj}) \right. \\ &\quad \left. + 2 \sum_{\substack{\ell=2 \\ \ell \neq j}}^q \left( \frac{1}{s_\ell} - 1 \right) \mathcal{G}_{1\ell} \mathcal{G}_{\ell j} - 2 \mathcal{G}_{1j} \sum_{i=2}^q \frac{1}{s_i} (\partial_t + (1 - s_i) \mathcal{G}_{ii}) \right) + O(\Delta x^3). \end{aligned} \quad (31)$$

In general, we have written, for the first row, the leading terms in  $\text{adj}(\zeta \mathbf{I} - \mathcal{A})$ . We shall take its product with  $\mathcal{B}$ . Thus, one has

$$\begin{aligned} \text{adj}(\zeta \mathbf{I} - \mathcal{A}) \mathcal{B} &= \text{adj}(\zeta \mathbf{I} - \mathcal{A})^{(0)} \mathcal{B}^{(0)} + \Delta x (\text{adj}(\zeta \mathbf{I} - \mathcal{A})^{(0)} \mathcal{B}^{(1)} + \text{adj}(\zeta \mathbf{I} - \mathcal{A})^{(1)} \mathcal{B}^{(0)}), \\ &+ \Delta x^2 (\text{adj}(\zeta \mathbf{I} - \mathcal{A})^{(0)} \mathcal{B}^{(2)} + \text{adj}(\zeta \mathbf{I} - \mathcal{A})^{(1)} \mathcal{B}^{(1)} + \text{adj}(\zeta \mathbf{I} - \mathcal{A})^{(2)} \mathcal{B}^{(0)}) + O(\Delta x^3), \end{aligned} \quad (32)$$

generating products of terms in the fashion of the Cauchy product. This completes the preliminary results needed to prove Theorem 4.4 and Theorem 4.5.

### 5.3 Overall computation

We now put all the previous calculations together to prove Theorem 4.4 and Theorem 4.5. As previously pointed out, we can assume, without loss of generality, that  $s_1 = 0$ , passing to the limit. This allows to deal with simpler expressions with less terms.

#### 5.3.1 First-order equations

Starting with Theorem 4.4, it is sufficient to truncate all the formal power series at  $O(\Delta x^2)$ . In particular, using the fact that the first column of  $\mathcal{B}$  is zero for  $s_1 = 0$ , we have that  $\lim_{s_1 \rightarrow 0} (\text{adj}(\zeta \mathbf{I} - \mathcal{A}) \mathcal{B})_{11} = 0$ . Observe that if the relaxation parameter corresponding to the conserved moment were not equal to zero, we would have had  $(\text{adj}(\zeta \mathbf{I} - \mathcal{A}) \mathcal{B})_{11} = O(1)$ . Still the matrix  $\mathcal{S}$  would not have been singular thus we would have some non vanishing zero-order term in  $\det(\zeta \mathbf{I} - \mathcal{A})$  to compensate the one from the adjugate.

On the other hand, for any  $j \in [2..q]$ , using Equation (31), Lemma 5.2 and Equation (32), entails

$$\begin{aligned} \lim_{s_1 \rightarrow 0} (\text{adj}(\zeta \mathbf{I} - \mathcal{A}) \mathcal{B})_{1j} &= \Delta x \left( (\text{adj}(\zeta \mathbf{I} - \mathcal{A})^{(0)} \mathcal{B}^{(1)})_{1j} + (\text{adj}(\zeta \mathbf{I} - \mathcal{A})^{(1)} \mathcal{B}^{(0)})_{1j} \right) + O(\Delta x^2), \\ &= \Delta x \left( -\Pi \frac{1}{\lambda} s_j \mathcal{G}_{1j} - \Pi \frac{1}{\lambda} s_j \left( \frac{1}{s_j} - 1 \right) \right) + O(\Delta x^2) = -\Pi \frac{\Delta x}{\lambda} \mathcal{G}_{1j} + O(\Delta x^2). \end{aligned}$$

On the other hand, Equation (23) directly yields

$$\lim_{s_1 \rightarrow 0} \det(\zeta \mathbf{I} - \mathcal{A}) = \Pi \frac{\Delta x}{\lambda} (\partial_t + \mathcal{G}_{11}) + O(\Delta x^2),$$

thus we obtain the modified equation (whatever the choice of  $s_1 \in \mathbb{R}$ )

$$\Pi \frac{\Delta x}{\lambda} \left( \partial_t m_1 + \mathcal{G}_{11} m_1 + \sum_{j=2}^q \mathcal{G}_{1j} m_j^{\text{eq}} \right) = O(\Delta x^2),$$

which is the desired result from Theorem 4.4 for  $N = 1$  upon dividing by the constant  $\Pi = O(1)$ . Observe that the term  $\Pi$  is never present in the computations by [15] because they are done on the original lattice Boltzmann scheme Equation (5) or Equation (8) which has only one time step. For instance, in [15], the multi-step

nature of the problem, generated by the non-conserved moments relaxing away from the equilibrium, is damped at the very beginning of the procedure by performing the Taylor expansions of the scheme on the non-conserved variables and then plugging them into the expansions for the conserved moments.

Before proceeding, let us utilize the previous equation to get rid of the time derivatives in the second order terms. This is in general not compulsory but constitutes the policy by [13, 15] and is common to all the approaches (Chapman-Enskog, equivalent equation, Maxwell iteration, *etc.*) in order to find the value of the diffusion coefficients from the second-order terms. Notice that in this case, where  $N = 1$ ,  $\gamma_1$ , is a scalar here denoted  $\gamma_1$  for brevity.

$$\partial_t m_1 = -\mathcal{G}_{11} m_1 - \sum_{j=2}^q \mathcal{G}_{1j} m_j^{\text{eq}} + O(\Delta x) = -\gamma_1 + O(\Delta x), \quad (33)$$

$$\partial_t m^{\text{eq}} = \frac{dm^{\text{eq}}}{dm_1} \partial_t m_1 = -\frac{dm^{\text{eq}}}{dm_1} \left( \mathcal{G}_{11} m_1 + \sum_{j=2}^q \mathcal{G}_{1j} m_j^{\text{eq}} \right) + O(\Delta x) = -\frac{dm^{\text{eq}}}{dm_1} \gamma_1 + O(\Delta x), \quad (34)$$

$$\partial_{tt} m_1 = -\partial_t \left( \mathcal{G}_{11} m_1 + \sum_{j=2}^q \mathcal{G}_{1j} m_j^{\text{eq}} \right) + O(\Delta x) = -\mathcal{G}_{11} \partial_t m_1 - \sum_{j=2}^q \mathcal{G}_{1j} \partial_t m_j^{\text{eq}} + O(\Delta x), \quad (35)$$

$$= \mathcal{G}_{11} \gamma_1 + \sum_{j=2}^q \mathcal{G}_{1j} \frac{dm_j^{\text{eq}}}{dm_1} \gamma_1 + O(\Delta x) = \mathcal{G}_{11} \mathcal{G}_{11} m_1 + \mathcal{G}_{11} \sum_{j=2}^q \mathcal{G}_{1j} m_j^{\text{eq}} + \sum_{j=2}^q \mathcal{G}_{1j} \frac{dm_j^{\text{eq}}}{dm_1} \gamma_1 + O(\Delta x). \quad (36)$$

These equalities are formal and obtained by taking advantage either of the chain rule, since the moments at equilibrium are functions of the conserved moments, or of the re-injection of Equation (33) by assuming that the differentiation preserves the asymptotic relations from the symbol  $O(\cdot)$ . This way of proceeding is fully understood in the framework of Finite Difference schemes for PDEs, see [45, 1].

### 5.3.2 Second-order equations

We can now go to the proof of Theorem 4.5, which is more involved due to the presence of more terms to estimate. To make the link with the findings of [15], the increased complexity comes from the more intricate and entangled block structure of  $\mathcal{G}^2$ . We have to treat the second-order term in Equation (32), made up of three products. For any  $j \in [2..q]$  (once again, the first component vanishes for  $s_1 = 0$ )

- Using Lemma 5.2 and the zero-order expansion of the adjugate gives

$$\lim_{s_1 \rightarrow 0} (\text{adj}(\zeta \mathbf{I} - \mathcal{A})^{(0)} \mathcal{B}^{(2)})_{1j} = \Pi \mathcal{B}_{1j}^{(2)} = \Pi \frac{1}{2\lambda^2} s_j \left( \mathcal{G}_{11} \mathcal{G}_{1j} + \sum_{\ell=2}^q \mathcal{G}_{1\ell} \mathcal{G}_{\ell j} \right).$$

- Using Lemma 5.2 with Equation (30) and Equation (31)

$$\begin{aligned} \lim_{s_1 \rightarrow 0} (\text{adj}(\zeta \mathbf{I} - \mathcal{A})^{(1)} \mathcal{B}^{(1)})_{1j} &= -\frac{1}{\lambda} s_j \sum_{\ell=1}^q (\text{adj}(\zeta \mathbf{I} - \mathcal{A})^{(1)})_{1\ell} \mathcal{G}_{\ell j}, \\ &= -\Pi \frac{1}{\lambda^2} s_j \left( \mathcal{G}_{1j} \left( \sum_{\ell=2}^q \frac{1}{s_\ell} \right) \partial_t + \mathcal{G}_{1j} \sum_{\ell=2}^q \left( \frac{1}{s_\ell} - 1 \right) \mathcal{G}_{\ell\ell} - \sum_{\ell=2}^q \left( \frac{1}{s_\ell} - 1 \right) \mathcal{G}_{1\ell} \mathcal{G}_{\ell j} \right). \end{aligned}$$

- Using Lemma 5.2 and Equation (31)

$$\begin{aligned} \lim_{s_1 \rightarrow 0} (\text{adj}(\zeta \mathbf{I} - \mathcal{A})^{(2)} \mathcal{B}^{(0)})_{1j} &= s_j (\text{adj}(\zeta \mathbf{I} - \mathcal{A})^{(2)})_{1j}, \\ &= \Pi \frac{1}{\lambda^2} \left( 1 - s_j \right) \left( \frac{1}{2} \mathcal{G}_{11} \mathcal{G}_{1j} + \sum_{\ell=2}^q \left( \frac{1}{s_\ell} - \frac{1}{2} \right) \mathcal{G}_{1\ell} \mathcal{G}_{\ell j} - \left( \frac{1}{s_j} - 1 \right) \mathcal{G}_{1j} \mathcal{G}_{jj} \right. \\ &\quad \left. + \frac{1}{s_j} \mathcal{G}_{1j} (\partial_t + (1 - s_j) \mathcal{G}_{jj}) - \mathcal{G}_{1j} \sum_{i=2}^q \frac{1}{s_i} (\partial_t + (1 - s_i) \mathcal{G}_{ii}) \right). \end{aligned}$$

Summing these three contributions and after some straightforward but tedious computations, the second-order term in Equation (32) is given by

$$\lim_{s_1 \rightarrow 0} ((\text{adj}(\zeta I - \mathcal{A})\mathcal{B})^{(2)})_{1j} = \Pi \frac{1}{\lambda^2} \left( \frac{1}{2} \mathcal{G}_{11} \mathcal{G}_{1j} + \sum_{\ell=2}^q \left( \frac{1}{s_\ell} - \frac{1}{2} \right) \mathcal{G}_{1\ell} \mathcal{G}_{\ell j} - \mathcal{G}_{1j} \left( 1 + \sum_{\substack{\ell=2 \\ \ell \neq j}}^q \frac{1}{s_\ell} \right) \partial_t - \mathcal{G}_{1j} \sum_{\ell=2}^q \left( \frac{1}{s_\ell} - 1 \right) \mathcal{G}_{\ell\ell} \right).$$

Hence, using Equation (34) to get rid of the time derivative of the equilibria, we have

$$\begin{aligned} \lim_{s_1 \rightarrow 0} \sum_{j=2}^q ((\text{adj}(\zeta I - \mathcal{A})\mathcal{B})^{(2)})_{1j} m_j^{\text{eq}} &= \Pi \frac{1}{\lambda^2} \sum_{j=2}^q \left( \frac{1}{2} \mathcal{G}_{11} \mathcal{G}_{1j} m_j^{\text{eq}} + \sum_{\ell=2}^q \left( \frac{1}{s_\ell} - \frac{1}{2} \right) \mathcal{G}_{1\ell} \mathcal{G}_{\ell j} m_j^{\text{eq}} \right. \\ &\quad \left. + \mathcal{G}_{1j} \left( 1 + \sum_{\substack{\ell=2 \\ \ell \neq j}}^q \frac{1}{s_\ell} \right) \frac{dm_j^{\text{eq}}}{dm_1} \gamma_1 - \mathcal{G}_{1j} \sum_{\ell=2}^q \left( \frac{1}{s_\ell} - 1 \right) \mathcal{G}_{\ell\ell} m_j^{\text{eq}} \right) + O(\Delta x). \end{aligned}$$

Notice that in this result, a reminder of order  $O(\Delta x)$  appears. Once more, using Equation (33) and Equation (36) to eliminate the time derivatives in the second-order terms from Equation (23) gives

$$\begin{aligned} \lim_{s_1 \rightarrow 0} (\det(\zeta I - \mathcal{A}))^{(2)} m_1 &= \Pi \frac{1}{\lambda^2} \left( \left( \frac{1}{2} + \sum_{\ell=2}^q \frac{1}{s_\ell} \right) \partial_{tt} m_1 + \left( \sum_{\ell=2}^q \frac{1}{s_\ell} \right) \mathcal{G}_{11} \partial_t m_1 + \left( \sum_{i=2}^q \left( \frac{1}{s_i} - 1 \right) \mathcal{G}_{ii} \right) \partial_t m_1 \right. \\ &\quad \left. - \frac{1}{2} \mathcal{G}_{11} \mathcal{G}_{11} m_1 - \sum_{\ell=2}^q \left( \frac{1}{s_\ell} - \frac{1}{2} \right) \mathcal{G}_{1\ell} \mathcal{G}_{\ell 1} m_1 + \mathcal{G}_{11} \left( \sum_{i=2}^q \left( \frac{1}{s_i} - 1 \right) \mathcal{G}_{ii} \right) m_1 \right), \\ &= \Pi \frac{1}{\lambda^2} \left( \left( \frac{1}{2} + \sum_{\ell=2}^q \frac{1}{s_\ell} \right) \left( \mathcal{G}_{11} \mathcal{G}_{11} + \mathcal{G}_{11} \sum_{j=2}^q \mathcal{G}_{1j} + \sum_{j=2}^q \mathcal{G}_{1j} \frac{dm_j^{\text{eq}}}{dm_1} \gamma_1 \right) \right. \\ &\quad \left. - \left( \sum_{\ell=2}^q \frac{1}{s_\ell} \right) \mathcal{G}_{11} \left( \mathcal{G}_{11} m_1 + \sum_{j=2}^q \mathcal{G}_{1j} m_j^{\text{eq}} \right) - \left( \sum_{i=2}^q \left( \frac{1}{s_i} - 1 \right) \mathcal{G}_{ii} \right) \left( \mathcal{G}_{11} m_1 + \sum_{j=2}^q \mathcal{G}_{1j} m_j^{\text{eq}} \right) \right. \\ &\quad \left. - \frac{1}{2} \mathcal{G}_{11} \mathcal{G}_{11} m_1 - \sum_{\ell=2}^q \left( \frac{1}{s_\ell} - \frac{1}{2} \right) \mathcal{G}_{1\ell} \mathcal{G}_{\ell 1} m_1 + \mathcal{G}_{11} \left( \sum_{i=2}^q \left( \frac{1}{s_i} - 1 \right) \mathcal{G}_{ii} \right) m_1 \right) + O(\Delta x). \end{aligned}$$

With this, after simplifications, we obtain the remaining term to master the second-order contributions in the expansion of the Finite Difference scheme Equation (9).

$$\begin{aligned} &(\det(\zeta I - \mathcal{A}))^{(2)} m_1 - \sum_{j=2}^q ((\text{adj}(\zeta I - \mathcal{A})\mathcal{B})^{(2)})_{1j} m_j^{\text{eq}} \\ &= \Pi \frac{1}{\lambda^2} \left( - \sum_{j=2}^q \left( \frac{1}{s_j} - \frac{1}{2} \right) \mathcal{G}_{1j} \mathcal{G}_{j1} m_1 - \sum_{j=2}^q \sum_{\ell=2}^q \left( \frac{1}{s_\ell} - \frac{1}{2} \right) \mathcal{G}_{1\ell} \mathcal{G}_{\ell j} m_j^{\text{eq}} + \sum_{j=2}^q \left( \frac{1}{s_j} - \frac{1}{2} \right) \mathcal{G}_{1j} \frac{dm_j^{\text{eq}}}{dm_1} \gamma_1 \right) + O(\Delta x). \end{aligned}$$

To wrap up, these computations yield, together with the ones from Section 5.3.1, the expected result for  $N = 1$ , which reads

$$\frac{\Delta x}{\lambda} \Pi \left( \partial_t m_1 + \gamma_1 + \frac{\Delta x}{\lambda} \sum_{j=2}^q \left( \frac{1}{s_j} - \frac{1}{2} \right) \mathcal{G}_{1j} \left( \frac{dm_j}{dm_1} \gamma_1 - \mathcal{G}_{j1} m_1 - \sum_{\ell=2}^q \mathcal{G}_{j\ell} m_\ell^{\text{eq}} \right) \right) = O(\Delta x^3),$$

and thus proves Theorem 4.5.

## 6 Extension of the proofs to several conserved moments: key ideas

In this Section, we sketch the demonstration of Theorem 4.4 and Theorem 4.5 for any  $N \geq 1$ . For the sake of providing a quick and effective presentation of this matter, we limit ourselves to first-order in  $\Delta x$ . Select a conserved moment, which shall be indexed by  $i \in [1..N]$ .

**Remark 6.1.** The operation selecting rows and columns to yield  $\mathbf{A}_i$  and  $\mathbf{A}_i^\diamond$  from Proposition 3.5 bis does not change the orders of the expansions. This is, let  $\mathbf{C} \in \mathcal{M}_q(\mathbb{R}[\mathbf{z}] \otimes_{\mathbb{R}} \mathbb{D})$  and  $\mathcal{C} = \sum_{r=0}^{+\infty} \Delta x^r \mathbf{C}^{(r)} \in \mathcal{M}_q(\mathcal{S})$  such that  $\mathbf{C} \asymp \mathcal{C}$  and  $I \subset [1..q]$  a set of indices, then

$$\mathbf{C}_I \asymp \left( \sum_{r=0}^{+\infty} \Delta x^r \mathbf{C}^{(r)} \right)_I = \sum_{r=0}^{+\infty} \Delta x^r (\mathbf{C}^{(r)})_I.$$

Thus we have the analogous of Lemma 5.3, where  $\mathbf{zI} - \mathbf{A}_i \asymp \zeta \mathbf{I} - \mathcal{A}_i$ , with

$$\zeta \mathbf{I} - \mathcal{A}_i = \sum_{r=0}^{+\infty} \frac{\Delta x^r}{\lambda^r r!} \left( \partial_t^r \mathbf{I} - (-1)^r (\mathcal{G}^r (\mathbf{I} - \mathbf{S}))_{\{i\} \cup [N+1..q]} \right).$$

The first two term in the expansion of the inverse of the resolvent are

$$(\zeta \mathbf{I} - \mathcal{A}_i)^{(0)} = \text{diag}(1, \dots, 1, s_i, 1, \dots, 1, s_{N+1}, \dots, s_q).$$

In the spirit of Remark 5.5 and the computations developed in Section 5, for the case  $s_i = 0$ , we introduce a regularization with  $s_i \neq 0$  and then we pass to the limit. Moreover

$$(\zeta \mathbf{I} - \mathcal{A}_i)^{(1)} = \frac{1}{\lambda} \partial_t \mathbf{I} + \frac{1}{\lambda} \begin{pmatrix} 0 & \cdots & 0 & | & 0 & & 0 & \cdots & 0 & || & 0 & & \cdots & 0 \\ \vdots & \ddots & \vdots & | & \vdots & & \vdots & \ddots & \vdots & || & \vdots & & \ddots & \vdots \\ 0 & \cdots & 0 & | & 0 & & 0 & \cdots & 0 & || & 0 & & \cdots & 0 \\ \hline 0 & \cdots & 0 & | & (1-s_i)\mathcal{G}_{ii} & & 0 & \cdots & 0 & || & (1-s_{N+1})\mathcal{G}_{i(N+1)} & & \cdots & (1-s_q)\mathcal{G}_{iq} \\ 0 & \cdots & 0 & | & 0 & & 0 & \cdots & 0 & || & 0 & & \cdots & 0 \\ \vdots & \ddots & \vdots & | & \vdots & & \vdots & \ddots & \vdots & || & \vdots & & \ddots & \vdots \\ 0 & \cdots & 0 & | & 0 & & 0 & \cdots & 0 & || & 0 & & \cdots & 0 \\ \hline 0 & \cdots & 0 & | & (1-s_i)\mathcal{G}_{(N+1)i} & & 0 & \cdots & 0 & || & (1-s_{N+1})\mathcal{G}_{(N+1)(N+1)} & & \cdots & (1-s_q)\mathcal{G}_{(N+1)q} \\ \vdots & \ddots & \vdots & | & \vdots & & \vdots & \ddots & \vdots & || & \vdots & & \ddots & \vdots \\ 0 & \cdots & 0 & | & (1-s_i)\mathcal{G}_{qi} & & 0 & \cdots & 0 & || & (1-s_{N+1})\mathcal{G}_{q(N+1)} & & \cdots & (1-s_q)\mathcal{G}_{qq} \end{pmatrix}.$$

We thus have

- As for the case  $N = 1$  treated in detail, we have that  $\lim_{s_i \rightarrow 0} \det(\zeta \mathbf{I} - \mathcal{A}_i)^{(0)} = 0$ . On the other hand, using the formula for the adjugate of an upper triangular matrix, see [22], we have  $\lim_{s_i \rightarrow 0} \text{adj}(\zeta \mathbf{I} - \mathcal{A}_i)^{(0)} = \Pi \mathbf{e}_i \otimes \mathbf{e}_i$ , where in this Section  $\Pi := \prod_{\ell=N+1}^{\ell=q} s_\ell$ .
- Taking  $\mathbf{C} = (\zeta \mathbf{I} - \mathcal{A}_i)^{(0)} \in \text{GL}_q(\mathbb{R}) \subset \text{GL}_q(\mathcal{S})$  and  $\mathbf{D} = \Delta x (\zeta \mathbf{I} - \mathcal{A}_i)^{(1)} + O(\Delta x^2) \in \mathcal{M}_q(\mathcal{S})$  in the Jacobi formula Equation (17)

$$\begin{aligned} \lim_{s_i \rightarrow 0} \text{D}_C(\det(\mathbf{C}))(\mathbf{D}) &= \lim_{s_i \rightarrow 0} s_i \Pi \left( (N-1) \partial_t + \frac{1}{s_i} (\partial_t + \mathcal{G}_{ii}) + \sum_{\ell=N+1}^q \frac{1}{s_\ell} (\partial_t + (1-s_\ell)\mathcal{G}_{\ell\ell}) \right) + O(\Delta x^2) \\ &= \Pi (\partial_t + \mathcal{G}_{ii}) + O(\Delta x^2). \end{aligned}$$

To handle the term with the adjugate, observe that the first-order term is made up of the terms

$$(\text{adj}(\zeta \mathbf{I} - \mathcal{A}_i) \mathcal{A}_i^\diamond)^{(1)} = (\text{adj}(\zeta \mathbf{I} - \mathcal{A}_i))^{(0)} (\mathcal{A}_i^\diamond)^{(1)} + (\text{adj}(\zeta \mathbf{I} - \mathcal{A}_i))^{(1)} (\mathcal{A}_i^\diamond)^{(0)}, \quad (37)$$

and in particular, we are interested in the  $i$ -th line of this matrix. Because of the fact that  $(\mathcal{A}_i^\diamond)^{(0)} = \text{diag}(1 - s_1, \dots, 1 - s_{i-1}, 0, 1 - s_{i+1}, \dots, 1 - s_N, 0, \dots, 0)$ , the  $i$ -th line of the second term on the right hand side of Equa-

tion (37) is zero, thus we do not have to study it. For the remaining term, it can be easily seen that

$$(\mathcal{A}_i^\diamond)^{(1)} = -\frac{1}{\lambda}(\mathbf{I} - \mathbf{S}) \begin{pmatrix} \mathcal{G}_{11} & \cdots & \mathcal{G}_{1(i-1)} & \mathcal{G}_{1i} & \mathcal{G}_{1(i+1)} & \cdots & \mathcal{G}_{1N} & \mathcal{G}_{1(N+1)} & \cdots & \mathcal{G}_{1q} \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \mathcal{G}_{(i-1)1} & \cdots & \mathcal{G}_{(i-1)(i-1)} & \mathcal{G}_{(i-1)i} & \mathcal{G}_{(i-1)(i+1)} & \cdots & \mathcal{G}_{(i-1)N} & \mathcal{G}_{(i-1)(N+1)} & \cdots & \mathcal{G}_{(i-1)q} \\ \mathcal{G}_{i1} & \cdots & \mathcal{G}_{ii} & 0 & \mathcal{G}_{i(i+1)} & \cdots & \mathcal{G}_{iN} & 0 & \cdots & 0 \\ \mathcal{G}_{(i+1)1} & \cdots & \mathcal{G}_{(i+1)(i-1)} & \mathcal{G}_{(i+1)i} & \mathcal{G}_{(i+1)(i+1)} & \cdots & \mathcal{G}_{(i+1)N} & \mathcal{G}_{(i+1)(N+1)} & \cdots & \mathcal{G}_{(i+1)q} \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \mathcal{G}_{N1} & \cdots & \mathcal{G}_{N(i-1)} & \mathcal{G}_{Ni} & \mathcal{G}_{N(i+1)} & \cdots & \mathcal{G}_{NN} & \mathcal{G}_{N(N+1)} & \cdots & \mathcal{G}_{Nq} \\ \mathcal{G}_{(N+1)1} & \cdots & \mathcal{G}_{(N+1)(i-1)} & 0 & \mathcal{G}_{(N+1)(i+1)} & \cdots & \mathcal{G}_{(N+1)N} & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \mathcal{G}_{q1} & \cdots & \mathcal{G}_{q(i-1)} & 0 & \mathcal{G}_{q(i+1)} & \cdots & \mathcal{G}_{qN} & 0 & \cdots & 0 \end{pmatrix},$$

thus we deduce that

$$((\text{adj}(\zeta \mathbf{I} - \mathcal{A}_i) \mathcal{A}_i^\diamond)^{(1)})_{i,\cdot} = -\Pi \frac{1}{\lambda} ((1 - s_1) \mathcal{G}_{i1}, \dots, (1 - s_{i-1}) \mathcal{G}_{i(i-1)}, 0, (1 - s_{i+1}) \mathcal{G}_{i(i+1)}, \dots, (1 - s_N) \mathcal{G}_{iN}, 0, \dots, 0).$$

Dealing with the zero and first order term in  $\text{adj}(\zeta \mathbf{I} - \mathcal{A}_i) \mathcal{B}$  works the same than  $N = 1$ , thus we do not repeat it. Moreover, these terms allow for the compensation of the dependence on the choice of the relaxation parameter of the other conserved moments  $s_1, \dots, s_{i-1}, s_{i+1}, \dots, s_N$  in the previous equation, as claimed in Section 3.4, thanks to Equation (1).

Putting all the previously discussed facts together into the truncated Equation (16) yields

$$\Pi \frac{\Delta x}{\lambda} \left( \partial_t m_i + \mathcal{G}_{ii} m_i + \sum_{\substack{j=1 \\ j \neq i}}^N \mathcal{G}_{ij} m_j + \sum_{j=N+1}^q \mathcal{G}_{ij} m_j^{\text{eq}} \right) = O(\Delta x^2),$$

which is the result from Theorem 4.4 for  $N \geq 1$ . Analogous reasonings yield Theorem 4.5.

## 7 Link with the existing approaches

To finish our contribution, we briefly sketch the links with previous works on the macroscopic equations like [46] and [13, 15]. A more complete study shall be the object of future investigations.

### 7.1 Equivalent equations

Our results Theorem 4.4 and Theorem 4.5 coincide with the analogous results in [15] up to second order. The substantial difference is that we apply the Taylor expansions to the solution of the corresponding Finite Difference scheme given either by Proposition 3.4 bis or Proposition 3.5 bis, where non-conserved moments have been removed. We therefore reasonably conjecture that the obtained macroscopic equations coincide at any order. The mathematical justification of this conjecture shall be the object of future investigations.

However, the quasi-equilibrium, which is extensively used in [15] can be somehow recovered in our previous discussion. Let  $N = 1$  to fix ideas. In the proof of Proposition 3.4 bis, nothing prevents us from selecting, instead of the first row, the  $i \in [2..q]$  row, corresponding to a non-conserved moment. This is

$$\det(\mathbf{zI} - \mathbf{A}) m_i = (\text{adj}(\mathbf{zI} - \mathbf{A}) \mathbf{B} \mathbf{m}^{\text{eq}})_i. \quad (38)$$

Let us stress that even if this could seem to be a viable Finite Difference scheme for the non-conserved variable  $m_i$ , it is not independent from the conserved moment  $m_1$  the equilibria depend on and furthermore, this formulation

certainly depends on the choice of  $s_1$ , the relaxation parameter of the conserved moment. This is somehow unwanted since  $s_1$  is *in fine* not present in the original lattice Boltzmann scheme. From the computations of Section 5, we see that

$$\det(\zeta \mathbf{I} - \mathcal{A}) = s_1 \Pi + O(\Delta x), \quad \text{adj}(\zeta \mathbf{I} - \mathcal{A}) = \Pi \text{diag}\left(1, \frac{s_1}{s_2}, \dots, \frac{s_1}{s_q}\right) + O(\Delta x), \quad \mathcal{B} = \mathbf{S} + O(\Delta x).$$

Using the asymptotic equivalents truncated at leading order in Equation (38) thus provides

$$s_1 \Pi m_i + O(\Delta x) = s_1 \Pi m_i^{\text{eq}} + O(\Delta x), \quad \text{hence also} \quad m_i = m_i^{\text{eq}} + O(\Delta x),$$

provided that  $s_1 \neq 0$ . This is the quasi-equilibrium of the non-conserved moments, which is re-injected in the lattice Boltzmann schemes to eliminate them in the procedure by [15]. However, in our framework, we are not really allowed to write Equation (38).

## 7.2 Maxwell iteration

In [46], the computations have been carried only for the  $D_2Q_9$  scheme by [31] with  $N = 3$ . In this part of our work, we are going to develop the computations until third-order for any lattice Boltzmann scheme under acoustic scaling, *i.e.* Assumptions 4.1. In this Section, it is crucial to assume that  $\mathbf{S} \in \text{GL}_q(\mathbb{R})$ . Observe that this assumption ensures that  $\det(\zeta \mathbf{I} - \mathcal{A})$  is a unit (invertible) in the ring  $\mathcal{S}$  or equivalently that  $\zeta \mathbf{I} - \mathcal{A}$  belongs to  $\text{GL}_q(\mathcal{S})$ . The Maxwell iteration [46] at step  $k \in \mathbb{N}$  reads, after simple computations

$$\mathbf{m}^{[k]} = \left( \sum_{r=0}^k (-\mathbf{S}^{-1}(\zeta \bar{\mathcal{T}} - \mathbf{I}))^r \right) \mathbf{m}^{\text{eq}}, \quad (39)$$

where the quasi-equilibrium is encoded in the choice  $\mathbf{m}^{[0]} = \mathbf{m}^{\text{eq}}$  and where we have taken, as for Equation (14)

$$\bar{\mathbf{T}} := \mathbf{M} \text{diag}(t_{-c_1}, \dots, t_{-c_q}) \mathbf{M}^{-1} \asymp \mathbf{M} \left( \sum_{|\mathbf{v}| \geq 0} \frac{\Delta x^{|\mathbf{v}|}}{\mathbf{v}!} \text{diag}(\mathbf{c}_1^{\mathbf{v}}, \dots, \mathbf{c}_q^{\mathbf{v}}) \partial^{\mathbf{v}} \right) \mathbf{M}^{-1} =: \bar{\mathcal{T}} \in \mathcal{M}_q(\mathcal{S}).$$

It is easy to see that  $\bar{\mathcal{T}} \bar{\mathcal{T}} = \bar{\mathcal{T}} \mathcal{T} = \mathbf{I}$  and moreover, in analogy with Lemma 5.3

$$\zeta \bar{\mathcal{T}} - \mathbf{I} = \frac{\Delta x}{\lambda} (\partial_t \mathbf{I} + \mathcal{G}) + \frac{\Delta x^2}{2\lambda^2} (\partial_{tt} \mathbf{I} + 2\mathcal{G} \partial_t + \mathcal{G}^2) + O(\Delta x^3). \quad (40)$$

The Maxwell iteration works by assuming that  $\mathbf{m} = \mathbf{m}^{[k]} + O(\Delta x^{k+1})$ . Taking  $k = 1$  in Equation (39) and using Equation (40), we have

$$\mathbf{m} = \mathbf{m}^{\text{eq}} - \mathbf{S}^{-1} \frac{\Delta x}{\lambda} (\partial_t \mathbf{I} + \mathcal{G}) \mathbf{m}^{\text{eq}} + O(\Delta x^2).$$

Let  $i \in [1..N]$ , then taking advantage of Equation (1)

$$m_i = m_i - \frac{\Delta x}{\lambda s_i} \left( \partial_t m_i + \sum_{j=1}^N \mathcal{G}_{ij} m_j + \sum_{j=N+1}^q \mathcal{G}_{1j} m_j^{\text{eq}} \right) + O(\Delta x^2),$$

which upon division, is the same result than Theorem 4.4. Going up to order two considering  $k = 2$ , we have

$$\begin{aligned} \mathbf{m} &= \mathbf{m}^{\text{eq}} - \mathbf{S}^{-1} \frac{\Delta x}{\lambda} (\partial_t \mathbf{I} + \mathcal{G}) \mathbf{m}^{\text{eq}} \\ &\quad + \frac{\Delta x^2}{2\lambda^2} \mathbf{S}^{-1} \left( (2\mathbf{S}^{-1} - \mathbf{I}) \partial_{tt} + 2(\mathbf{S}^{-1} \mathcal{G} + \mathcal{G} \mathbf{S}^{-1} - \mathcal{G}) \partial_t + \mathcal{G} (2\mathbf{S}^{-1} - \mathbf{I}) \mathcal{G} \right) \mathbf{m}^{\text{eq}} + O(\Delta x^3). \end{aligned}$$

Once more, selecting the  $i$ -th row provides

$$\begin{aligned} m_i &= m_i - \frac{\Delta x}{\lambda s_i} \left( \partial_t m_i + \gamma_{1,i} + \frac{\Delta x}{\lambda} \left( \left( \frac{1}{s_i} - \frac{1}{2} \right) \partial_{tt} m_i + \sum_{j=1}^q \left( \frac{1}{s_i} + \frac{1}{s_j} - 1 \right) \mathcal{G}_{ij} \partial_t m_j^{\text{eq}} + \sum_{j=1}^q \sum_{\ell=1}^q \left( \frac{1}{s_\ell} - \frac{1}{2} \right) \mathcal{G}_{i\ell} \mathcal{G}_{\ell j} m_j^{\text{eq}} \right) \right) \\ &\quad + O(\Delta x^3). \end{aligned}$$

Using relations analogous to Equation (33), Equation (34) and Equation (36) for  $N \geq 1$ , formally obtained by differentiating the result from Theorem 4.4, we finally obtain, after tedious but elementary computations

$$m_i = m_i - \frac{\Delta x}{\lambda s_i} \left( \partial_t m_i + \gamma_{1,i} + \frac{\Delta x}{\lambda} \sum_{j=N+1}^q \left( \frac{1}{s_j} - \frac{1}{2} \right) \mathcal{G}_{ij} \left( \sum_{\ell=1}^N \frac{dm_j^{\text{eq}}}{dm_\ell} \gamma_{1,\ell} - \sum_{\ell=1}^N \mathcal{G}_{j\ell} m_\ell - \sum_{\ell=N+1}^q \mathcal{G}_{j\ell} m_\ell^{\text{eq}} \right) \right) + O(\Delta x^3),$$

which coincides with the result from Theorem 4.5. Therefore, up to order two, our approach yields results consistent with those from [46].

To intuitively illustrate why it is reasonable to believe that we recover the same result at any order, let us assume  $N = 1$ . Then we have, using that  $\mathbf{S} \in \text{GL}_q(\mathbb{R})$ ,  $\mathcal{T}\overline{\mathcal{T}} = \overline{\mathcal{T}}\mathcal{T} = \mathbf{I}$ , the rule for the inverse of a product of matrix and the identity relative to geometric series in the context of formal power series, that

$$\begin{aligned} \mathbf{0} &= \det(\zeta \mathbf{I} - \mathcal{A}) \mathbf{m} - \text{adj}(\zeta \mathbf{I} - \mathcal{A}) \mathcal{B} \mathbf{m}^{\text{eq}} = \det(\zeta \mathbf{I} - \mathcal{A}) \left( \mathbf{m} - (\zeta \mathbf{I} - \mathcal{T}(\mathbf{I} - \mathbf{S}))^{-1} \mathcal{T} \mathbf{S} \mathbf{m}^{\text{eq}} \right), \\ &= \det(\zeta \mathbf{I} - \mathcal{A}) \left( \mathbf{m} - (\mathbf{S}^{-1} \overline{\mathcal{T}} (\zeta \mathbf{I} - \mathcal{T}(\mathbf{I} - \mathbf{S})))^{-1} \mathbf{m}^{\text{eq}} \right) = \det(\zeta \mathbf{I} - \mathcal{A}) \left( \mathbf{m} - (\mathbf{I} + \mathbf{S}^{-1} (\zeta \overline{\mathcal{T}} - \mathbf{I}))^{-1} \mathbf{m}^{\text{eq}} \right), \\ &= \det(\zeta \mathbf{I} - \mathcal{A}) \left( \mathbf{m} - \left( \sum_{r=0}^{+\infty} (-\mathbf{S}^{-1} (\zeta \overline{\mathcal{T}} - \mathbf{I}))^r \right) \mathbf{m}^{\text{eq}} \right) = \det(\zeta \mathbf{I} - \mathcal{A}) \left( \mathbf{m} - \lim_{k \rightarrow +\infty} \mathbf{m}^{[k]} \right). \end{aligned}$$

Therefore the expansion of the Finite Difference scheme from Proposition 3.4 bis and the non-truncated Maxwell iteration method on the lattice Boltzmann scheme coincide up to a multiplication by a formal power series of time-space differential operators, *i.e.*  $\det(\zeta \mathbf{I} - \mathcal{A}) \in \mathcal{S}$ . *A priori*, the resulting macroscopic equations are not the same, but since  $\det(\zeta \mathbf{I} - \mathcal{A}) = \det(\mathbf{S}) + O(\Delta x) = s_1 \Pi + O(\Delta x)$ , thus we “pay” only a constant factor we can divide by at dominant order, the macroscopic equations at leading order are the same. Then, at each order, the result must be the same because we re-inject, in a recursive fashion, the solution truncated at the previous order to eliminate the higher-order time derivatives, see for instance Equation (34) and Equation (36).

## 8 Conclusions and perspectives

In this original paper, we have rigorously derived the macroscopic equations up to order two for any lattice Boltzmann scheme under acoustic scaling by restating it as a multi-step macroscopic Finite Difference scheme on the conserved moments [3]. Since the passage from the kinetic to the macroscopic standpoint is fully discrete, our analysis can handle any type of time-space scaling and be pushed forward to reach higher orders of accuracy in the discretization parameters. Contrarily to the existing techniques, the quasi-equilibrium of the non-conserved moments in the limit of small discretization parameters is not the key to eliminate the non-conserved variables from the macroscopic equations. The obtained results confirm, going beyond empirical evidence, that the formal Taylor expansion by [13, 15] and the Maxwell iteration by [46] are well-grounded from the perspective of numerical analysts and traditional numerical methods for PDEs, such as Finite Difference.

An improvement of the present work could be the establishment of the equivalence between different consistency analyses for higher orders and ideally for any order. Even if more involved from the standpoint of computations, the extension can be easily done by considering derivatives of higher order for the determinant and adjugate functions, in the spirit of Lemma 5.4 and Lemma 5.6. In this work, all the computations have been done by hand but one could envision to seek some help from symbolic computations. This is a current path of investigation which final aim is to provide the computation – inside the package `pylbn`<sup>5</sup> – of the equivalent equations of any lattice Boltzmann scheme either by the corresponding Finite Difference scheme or using the Maxwell iteration. Another interesting track for future researches could be the re-state of our results using a parabolic scaling between time and space.

## Acknowledgments

The author deeply thanks his PhD advisors, M. Massot and B. Graille, for the fruitful discussions and advice on the subject, as well as his brother P. Bellotti for the useful tips to improve the style of manuscript. The author is supported by a PhD funding (year 2019) from the Ecole polytechnique.

<sup>5</sup><https://pylbn.readthedocs.io>



## References

- [1] ALLAIRE, G. *Numerical analysis and optimization: an introduction to mathematical modelling and numerical simulation*. Oxford University Press, 2007.
- [2] BELLOTTI, T., GOUARIN, L., GRAILLE, B., AND MASSOT, M. High accuracy analysis of adaptive multiresolution-based lattice Boltzmann schemes via the equivalent equations. *arXiv preprint arXiv:2105.13816, Submitted to SMAI Journal of Computational Mathematics* (2021).
- [3] BELLOTTI, T., GRAILLE, B., AND MASSOT, M. Finite difference formulation of any lattice Boltzmann scheme. *arXiv preprint arXiv:2201.05354* (2021).
- [4] BOUCHUT, F., GUARGUAGLINI, F. R., AND NATALINI, R. Diffusive BGK approximations for nonlinear multidimensional parabolic equations. *Indiana University Mathematics Journal* (2000), 723–749.
- [5] BREWER, J. W., BUNCE, J. W., AND VAN VLECK, F. S. *Linear systems over commutative rings*. CRC Press, 1986.
- [6] CAIAZZO, A., JUNK, M., AND RHEINLÄNDER, M. Comparison of analysis techniques for the lattice Boltzmann method. *Computers & Mathematics with Applications* 58, 5 (2009), 883–897.
- [7] CHAPMAN, S., AND COWLING, T. G. *The mathematical theory of non-uniform gases: an account of the kinetic theory of viscosity, thermal conduction and diffusion in gases*. Cambridge university press, 1990.
- [8] CHEN, S., AND DOOLEN, G. D. Lattice Boltzmann method for fluid flows. *Annual Review of Fluid Mechanics* 30, 1 (1998), 329–364.
- [9] CHENG, S. S. *Partial difference equations*, vol. 3. CRC Press, 2003.
- [10] DELLACHERIE, S. Construction and analysis of lattice Boltzmann methods applied to a 1D convection-diffusion equation. *Acta Applicandae Mathematicae* 131, 1 (2014), 69–140.
- [11] D’HUMIÈRES, D. *Generalized Lattice-Boltzmann Equations*. American Institute of Aeronautics and Astronautics, Inc., 1992, pp. 450–458.
- [12] DING, J., AND ZHOU, A. Eigenvalues of rank-one updated matrices with some applications. *Applied Mathematics Letters* 20, 12 (2007), 1223–1226.
- [13] DUBOIS, F. Equivalent partial differential equations of a lattice Boltzmann scheme. *Computers & Mathematics with Applications* 55, 7 (2008), 1441–1449.
- [14] DUBOIS, F. General third order Chapman-Enskog expansion of lattice Boltzmann schemes. In *16th International Conference for Mesoscopic Methods in Engineering and Science, Edinburgh, 22–26 July 2019*. (Edinburgh, United Kingdom, July 2019).
- [15] DUBOIS, F. Nonlinear fourth order Taylor expansion of lattice Boltzmann schemes. *Asymptotic Analysis*, 1 Jan. 2021 (2021), 1–41.
- [16] DUBOIS, F., AND LALLEMAND, P. Towards higher order lattice Boltzmann schemes. *Journal of Statistical Mechanics: Theory and Experiment* 2009, 06 (2009), P06006.
- [17] DUBOIS, F., AND LALLEMAND, P. Quartic parameters for acoustic applications of lattice Boltzmann scheme. *Computers & Mathematics with Applications* 61, 12 (2011), 3404–3416.
- [18] FUČÍK, R., AND STRAKA, R. Equivalent finite difference and partial differential equations for the lattice Boltzmann method. *Computers & Mathematics with Applications* 90 (2021), 96–103.
- [19] GUO, Z., AND SHU, C. *Lattice Boltzmann method and its application in engineering*, vol. 3. World Scientific, 2013.
- [20] HÉNON, M. Viscosity of a lattice gas. *Lattice Gas Methods for Partial Differential Equations* (1987), 179–207.

- [21] HIGUERA, F. J., AND JIMÉNEZ, J. Boltzmann approach to lattice gas simulations. *EPL (Europhysics Letters)* 9, 7 (1989), 663.
- [22] HORN, R. A., AND JOHNSON, C. R. *Matrix analysis*. Cambridge university press, 2012.
- [23] HUANG, K. *Statistical Mechanics*, 2 ed. John Wiley & Sons, 1987.
- [24] JOHNSON, W. P. The curious history of Faà di Bruno’s formula. *The American Mathematical Monthly* 109, 3 (2002), 217–234.
- [25] JUNK, M., KLAR, A., AND LUO, L.-S. Asymptotic analysis of the lattice Boltzmann equation. *Journal of Computational Physics* 210, 2 (2005), 676–704.
- [26] JUNK, M., AND YANG, Z. Convergence of lattice Boltzmann methods for Navier–Stokes flows in periodic and bounded domains. *Numerische Mathematik* 112, 1 (2009), 65–87.
- [27] JUNK, M., AND YONG, W.-A. Rigorous Navier–Stokes limit of the lattice Boltzmann equation. *Asymptotic Analysis* 35, 2 (2003), 165–185.
- [28] JURY, E. I. *Theory and Application of the z-Transform Method*. Krieger Publishing Co., 1964.
- [29] KASSEL, C. *Quantum Groups*, 1 ed. Graduate Texts in Mathematics. Springer-Verlag New York, 1995.
- [30] KRÜGER, T., KUSUMAATMAJA, H., KUZMIN, A., SHARDT, O., SILVA, G., AND VIGGEN, E. M. The lattice Boltzmann method. *Springer International Publishing* 10, 978-3 (2017).
- [31] LALLEMAND, P., AND LUO, L.-S. Theory of the lattice Boltzmann method: Dispersion, dissipation, isotropy, Galilean invariance, and stability. *Physical Review E* 61, 6 (2000), 6546.
- [32] LANG, S. *Algebra*, 3 ed. Graduate Texts in Mathematics. Springer-Verlag New York, 2002.
- [33] MCNAMARA, G. R., AND ZANETTI, G. Use of the Boltzmann equation to simulate lattice-gas automata. *Physical Review Letters* 61, 20 (1988), 2332.
- [34] MILLER, K. S. *An Introduction to the Calculus of Finite Differences and Difference Equations*. Dover Publications, 1960.
- [35] MONFORTE, A. A., AND KAUERS, M. Formal Laurent series in several variables. *Expositiones Mathematicae* 31, 4 (2013), 350–367.
- [36] NIVEN, I. Formal power series. *The American Mathematical Monthly* 76, 8 (1969), 871–889.
- [37] QIAN, Y.-H., AND ZHOU, Y. Higher-order dynamics in lattice-based models using the Chapman-Enskog method. *Physical Review E* 61, 2 (2000), 2103.
- [38] RHEINLÄNDER, M. K. *Analysis of lattice-Boltzmann methods: asymptotic and numeric investigation of a singularly perturbed system*. PhD thesis, 2007.
- [39] ROMAN, S. *The Umbral Calculus*. Dover Publications, 2005.
- [40] STEWART, G. On the adjugate matrix. *Linear Algebra and its Applications* 283, 1-3 (1998), 151–164.
- [41] STRIKWERDA, J. C. *Finite difference schemes and partial differential equations*. SIAM, 2004.
- [42] SUCCI, S. *The lattice Boltzmann equation: for fluid dynamics and beyond*. Oxford University Press, 2001.
- [43] SUGA, S. An accurate multi-level finite difference scheme for 1D diffusion equations derived from the lattice Boltzmann method. *Journal of Statistical Physics* 140, 3 (2010), 494–503.
- [44] VAN LEEMPUT, P., RHEINLÄNDER, M., AND JUNK, M. Smooth initialization of lattice Boltzmann schemes. *Computers & Mathematics with Applications* 58, 5 (2009), 867–882.

- [45] WARMING, R. F., AND HYETT, B. The modified equation approach to the stability and accuracy analysis of finite-difference methods. *Journal of Computational Physics* 14, 2 (1974), 159–179.
- [46] YONG, W.-A., ZHAO, W., LUO, L.-S., ET AL. Theory of the lattice Boltzmann method: Derivation of macroscopic equations via the Maxwell iteration. *Physical Review E* 93, 3 (2016), 033310.
- [47] ZHAO, W., AND YONG, W.-A. Maxwell iteration for the lattice Boltzmann method with diffusive scaling. *Physical Review E* 95, 3 (2017), 033311.
- [48] ZWILLINGER, D. *CRC standard mathematical tables and formulas*. Chapman and Hall - CRC, 2018.