

# Derivation and analysis of a class of relaxation operators in kinetic theory

Stéphane Brull<sup>1\*</sup>, Vincent Pavan<sup>2</sup> and Jacques Schneider<sup>3</sup>

<sup>1\*</sup>Institut de Mathématiques de Bordeaux, Bordeaux INP, Univ. Bordeaux, CNRS, 351, cours de la Libération, Talence, F-33400, France.

<sup>2</sup>IUSTI UMR 7343, Aix-Marseille Université, 5 rue Enrico Fermi,, Marseille, F-13013 Marseille, France.

<sup>3</sup>IMATH , Université de Toulon, Campus La Garde CS, Toulon, 83041, France.

\*Corresponding author(s). E-mail(s):

[stephane.brull@math.u-bordeaux.fr](mailto:stephane.brull@math.u-bordeaux.fr);

Contributing authors: [vincent.pavan@univ-amu.fr](mailto:vincent.pavan@univ-amu.fr);

[jacques.schneider@univ-tln.fr](mailto:jacques.schneider@univ-tln.fr);

## Abstract

We aim to present a theory for the derivation of a class of relaxation operators approximating the Boltzmann collision operator. The construction is based on an approximation of the inverse Boltzmann linearized operator, on relaxation equations on the moments of the distribution function and finally on a variational problem to be solved. The theory comprises a characterization of the set of moments of non negative integrable functions, a study of those linear application whose range lies in this set and a generalization of the functional to be minimized under moment constraints. In particular we clarify but also modify some steps in the proof of Junk's theorem on the characterization of moments of non negative functions [30]. We also reestablish a theorem of Csiszar's [20] by different means on a class of functionals leading to well-posed variational problems. The present theory encompasses the derivation of known models and that of new models.

**Keywords:** relaxation operator, truncated moment problem, positive polynomials,  $\phi$ -divergence

# 1 Introduction

Kinetic models are used to simulate rarefied gases in the context of atmospheric reentry, CVD deposition, micro-channels and other processes. The question arises whether the very fine description of the gas that is given by the original Boltzmann equation is required or not for such simulations. Depending on the Knudsen number, the collision operator may be replaced by simpler models that are easier to handle such as relaxation operators which seminal model is the BGK one [7]. Another reason for considering such operators is the difficulty to obtain the physical parameters that rule the interaction between molecules beyond the case of mono-atomic molecules. Contrarily, relaxation operators are often scaled by measurements obtained at the macroscopic level such as the diffusion coefficients. The idea that consists in taking some information available at this level is also used for the Boltzmann equation with continuous internal energy states - see e.g [10] - but this will not be our concern here. So, there exists plenty relaxation operators that are used in different context but no unified approach while many feature the same patterns. A linear or almost linear behavior with respect to the moments of the distribution function, most properties that are satisfied by the original equation such as positivity of solution or the existence of an entropy, and finally a correct hydrodynamic limit up to the Navier- Stokes equation.

In the present article, we aim to develop a theoretical approach that applies to existing models such as the BGK [7], ESBGK [29] or Shakhov model [42] and serves as a ground from which ongoing models can be constructed. The theory is presented only in the case of mono-atomic molecules for it requires much technical matters. The paper goes along the steps of what we name the method of moments relaxation which is a generalization of the work presented in [12, 13] together with some applications for modeling multicomponent gases (see for example [11, 14, 15, 41]). For short, the construction is based on relaxation equations that are relations between moments of the operator and moments of the distribution function with respect to a vector of weight functions  $\mathbf{m}$ .  $\mathbf{m}$  together with relaxation parameters are suitably chosen in order to obtain for example the right transport coefficients in the hydrodynamic limit. The relaxation equations are restated in term of linear relations between the moments of the probability density function  $f$  and those of the target function to be found. So the question arises whether those relations have a range into the set of realizable moments  $\mathcal{R}_{\mathbf{m}}^+$ , that is the set of vectors which are moments of nonnegative functions with respect to the weight functions  $\mathbf{m}$ . So we will characterize  $\mathcal{R}_{\mathbf{m}}^+$  and next specify which are the admissible relaxation equations. In one dimension and when the domain of velocity is  $\mathbb{R}$  and  $\mathbf{m} = \{1, \dots, v^{2N}\}^T$ , the characterization of  $\mathcal{R}_{\mathbf{m}}^+$  is known as the Hamburger moment problem. In this case, a moment  $\boldsymbol{\rho}$  is realizable w.r.t.  $\mathbf{m}$  iff the moment matrix  $(\rho_{i+j})_{i=0, \dots, N; j=0, \dots, N}$  is symmetric positive definite (SPD) (see for example [3]). In higher dimension which is the case of the usual velocity domain, the problem has given rise to many research (see for example [21, 23, 24, 27, 28, 33] and also the survey in [25]). Most results deal with moments of positive Borel measures and eventually of moments of atomic measures. At this moment, we must cite the contribution of Pichard [37, 38] to moments closure in kinetic theory for his analysis of such difficult results. So in dimension bigger than one, the characterization of realizability by positive Borel measure is rather abstract and do not lead to tractable

conditions such as in dimension 1. Nevertheless the main idea which consists in relating realizability with nonnegative polynomials provides more information when the measures are represented by integrable functions [30]. Then characterizing realizable moments amounts to characterizing nonnegative polynomials. This problem is known as the 17<sup>th</sup> Hilbert problem (see [6] for a survey).

The next problem is the way to define  $G$  once its moments are known. For most relaxation operators, it is usually done by minimizing the natural entropy under moment constraints. Unfortunately, this functional is not suited neither to the simple case of Grad thirteen moments nor to moments of higher degree. Junk was the first to raise and characterize the problem that one may face with Levermore's closure [30, 34]. He was then followed by different authors giving their own insight on the problem [26, 36, 40]. One way to bypass the problem was proposed by Abdelmalik and Van Brummelen in the context of moments closures [2] by approximating the relative entropy using  $\phi$ -divergence [18]. More precisely, they have proposed an approximation of the relative entropy that satisfies the condition established by Csiszar [20] in order to set a well-posed variational problem. That is a minimization problem which solution is unique and qualifies all moments constraints. The generalization to different approximations of the relative entropy by  $\phi$ -divergence can finally be found in [1].

This article is organized as follows. We display in section 2 a set of properties that should be satisfied by the relaxation operators we are going to construct. Then we present the method of moments relaxation together with the mathematical problems that are related with this construction. Section 3 is devoted to the characterization of realizable moments and of admissible linear relaxations. Firstly, we revisit and try to clarify Junk's theorem [30] on the characterization of  $\mathcal{R}_m^+$ . Secondly, we display the known results on the 17<sup>th</sup> Hilbert problems, how they can be used but also what is the limitation of the characterization of  $\mathcal{R}_m^+$  by using nonnegative polynomials. Thirdly we consider the case of Grad thirteen moments and fully describe admissible relaxations whose range are in  $\mathcal{R}_m^+$ . We close this section by a studying a class of linear operators that let  $\mathcal{R}_m^+$  invariant and that are related to Galilean invariance of the models to be constructed. Section 4 deals with the optimization problem. Since this problem is central in the derivation of BGK-type models but also of moments closure in kinetic theory, we study the existence theorems related to different  $\phi$ -divergences. In particular we focus on one theorem by Csiszar [20] and reestablish it by using convex analysis (see in particular [8]). This framework is more suited to derive the shape of the solution in the problems we are dealing with. In section 5, we show that the model which is constructed just basing on relaxations on the Grad thirteen moments is well-posed in the sense of section 2. We also address the general case and point out some problems to be solved. We then display some known models that can be derived in this framework just by using different functional in the variational problem [7, 29, 42]. We compare the present approach to Levermore's sum of relaxation operators [34]. Velocity dependent relaxation frequency models are not contained within the present approach even if they share a lot of common points with the present study [9, 35, 43, 44]. Finally, all the proofs are presented separately in section 6.

## 2 Method of moments relaxation

In this section, we present the general framework for the derivation of BGK type models that will be studied in this paper. We start by considering the Boltzmann equation

$$\partial_t f + \mathbf{v} \cdot \nabla_{\mathbf{x}} f = Q(f, f), \quad (1)$$

where  $Q(f, f)$  is the classical Boltzmann operator [16, 17]. Our aim is to construct a relaxation operator  $K(f)$  that may approximate  $Q(f, f)$  in some physical regime. In a first time we display some expected properties on  $K(f)$  that are inherited from those of  $Q(f, f)$  and in a second one we describe the framework for the derivation of a class of relaxation models and set some mathematical problems that will be addressed in the two next sections.

### 2.1 Well-defined operator

We expect that an operator of the form

$$K(f) = \nu(G - f),$$

satisfies properties that are similar to those of the Boltzmann collision operator. In some sense, such an operator is well-defined as soon as there hold the following properties.

1. Conservation laws :

$$\forall f, \int K(f) \phi(\mathbf{v}) d\mathbf{v} = 0 \Leftrightarrow \phi \in \mathbb{K},$$

with

$$\mathbb{K} = \text{span} \{ \mathbf{1}, \mathbf{v}, \mathbf{v}^2 \} \quad (2)$$

2. Preservation of positivity. Starting from a nonnegative initial condition  $f(0, \mathbf{x}, \mathbf{v}) \geq 0$  the solution must remain nonnegative and there must hold  $\nu(t, \mathbf{x}) > 0$  and  $G[f](t, \mathbf{x}, \mathbf{v}) \geq 0$  at any  $(t, \mathbf{x}, \mathbf{v}) \in \mathbb{R}^+ \times \mathbb{R}^3 \times \mathbb{R}^3$ .
3. There exists a strictly convex function  $\eta$  with  $\text{dom}(\eta) = \mathbb{R}^+$  (more condition on  $\eta$  are studied in section 4.1) such that the functional  $\mathcal{H}(f) = \int \eta(f) d\mathbf{v}$  satisfies

$$\int K(f) \eta'(f) d\mathbf{v} \leq 0. \quad (3)$$

Equation (3) must then be completed with the usual characterization of local equilibrium.

$$K(f) = 0 \Leftrightarrow \int \eta'(f) K(f) = 0 \Leftrightarrow \eta'(f) \in \mathbb{K}, \quad (4)$$

with

$$K(f) = 0 \Leftrightarrow \mathcal{M}(\mathbf{v}) = \frac{n}{(2\pi k_B T/m)^{3/2}} \exp\left(-\frac{m(\mathbf{v} - \mathbf{u})^2}{2k_B T}\right) \quad (5)$$

where  $n, \mathbf{u}, T$  are the density, velocity and temperature. We refer to the properties (3), (4) and (5) as the extended H-theorem or H-theorem in short.

4. Galilean invariance. For any translation:  $\tau\mathbf{v} = \mathbf{v} - \mathbf{u}$  and rotation  $\tau\mathbf{v} = \Theta\mathbf{v}$  in the velocity space there holds

$$K([\tau f]) = [\tau K(f)]$$

where by definition

$$\forall \mathbf{v}, \quad [\tau f](\mathbf{v}) := f(\tau\mathbf{v}).$$

5.  $\mathcal{M}$  being fixed one defines as usual the linearized operator  $\mathcal{L}$  with

$$\mathcal{L}(g) := \lim_{\epsilon \rightarrow 0^+} \frac{1}{\epsilon \mathcal{M}} K(\mathcal{M}(1 + \epsilon g)), \quad (6)$$

whose domain is  $L^2(\mathcal{M})$ . Then there must hold

- (a)  $\ker \mathcal{L} = \mathbb{K}$
- (b)  $\mathcal{L}$  is Fredholm, with pseudo-inverse  $\mathcal{L}^{-1}$  defined on the orthogonal of  $\ker \mathcal{L}$  in  $L^2(\mathcal{M})$ .
- (c) It is symmetrical negative on  $(\ker \mathcal{L})^\perp$
- (d) The viscosity and thermal conductivity computed from  $\mathcal{L}^{-1}$  must be the same as the ones derived from the Boltzmann equation (see (21)).

**Remark 1.** *The moment approach of relaxation operators is not suited to the Boltzmann entropy  $\mathcal{H}(f) = \int f \ln(f) d\mathbf{v}$  in most cases. We do not ask  $\mathcal{H}(f) = \int \eta(f) d\mathbf{v}$  to be a Lyapounov functional for the non homogeneous equation. The weakened property (3) that is required here enlarges the choice in  $\eta$ . This may give some stability if the variations of  $f$  due to the transport of particles are smaller than those due to the relaxation operator.*

## 2.2 Method of moments relaxation

Let us now consider the kinetic equation

$$\partial_t f + \mathbf{v} \cdot \nabla_{\mathbf{x}} f = K(f) \quad (7)$$

with

$$K(f) = \nu(G - f).$$

The question is: how to define  $\nu$  and  $G$  in such a way that the solution of equation (7) behaves as that of the Boltzmann equation (1). The range of validity in this comparison is understood in the following way.  $\mathcal{M}$  being a local equilibrium function,  $f = \mathcal{M} + \mathcal{M}g$  and the Boltzmann operator reads

$$Q(f, f) = \mathcal{M}\mathcal{L}_B(g) + Q(\mathcal{M}g, \mathcal{M}g),$$

where  $\mathcal{L}_B$  is the linearized Boltzmann operator. So we assume that  $Q(Mg, Mg)$  is negligible compared to  $\mathcal{M}\mathcal{L}_B(g)$  and we are looking for a relaxation operator  $K(f)$  approximating  $\mathcal{M}\mathcal{L}_B(g)$ .

### 2.2.1 Setting moments relaxation

The above approximation writes in the weak form

$$\begin{aligned} \int_{\mathbb{R}^3} K(f)\phi d\mathbf{v} &\approx \int_{\mathbb{R}^3} \mathcal{M}\mathcal{L}_B(g)\phi d\mathbf{v} = \int_{\mathbb{R}^3} \mathcal{M}g\mathcal{L}_B(\phi) d\mathbf{v} \\ &= \int_{\mathbb{R}^3} (\mathcal{M} + \mathcal{M}g)\mathcal{L}_B(\phi) d\mathbf{v} = \int_{\mathbb{R}^3} f\mathcal{L}_B(\phi) d\mathbf{v}. \end{aligned}$$

This approximation is of course impossible if one considers all test functions  $\phi$  in  $L^2(\mathcal{M})$  since this would imply  $K(f) = \mathcal{M}\mathcal{L}_B(g)$ . So we must restrict ourselves to a space of finite dimension  $\mathbb{P}$ . In the case of Maxwell molecules, one just has to consider the space  $\mathbb{P}$  spanned by the  $q$  first eigenfunctions of  $\mathcal{L}_B : (m_i(\mathbf{v}))_{i=1, \dots, q}$ . One then sets the following relaxation equations

$$\int_{\mathbb{R}^3} K(f) m_i d\mathbf{v} = \int_{\mathbb{R}^3} f\mathcal{L}_B(m_i) d\mathbf{v} = -\nu_i \int_{\mathbb{R}^3} f m_i d\mathbf{v}. \quad (8)$$

Such relations are constitutive properties of the models to be constructed. Their generalization to other types of molecular interaction can be done in two ways. Let  $\mathbb{P}$  be a (polynomial) space of dimension  $q$  with  $\mathbb{K} \subset \mathbb{P}$  and  $\mathcal{P}$  be the projection onto  $\mathbb{P}$  in  $L^2(\mathcal{M})$ . One may then either approximate the restriction of  $\mathcal{L}_B$  on  $\mathbb{P}$  or that of  $\mathcal{L}_B^{-1}$  on  $\mathbb{P} \cap \mathbb{K}^\perp$ . In the first case, one considers the linear operator  $\tilde{\mathcal{L}}$  which restriction to  $\mathbb{P}$  is  $\mathcal{P}\mathcal{L}_B\mathcal{P}$ . With this, the restriction of  $\tilde{\mathcal{L}}$  on  $\mathbb{P}$  has a range equal to  $\mathbb{P}$  and this restriction is self-adjoint. Thus, there exists an orthogonal basis  $(\tilde{m}_i(\mathbf{v}))_{i=1, \dots, q}$  of  $\mathbb{P}$  such that

$$\forall g \in \mathbb{P}, \quad \tilde{\mathcal{L}}(g) = \mathcal{P}\mathcal{L}_B\mathcal{P}(g) = -\sum_{i=1}^q \tilde{\nu}_i \mathcal{P}\tilde{m}_i(g), \quad (9)$$

where  $\mathcal{P}_{\tilde{m}_i}$  denotes the orthogonal projection onto the eigenfunction  $\tilde{m}_i$  in  $L^2(\mathcal{M})$  with  $\tilde{\nu}_i = 0$  for  $\tilde{m}_i \in \mathbb{K}$  and  $\tilde{\nu}_i > 0$  in  $\mathbb{K}^\perp$ . Notice that this approximation corresponds to the one performed in [16] for the linearized and linear Boltzmann equation. One then replace  $\mathcal{L}_B$  with  $\tilde{\mathcal{L}}$  in (8). Unfortunately such a direct approximation of  $\mathcal{L}_B$  does not give the right transport coefficients in the hydrodynamic limit as will be shown in the example below. Instead, the idea consists in approximating  $\mathcal{L}_B^{-1}$  on  $\mathbb{P} \cap \mathbb{K}^\perp$  by stating

$$\forall g \in \mathbb{P} \cap \mathbb{K}^\perp, \quad \mathcal{L}^{-1}(g) = \mathcal{P}_{\mathbb{K}^\perp} \mathcal{L}_B^{-1} \mathcal{P}_{\mathbb{K}^\perp}(g) = -\sum_{i=6}^q \nu_i^{-1} \mathcal{P}_{m_i}(g), \quad (10)$$

while letting the kernel of the operator being  $\mathbb{K}$ . Here  $\mathcal{P}_{\mathbb{K}^\perp}$  is the restriction of  $\mathcal{P}$  to  $\mathbb{K}^\perp$ ,  $\mathcal{P}_{m_i}$  is the projection onto the eigenfunction  $m_i$  for all  $i = 6, \dots, q$  and the

eigenvalues  $(-\nu_i^{-1})_{i=6,\dots,q}$  are strictly negative. Thus one has

$$\text{Ker}\mathcal{L} = \mathbb{K}, \quad (11)$$

$$\forall g \in \mathbb{P} \cap \mathbb{K}^\perp, \quad \mathcal{L}(g) = (\mathcal{P}_{\mathbb{K}^\perp} \mathcal{L}_B^{-1} \mathcal{P}_{\mathbb{K}^\perp})^{-1}(g) = - \sum_{i=6}^q \nu_i \mathcal{P}_{m_i}(g) \quad (12)$$

which we plug again into (8) and obtain the relaxation equations

$$\int_{\mathbb{R}^3} K(f) (1, \mathbf{v}, |\mathbf{v}|^2) d\mathbf{v} = 0 \quad (13)$$

and for all eigenfunctions  $m_i \in \mathbb{P} \cap \mathbb{K}^\perp$ ,

$$\int_{\mathbb{R}^3} K(f) m_i(\mathbf{v}) d\mathbf{v} = -\nu_i \int_{\mathbb{R}^3} f m_i(\mathbf{v}) d\mathbf{v}. \quad (14)$$

It must be noted that  $\mathcal{L}^{-1}$  is not the pseudo-inverse of  $\tilde{\mathcal{L}}$  onto  $\mathbb{K}^\perp$  except in the case of Maxwell molecules.

At the moment, it is not clear that a BGK-type model satisfying (13) and (14) will provide us with the right linearized operator in the hydrodynamic limit and one formally proves the following proposition.

**Proposition 1.** *Assume (4) and (5) hold. Assume moreover that  $f \rightarrow G(f)$  is smooth, then if  $K(f)$  satisfies (13, 14), its linearized operator defined in (6) is*

$$\mathcal{L} = \nu \left( (\mathcal{P}_{\mathbb{K}} - \mathcal{I}) + \sum_i \left( 1 - \frac{\nu_i}{\nu} \right) \mathcal{P}_{m_i} \right). \quad (15)$$

*Details of the proof are let to the last section of this article in the framework of the example below - proposition 19 - which can be easily extended to the general case.*

The above considerations are finally made explicit in the most meaningful case of relaxation on the Grad thirteen moments.

**Example 1.** *For  $\mathcal{M}$  defined in (5) we consider the polynomial space*

$$\mathbb{P} = \mathbb{K} \oplus^\perp (\mathbb{A})_{i,j} \oplus^\perp (\mathbf{b})_i, \quad i, j = 1, 2, 3 \quad (16)$$

where  $\mathbb{A}$  and  $\mathbf{b}$  are the Sonine "polynomials"

$$\mathbb{A}(\mathbf{v} - \mathbf{u}) = \frac{m}{k_B T} \left[ (\mathbf{v} - \mathbf{u}) \otimes (\mathbf{v} - \mathbf{u}) - \frac{1}{3} \|\mathbf{v} - \mathbf{u}\|^2 \mathbb{I} \right], \quad (17)$$

$$\mathbf{b}(\mathbf{v} - \mathbf{u}) = (\mathbf{v} - \mathbf{u}) \left[ \frac{1}{2} m (\mathbf{v} - \mathbf{u})^2 - \frac{5}{2} k_B T \right], \quad (18)$$

and orthogonality in (16) holds for the usual scalar product in  $L^2(\mathcal{M})$ . Then there exists functions  $a(|\mathbf{V}|, T) > 0$  and  $b(|\mathbf{V}|, T) > 0$  with  $\mathbf{V} = (\mathbf{v} - \mathbf{u}) / \sqrt{k_B T / m}$  (see e.g

[17]) such that

$$\mathcal{L}_B^{-1}(\mathbb{A}) = -a(|\mathbf{V}|, T)\mathbb{A}, \quad \mathcal{L}_B^{-1}(\mathbf{b}) = -b(|\mathbf{V}|, T)\mathbf{b} \quad (19)$$

so that  $\mathcal{L}_B^{-1}(\mathbb{A}) \perp \mathbf{b}$ ,  $\mathcal{L}_B^{-1}(\mathbf{b}) \perp \mathbb{A}$ . As a consequence  $\mathcal{L}^{-1}$  in (10) satisfies

$$\mathcal{L}^{-1}(\mathbb{A}) = -\nu_{\mathbb{A}}^{-1}\mathbb{A}, \quad \mathcal{L}^{-1}(\mathbf{b}) = -\nu_{\mathbf{b}}^{-1}\mathbf{b}, \quad (20)$$

for some positive values  $\nu_{\mathbb{A}}$  and  $\nu_{\mathbf{b}}$ . Those values are then related to the viscosity  $\mu_B$  and the thermal conductivity  $\kappa_B$  obtained in the Navier-Stokes limit of the Boltzmann equation

$$\mu_B = -\frac{k_B T}{10} \langle \mathcal{L}_B^{-1}(\mathbb{A}), \mathbb{A} \rangle = -\frac{k_B T}{10} \langle \mathcal{L}^{-1}(\mathbb{A}), \mathbb{A} \rangle = \frac{nk_B T}{\nu_{\mathbb{A}}}, \quad (21)$$

$$\kappa_B = -\frac{1}{3k_B T^2} \langle \mathcal{L}_B^{-1}(\mathbf{b}), \mathbf{b} \rangle = -\frac{1}{3k_B T^2} \langle \mathcal{L}^{-1}(\mathbf{b}), \mathbf{b} \rangle = \frac{5}{2} \frac{nk_B^2 T}{m\nu_{\mathbf{b}}}, \quad (22)$$

where the scalar product is extended to tensors by performing full contraction. With those eigenfunctions and eigenvalues at hand one set two relaxation constraints in addition to the conservation laws (13)

$$\int_{\mathbb{R}^3} \nu(G - f)\mathbb{A} \, d\mathbf{v} = -\nu_{\mathbb{A}} \int_{\mathbb{R}^3} f\mathbb{A} \, d\mathbf{v}, \quad (23)$$

$$\int_{\mathbb{R}^3} \nu(G - f)\mathbf{b} \, d\mathbf{v} = -\nu_{\mathbf{b}} \int_{\mathbb{R}^3} f\mathbf{b} \, d\mathbf{v}. \quad (24)$$

Under the assumptions of proposition 1, one can perform a Chapman-Enskog expansion in

$$\partial_t f + \mathbf{v} \cdot \nabla_{\mathbf{x}} f = \frac{1}{\varepsilon} K(f). \quad (25)$$

$f$  then writes as  $f = \mathcal{M} + \varepsilon \mathcal{M}g + \mathcal{O}(\varepsilon^2)$ . This provides us with the Euler equation up to an order  $\mathcal{O}(\varepsilon)$  while at the next order  $g$  satisfies

$$\mathcal{L}(g) = \mathbb{A} : \mathbb{D}(\mathbf{u}) + \mathbf{b} \cdot \nabla_{\mathbf{x}} \left( -\frac{1}{k_B T} \right), \quad (26)$$

where  $\mathbb{D}(\mathbf{u})$  is the Reynolds tensor

$$\mathbb{D}(\mathbf{u}) = \left[ \nabla_{\mathbf{x}} \mathbf{u} + (\nabla_{\mathbf{x}} \mathbf{u})^T \right] - \frac{2}{3} (\nabla_{\mathbf{x}} \cdot \mathbf{u}) \mathbb{I}$$

Proposition 1 allows to solve (26) and one finds that the corresponding diffusion coefficients are those given in (21). Contrarily, if one considers the operator defined in (9) the diffusion coefficients in the Navier-Stokes equations are not correct. One first



remarks that  $\langle \mathcal{L}(\mathbb{A}), \mathbf{b} \rangle = 0$  thanks to the even/odd symmetry. Thus  $\mathbb{A}$  and  $\mathbf{b}$  are eigenfunctions of  $\tilde{\mathcal{L}}$  with corresponding eigenvalues  $-\tilde{\nu}_{\mathbb{A}}$  and  $-\tilde{\nu}_{\mathbf{b}}$  defined by

$$\tilde{\nu}_{\mathbb{A}} = -\frac{\langle \mathcal{L}_B(\mathbb{A}), \mathbb{A} \rangle}{\langle \mathbb{A}, \mathbb{A} \rangle}, \quad \tilde{\nu}_{\mathbf{b}} = -\frac{\langle \mathcal{L}_B(\mathbf{b}), \mathbf{b} \rangle}{\langle \mathbf{b}, \mathbf{b} \rangle}.$$

So, in the Chapman-Enskog expansion, the Navier-Stokes equations are obtained with the following diffusion coefficients

$$\begin{aligned} \mu_K &= \frac{nk_B T}{\tilde{\nu}_{\mathbb{A}}} = -nk_B T \frac{\langle \mathbb{A}, \mathbb{A} \rangle}{\langle \mathcal{L}_B(\mathbb{A}), \mathbb{A} \rangle}, \\ \kappa_K &= \frac{5}{2} \frac{nk_B^2 T}{m\tilde{\nu}_{\mathbf{b}}} = -\frac{5}{2} \frac{nk_B^2 T}{m} \frac{\langle \mathbf{b}, \mathbf{b} \rangle}{\langle \mathcal{L}_B(\mathbf{b}), \mathbf{b} \rangle}. \end{aligned}$$

By comparison with (21) one finds that  $\mu_K = \mu_B$  and  $\kappa_K = \kappa_B$  for Maxwell molecules. In the other cases,  $\mu_K = \mu_B$  reads also  $\langle \mathcal{L}_B^{-1}(\mathbb{A}), \mathbb{A} \rangle = C(n, T) \langle \mathcal{L}_B(\mathbb{A}), \mathbb{A} \rangle^{-1}$  where  $C(n, T)$  does not depend on the interaction potential which is untrue (and likewise for the heat conductivity). Remark finally that such a problem arises in the hydrodynamic limit of moment system for the Boltzmann equation [34]. It is then found that  $\mu_K < \mu_B$  and  $\kappa_K < \kappa_B$ .

### 2.2.2 Definition of $\nu$ and $G$ and related mathematical problems

Remark that the considerations in the above section do not require  $K(f)$  to be a relaxation operator as soon as the solution to (7) exists and  $K(f)$  satisfies the properties displayed in section 2.1 except property 2. But the natural way to satisfy this property is to let  $K(f)$  be a relaxation operator with  $\nu > 0$  and  $G \geq 0$ . Such conditions are part of the following mathematical problems we are going to study :

1. **Definition of  $\nu$**  : the set of constraints writes

$$\int_{\mathbb{R}^3} G m_i d\mathbf{v} = \left(1 - \frac{\nu_i}{\nu}\right) \int_{\mathbb{R}^3} f m_i d\mathbf{v}, \quad i \in \{1, \dots, q\}, \quad (27)$$

The relation (27) can be rewritten in a compact shape

$$\boldsymbol{\rho}_G := L(\boldsymbol{\rho}_f) \quad (28)$$

where the vectors of moments  $\boldsymbol{\rho}_G$  and  $\boldsymbol{\rho}_f$  are defined with

$$\boldsymbol{\rho}_G = \int_{\mathbb{R}^3} G \mathbf{m}(\mathbf{v}) d\mathbf{v}, \quad \boldsymbol{\rho}_f = \int_{\mathbb{R}^3} f \mathbf{m}(\mathbf{v}) d\mathbf{v} \quad (29)$$

with  $\mathbf{m} := (m_1, \dots, m_k, \dots, m_q)^T$ . If one assumes that the moments  $\boldsymbol{\rho}_f$  of the nonnegative function are bounded, then the moments of  $G$  must be those of a

nonnegative function as well. That is

$$C_f = \left\{ h \geq 0, \int h(\mathbf{v}) \mathbf{m}(\mathbf{v}) d\mathbf{v} = L(\boldsymbol{\rho}_f) \right\} \neq \emptyset. \quad (30)$$

Clearly, this condition just relies on the value of  $\nu$  since the relation frequencies are defined in a unique way according to (10). In order to state this question more clearly, let us give some few definitions.

**Definition 1.** We note with  $\mathbb{L}^1(\mathbf{m})$  the set of integrable functions  $f : \mathbb{R}^3 \mapsto \mathbb{R}$  such that

$$\forall k \in [1, q], \int \|\mathbf{m}_k(\mathbf{v})\| |f(\mathbf{v})| d\mathbf{v} < +\infty.$$

When  $f \in \mathbb{L}^1(\mathbf{m})$  is non negative, we note  $f \in \mathbb{L}^{1,+}(\mathbf{m})$ , and when it is not zero we note  $f \in \mathbb{L}^{1,*,+}(\mathbf{m})$ .

Now we introduce the moment map

**Definition 2** (Realizable moments). Let  $R : \mathbb{L}^1(\mathbf{m}) \mapsto \mathbb{R}^q$  be defined as

$$\forall f \in \mathbb{L}^1(\mathbf{m}), \quad R[f] = \int \mathbf{m}(\mathbf{v}) f(\mathbf{v}) d\mathbf{v}$$

we adopt for the sequel the following notations

$$\mathcal{R}_{\mathbf{m}}^+ = \{R[f], f \in \mathbb{L}^{1,+}(\mathbf{m})\}, \quad \mathcal{R}_{\mathbf{m}}^{*,+} = \{R[f], f \in \mathbb{L}^{1,*,+}(\mathbf{m})\}. \quad (31)$$

$\mathcal{R}_{\mathbf{m}}^+$  is named the set of realizable moments.

Coming back to the definition of  $\nu$ , one needs to have a tractable way to characterize  $\mathcal{R}_{\mathbf{m}}^+$  and to study the range of  $L$  as a linear operator depending on  $\nu$ . In other words, there must hold

$$L(\mathcal{R}_{\mathbf{m}}^+) \subset \mathcal{R}_{\mathbf{m}}^+. \quad (32)$$

Those questions will be studied in the next section.

2. **Definition of  $G(f)$ :** Assume that

$$C_f = \left\{ h \geq 0, \int h(\mathbf{v}) \mathbf{m}(\mathbf{v}) d\mathbf{v} = L(\boldsymbol{\rho}_f) \right\} \neq \emptyset. \quad (33)$$

The final step of the method is performed by solving a variational problem and defining  $G$  as

$$G = \arg \min_{h \in C_f} \int \eta(h) d\mathbf{v}, \quad (34)$$

where  $\eta$  is a strictly convex function with  $\text{dom}(h) = \mathbb{R}^+$  (more suitable conditions will be given in section 4). Statistical physics states that  $\eta$  must be defined as  $\eta_B(\cdot) = \cdot \ln(\cdot)$ . Unfortunately theoretical and numerical studies have revealed

numerous problems as soon as there are more constraints than (13) and (23) [30]. In particular there is no solution to the variational problem under the constraints (13, 14) in the example of Grad relaxation on thirteen moments. Adding more constraints neither provides us with a well-posed variational problem in general. This explains why other choices of  $\eta$  must be found and why we only ask for some local stability through (3), (4) and (5). In section 4, we will address this problem in the general framework of  $\phi$ -divergence [18] following the path in [2]. The well-posedness of the related variational problem will be studied in details by recalling known results and also establishing a Theorem by means of convex analysis.

### 3 The moment problem

Our concern is now to define a relaxation frequency  $\nu$  such that (33) holds for  $f \in \mathbb{L}^{1,+}(\mathbf{m})$  (see definition 2). This requires on one hand to have a tractable criteria to characterize the set of realizable moments  $\mathcal{R}_{\mathbf{m}}^+$  (31) and on the other hand to study the range of  $L$  (28). The characterization of  $\mathcal{R}_{\mathbf{m}}^+$  is part of a broader program which addresses the question of realizability by nonnegative Borel measures. This problem is known as the truncated moment problem and we recall in section 3.1 the most significant results on this topic. Since those characterization fail to be practical but also too general in our context, we focus in section 3.2 on the specific case where nonnegative Borel measures are represented by nonnegative integrable functions. In this case, we recall a theorem by Junk [30], clarify and somehow complete his proof. We then derive in section 3.3 a tractable necessary and sometime sufficient condition for a moment to be in  $\mathcal{R}_{\mathbf{m}}^+$  for a certain class of polynomial space  $\mathbb{P} = \text{span}(m_1, \dots, m_q)$ . In section 3.4 This condition is used in the case of Grad relaxation (example 1) to define the set of admissible relaxation frequencies  $\nu$ , i.e those frequencies for which (33) holds. Finally, we study in section 3.5 the relation between  $\mathcal{R}_{\mathbf{m}}^+$  and Galilean invariance.

#### 3.1 The truncated moment problem

There exists an extensive literature dedicated to the realizability of moments by nonnegative Borel measure. Here we want to recall the most significant ones. In order to relate them with the previous section we let  $\mathbb{P} = \mathbb{P}_{2n}$  with  $q = \dim(\mathbb{P}_{2n})$ . We also consider the general case where  $\mathbf{v} \in \mathbb{R}^d$ .  $\mathbf{i}$  being a multiindex in  $\mathbb{N}^d$ , with  $|\mathbf{i}| = i_1 + i_2 + \dots + i_d \leq 2n$ , one considers the d-dimensional real multisequence  $\boldsymbol{\rho} = \{\rho_{\mathbf{i}}\}$  ( $\boldsymbol{\rho}$  is just another notation for a vector in  $\mathbb{R}^q$ ). The so-called truncated moments problem asks for condition on  $\boldsymbol{\rho}$  such that there exists a nonnegative Borel measure satisfying

$$\rho_{\mathbf{i}} = \int \mathbf{v}^{\mathbf{i}} d\mu, \mathbf{i} \in \mathbb{N}^d,$$

with the standard notation  $\mathbf{v}^{\mathbf{i}} = v_1^{i_1} \dots v_d^{i_d}$ .

**Definition 3.** For a polynomial  $P(\mathbf{v}) = \sum_{\mathbf{i} \in \mathbb{N}^d} \alpha_{\mathbf{i}} \mathbf{v}^{\mathbf{i}}$  one defines the Riesz functional  $L_{\boldsymbol{\rho}}(P)$  at  $P$  with

$$L_{\boldsymbol{\rho}}(P) = \sum_{\mathbf{i} \in \mathbb{N}^d} \rho_{\mathbf{i}} \alpha_{\mathbf{i}}.$$

It is positive when  $L_{\boldsymbol{\rho}}(P) \geq 0$  for all  $P \geq 0$  in  $\mathbb{P}_{2n}$ .

**Theorem 2.** (Analogue of the Riesz-Haviland theorem [24]) Let  $\boldsymbol{\rho} = \boldsymbol{\rho}^{2n}$  be a  $d$ -dimensional multisequence, then  $\boldsymbol{\rho}$  has a representing measure if and only if  $\boldsymbol{\rho}$  admits an extension  $\tilde{\boldsymbol{\rho}} = \tilde{\boldsymbol{\rho}}^{2n+2}$  such that  $L_{\tilde{\boldsymbol{\rho}}}$  is positive.

As we will see in the next section, of particular interest are the measures that can be represented by atomic measure. One then associates to a  $d$ -dimensional multisequence  $\boldsymbol{\rho} = \boldsymbol{\rho}^{2n}$  a moment matrix  $M_n = M_n(\boldsymbol{\rho})$  in the following way. Considering again a polynomial  $P(\mathbf{v}) = \sum_{\mathbf{i} \in \mathbb{N}^d} \alpha_{\mathbf{i}} \mathbf{v}^{\mathbf{i}}$ ,  $(\alpha_{\mathbf{i}})$  denotes the coefficient vector of  $P$  relative to the basis of  $\mathbb{P}_{2n}$  of monomials in degree-lexicographic order. The rows and columns of  $M_n$  are then indexed by the elements  $\mathbf{v}^{\mathbf{i}}$  of  $\mathbb{P}_{2n}$  so that the entry of  $M_n$  in row  $X^{\mathbf{i}}$  and column  $X^{\mathbf{j}}$  is  $\rho_{\mathbf{i}+\mathbf{j}}$ .

**Theorem 3.** (Flat extension theorem [25])  $\boldsymbol{\rho} = \boldsymbol{\rho}^{2n}$  has a rank  $M_n$ -atomic representing measure if only if  $M_n$  is semi definite positive and  $M_n$  admits a flat extension i.e. a moment matrix extension  $M_{n+1}$  satisfying  $\text{rank}(M_{n+1}) = \text{rank}(M_n)$ . In this case,  $\boldsymbol{\rho}^{2n+2}$  admits a unique representing measure  $\mu = \mu_{M_{n+1}}$ , satisfying  $\text{supp } \mu = V(M_{n+1})$  and  $\text{card } \text{supp } \mu = \text{rank}(M_n)$ . Further,  $M_{n+1}$  admits a unique successive positive moment matrix extensions  $M_{n+2}, M_{n+3}, \dots$ , and these are flat extensions.

**Remark 2.** Clearly, none of the above theorems provides us with a criteria of practical use to know whether a moment  $\boldsymbol{\rho}$  admits a representing nonnegative Borel measure since extensions are not explicit. More generally and as can be seen for example in [25], the theorems addressing this problem always give abstract characterizations by using either the Riesz functional or the moment matrices. Nevertheless, as we will see in the next sections, this is not the case when the problem is restricted to realizability by nonnegative functions.

### 3.2 Characterization of $\mathcal{R}_{\mathbf{m}}^+$

We again assume that  $\mathbf{m}^T := (m_1, \dots, m_k, \dots, m_q)$  is a basis of a polynomial space  $\mathbb{P}$ , where  $\mathbb{P}$  is not necessarily  $\mathbb{P}_{2n}$ . The problem of realizability is now stated as follows : for a given vector  $\boldsymbol{\rho} \in \mathbb{R}^q$ , is there an integrable non negative function  $f : \mathbb{R}^3 \mapsto \mathbb{R}^+$  such that

$$\forall k \in [1, q], \quad \int m_k(\mathbf{v}) f(\mathbf{v}) d\mathbf{v} = \rho_k. \quad (35)$$

$\boldsymbol{\rho} := (\rho_1, \dots, \rho_q)$  is said realizable when (35) holds. Then there is the following theorem

**Theorem 4.** (Junk [30]) A vector  $\boldsymbol{\rho} \in \mathcal{R}_{\mathbf{m}}^{*+}$  iff for all  $\boldsymbol{\alpha} \neq 0 \in \mathbb{R}^q$  which satisfies  $\boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v}) \leq 0$  on  $\mathbb{R}^3$ , the relation  $\boldsymbol{\alpha} \cdot \boldsymbol{\rho} < 0$  holds.

In particular,  $\mathcal{R}_{\mathbf{m}}^{*+}$  is an open set. Moreover each  $\boldsymbol{\rho} \in \mathcal{R}_{\mathbf{m}}^{*+}$  is a moment vector of bounded  $f \in \mathbb{L}^{1,*}(\mathbf{m})$  which is compactly supported.

The proof relies essentially on a duality argument. The right implication is quite obvious. Consider indeed a moment  $\boldsymbol{\rho} \in \mathcal{R}_{\mathbf{m}}^{*+}$ , that is a vector such that there exists a

nonnegative and non null function  $f \in \mathbb{L}^{1,+}(\mathbf{m})$  satisfying  $\boldsymbol{\rho} = \int \mathbf{m}(\mathbf{v}) f(\mathbf{v}) d\mathbf{v}$ . Let  $\boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v})$  be a non positive and non null polynomial, then

$$\boldsymbol{\alpha} \cdot \boldsymbol{\rho} = \boldsymbol{\alpha} \cdot \int \mathbf{m}(\mathbf{v}) f(\mathbf{v}) d\mathbf{v} = \int \boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v}) f(\mathbf{v}) d\mathbf{v} < 0 \quad (36)$$

since  $f \neq 0$  on a set of non zero measure and  $\boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v}) \neq 0$  almost everywhere.  $\mathcal{R}_{\mathbf{m}}^+$  is a positive cone and the above computation proves that its polar cone is the set of coefficients of non positive polynomials :

$$\begin{aligned} (\mathcal{R}_{\mathbf{m}}^+)^{\circ} &= \{ \boldsymbol{\alpha} \in \mathbb{R}^q, \forall \boldsymbol{\rho} \in \mathcal{R}_{\mathbf{m}}^+, \boldsymbol{\alpha} \cdot \boldsymbol{\rho} \leq 0 \} \\ &= \left\{ \boldsymbol{\alpha} \in \mathbb{R}^q, \forall f \in \mathbb{L}^{1,+}(\mathbf{m}), \int \boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v}) f(\mathbf{v}) d\mathbf{v} \leq 0 \right\} \\ &= \{ \boldsymbol{\alpha} \in \mathbb{R}^q, \boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v}) \leq 0 \}. \end{aligned}$$

Note by  $C^{\circ}$  the set  $\{ \boldsymbol{\alpha} \in \mathbb{R}^q, \boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v}) \leq 0 \}$ . If the convex cone  $\mathcal{R}_{\mathbf{m}}^+$  is of non empty interior, one finds that  $\mathcal{R}_{\mathbf{m}}^{+*}$  is characterized by the right statement in theorem 4 according to (36) and the following characterization of cone's interior

**Theorem 5.** *Assume  $C$  is a convex cone of  $\mathbb{R}^q$  with non empty interior, then*

$$\forall \mathbf{y} \in \mathbb{R}^q, [\mathbf{y} \in \text{int}(C) \Leftrightarrow \forall \boldsymbol{\alpha} \in C^{\circ}, \boldsymbol{\alpha} \neq \mathbf{0} \Rightarrow \boldsymbol{\alpha} \cdot \mathbf{y} < 0].$$

It is possible to prove that  $\mathcal{R}_{\mathbf{m}}^{+*}$  is an open set in the case where the space of velocity is a bounded set in one dimension [37] but such a construction is not possible in other cases. A more general argument consists in considering the set

$$C = \left\{ \sum_i \lambda_i \mathbf{m}(\mathbf{v}_i), \lambda_i \geq 0, \forall \mathbf{v}_i \in \mathbb{R}^3 \right\} \quad (37)$$

which is the set of moments of all atomic measures. In order to understand why we consider this set, we first recall the following theorem.

**Theorem 6.** *(Tchakaloff theorem [25]) Let  $\mu$  be a nonnegative Borel measure compactly supported in  $K \subset \mathbb{R}^3$ . Then there exist finite many points in  $K$ ,  $w_1, \dots, w_N$  with  $N \leq q$ , and positive weights  $\alpha_1, \dots, \alpha_N$ , such that for each  $p \in \mathbb{P}$*

$$\forall p \in \mathbb{P}, \quad \int_{\mathbb{R}^d} p(\mathbf{v}) d\mu = \sum_{i=1}^N \alpha_i p(w_i).$$

This theorem implies that compactly supported  $L^1$  nonnegative function have the same moments as compactly supported atomic measures. For non compactly supported  $L^1$  nonnegative function, one proves that  $\mathcal{R}_{\mathbf{m}}^{+*} \subset C$  in the following way. One first remark that the polar cone of  $C$  is the same as the one of  $\mathcal{R}_{\mathbf{m}}^+$

$$C^{\circ} = \{ \boldsymbol{\alpha} \in \mathbb{R}^q, \boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v}) \leq 0, \forall \mathbf{v} \in \mathbb{R}^3 \} \quad (38)$$

$$= (\mathcal{R}_m^+)^{\circ}. \quad (39)$$

Then there is

**Proposition 7.** *For  $C$  defined in (37), there is  $\text{int}(C) \neq \emptyset$ .*

Thus  $\mathcal{R}_m^{+*} \subset \text{int } C$  according to theorem 5. To finish the proof, remark that for  $\Psi_\varepsilon \in C_c^\infty(\mathbb{R}^3)$ ,  $\Psi_\varepsilon \geq 0$  such that  $\Psi_\varepsilon \rightarrow \delta_0$  as  $\varepsilon \rightarrow 0$ , one has

$$\forall \mathbf{v} \in \mathbb{R}^d, \quad \int \mathbf{m}(x) \Psi_\varepsilon(\mathbf{v} - \mathbf{x}) d\mathbf{x} \longrightarrow \mathbf{m}(\mathbf{v})$$

This proves that the set of moments of nonnegative functions in  $C_c^\infty(\mathbb{R}^d)$  is dense in  $C$  ( $C \subset \text{cl } \mathcal{R}_m^+$ ). To summarize, one has

$$\mathcal{R}_m^{+*} \subset \text{int } C \subset \text{int}(\text{cl } \mathcal{R}_m^+)$$

$\mathcal{R}_m^{+*}$  is a convex set so that  $\mathcal{R}_m^{+*} = \text{int}(\text{cl } \mathcal{R}_m^+)$  according to a Caratheodory theorem in finite dimension [39] : each  $\boldsymbol{\rho} \in \text{int}(\text{cl } \mathcal{R}_m^+)$  is the convex combination of  $q + 1$  affinely independent points in  $\text{int}(\text{cl } \mathcal{R}_m^+)$  and by density of  $q + 1$  affinely independent points in  $\mathcal{R}_m^{+*}$ . So finally there is

$$\mathcal{R}_m^{+*} = \text{int } C = \text{int}(\text{cl } \mathcal{R}_m^+)$$

which proves the first assertion in theorem 4 together with " $\mathcal{R}_m^{+*}$  is an open set". Finally, the  $q + 1$  affinely independent points in the above reasoning can be chosen as moments of nonnegative functions in  $C_c^\infty(\mathbb{R}^d)$  by density which ends the proof.

**Remark 3.** *Then first statement of theorem 4 reads:*

$$\boldsymbol{\rho} \in \mathcal{R}_m^{+*} \iff [L_{\boldsymbol{\rho}}(P) > 0 \quad \forall P \in \mathbb{P} \quad \text{with} \quad P(\mathbf{v}) \geq 0 \quad (P(\mathbf{v}) \neq 0)].$$

Thus, when  $\mathbb{P} = \mathbb{P}_{2n}$ , the characterization of a realizable moment does not require any property related to polynomials of higher degree contrarily to Theorem 2.

### 3.3 Characterization of $\mathcal{R}_m^+$ in quadratic structured polynomial spaces

According to Theorem 4 and Remark 3, characterizing realizable moments is equivalent to characterizing those nonnegative polynomials  $p(\mathbf{v}) = \boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v}) \in \mathbb{P}$  on  $\mathbb{R}^3$ . This problem is related in some way to the 17<sup>th</sup> Hilbert problem which asks for the characterization of nonnegative rational function with real coefficients as the sum of square of nonnegative rational function with real coefficients. Whilst this result was proved by Artin in 1927, it is not of practical use and one should rather investigate the nature of nonnegative polynomials. Let us first give a definition.

**Definition 4.** *Let  $\mathbb{P}$  be the polynomial space generated by the basis  $\{1, \widetilde{m}_2, \dots, \widetilde{m}_p\}$ . We denote with*

$$\mathbb{P} = \text{span}(\widetilde{m}_i \widetilde{m}_j), \quad 1 \leq i \leq j \leq p$$

the space generated by the family of pairwise multiplication of any element in  $\tilde{\mathbf{m}}(\mathbf{v})$  and with  $q$  its dimension. Since the value of  $\tilde{m}_i \tilde{m}_j$  may appear for different pair of indexes  $(i, j)$ , one may extract a polynomial basis of  $\mathbb{P}$  which is composed as before of the components of  $\mathbf{m}(\mathbf{v})$ .

We say that  $\mathbb{P}$  is the quadratic space over  $\tilde{\mathbb{P}}$ .

Example 1: Each polynomial space  $\mathbb{P}_{2n}$ ,  $n \geq 1$  is quadratic over  $\mathbb{P}_n$ .

Example 2: The space of collisional invariant  $\mathbb{K}$  is not quadratic over another space if  $d > 1$  but the space generated by  $\{\mathbf{1}, \mathbf{v}, \mathbf{v} \otimes \mathbf{v}, \mathbf{v}\mathbf{v}^2, \mathbf{v}^4\}$  is quadratic over  $\mathbb{K}$  (here we have used the tensorial notation for convenience.) It is a space strictly contained in  $\mathbb{P}_4$ .

For  $\boldsymbol{\rho} \in \mathbb{R}^q$ , one can construct a  $p \times p$  moment matrix  $\mathbf{H}$  in the same way as in section 3.1. In particular,  $\mathbf{H} = M_n(\boldsymbol{\rho})$  when  $\mathbb{P} = \mathbb{P}_{2n}$ . Then there is the following characterisation of  $\mathcal{R}_{\mathbf{m}}^+$  by using theorem 4.

**Corollary 1.** *Let  $\mathbb{P}$  be a quadratic structured space. Then for  $\boldsymbol{\rho} \in \mathcal{R}_{\mathbf{m}}^{+*}$  the matrix  $\mathbf{H}$  is SPD. Moreover if all nonnegative polynomials in  $\mathbb{P}$  are sum of square polynomials then there is*

$$\boldsymbol{\rho} \in \mathcal{R}_{\mathbf{m}}^+ \iff \mathbf{H} \text{ is SPD.} \quad (40)$$

*Proof.* For  $\boldsymbol{\rho} \in \mathcal{R}_{\mathbf{m}}^{+*}$ , there exists a nonnegative function  $f \in \mathbb{L}^1(\mathbf{m})$  such that  $\int f \mathbf{m}(\mathbf{v}) d\mathbf{v} = \boldsymbol{\rho}$ . For  $\boldsymbol{\beta} \in \mathbb{R}^p \setminus \{0\}$ , we consider the square polynomial  $p(\mathbf{v}) = \boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v}) = (\boldsymbol{\beta} \cdot \tilde{\mathbf{m}}(\mathbf{v}))^2$  where  $\tilde{\mathbf{m}} = (1, \tilde{m}_2, \dots, \tilde{m}_p)^T$ . Then there is

$$\boldsymbol{\rho} \cdot \boldsymbol{\alpha} = \int p(\mathbf{v}) f(\mathbf{v}) d\mathbf{v} = \boldsymbol{\beta}^T \cdot \left( \int \tilde{\mathbf{m}} \cdot \tilde{\mathbf{m}}^T f(\mathbf{v}) d\mathbf{v} \right) \cdot \boldsymbol{\beta} = \boldsymbol{\beta}^T \cdot \mathbf{H} \cdot \boldsymbol{\beta} > 0. \quad (41)$$

So if we write the condition  $\boldsymbol{\rho} \cdot \boldsymbol{\alpha} > 0$  for all square polynomials (and thus for all sum of square polynomials), we end up with an equivalent condition which is: the moment matrix  $\mathbf{H} = (\rho_{ij})$  is symmetric positive definite (SPD).  $\square$

In one dimension, (40) holds true when  $\mathbb{P} = \mathbb{P}_{2n}$ . This is the known characterization of realizable moments for the Hamburger moment problem [3]. In higher dimension, especially when  $d = 3$ , positive polynomials in  $\mathbb{P} = \mathbb{P}_{2n}$  with  $n \geq 2$  are not necessary sum of square polynomials. We refer to [6] for a bibliography on the topic. Thus the condition  $\mathbf{H} = M_n(\boldsymbol{\rho})$  (see notations before theorem 3) is SPD is only a necessary condition.

### 3.4 Application to relaxation of the Grad moments

So far, we have focused on pretty general consideration on realizable moment. In practice, some usual basis in kinetic theory are

- "Euler" basis  $\{\mathbf{1}, \mathbf{v}, \mathbf{v}^2\}$
- Gauss basis  $\{\mathbf{1}, \mathbf{v}, \mathbf{v} \otimes \mathbf{v}\}$
- Grad basis  $\{\mathbf{1}, \mathbf{v}, \mathbf{v} \otimes \mathbf{v}, \mathbf{v} \mathbf{v}^2\}$
- Levermore basis  $\{\mathbf{1}, \mathbf{v}, \mathbf{v} \otimes \mathbf{v}, \mathbf{v}\mathbf{v}^2, \mathbf{v}^4\}$

where we have used the tensor notation.

The last one has no particular physical interpretation and was just introduced for solving variational problems with the usual entropy  $\int f \ln(f) d\mathbf{v}$ . Among those basis, Grad basis is the most important because it contains all physically meaningful moments (mass, momentum, energy, pressure tensor, heat flux). In the sequel, we use the decomposition of the Grad space defined in (16). The polynomial family of interest is the polynomial families noted as:

$$\mathbf{a}(\mathbf{v} - \mathbf{u}) = \left( \mathbf{1}, (\mathbf{v} - \mathbf{u}), (\mathbf{v} - \mathbf{u})^2 - 3 \frac{k_B T}{m}, \frac{k_B T}{m} \mathbb{A}(\mathbf{v} - \mathbf{u}), \mathbf{b}(\mathbf{v} - \mathbf{u}) \right) \quad (42)$$

$$:= (a_0(\mathbf{v} - \mathbf{u}), \mathbf{a}_1(\mathbf{v} - \mathbf{u}), a_2(\mathbf{v} - \mathbf{u}), a_3(\mathbf{v} - \mathbf{u}), \mathbf{a}_4(\mathbf{v} - \mathbf{u})). \quad (43)$$

It is still composed of tensors of even or odd ranks in the variable  $\mathbf{v} - \mathbf{u}$ . The polar cone of  $\mathcal{R}_{\mathbf{a}}^{*+}$  is given by:

$$\mathcal{C}_{\mathbf{a}}^{\circ} = \{ \boldsymbol{\alpha} \in \mathbb{R}^q, \forall \mathbf{v} \in \mathbb{R}^d, \boldsymbol{\alpha} \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}) \leq 0 \}.$$

It is immediate that any  $\boldsymbol{\alpha} \in \mathcal{C}_{\mathbf{a}}^{\circ}$  has a null component on  $(\mathbf{v} - \mathbf{u})(\mathbf{v} - \mathbf{u})^2$  and as a consequence on  $\mathbf{b}(\mathbf{v} - \mathbf{u})$ . Then we have

$$\mathcal{C}_{\mathbf{a}}^{\circ} = \left\{ (\boldsymbol{\beta}, \mathbf{0}), \forall \mathbf{v} \in \mathbb{R}^d, \boldsymbol{\beta} \cdot \left( \mathbf{1}, (\mathbf{v} - \mathbf{u}), (\mathbf{v} - \mathbf{u})^2 - 3 \frac{k_B T}{m}, a_3(\mathbf{v} - \mathbf{u}) \right) \leq 0 \right\}.$$

The study of realizable moments on  $\mathbf{a}(\mathbf{v} - \mathbf{u})$  then simplifies dramatically since we just need to know which moments are realizable in the Gauss basis defined as  $\mathbb{P}_{\text{Gauss}} = \text{span}\{\mathbf{1}, \mathbf{v}, \mathbf{v}^2, \mathbb{A}(\mathbf{v})\}$ . Using Theorem 4, the characterization of  $\mathcal{R}_{\mathbf{a}}^{*+}$  is equivalent to the characterization of the positive polynomials in  $\mathbb{P} = \text{span}\{\mathbf{1}, \mathbf{v}, \mathbf{v} \otimes \mathbf{v}\}$ . From Hilbert's theorems, any positive polynomial in  $\mathbb{P}$  can be written as a sum of square of polynomials, that is

$$g(\mathbf{v}) = \sum_i (p_i + \mathbf{q}_i \cdot \mathbf{v})^2,$$

with  $p_i \in \mathbb{R}$  and  $q_i \in \mathbb{R}^3$ . Let  $\boldsymbol{\rho} = (n, n\mathbf{u}, n\mathbb{D})$  be a moment w.r.t.  $(\mathbf{1}, \mathbf{v}, \mathbf{v} \otimes \mathbf{v})$ . According to Theorem 4,  $\boldsymbol{\rho}$  is realizable if and only if for any non negative, non null polynomial  $g(\mathbf{v})$  defined by its list of coefficients  $\boldsymbol{\gamma}$  there holds  $\boldsymbol{\gamma} \cdot \boldsymbol{\rho} > 0$ . Rearranging the components of  $g(\mathbf{v})$  in the Gauss basis, we find after some algebra that this condition can be written as

$$n \sum_i \left( [p_i \ \mathbf{q}_i] \begin{bmatrix} \mathbf{1} & \mathbf{u}^t \\ \mathbf{u} & \mathbb{D} \end{bmatrix} \begin{bmatrix} p_i \\ \mathbf{q}_i \end{bmatrix} \right) > 0 := n \sum_i \left( [p_i \ \mathbf{q}_i] \mathbf{H} \begin{bmatrix} p_i \\ \mathbf{q}_i \end{bmatrix} \right).$$

This is equivalent to  $n > 0$  and to the positivity of the moment matrix  $\mathbf{H}$ . But  $\mathbf{H}$  is positive if and only  $\mathbb{D} - \mathbf{u} \otimes \mathbf{u}$  is positive as proved in the following lemma:



**Lemma 1.** *The matrix*

$$\begin{bmatrix} 1 & \mathbf{u}^t \\ \mathbf{u} & \mathbb{D} \end{bmatrix}$$

*is positive if and only if the matrix  $\mathbb{D} - \mathbf{u} \otimes \mathbf{u}$  is positive.*

It is interesting to remark that for a given realizable moment  $\boldsymbol{\rho} := \{n, n\mathbf{u}, 3ne, n\Pi, n\mathbf{Q}\}$ , the heat flux  $n\mathbf{Q}$  can take any value in  $\mathbb{R}^3$  because of the characterization given by theorem 4 which does not include moments of order 3. This is due to the fact that one can always add to a positive function having this moment a perturbation that let it nonnegative, keeps all of its moments in the Gauss basis but not that with respect to  $\mathbf{b}(\mathbf{v} - \mathbf{u})$ . The Grad basis being conveniently written in (42), that is starting from a nonnegative function  $f$ , there is

$$\int_{\mathbb{R}^3} f \mathbf{a}(\mathbf{v} - \mathbf{u}) d\mathbf{v} = (n, 0, 0, n\Pi, n\mathbf{Q}),$$

where  $n\Pi$  is the traceless pressure tensor and  $n\mathbf{Q}$  is the heat flux. A straightforward consequence of lemma 1 is the following proposition

**Proposition 8.** *[Grad relaxation] If  $(n, \mathbf{0}, 0, n\Pi, n\mathbf{Q})$  is realizable, then for any  $\lambda_{\mathbb{A}} \in [-\frac{1}{2}, 1]$ , and  $\lambda_{\mathbb{B}} \in \mathbb{R}$  the moment  $(n, \mathbf{0}, 0, n\lambda_{\mathbb{A}}\Pi, n\lambda_{\mathbb{B}}\mathbf{Q})$  is still realizable.*

Coming back to the relaxation constraints, the above result states that the relaxed moment  $(n, 0, 0, (1 - \frac{\nu_{\mathbb{A}}}{\nu})n\Pi, (1 - \frac{\nu_{\mathbb{B}}}{\nu})n\mathbf{Q})$  is realizable when  $0 \leq 1 - \frac{\nu_{\mathbb{A}}}{\nu} \leq 1$  and for all  $\nu_{\mathbb{B}} \in \mathbb{R}$ . Also, the admissible relaxation on  $n\Pi$  is just the one that is found in the study of the ESBGK model [4, 12].

### 3.5 Galilean invariance

We may now address the problem of finding linear application that let  $\mathcal{R}_{\mathbf{m}}^+$  invariant. In particular we are going to focus on the relation between such linear maps and Galilean invariance. The method developed in section 2.2 for constructing relaxation operators  $K(f)$  must satisfy

$$\tau_{\mathbf{u}}(K(f)) = K(\tau_{\mathbf{u}}f) \quad \forall \mathbf{u} \in \mathbb{R}^3 \quad \text{and} \quad \tau_{\Theta}K(f) = K(\tau_{\Theta}f) \quad \forall \Theta \in SO(3), \quad (44)$$

where

$$\tau_{\mathbf{u}}f(\mathbf{v}) = f(\mathbf{v} + \mathbf{u}) \quad \text{and} \quad \tau_{\Theta}f(\mathbf{v}) = f(\tau_{\Theta}\mathbf{v}).$$

Recall that  $K(f) = \nu(G - f)$ . The setting of  $\nu$  and construction of  $G$  just depends on the moments of  $f$  -  $\boldsymbol{\rho} = \int f(\mathbf{v})\mathbf{m}d\mathbf{v} \in \mathcal{R}_{\mathbf{m}}^{*+}$  - and not on  $f$  itself. We may write  $\nu = \nu(\boldsymbol{\rho}(f))$  and  $G = G(\boldsymbol{\rho}(f))$ . So, starting from a function  $f \in \mathbb{L}^{1,+}(\mathbf{m})$ ,  $\tau_{\mathbf{u}}f$  and  $\tau_{\Theta}f$  must themselves be function of  $\mathbb{L}^{1,+}(\mathbf{m})$  to make the construction  $G(\boldsymbol{\rho}(\tau_{\mathbf{u}}f))$  (likewise  $G(\boldsymbol{\rho}(\tau_{\Theta}f))$ ) possible. This writes

$$\boldsymbol{\rho}(\tau_{\mathbf{u}}f) = \int (\tau_{\mathbf{u}}f)\mathbf{m}(\mathbf{v})d\mathbf{v} = \int f(\mathbf{w})\mathbf{m}(\tau_{-\mathbf{u}}(\mathbf{w}))d\mathbf{w} \in \mathcal{R}_{\mathbf{m}}^{*+}.$$

A sufficient condition for integrability of  $\tau_{\mathbf{u}}f$  in  $\mathbb{L}^1(\mathbf{m})$  is:  $\mathbb{P} = \text{span}(\mathbf{m})$  is invariant under the action of  $\tau_{-\mathbf{u}}$ , which can be expressed as

$$\exists \Lambda(-\mathbf{u}) \in \mathbb{R}^q \times \mathbb{R}^q \text{ such that } \mathbf{m}(\tau_{-\mathbf{u}}(\mathbf{w})) = \Lambda(-\mathbf{u})\mathbf{m}(\mathbf{w}). \quad (45)$$

Positivity is then satisfied since  $\tau_{\mathbf{u}}f \geq 0$  and there is  $\boldsymbol{\rho}(\tau_{\mathbf{u}}f) \in \mathcal{R}_{\mathbf{m}}^{*+}$ . As a consequence  $\mathcal{R}_{\mathbf{m}}^{*+}$  is stable under the action of the linear map  $\Lambda(-\mathbf{u})$  which inverse is  $\Lambda(\mathbf{u})$ . Likewise, if there exists  $\Lambda(\Theta^t) \in \mathbb{R}^q \times \mathbb{R}^q$  such that

$$\mathbf{m}(\Theta^t \mathbf{w}) = \Lambda(\Theta^t)\mathbf{m}(\mathbf{w}), \quad (46)$$

the same conclusion holds. The above conditions relate polynomial spaces  $\mathbb{P}$  which are invariant under the action of translations and rotations to set of moments  $\mathcal{R}_{\mathbf{m}}^{*+}$  which are invariant under the above mapping. Such polynomial spaces are named Galilean invariant and are necessarily polynomial as proved by Junk and Unterreiter [31].

**Proposition 9.** *Assume that the space  $\mathbb{P} = \text{span}\{\mathbf{m}\}$  is invariant under the translations and the rotations. Then  $\mathcal{R}_{\mathbf{m}}^{*+}$  is invariant under the action of  $\Lambda(\mathbf{u})$  and  $\Lambda(\Theta)$  for any  $\mathbf{u} \in \mathbb{R}^3$  and any  $\Theta \in SO(3)$ , where  $\Lambda(\mathbf{u})$  and  $\Lambda(\Theta)$  are the matrices defined in (45) and (46).*

Instead of proving the equivalence in the above conclusion, we prefer to look at the construction of the model and state a necessary condition on  $\mathbb{P}$ .

**Proposition 10.** *Let  $G : \mathcal{R}_{\mathbf{m}}^{*+} \rightarrow \mathbb{L}^{1,+}(\mathbf{m})$  s.t.  $G(\boldsymbol{\rho}_1) \neq G(\boldsymbol{\rho}_2)$  for any  $\boldsymbol{\rho}_1 \neq \boldsymbol{\rho}_2 \in \mathcal{R}_{\mathbf{m}}^{*+}$ . Then if  $G$  satisfies*

$$\begin{aligned} (\forall f \in \mathbb{L}^{1,+}(\mathbf{m})), (\forall \mathbf{u} \in \mathbb{R}^3), (\forall \Theta \in SO(3)), G(R[\tau_{\mathbf{u}}f]) &= \tau_{\mathbf{u}}G(R(f)), \\ G(R[\tau_{\Theta}f]) &= \tau_{\Theta}G(R[f]), \end{aligned}$$

the space  $\mathbb{P} = \text{span}\{\mathbf{m}\}$  is invariant under Galilean transforms.

## 4 Solving the variational problem

In this section, we address the problem of finding the target function  $G$  which is the last step in the method of moments relaxation (see section 2.2.2). In a first time (section 4.1), we state the variational problem in a quite general form and display classical properties for the problem to be well-posed. The variational problem is then restated in the framework of  $\phi$ -divergence [18]. The (relative) entropy being also a  $\phi$ -divergence, the corresponding minimization problem is addressed in section 4.2.1. We summarize known results in this case and eventually outline a lacking property in the set listed in section 4.1. In the next section, we recall a theorem by Csiszar [20] together with some ideas in his proof. We focus on the reason why a specific condition in his theorem somehow allows to satisfy this lacking property. In the last section, we reestablish the previous theorem by using convex analysis and recall the results obtained in [2]. All along this section, we keep the notations introduced in section 2.2.

## 4.1 The variational problem

The problem we want to address is the following.

**Problem 1.** Let  $\boldsymbol{\rho} \in \mathbb{R}^q$  and consider the convex domain

$$D(\boldsymbol{\rho}) = \left\{ g \in \mathbb{L}^1(\mathbf{m}), \int \mathbf{m} g = \boldsymbol{\rho} \right\}. \quad (47)$$

For a strictly convex function  $\eta$  with  $\text{dom}(\eta) = \mathbb{R}^+$ , we consider the functional

$$\mathcal{H}(g) = \int \eta(g) d\mathbf{v}. \quad (48)$$

Then the problem is to find if possible a nonnegative function  $G \in D(\boldsymbol{\rho})$  such that

$$\mathcal{H}(G) = \inf_{g \in D(\boldsymbol{\rho})} \mathcal{H}(g). \quad (49)$$

Clearly, such a problem much depends on the property of  $\eta$ . The above ones being too general, we display conditions that ensure that this variational problem is well-posed [8].

### Ideal properties

1. Together with  $\eta$  is strictly convex with  $\text{dom}(\eta) = \mathbb{R}^+$  (and thus closed on  $\mathbb{R}^+$ ), one requires  $\eta$  to be proper and super linear at infinity.
2. There exists a weak topology for which the functional  $g \mapsto \int \eta(g) d\mathbf{v}$  is semi-lower continuous and which is suited to the continuity of the constraints with respect to  $g$ .
3. The next condition may be stated in simple words as follows. There exists a feasible function  $g$ , that is a function satisfying at the same time  $\int \eta(g) d\mathbf{v} \in \mathbb{R}$  and  $\int g \mathbf{m} d\mathbf{v} = \boldsymbol{\rho}$ .

With this set of conditions, it may be proved that there is

$$G(\mathbf{v}) = \eta^{*'}(\boldsymbol{\alpha}[\boldsymbol{\rho}] \cdot \mathbf{m}(\mathbf{v})). \quad (50)$$

for some  $\boldsymbol{\alpha} \in \mathbb{R}^q$ . As we will see below, the second condition is hardly compatible with kinetic theory when one considers the entropy:  $g \mapsto \int g \ln(g) d\mathbf{v}$  in an unbounded domain.

## 4.2 $\phi$ -divergence

We now want to restate the above problem in term of  $\phi$ -divergence that were introduced by Csiszar in 1970 [18].

**Definition 5.** Let  $\phi$  be a strictly convex function with domain on  $(0, +\infty)$  and minimum at  $x = 1$ . The  $\phi$ -divergence of two distribution function is defined as

$$I(p||r) = \int_{\mathbb{R}^3} r(\mathbf{v}) \phi\left(\frac{p(\mathbf{v})}{r(\mathbf{v})}\right) d\mathbf{v}. \quad (51)$$

In our case, we are much concerned with a kinetic equation whose solution should be closed to the local equilibrium states. It is then natural to let  $r = \mathcal{M}$  into (51) since then minimizing  $I(g||\mathcal{M})$  under moment constraints amounts to look for a function  $G$  which is as closed as possible to  $\mathcal{M}$  relatively to the divergence  $\phi$ . This approach is similar to that of Abdelmalik and Van Brummelen and their analysis of moments closure [2]. Then, the variational problem is stated as follows.

**Problem 2.** [primal problem] Let  $\phi$  be a  $\phi$ -divergence. Define  $\mathcal{H}$  with

$$\mathcal{H}(g) = \int \mathcal{M}\phi\left(\frac{g}{\mathcal{M}}\right) d\mathbf{v}, \quad (52)$$

and the real extended value function  $h: \mathbb{R}^q \rightarrow \overline{\mathbb{R}}$  for any  $\boldsymbol{\rho} \in \mathbb{R}^q$  by

$$h(\boldsymbol{\rho}) = \inf_{g \in D(\boldsymbol{\rho})} \mathcal{H}(g).$$

The primal problem consists in finding if possible a function  $G$  s.t.

1.  $G \in D(\boldsymbol{\rho})$
2.  $\mathcal{H}(G) = h(\boldsymbol{\rho})$

**Remark 4.** The above problem amounts to let

$$\eta(f) = \mathcal{M}\phi\left(\frac{f}{\mathcal{M}}\right), \quad (53)$$

in problem 1. Yet, the problems 1 and 2 are very much similar if one substitutes to the Lebesgue measure the measure  $d\mu = \mathcal{M}d\mathbf{v}$ . Then one may define the functional

$$\mathcal{F}(u) = \int \phi(u) d\mu, \quad (54)$$

and consider the set

$$\tilde{D}(\boldsymbol{\rho}) = \left\{ u \in \mathbb{L}^1(\mathbf{m})(d\mu), \int \mathbf{m}(\mathbf{v}) u(\mathbf{v}) d\mu = \boldsymbol{\rho} \right\}. \quad (55)$$

Then problem 1 is replaced by : find if possible a nonnegative function  $P \in \tilde{D}(\boldsymbol{\rho})$  such that

$$\mathcal{F}(P) = \inf_{u \in \tilde{D}(\boldsymbol{\rho})} \mathcal{F}(u), \quad (56)$$

and the solution, if it exists, gives the solution to problem 2 as

$$G(\mathbf{v}) = \mathcal{M}(\mathbf{v}) P(\mathbf{v}). \quad (57)$$

So the main problem amounts to set the right conditions on  $\phi$  in order to have a well-posed problem. Such conditions should be in some way the same as the ones referred as "ideal properties" in section 4.1.

#### 4.2.1 Analysis of the case $\phi(x) = x \ln(x)$

We now want to recall some key points related to the minimization of the natural entropy  $\mathcal{H}(g) = \int g \ln(g) d\mathbf{v}$  under moment constraints. This problem has been analysed in a serie of articles [26, 30, 36, 40]. In order to recast this functional within the framework of  $\phi$ -divergence, one just let  $\phi(x) = x \ln(x)$  into (52) in which case  $\mathcal{H}(g) = \int g \ln(g/\mathcal{M}) d\mathbf{v}$ . In the case where  $\mathbb{K} \subset \mathbb{P}$  and  $\int g \psi d\mathbf{v} = \int \mathcal{M} \psi d\mathbf{v}$  for all  $\psi \in \mathbb{K}$ ,  $\mathcal{H}(g)$  is the classical relative entropy.

Let us first remark that  $\phi$  satisfies the first set of ideal properties displayed in section 4.1. The condition 3 is addressed in the following proposition.

**Proposition 11.** *The domain of  $h$  is  $\mathcal{R}_m^+$ .*

*Proof.* ( $\forall f \geq 0$ ),  $f \ln(f/\mathcal{M}) \geq -\frac{1}{e}\mathcal{M}$  as a simple consequence of the inequality  $x \ln(x) \geq -\frac{1}{e}$ . Thus ( $\forall f \geq 0$ ),  $\mathcal{H}(f) \geq -\frac{n}{e}$ .

Next,  $\forall \boldsymbol{\rho} \in \mathcal{R}_m^+$ ,  $\exists \Psi_{\boldsymbol{\rho}} \in \mathcal{C}_c^\infty(\mathbb{R}^3)$  such that  $\Psi_{\boldsymbol{\rho}} \in D(\boldsymbol{\rho})$ . Hence  $f/\mathcal{M}$  is compactly supported, bounded and  $\mathcal{H}(\Psi_{\boldsymbol{\rho}}) < +\infty$ . And there is

$$D^+(\boldsymbol{\rho}) = \{g \text{ s.t. } \int_{\mathbb{R}^3} \mathbf{m} g d\mathbf{v} = \boldsymbol{\rho} \text{ and } \mathcal{H}(g) < +\infty\} \neq \emptyset.$$

As a consequence  $h(\boldsymbol{\rho})$  is well defined for all  $\boldsymbol{\rho} \in \mathcal{R}_m^+$ .

Assume now that  $\boldsymbol{\rho} \notin \mathcal{R}_m^+$ . Then for any  $f$  s.t.  $\int f \mathbf{m} d\mathbf{v} = \boldsymbol{\rho}$ , there is an open set  $\omega_f$  of non zero measure s.t.  $f < 0$  on  $\omega_f$ .  $\phi$  being equal to  $+\infty$  when  $x < 0$ , we have  $h(\boldsymbol{\rho}) = +\infty$ .  $\square$

Utilizing theorem 4, for any  $\boldsymbol{\rho} \in \mathcal{R}_m^+$  we may restrict  $D^+(\boldsymbol{\rho})$  to

$$D_{\Psi}^+(\boldsymbol{\rho}) = \{g \text{ s.t. } \int_{\mathbb{R}^3} \mathbf{m} g d\mathbf{v} = \boldsymbol{\rho} \text{ and } \mathcal{H}(g) \leq \mathcal{H}(\Psi_{\boldsymbol{\rho}})\} \neq \emptyset.$$

The superlinearity of  $\phi$  together with the boundedness of moments of order more than 1 show that  $D_{\Psi}^+$  is weakly relatively compact in  $\mathbb{L}^1$  according to Dunford-Pettis lemma. Thus, a minimizing sequence  $g_n \in D_{\Psi}^+(\boldsymbol{\rho})$  (which is also a minimizing sequence in  $D^+(\boldsymbol{\rho})$ ) converges weakly in  $\mathbb{L}^1$  to a function  $G$ . But this does not prove more and property 2 in section 4.1 is not satisfied. A consequence is that  $G$  might not satisfy  $\int G \mathbf{m} d\mathbf{v} = \boldsymbol{\rho}$ . Junk has shown in a famous paper [30] that the constraint of highest degree might drop in when looking at the infimum of the natural entropy in  $D^+(\boldsymbol{\rho})$ . We want here to give a rapid hint into that problem. Consider the dual function  $h^*$  of  $h$  defined on its domain  $\Lambda$  by

$$h^*(\boldsymbol{\alpha}) = \int \exp(\boldsymbol{\alpha} \cdot \mathbf{m}) d\mathbf{v} \quad (58)$$

and assume that  $\Lambda \cap \partial\Lambda \neq \emptyset$ . For  $\boldsymbol{\alpha} \in \Lambda \cap \partial\Lambda$ ,  $h^*$  has only sided derivative at  $\boldsymbol{\alpha}$  on  $\Lambda \cap \partial\Lambda$  and in the directions pointing into the domain and there is only a subdifferential at  $\boldsymbol{\alpha}$ . One can prove the following. Firstly, each moment

$$\boldsymbol{\rho}_+ = (\rho_1, \dots, \rho_q + t), \quad t > 0, \quad (59)$$

where

$$\boldsymbol{\rho} = \int \mathbf{m}(v) \exp(\boldsymbol{\alpha} \cdot \mathbf{m}) dv$$

belongs to  $\mathcal{R}_{\mathbf{m}}^+$  when  $\mathbf{m} = \{\mathbf{1}, \mathbf{v}, \mathbf{v}^2, \dots, |\mathbf{v}|^N\}$  (here  $N$  is the maximal degree of the component). Indeed,  $\boldsymbol{\rho}$  being a realizable moment, the necessary and sufficient condition  $\boldsymbol{\rho} \cdot \boldsymbol{\alpha} > 0$  for any positive polynomial implies  $\boldsymbol{\rho}_+ \cdot \boldsymbol{\alpha} > 0$  for any  $\boldsymbol{\rho}_+$ . One can then prove that [26, 36]

$$h^{**}(\boldsymbol{\rho}_+) = \max_{\tilde{\boldsymbol{\alpha}}} \{\tilde{\boldsymbol{\alpha}} \cdot \boldsymbol{\rho}_+ - h^*(\tilde{\boldsymbol{\alpha}})\} = \boldsymbol{\alpha} \cdot \boldsymbol{\rho} - h^*(\boldsymbol{\alpha}) = h(\boldsymbol{\rho}).$$

$h$  being semi lower continuous in  $\mathcal{R}_{\mathbf{m}}^{+*}$  (implying  $h^{**} = h$ ) this proves that the subdifferential  $\partial h^*(\boldsymbol{\alpha})$  is the whole half-line  $\boldsymbol{\rho}_+$ . Therefore  $\inf_{g \in D(\boldsymbol{\rho}_+)} \mathcal{H}(g) = h(\boldsymbol{\rho})$  is attained at the function  $\exp(\boldsymbol{\alpha} \cdot \mathbf{m}) \notin D(\boldsymbol{\rho}_+)$ . To summarize, the existence of a solution to the primal problem is subjected to the shape of the domain of definition of the dual function  $h^*$ .

#### 4.2.2 Csiszar theorem

In a remarkable paper, Csizsár [20] has shown existence of solution to variational problem for a wide variety of  $\phi$ -divergence (see definition 5). The main concern in the present context and to that of moments closure is to obtain a unique nonnegative minimizer qualifying all constraints. In order to do so, Csiszar adds an assumption and his results reads as follows.

**Theorem 12.** *Let  $\phi$  be a strictly convex differentiable function defined on  $]0, +\infty[$ . For  $g \in L^1(\mathbb{R}^3)$  define  $\mathcal{H}(g)$  with*

$$\mathcal{H}(g) = \int_{\mathbb{R}^3} \mathcal{M} \phi\left(\frac{g}{\mathcal{M}}\right) d\mathbf{v}.$$

Assume that

$$\phi(1) = \phi'(1) = 0, \quad \lim_{p \rightarrow +\infty} \phi'(p) = +\infty, \quad (60)$$

$$(\forall \lambda > 0), \quad \int_{\mathbb{R}^3} \phi^*(\lambda |m_i(\mathbf{v})|) \mathcal{M} d\mathbf{v} < +\infty. \quad (61)$$

Let  $\boldsymbol{\rho} \in \mathbb{R}^q$  and  $D(\boldsymbol{\rho})$  be the set of constraints defined in (55). Then if

$$\inf_{f \in D(\boldsymbol{\rho})} \mathcal{H}(f) \in \mathbb{R}, \quad (62)$$

there exists a unique function  $G \geq 0$  in  $D(\boldsymbol{\rho})$  such that

$$\mathcal{H}(G) = \inf_{g \in D(\boldsymbol{\rho})} \mathcal{H}(g).$$

We present here some ideas of the proof. First of all (62) is equivalent to state that there exists a feasible function  $g$  with the following meaning:  $g \in D(\boldsymbol{\rho})$  and  $\mathcal{H}(g) < +\infty$ . The first part of the proof is rather technical and requires some knowledge in measure theory. In the previous section, we have seen that under the conditions (60) together with  $\phi$  is strictly convex and differentiable on  $]0, +\infty[$  - a minimizing sequence  $g_n$  in  $D^+(\boldsymbol{\rho})$  converges weakly to a function  $G$  in  $L^1$  by using the Dunford-Pettis lemma. Csiszar approach is stronger and allows to prove that under the same conditions  $g_n \rightarrow G$  in  $L^1$ .  $G$  is then named as the general projection of  $\mathcal{M}$  onto  $D^+(\boldsymbol{\rho})$ . The second part of the proof deals somehow with the properties 2 in section 4.1 in a different way. Instead of considering the functional  $\mathcal{H}(g)$ , Csiszar proposes to use the theory of Orlicz space whose main purpose is to extend to convex functional the usual norm based on function such as  $\psi(x) = x^p$  for  $1 \leq p$ . One considers the function  $\psi$  defined by  $\psi(x) = \phi(x+1)$ .  $\psi$  is a strictly convex and differentiable function on  $] -1, +\infty[$ ,  $\psi(0) = \psi'(0) = 0$  and  $\lim_{p \rightarrow +\infty} \psi'(p) = +\infty$ . The (Banach) Orlicz space  $L_\psi$  is defined as the space of those functions  $u$  such that there exists a  $\lambda > 0$  satisfying

$$\int_{\mathbb{R}^3} \psi(\lambda|u|) d\mu < +\infty,$$

where  $d\mu = \mathcal{M}d\mathbf{v}$ . It is equipped with the norm

$$\|g\|_{L_\psi} = \inf_{\lambda > 0} \frac{1}{\lambda} \left( 1 + \int_{\mathbb{R}^3} \psi(\lambda|g|) d\mu \right)$$

(different types of norm exist). The condition (61) on the convex conjugate  $\phi^*$  is related to a weak topology on  $L_\psi$  which is set as follows. The convex conjugate  $\psi^*$  shares the same properties as  $\psi$  on  $\mathbb{R}^+$  (see lemma 2 in section 4.3) so that one may as well consider the Orlicz space  $L_{\psi^*}$  and its subspace  $N_{\psi^*}$  - named the "small" Orlicz space - of those function  $u^*$  such that

$$(\forall \lambda > 0), \quad \int_{\mathbb{R}^3} \psi^*(\lambda|u^*|) d\mu < +\infty.$$

$N_{\psi^*}$  and  $L_\psi$  are dual spaces and one defines on  $L_\psi$  the  $N_{\psi^*}$ -weak topology with  $u_n \rightharpoonup u$  in  $L_\psi$  iff  $\forall \theta \in N_{\psi^*}$ ,  $\int u_n \theta d\mu \rightarrow \int u \theta d\mu$ . In order to relate (61) to this weak topology, we prefer to use the classical framework of convex analysis by letting  $\phi(x) = +\infty$ ,  $\forall x < 0$  and  $\phi(0) = \lim_{x \rightarrow 0^+} \phi(x)$  (which is either a positive finite value or  $+\infty$ ). Then there is

$$\begin{aligned} \psi^*(x^*) &= \max_{x \in \mathbb{R}} (xx^* - \phi^*(x+1)) = -x^* + \max_{z \in \mathbb{R}} (zx^* - \phi^*(z)) = \phi^*(x^*) - x^*, \\ \phi^*(0) &= \max_{x \in \mathbb{R}} (-\phi^*(0)) = 0 = \psi^*(0), \end{aligned}$$

so that  $\phi^*(x^*) \geq \psi^*(x^*)$  for  $x^* \geq 0$ . As a consequence, there is  $m_i \in N_{\phi^*}$  for all  $i$ . The reason why the minimizing sequence  $g_n$  converges to a function  $G \in D(\boldsymbol{\rho})$  lies on the following theorem.

**Theorem 13.** [32] *If a sequence  $(u_n)_n \subset L_\psi$  is bounded in Orlicz norm and converges in  $\mu$ -measure to some function  $u$  then  $u \in L_\psi$  and  $u_n \rightarrow u$   $N_{\psi^*}$ -weakly.*

Let  $u_n = g_n/\mathcal{M}$  where  $g_n$  is the minimizing sequence.  $u_n$  converges in  $\mu$ -measure to  $G/\mathcal{M}$  since  $g_n \rightarrow G$  in  $L^1$  while

$$\begin{aligned} \|u_n\|_{L_\psi} &\leq 2\left(1 + \int_{\mathbb{R}^3} \phi\left(1 + \frac{u_n}{2}\right) d\mu\right) \leq 2 + n\phi(2) + \int_{\mathbb{R}^3} \phi(u_n) d\mu \\ &= 2 + n\phi(2) + \mathcal{H}(g_n), \end{aligned}$$

by convexity which ends the proof (here  $n = \int d\mu$ ).

### 4.3 Assumption and main result

The strength of Csiszar approach is to consider property 2 in section 4.1 by setting a property on  $\phi^*$  and utilizing the theory of Orlicz space. Yet his proof should be completed with the characterization of the minimizing function  $G$ . So our purpose is to establish a result similar to Theorem 12 by using standard convex analysis which allows to characterize  $G$ . As usual in convex analysis, this will be done by considering the dual problem to problem 2.

**Problem 3** (dual problem). *The dual problem consists in defining the Legendre dual function  $h_{\mathbf{m}}^*$  from  $\mathbb{R}^q$  to the extended reals  $\overline{\mathbb{R}}$  - that is  $\mathbb{R} \cup \{+\infty, -\infty\}$  - as follows:*

$$\forall \boldsymbol{\alpha} \in \mathbb{R}^q, \quad h_{\mathbf{m}}^*(\boldsymbol{\alpha}) = \sup_{\boldsymbol{\rho} \in \mathbb{R}^q} (\boldsymbol{\rho} \cdot \boldsymbol{\alpha} - h_{\mathbf{m}}(\boldsymbol{\rho})). \quad (63)$$

Then the main theorem reads as.

**Theorem 14.** *Let  $\phi : \mathbb{R} \mapsto \mathbb{R} \cup \{+\infty\}$  be a strictly convex and differentiable function on its domain  $\text{dom}(\phi) = [0, +\infty)$ . Assume moreover that*

1. *There holds the following properties:*

$$\phi(0) = 0, \quad p_0 := \inf_{y>0} \frac{\phi(y)}{y} \in \mathbb{R}, \quad \sup_{y>0} \frac{\phi(y)}{y} = +\infty$$

2. *For any polynomial  $\pi(\mathbf{v}) = \boldsymbol{\gamma} \cdot \mathbf{m}(\mathbf{v})$  then  $\phi^*(\pi) \in \mathbb{L}^1(\mathcal{M}(\mathbf{v}) d\mathbf{v})$  where  $\phi^*$  is the Legendre transform of  $\phi$*

*Then there hold*

1. *For any  $\boldsymbol{\rho} \in \mathcal{R}_{\mathbf{m}}^{+*}$  there exists a unique  $\boldsymbol{\alpha} \in \mathbb{R}^q$  such that*

$$\boldsymbol{\rho} = \int \phi'^*(\boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v})) \mathbf{m}(\mathbf{v}) \mathcal{M} d\mathbf{v}$$



and the moments  $\boldsymbol{\rho}$  and its conjugate moment  $\boldsymbol{\alpha}$  are linked thanks to the sub-differential equation:

$$h_{\mathbf{m}}(\boldsymbol{\rho}) + h_{\mathbf{m}}^*(\boldsymbol{\alpha}) = \boldsymbol{\alpha} \cdot \boldsymbol{\rho} \quad (64)$$

2. Moreover the function  $G = \mathcal{M}\phi^{*'}(\boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v}))$  is the unique minimizer of the primal problem and satisfies

$$h_{\mathbf{m}}(\boldsymbol{\rho}) = \int \phi\left(\frac{G}{\mathcal{M}}\right) \mathcal{M}d\mathbf{v} \quad (65)$$

3.  $h_{\mathbf{m}}$  is strictly convex in its domain and  $\nabla h^*$  is a bijection from  $C^\circ$  to  $\mathcal{R}_{\mathbf{m}}^+ \setminus \{0\}$  where

$$C^\circ = \{\boldsymbol{\alpha} \in \mathbb{R}^q, \boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v}) > p_0, \text{ on a set } \omega_{\boldsymbol{\alpha}} \text{ of non 0 measure}\}. \quad (66)$$

Let us address the following remarks:

1. The first assumption -  $\text{dom}(\phi) = [0, +\infty)$  - enables to grant that the primal problem will have solutions that will be non negative. Also, it is quite natural to assume that 0 is in the domain of  $\phi$  unless the p.d.f 0 cannot appear as a solution. So we will be able to produce non negative modeling function for the BGK construction. Differentiation on  $[0, +\infty)$  implies semi-lower continuity on  $[0, +\infty)$  (which is essential in convex optimization) and will enable one to one property of  $\phi'$ . Strict convexity will have, as often, many implications in uniqueness consideration.
2. The third assumption is essential for integrations properties. In particular, it breaks when  $\phi(x) = x \ln(x)$  since there holds then  $\phi^*(y) = \exp(y - 1)$  and only few polynomial  $\exp(\pi(\mathbf{v}) - 1)$  have proper integration properties. This is one of the most difficult aspect of the Levermore program which was based first on a function like  $x \ln(x)$ .

The proof goes along different steps we are going to develop. The first one consists in generalizing Proposition 11 concerning the primal problem.

**Proposition 15.** *The following properties are satisfied*

1. The moment entropy function  $h_{\mathbf{m}}(\boldsymbol{\rho})$  has returned values in  $\mathbb{R} \cup \{+\infty\}$
2. Its domain is exactly  $\text{dom}(h_{\mathbf{m}}) = \mathcal{R}_{\mathbf{m}}^+ = \{\int \mathbf{m}g, g \in \mathbb{L}^{1+}(\mathbf{m})\}$ .
3. The interior of the domain is given by  $\text{int}(\text{dom}(h_{\mathbf{m}})) = \mathcal{R}_{\mathbf{m}}^{+*} = \text{dom}(h_{\mathbf{m}}) \setminus \{0\}$
4. The function  $h_{\mathbf{m}}$  is convex.

Then the Legendre conjugate  $\phi^*$  of  $\phi$  has the following properties.

**Lemma 2.** *Consider a function  $\phi$  such as in Theorem 14. Then the following properties hold*

1.  $\forall x \in (-\infty, p_0], \phi^*(x) = 0$
2.  $\forall x \in (p_0, +\infty), \phi^*(x) > 0$
3. The function  $\phi^* : \mathbb{R} \mapsto \mathbb{R}^+$  is  $C^1$  smooth and thus semi-lower continuous.

The link between strict convexity of a given function and smoothness of its conjugate is well known in convex analysis (see [39]) but we prefer to detail the whole proof of this lemma for the sake of consistency. Smoothness of the function  $\phi^*$  has many important implications on the smoothness of the dual function  $h^*$ , which, in turn, has

important consequences on the primal minimization problem and on existence and uniqueness of the minimizer.

We can now give an explicit expression of the dual function  $h^*$ .

**Theorem 16.** 1. For any  $\boldsymbol{\alpha} \in \mathbb{R}^q$  there is

$$h_{\mathbf{m}}^*(\boldsymbol{\alpha}) = \sup_{g \in \mathbb{L}^1(\mathbf{m})} \int \left[ \boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v}) g(\mathbf{v}) - \phi\left(\frac{g}{\mathcal{M}}\right) \mathcal{M}(\mathbf{v}) \right] d\mathbf{v} \in \mathbb{R}. \quad (67)$$

2. Moreover one can also compute for any  $\boldsymbol{\alpha} \in \mathbb{R}^q$

$$h_{\mathbf{m}}^*(\boldsymbol{\alpha}) = \int \phi^*(\boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v})) \mathcal{M}(\mathbf{v}) d\mathbf{v}. \quad (68)$$

3. The function  $h_{\mathbf{m}}^*$  is continuously differentiable on  $\mathbb{R}^q$  and there holds for any  $\boldsymbol{\alpha} \in \mathbb{R}^q$

$$h_{\mathbf{m}}^{*'}(\boldsymbol{\alpha}) = \int \phi^{*'}(\boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v})) \mathbf{m}(\mathbf{v}) \mathcal{M}(\mathbf{v}) d\mathbf{v}. \quad (69)$$

Now we are able to finish the proof of Theorem 14. In few words, it is based on well known properties of convex analysis.

1. A convex function  $h : \mathbb{R}^q \rightarrow \mathbb{R} \cup \{+\infty\}$  is continuous on the interior of its domain
2. A convex function  $h : X \rightarrow \mathbb{R} \cup \{+\infty\}$  which is continuous at  $\mathbf{x} \in \text{dom}()$  has a non void sub-differential at  $\mathbf{x}$ , that is  $\partial h(\mathbf{v}) \neq \emptyset$ .
3. For a proper convex function  $h$ , closed at  $\boldsymbol{\rho}$ , there is  $\boldsymbol{\alpha} \in \partial h(\boldsymbol{\rho}) \Leftrightarrow \boldsymbol{\rho} \in \partial h^*(\boldsymbol{\alpha})$ .

#### 4.3.1 Abdelmalik and Van Brummelen approximation of the Boltzmann entropy

In their article, the authors consider different  $\phi$ -divergence to analyze closures in moment systems for the Boltzmann equation. In particular they focus on the problem related to the case  $\phi(x) = x \ln(x)$  (see section 4.2.1) and on the expected solution to the corresponding variational problem 2. Recall that for  $\boldsymbol{\alpha} \in \Lambda \cap \partial\Lambda$  where  $\Lambda$  is the domain of definition of  $h^*$  (58), a minimizing sequence in  $D(\boldsymbol{\rho}_+)$  (see (59) and (60)) converges to  $\exp(\boldsymbol{\alpha} \cdot \mathbf{m}) \notin D(\boldsymbol{\rho}_+)$ . Instead, they propose to look for a solution of the form

$$G = \mathcal{M}^{\phi^{*'}}(\boldsymbol{\alpha}_N \cdot \mathbf{m}(\mathbf{v})) \quad \text{with} \quad \phi_N^{*'}(y) = \left(1 + \frac{y}{N}\right)_+^N, \quad (70)$$

where  $(x)_+ = \frac{1}{2}(x + |x|)$  (here  $N \in \mathbb{N}^*$ ). This amounts to let

$$\forall y \in \mathbb{R}, \quad \phi_N^*(y) = \frac{N}{N+1} \left(1 + \frac{y}{N}\right)_+^{N+1}, \quad (71)$$

which is the Legendre dual function of  $\phi_N$  defined with

$$\forall x \geq 0, \quad \phi_N(x) = \frac{N}{N+1} \left(x \widetilde{\log}_N(x) - x\right), \quad (72)$$

where

$$\forall x \geq 0, \quad \widetilde{\log}_N(x) = (\phi_N^{*'})^{-1}(x) = Nx^{1/N} - N.$$

The motivation of the definition of  $\phi_N^{*'}(70)$  comes from the Stirling-Tsallis approximation of the exponential function since there is the pointwise limits

$$\forall y \in \mathbb{R}, \quad \lim_N (\phi_N^{*'})'(y) = \exp(y).$$

together with

$$\forall x \geq 0, \quad \lim_{N \rightarrow +\infty} \phi_N(x) = x \ln(x) - x, \quad \forall y \in \mathbb{R}, \quad \lim_N \phi_N^*(y) = \exp(y).$$

So while the variational problem 2 with  $\phi(x) = x \ln x$  is not well-posed for every  $\rho \in \mathcal{R}_{\mathbf{m}}^{+*}$ , its approximation  $\phi_N$  satisfies all assumptions in theorems 14. Thus for  $\rho \in \mathcal{R}_{\mathbf{m}}^{+*}$  there exists a unique function  $G \in D(\rho)$  such that

$$\mathcal{H}(G) = \inf_{g \in D(\rho)} \int_{\mathbb{R}^3} \mathcal{M} \phi_N \left( \frac{g}{\mathcal{M}} \right) d\mathbf{v}.$$

Moreover  $G$  reads as (70) according to theorem 14.

## 5 Application to the construction of BGK models using $\phi$ -divergence

### 5.1 Relaxation on the Grad thirteen moments

Let  $f(t, x, v)$  be a nonnegative function at  $(t, x)$ . Denote with  $n$ ,  $\mathbf{u}$  and  $T$  the corresponding density, velocity temperature and with  $\mathcal{M}$  the local Maxwellian associated to  $f$ . Let  $\mathbf{a}(\mathbf{v} - \mathbf{u})$  be the local Grad basis and denote

$$\rho_f = (n, \mathbf{0}, 0, n\Pi, n\mathbf{Q}) \tag{73}$$

$$= \int \left( \mathbf{1}, (\mathbf{v} - \mathbf{u}), (\mathbf{v} - \mathbf{u})^2 - 3 \frac{k_B T}{m}, \mathbb{A}(\mathbf{v} - \mathbf{u}), \mathbf{b}(\mathbf{v} - \mathbf{u}) \right) f(\mathbf{v}) d\mathbf{v}. \tag{74}$$

Remark that there is by definition  $\rho_{\mathcal{M}} = (n, \mathbf{0}, 0, 0, 0)$ .

#### 5.1.1 Principle of construction

We just recall here the steps in the derivation of a relaxation operator in the framework of example 1.

1. The relaxation frequencies  $\nu_{\mathbb{A}}$  and  $\nu_{\mathbf{b}}$  being defined in (21), one may take any value for  $\nu$  with the condition  $\nu > \nu_{\mathbb{A}}, \nu_{\mathbf{b}}$ .

2. The relaxed moment  $\mathbf{L}(\boldsymbol{\rho}_f)$  writes

$$\mathbf{L}(\boldsymbol{\rho}_f) = (n, \mathbf{0}, 0, \lambda_{\mathbb{A}} n \mathbb{A}, \lambda_{\mathbb{B}} n \mathbf{b}), \quad \lambda_{\mathbb{A}} = 1 - \frac{\nu_{\mathbb{A}}}{\nu}, \quad \lambda_{\mathbb{B}} = 1 - \frac{\nu_{\mathbb{B}}}{\nu}. \quad (75)$$

With the above value of  $\nu$ ,  $\mathbf{L}(\boldsymbol{\rho}_f)$  is still realizable from proposition 8

3. Choose a  $\phi$ -divergence function satisfying the properties of Theorem 14. Then replace in Theorem 14,  $\boldsymbol{\rho}$  with  $\mathbf{L}(\boldsymbol{\rho}_f)$ . Then  $G$  is defined as  $\mathcal{M}(\phi^*)(\boldsymbol{\alpha} \cdot \mathbf{m})$
4. The BGK operator reads as  $K(f) = \nu(G - f)$ .

### 5.1.2 Properties of the model

Remark firstly that  $G = \mathcal{M}(\phi^*)(\boldsymbol{\alpha} \cdot \mathbf{m})$  is nonnegative. Thus the solution  $f$  to (7), if it exists, is nonnegative as well. Next the Grad space satisfies the necessary conditions of proposition 10 as concerns Galilean invariance of the modeling equation (7). Then Galilean invariance holds according to

**Proposition 17.** For  $\tau \in \{\tau_{\mathbf{u}}; \tau_{\theta}\}$   $\tau G(f) = G(\tau f)$

Let us now prove the (full) H-theorem.

**Theorem 18.** [H-theorem] Recall that

$$\mathcal{H}(f) = \int_{\mathbb{R}^3} \mathcal{M} \phi(f/\mathcal{M}) d\mathbf{v}.$$

Then there hold

$$\forall f \geq 0 \in \mathbb{L}^1(\mathbf{a}), \quad \left\langle K(f) \phi' \left( \frac{f}{\mathcal{M}} \right) \right\rangle \leq 0.$$

together with the characterization of equilibrium

$$K(f) = 0 \Leftrightarrow \left\langle K(f) \phi' \left( \frac{f}{\mathcal{M}} \right) \right\rangle = 0 \Leftrightarrow f = \mathcal{M}.$$

It must be denoted that the condition  $\nu > \nu_{\mathbb{A}}, \nu_{\mathbb{B}}$  is necessary in order to obtain the above results. This may be easily understood if one considers the relations between the moments of  $G$  and those of  $f$  (27) which only equate at  $\boldsymbol{\rho}_G = \boldsymbol{\rho}_f = \boldsymbol{\rho}_{\mathcal{M}}$ . Let us now derive proposition 1 in the case of relaxation in Grad space together with other properties stated in section 2.1 (the proof easily extends to the general case).

**Proposition 19.** For  $K(f)$  derived in section 5.1.1, the linearized operator  $\mathcal{L}$  defined in (6) reads as

$$\mathcal{L}(g) = \nu \left( \mathcal{P}_{\mathbb{K}} - \mathcal{I} + \left(1 - \frac{\nu_{\mathbb{A}}}{\nu}\right) \mathcal{P}_{\mathbb{A}} + \left(1 - \frac{\nu_{\mathbb{B}}}{\nu}\right) \mathcal{P}_{\mathbb{B}} \right). \quad (76)$$

As a consequence there holds

1. The kernel of the operator  $\mathcal{L}$  is exactly  $\mathbb{K}$  and there is also

$$\forall f, \quad \left[ \int K(f) \phi = 0 \right] \Leftrightarrow \phi \in \mathbb{K}$$

2. The operator is Fredholm, self-adjoint and negative on  $\mathbb{K}^\perp$
3. The diffusion coefficients in the Navier-Stokes limit of (7) are given by (21).

Remark that the second item is easily seen from (76). The third one is already proved in example 1.

## 5.2 The general case

We now consider a polynomial space  $\mathbb{P}$  satisfying the condition in proposition 10 and containing  $\mathbb{P}_{Grad} \subset \mathbb{P}$ .  $\mathbb{P}$  being invariant under translations and rotations, we may write

$$\mathbb{P} = \mathbb{K} \oplus^\perp m_6 \oplus^\perp \dots \oplus^\perp m_q. \quad (77)$$

where the polynomials  $(m_i)_i$  are defined in (10). If the collision invariants are written in an orthogonal basis for the  $L^2(\mathcal{M})$  dot product, then the vectors  $\boldsymbol{\rho}_f$ ,  $\boldsymbol{\rho}_{\mathcal{M}}$  and  $\boldsymbol{\rho}_G = L(\boldsymbol{\rho}_f)$  in the above decomposition read respectively

$$\begin{aligned} \boldsymbol{\rho}_f &= (n, \mathbf{0}, 0, \rho_6, \dots, \rho_q)^T, \\ \boldsymbol{\rho}_{\mathcal{M}} &= (n, \mathbf{0}, 0, 0, \dots, 0)^T, \\ \boldsymbol{\rho}_G &= \left( n, \mathbf{0}, 0, \left(1 - \frac{\nu_6}{\nu}\right)\rho_6, \dots, \left(1 - \frac{\nu_q}{\nu}\right)\rho_q \right)^T. \end{aligned}$$

The criteria of realizability of  $\boldsymbol{\rho}_G$  through symmetric positive definite moment matrix (see section 3.3) is unsatisfactory for two reasons : 1 - if  $\mathbb{P}$  is a quadratic space (definition 4), it is difficult to express the moment matrix corresponding to  $\boldsymbol{\rho}_G$  without knowing explicitly the eigenfunctions  $(m_i)_i$  and the corresponding eigenvalues  $(\nu_i)_i$ , 2 - even if it would be possible this criteria is not sufficient if all positive polynomials in  $\mathbb{P}$  are not sum of square polynomials.

So, for want of anything better we make the following assumption : the solution to (7) is such that the ball  $B(\boldsymbol{\rho}_{\mathcal{M}}, r)$  of radius  $r = \|\boldsymbol{\rho}_f - \boldsymbol{\rho}_{\mathcal{M}}\|$  stays in  $\mathcal{R}_{\mathbf{m}}^{*+}$ . In this case, it is easily seen that  $\boldsymbol{\rho}_G \in \mathcal{R}_{\mathbf{m}}^{*+}$  if  $\nu > \nu_i, \forall i$ .

Let us consider again a  $\phi$ -divergence function satisfying the properties of Theorem 14. The question whether

$$\int K(f)\phi' \left( \frac{f}{\mathcal{M}} \right) d\mathbf{v} \leq 0 \quad (78)$$

holds or not is for the moment an open problem. However the characterization of local equilibrium is easily found since  $K(f) = 0$  if and only if  $\boldsymbol{\rho}_f = \boldsymbol{\rho}_G$  which occurs only at  $G = \mathcal{M}$ . All other properties in the preceding section are satisfied. In particular, in the Chapman-Enskog expansion, one still finds that the solution satisfies the Euler equation in  $O(\varepsilon)$  while the Navier-Stokes equation is obtained with the right viscosity and heat conductivity just by using the definition of (10) in (21).

## 5.3 Some known models

### 5.3.1 BGK and ESBGK models

Let  $\phi(x) = x \ln(x)$ . If one just considers the conservation laws (13) together with the relaxation equation (23), then the variational problem is well-posed for  $-\frac{1}{2} \leq 1 - \frac{\nu_\Delta}{\nu} \leq$

1. Indeed on one hand the of constraints (33) is non empty according to proposition 8 and on the other hand the domain of  $h^*$  is non empty and open [30].  $\nu = \nu_{\mathbb{A}}$  with  $\nu_{\mathbb{A}}$  defined in (21) gives the well-known BGK operator [7]. Remark in this case that (24) is satisfied for  $\nu_{\mathbf{b},BGK} = \nu$ . So while the right viscosity is recovered in the hydrodynamic limit, the heat conductivity  $\kappa_{BGK}$  is such that the Prandtl number

$$Pr = \frac{5}{2} R \frac{\mu_B}{\kappa_{BGK}} = \frac{\nu_{\mathbf{b},BGK}}{\nu_{\mathbb{A}}} = 1.$$

More generally, for  $0 \leq \frac{\nu_{\mathbb{A}}}{\nu} \leq \frac{3}{2}$ , the solution to the variational problem always satisfies  $\int G \mathbf{b}(\mathbf{v} - \mathbf{u}) d\mathbf{v} = 0$ . The ESBGK model is then found by letting  $\nu = \nu_{\mathbf{b}}$  with  $\nu_{\mathbf{b}}$  defined in (21) which corresponds to the limit  $\frac{\nu_{\mathbb{A}}}{\nu} = \frac{3}{2}$  and  $Pr = \frac{2}{3}$ .

### 5.3.2 Shakhov model

Let now  $\phi(x) = \frac{1}{2}(x - 1)^2$ . In the Grad space, one considers the system (13, 23, 24). Assume that  $\mu_B$  and  $\kappa_B$  are either given by the exact computations in (21, 22) or by using some approximations of them. Let  $\nu = \nu_{\mathbb{A}} = \frac{n k_B T}{\mu_B}$  and  $\nu_{\mathbf{b}} = \frac{5}{2} \frac{n k_B^2 T}{m \kappa_B}$ . Remark that

$$Pr = \frac{5}{2} R \frac{\mu_B}{\kappa_B} = \frac{\nu_{\mathbf{b}}}{\nu_{\mathbb{A}}} = \frac{\nu_{\mathbf{b}}}{\nu}.$$

Then the system (13, 23, 24) together with the minimization problem give

$$G_S = \mathcal{M} \left( 1 + \frac{1 - Pr}{5} \frac{m}{n(k_B T)^2} \mathbf{q} \cdot (\mathbf{v} - \mathbf{u}) \left( m \frac{(\mathbf{v} - \mathbf{u})^2}{k_B T} - 5 \right) \right),$$

where  $\mathbf{q}$  is the heat flux defined by

$$\mathbf{q} = \frac{1}{2} m \int_{\mathbb{R}^3} f(\mathbf{v} - \mathbf{u})(\mathbf{v} - \mathbf{u})^2 d\mathbf{v}.$$

Originally  $G_S$  was computed in such a way that

$$\int_{\mathbb{R}^3} \mathbf{a}(\mathbf{v} - \mathbf{u}) \nu (G_S - f) d\mathbf{v} = \int_{\mathbb{R}^3} \mathbf{a}(\mathbf{v} - \mathbf{u}) Q(f, f) d\mathbf{v} \quad (79)$$

for Maxwell molecules and then adapted to other types of molecular interaction by introducing  $Pr$  into the definition of  $G_S$ . In the later case, the above equation is not valid.

The generalization through the diagonalization in (10) is easily performed by letting

$$G = \mathcal{M} \left( 1 + \sum_i \left( 1 - \frac{\nu_i}{\nu} \right) \mathcal{P}_{m_i}(g) \right),$$

where  $g = f/\mathcal{M} - 1$  and  $\nu > \nu_i, \forall i$ . In both cases  $G$  is not nonnegative. However one must point out that  $H(f) = \int \mathcal{M}\phi(f/\mathcal{M}) d\mathbf{v}$  is the natural entropy related to the whole method since

$$\int_{\mathbb{R}^3} \nu(G_S - f)\phi'\left(\frac{f}{\mathcal{M}} - 1\right) d\mathbf{v} = \langle \mathcal{L}(g), g \rangle \leq 0$$

where  $g = f/\mathcal{M} - 1$  and  $\mathcal{L}$  is defined in (15). It might happen that  $g \notin L^2(\mathcal{M})$  in which case the above value is  $-\infty$ . However  $\langle \mathcal{L}(g), g \rangle = 0$  only for  $g = 0$  or equivalently  $f = \mathcal{M}$ . Every other properties of section 2.1 are satisfied except the nonnegativity of  $G$ .

### 5.3.3 Levermore's operator

The analysis of the Chapman-Enskog expansion for moment system of the Boltzmann equation shows that wrong diffusion coefficients are obtained at the Navier-Stokes level [34]. Levermore has then proposed to substitute to the collision operator  $Q(f, f)$  a sum of relaxation operators constructed as follows. Let  $\mathbb{K} = \mathbb{M}_1 \subset \mathbb{M}_2 \subset \dots \subset \mathbb{M}_N$  and  $0 < \eta_1 < \eta_2 < \dots < \eta_N$ . Set

$$\mathcal{M}_k = \text{Argmin} \left\{ \int g \ln(g) / \int gp(\mathbf{v}) d\mathbf{v} = \int fp(\mathbf{v}) d\mathbf{v}, \forall p \in \mathbb{M}_k \right\}, \quad (80)$$

then  $K_{Lev}(f)$  writes

$$K_{Lev}(f) = \eta_1(\mathcal{M} - f) + \sum_{k=2}^N (\eta_k - \eta_{k-1})(\mathcal{M}_k - f).$$

Due to the assumption on each relaxation frequencies  $\nu_i$ , it is clear that  $K_{Lev}(f)$  preserves positivity, together with conservation laws. Also  $\int f \ln(f) d\mathbf{v}$  is the entropy in the non homogeneous equation (7). The linearized operator reads as

$$\mathcal{L}_{Lev} = - \sum_{k=1}^{N-1} \eta_k (\mathcal{P}_{k+1} - \mathcal{P}_k) + \eta_N (\mathcal{P}_N - \mathcal{I}),$$

where  $\mathcal{P}_k$  is the orthogonal projection onto  $\mathbb{M}_k$  in  $L^2(\mathcal{M})$ . Denoting with  $(m_{i,k})_{i=1, \dots, D_k}$  an orthogonal basis of  $\mathbb{M}_k \cap \mathbb{M}_{k-1}^\perp$ , there is

$$\begin{aligned} \mathcal{L}_{Lev} &= \eta_N (\mathcal{P}_N - \mathcal{I}) - \sum_{k=1}^{N-1} \eta_k \sum_{i=1}^{D_{k+1}} \mathcal{P}_{m_{i,k+1}} \\ &= \eta_N \left( (\mathcal{P}_{\mathbb{K}} - \mathcal{I}) + \sum_{k=1}^{N-1} \left( 1 - \frac{\eta_k}{\eta_N} \right) \sum_{i=1}^{D_{k+1}} \mathcal{P}_{m_{i,k+1}} \right) \end{aligned}$$

which has a form similar to (15). However there are many problems related to this construction. Junk was the first to point out that the solution to the variational problem might not satisfy all constraints [30]. Also, if  $\mathcal{L}_{Lev} = \mathcal{L}$  in (15), the solution to the variational problem in (80) may not exist as shows the simple case of Maxwell molecules. Indeed, some spaces  $\mathbb{M}_k$  have a maximal odd degree in which case there exists no solution in (80). Remark finally that  $K_{Lev}(f)$  does not satisfy relaxation equations such as (14). Indeed, for  $p \geq 3$  and  $1 \leq i \leq D_p$ , there is

$$\begin{aligned} \int K_{Lev}(f) m_{i,p} d\mathbf{v} &= \eta_1 \int (\mathcal{M} - f) m_{i,p} d\mathbf{v} + \sum_{k=2}^{p-1} (\eta_k - \eta_{k-1}) \int (\mathcal{M}_k - f) m_{i,p} d\mathbf{v} \\ &= -\eta_{p-1} \int f m_{i,p} d\mathbf{v} + \sum_{k=2}^{p-1} (\eta_k - \eta_{k-1}) \int \mathcal{M}_k m_{i,p} d\mathbf{v}. \end{aligned}$$

But for  $2 \leq k \leq p-1$ ,  $\int \mathcal{M}_k m_{i,p} d\mathbf{v}$  is not related to  $\int f m_{i,p} d\mathbf{v}$  in the minimization problem (80). Also, it does not vanish except if  $\mathcal{M}_k = \mathcal{M}q(\mathbf{v})$  for some polynomial  $q(\mathbf{v}) \in \mathbb{M}_k$  in which case the functional to be minimized in (80) is the one of the previous section.

The minimization problem in (80) can be fixed by using  $\phi$ -divergence such as in Theorem 14 since then the solution exists whatever the parity of the highest degree of the polynomials in the constraints. The operator satisfies by construction (78) and the characterization (4) follows under the sufficient condition  $\mathbb{M}_1 = \mathbb{K}$ . Thus, in the general case of section 5.2, the model is well-defined. Notice again that relations (14) still not hold so that it cannot be used in practice, especially in the context of moment systems for which it was originally designed.

## 6 Proofs

### 6.1 Proofs of the section 3

*Proof.* (Theorem 5) Let  $x \in \text{int}(C)$ . Then there exists  $\varepsilon > 0$  s.t.  $B(x, \varepsilon) \subset \text{int}(C)$ . So  $\forall y \in C^0$ ,  $y \neq 0$ ,  $x \cdot y \leq 0$ . But if there exists  $y \neq 0$  s.t.  $x \cdot y = 0$ , then by introducing  $z = x + \varepsilon \frac{y}{\|y\|}$ , we get  $z \cdot y > 0$ . But as  $z \in B(x, \varepsilon)$ , we get a contradiction.

Conversely, let  $x_0 \in C$  s.t.  $(\forall y \in C^0, y \neq 0)$ ,  $x_0 \cdot y < 0$ . Consider the linear form:  $y \mapsto x_0 \cdot y$ . Hence, by compactness of the unit sphere, we get

$$\sup_{y \in C^0, \|y\|=1} x_0 \cdot y = -\alpha < 0.$$

Therefore,  $(\forall y \in C^0)$ ,  $x_0 \cdot y \leq -\alpha \|y\|$ . Then  $\forall x \in B(x_0, \frac{\alpha}{2})$  and  $\forall y \in C^0$ ,

$$x \cdot y \leq (x - x_0) \cdot y + x_0 \cdot y \leq \frac{\alpha}{2} \|y\| - \alpha \|y\| \leq -\frac{\alpha}{2} \|y\|.$$



Then  $\forall x \in B(x_0, \frac{\alpha}{2}), \forall y \in C^0 \setminus \{0\}, x \cdot y < 0$ . By recalling that  $(C^0)^0 = \overline{C}$  and that  $(C^0)^0 = \{x \in \mathbb{R}^q, \text{ s.t. } \forall y \in C^0, x \cdot y \leq 0\}$ , we deduce that  $B(x_0, \frac{\alpha}{2}) \subset \overline{C}$  i.e.  $x_0 \in \text{int}(C)$ .  $\square$

*Proof.* (Proposition 7) Let  $q$  be the dimension of the space generated by  $\mathbf{m}$ . Then let us prove that there exists  $\mathbf{x}_1, \dots, \mathbf{x}_q$  such that the family  $\mathbf{m}(\mathbf{x}_k), k \in [1, q]$  is independent. It is obvious first that there exists  $\mathbf{x}_1$  such that  $\mathbf{m}(\mathbf{x}_1) \neq \mathbf{0}$ . Otherwise, for any  $\boldsymbol{\gamma}$  and for any  $\mathbf{x}$  there is  $\boldsymbol{\gamma} \cdot \mathbf{m}(\mathbf{v}) = 0$  and the family  $\mathbf{m}(\mathbf{v})$  cannot be a basis. That being said, there exists  $\mathbf{x}_2$  such that  $\mathbf{m}(\mathbf{x}_1), \mathbf{m}(\mathbf{x}_2)$  is independent. If we assume the contrary, this means that for any  $\mathbf{x}$  there is  $\lambda(\mathbf{x})$  such that  $\mathbf{m}(\mathbf{x}) = \lambda(\mathbf{x}) \mathbf{m}(\mathbf{x}_1)$ . Choose now  $\boldsymbol{\gamma}$  not zero orthogonal to  $\mathbf{m}(\mathbf{x}_1)$ . Then we have for any  $\mathbf{x}, \boldsymbol{\gamma} \cdot \mathbf{m}(\mathbf{x}) = \lambda(\mathbf{x}) \boldsymbol{\gamma} \cdot \mathbf{m}(\mathbf{x}_1) = 0$ . But this is a contradiction. By induction, if  $\mathbf{m}(\mathbf{x}_1), \dots, \mathbf{m}(\mathbf{x}_k)$  is independent such that  $k < q$ , we can always find  $\boldsymbol{\gamma}$  not zero in the orthogonal part. So there exists a  $\mathbf{x}_{k+1}$  such that the family  $\mathbf{m}(\mathbf{x}_1), \dots, \mathbf{m}(\mathbf{x}_{k+1})$  is independent. If  $C$  has an empty topological interior, then it is contained in an hyperplane and the subtraction of 2 elements of  $C$  is contained in an hyperplane containing 0. Here we just prove that  $\mathbb{R}^q = C - C$ . For this let us consider the family  $\mathbf{m}(\mathbf{x}_k), k \in [1, q]$ . By the former lemma it is (independent so also) generating. Then any  $\boldsymbol{\rho} \in \mathbb{R}^q$  can be written by a linear combination:

$$\boldsymbol{\rho} = \sum_{k=1}^{k=q} \lambda_k \mathbf{m}(\mathbf{x}_k).$$

Just let then write  $\lambda_k = \lambda_k^+ - \lambda_k^-$  with  $x^+, x^-$  the positive and negative part of any real  $x$ . Then we have immediately:

$$\boldsymbol{\rho} = \sum_{k=1}^{k=q} \lambda_k^+ \mathbf{m}(\mathbf{x}_k) - \sum_{k=1}^{k=q} \lambda_k^- \mathbf{m}(\mathbf{x}_k)$$

so proving  $\mathbb{R}^q = C - C$ .  $\square$

*Proof.* (Lemma 1) Let  $\mathbf{v} \neq \mathbf{0}$  and  $\alpha$  such that  $\alpha + \mathbf{u} \cdot \mathbf{v} = 0$ . A direct calculation gives

$$[\alpha \ \mathbf{v}] \begin{bmatrix} 1 & \mathbf{u} \\ \mathbf{u} & \mathbb{D} \end{bmatrix} \begin{bmatrix} \alpha \\ \mathbf{v} \end{bmatrix} = \alpha^2 + 2\alpha \mathbf{u} \cdot \mathbf{v} + \mathbb{D} : \mathbf{v} \otimes \mathbf{v} > 0$$

By factorization there holds

$$\alpha^2 + 2\alpha \mathbf{u} \cdot \mathbf{v} + \mathbb{D} : \mathbf{v} \otimes \mathbf{v} = (\alpha + \mathbf{u} \cdot \mathbf{v})^2 + (\mathbb{D} - \mathbf{u} \otimes \mathbf{u}) : \mathbf{v} \otimes \mathbf{v} > 0. \quad (81)$$

Then as  $\alpha + \mathbf{u} \cdot \mathbf{v} = 0$ , it comes that  $\mathbb{D} - \mathbf{u} \otimes \mathbf{u}$  is positive.

The converse statement is straightforward. If  $\mathbb{D} - \mathbf{u} \otimes \mathbf{u}$  is positive, choose  $\alpha, \mathbf{v} \neq \mathbf{0}$ . Since  $(\mathbb{D} - \mathbf{u} \otimes \mathbf{u})$  is positive, then RHS of (81) is always non negative. It is zero if and only if both  $\alpha + \mathbf{u} \cdot \mathbf{v}$  and  $(\mathbb{D} - \mathbf{u} \otimes \mathbf{u}) : \mathbf{v} \otimes \mathbf{v}$  are zero. From positiveness of  $\mathbb{D} - \mathbf{u} \otimes \mathbf{u}$  we get  $\mathbf{v} = \mathbf{0}$ . So  $\alpha = 0$ .  $\square$

Finally we prove quickly Proposition 8.

*Proof.* (Proposition 8). If  $(n, \mathbf{0}, ne, n\Pi, n\mathbf{Q})$  is realizable, then

$$\int_{\mathbb{R}^3} (\mathbf{v} - \mathbf{u}) \otimes (\mathbf{v} - \mathbf{u}) f \, d\mathbf{v} = n\Pi + n \frac{k_B T}{m} I_d$$

is SPD. We proceed as in [12]. Let  $\Theta_1, \Theta_2, \Theta_3$  the eigenvalues of  $\Pi$ . Then the eigenvalues of  $\lambda_{\mathbb{A}}\Pi + (1 - \lambda_{\mathbb{A}})n \frac{k_B T}{m} I_d$  are

$$\begin{aligned} \frac{1 + 2\lambda_{\mathbb{A}}}{3}\Theta_1 + (1 - \lambda_{\mathbb{A}})\frac{\Theta_2}{3} + (1 - \lambda_{\mathbb{A}})\frac{\Theta_3}{3}, & \quad (1 - \lambda_{\mathbb{A}})\frac{\Theta_1}{3} + \frac{1 + 2\lambda_{\mathbb{A}}}{3}\Theta_2 + (1 - \lambda_{\mathbb{A}})\frac{\Theta_3}{3}, \\ & \quad (1 - \lambda_{\mathbb{A}})\frac{\Theta_1}{3} + (1 - \lambda_{\mathbb{A}})\frac{\Theta_2}{3} + \frac{1 + 2\lambda_{\mathbb{A}}}{3}\Theta_3. \end{aligned}$$

So  $\lambda_{\mathbb{A}}\Pi + (1 - \lambda_{\mathbb{A}})n \frac{k_B T}{m} I_d$  is SPD for  $\lambda_{\mathbb{A}} \in [-\frac{1}{2}, 1]$  Hence, the relaxed moment  $(n, \mathbf{0}, ne, n\lambda_{\mathbb{A}}\Pi, n\lambda_{\mathbb{B}}\mathbf{Q})$  is realizable.  $\square$

*Proof.* (Proposition 10). Let  $\boldsymbol{\rho} \in \mathcal{R}_{\mathbf{m}}^+$  and  $f \in \mathbb{L}^{1,*}(\mathbf{m})$ ,  $f \geq 0$  s.t.  $R[f] = \boldsymbol{\rho}$ . Let  $\mathbf{u} \in \mathbb{R}^3$  and assume that

$$G(R[\tau_{\mathbf{u}}f]) = \tau_{\mathbf{u}}G(R[f]) = \tau_{\mathbf{u}}G(\boldsymbol{\rho}). \quad (82)$$

Hence there exists a relation between  $\boldsymbol{\rho}$  and  $R[\tau_{\mathbf{u}}f]$ . Remark that the relation defined by (82) does not depend on  $f$  as soon  $f \in R^{-1}(\boldsymbol{\rho})$ . Thus, the application  $L = R \circ \tau_{\mathbf{u}} \circ R^{-1} : \mathcal{R}_{\mathbf{m}}^+ \rightarrow \mathcal{R}_{\mathbf{m}}^+$  is well defined as soon as  $\tau_{\mathbf{u}}f \in \mathbb{L}^{1,+}(\mathbf{m})$  if  $f \in \mathbb{L}^{1,+}(\mathbf{m})$ .

Under this condition,  $R^{-1}$  defines a linear map from  $\mathcal{R}_{\mathbf{m}}^+$  into subsets of  $\mathbb{L}^{1,*}(\mathbf{m})$  as follows:

$$(\forall \lambda \geq 0), R(R^{-1}(\boldsymbol{\rho}_1) + \lambda R^{-1}(\boldsymbol{\rho}_2)) = \int (f + \lambda g)\mathbf{m}(\mathbf{v})d\mathbf{v} = \boldsymbol{\rho}_1 + \lambda\boldsymbol{\rho}_2$$

which is equivalent to  $R^{-1}(\boldsymbol{\rho}_1 + \lambda\boldsymbol{\rho}_2) = R^{-1}(\boldsymbol{\rho}_1) + \lambda R^{-1}(\boldsymbol{\rho}_2)$ . Hence  $L = R \circ \tau_{\mathbf{u}} \circ R^{-1} : \mathcal{R}_{\mathbf{m}}^+ \rightarrow \mathcal{R}_{\mathbf{m}}^+$  is linear.

Now remark that this relation can be extended to  $\mathbb{R}^q$ . Indeed,  $\mathcal{R}_{\mathbf{m}}^{+*}$  is an open and non void set in  $\mathbb{R}^q$  and contains  $q$  independent vectors  $(\boldsymbol{\rho}_1, \dots, \boldsymbol{\rho}_q)$  which form a basis of  $\mathbb{R}^q$ . Otherwise  $\mathcal{R}_{\mathbf{m}}^{+*}$  would be contained in an hyperplane. Thus, the linear application  $L$  is well defined on  $\mathbb{R}^q$  entirely.

$\forall \boldsymbol{\rho} \in \mathbb{R}^q, \forall f \in R^{-1}(\boldsymbol{\rho})$ , there is  $L\boldsymbol{\rho} = \int f L\mathbf{m}(\mathbf{v}) \, d\mathbf{v}$  and  $L\boldsymbol{\rho} = R \circ \tau_{\mathbf{u}}f$  at the same time i.e.

$$\int_{\mathbb{R}^3} \tau_{\mathbf{u}}f(\mathbf{v})\mathbf{m}(\mathbf{v}) \, d\mathbf{v} = \int_{\mathbb{R}^3} f(\mathbf{w}) \mathbf{m}(\mathbf{w} - \mathbf{u}) \, d\mathbf{w} = \int_{\mathbb{R}^3} f(\mathbf{w}) L\mathbf{m}(\mathbf{v}) \, d\mathbf{v}.$$

As  $Im(R^{-1}(\mathbb{R}^q)) = \mathbb{L}^1(\mathbf{m})$ , the previous relation must be true for any  $f \in \mathbb{L}^1(\mathbf{m})$ . Hence  $\mathbb{P} = \text{span}(\mathbf{m})$  must be invariant under the action of  $\tau_{-\mathbf{u}}$  and  $L = \Lambda(-\mathbf{u})$  as defined in (45). We can proceed in the same way for any  $\mathbf{u} \in \mathbb{R}^3$  and  $\Theta \in SO(3)$ , we deduce that  $\mathbb{P}$  is invariant by Galilean transforms.  $\square$

## 6.2 Proof of section 4

*Proof.* (Lemma 2). Let us prove any of the claimings:

1. First let us remark that from the definition of  $\phi$ -divergence, there holds the following property: since  $\phi$  is strictly convex and that there is  $\phi(0) = 0$  then the following function

$$\forall y \in (0, +\infty), \quad y \mapsto \frac{\phi(y)}{y} \in (p_0, +\infty)$$

is strictly increasing and one to one from  $(0, +\infty)$  onto  $(p_0, +\infty)$ . For any  $p \leq p_0$  and for any  $y > 0$  we have

$$p \leq p_0 < \frac{\phi(y)}{y}.$$

Then using  $\phi(0)$  we have  $\forall p \leq p_0, \forall y \geq 0, \quad py - \phi(y) \leq 0$ . Hence, taking the supremum on  $y \in \text{dom}(\phi) = [0, +\infty)$  there holds  $\forall p \leq p_0, \phi^*(p) \leq 0$ . But on the other hand, since the function  $\phi$  is convex and semi lower continuous, there holds  $\phi^{**} = \phi$ . As a consequence,

$$\phi(0) = - \inf_{p \in \mathbb{R}} \phi^*(p) = 0.$$

This means that  $\inf_{p \in \mathbb{R}} \phi^*(p) = 0$ . So we have  $\forall p \in \mathbb{R}, \phi^*(p) \geq 0, \forall p \leq p_0, \phi^*(p) \leq 0$ . So  $\forall p \leq p_0, \phi^*(p) = 0$ .

2. Now let us prove the non negativity of  $\phi^*$ . Since the function  $\phi(y)/y$  is strictly increasing and one to one, it is also continuous (characterization of bijection on intervals). For  $p \in (p_0, +\infty)$  there is just one element  $y_p \in (0, +\infty)$  such that  $py_p - \phi(y_p) = 0$ . Then for any  $y \in (0, y_p)$  we have  $py_p - \phi(y_p) > 0$  and for any  $y > y_p$  there holds  $py_p - \phi(y_p) < 0$ . In particular, there holds

$$\forall p > p_0, \quad \sup_{y \in \text{dom}(\phi)} (yp - \phi(y)) = \sup_{y \in [0, y_p]} (yp - \phi(y)) \geq 0.$$

But any semi-upper continuous function gets its supremum on a compact set (and  $yp - \phi(y)$  is semi-upper continuous). Then there exists  $z_p \in [0, y_p]$  such that

$$\forall p > p_0, \quad \phi^*(p) = z_p p - \phi(z_p) \geq 0, \quad \phi^*(p) = z_p p - \phi(z_p) \in \mathbb{R}.$$

Finally, since for all  $y \in (0, y_p), yp - \phi(y) > 0$  the supremum is of course  $> 0$ .

3. From the former property,  $\text{dom}(\phi^*) = \mathbb{R}$ . Then  $\phi^*$  is continuous on  $\mathbb{R}$  and at any point  $p \in \mathbb{R}$ , its sub-differential is not void. Assume that  $y_1 < y_2$  are in its sub-differential. Then, by the characterization of the sub-differential for  $\phi^*$ , and since  $\phi = \phi^{**}$  (semi-lower continuity)

$$\phi(y_1) = py_1 - \phi^*(p), \quad \phi(y_2) = py_2 - \phi^*(p)$$

this also means that  $p \in \partial\phi(y_1)$  and  $p \in \partial\phi(y_2)$ . So,  $\phi(y) - \phi(y_i) \geq p(y - y_i)$ . So

$$p(y_2 - y) \geq \phi(y_2) - \phi(y), \quad \phi(z) - \phi(y_1) \geq p(z - y_1).$$

For  $y = y_1$  and  $z = y_2$ , there holds  $p(y_2 - y_1) \geq \phi(y_2) - \phi(y_1) \geq p(y_2 - y_1)$ . So  $\phi(y_2) - \phi(y_1) = p(y_2 - y_1)$ . In particular, since  $\phi$  is convex, then for any  $\alpha \in [0, 1]$  and  $y = \alpha y_1 + (1 - \alpha) y_2$

$$\phi(y) \leq \alpha \phi(y_1) + (1 - \alpha) \phi(y_2) = \phi(y_2) + \alpha p(y_1 - y_2)$$

having  $\phi(y_2) = -\phi^*(p) + py_2$  we get  $\phi(y) \leq -\phi^*(p) + (1 - \alpha)py_2 + \alpha py_1 = -\phi^*(p) + py$ . This proves that  $p \in \partial\phi(y)$ . Using the same inequalities as above we have then

$$\forall y \in [y_1, y_2], \quad \phi(y) - \phi(y_1) = p(y - y_1).$$

So  $\phi$  is affine on  $[y_1, y_2]$  (with  $y_1 < y_2$ ) which contradicts that  $\phi$  it is strictly convex. Then for any  $p$  the sub-differential  $\partial\phi^*(p)$  has only one element. Then  $\phi^*$  is differentiable. Finally any convex function on  $\mathbb{R}$  which is differentiable is  $C^1$  smooth. □

*Proof.* (Proposition 15)

1. Define the entropy by

$$\forall g \in \mathbb{L}^1(\mathbf{m}), \quad \mathcal{H}(g) = \int \phi\left(\frac{g}{\mathcal{M}}\right) \mathcal{M} \in \mathbb{R} \cup \{+\infty\}.$$

First let us begin by assuming that  $g \geq 0$  almost everywhere.  $\phi$  is differentiable on  $[0, +\infty[$  strictly,  $\phi(0) = 0$  and  $\phi(p) \rightarrow +\infty$  imply that  $\phi$  is bounded from below. So

$$\mathcal{H}(g) \geq \min(\phi) \int \mathcal{M} d\mathbf{v}.$$

This proves that  $\mathcal{H}$  is bounded form below independently of  $g$ .

Now assume that  $g$  takes negative values on a non zero measure set of  $\mathbb{R}_+$  noted by  $\Omega$ . So  $\phi = +\infty$  on  $\Omega$ . Then

$$\mathcal{H}(g) = \int \phi\left(\frac{g}{\mathcal{M}}\right) \mathcal{M} = +\infty.$$

2. Finally  $\mathcal{H}$  is strictly convex on its domain comes thanks to the strict convexity of  $\phi$ .  
3. Now let us prove the rest of the proposition.

- Let  $\boldsymbol{\rho} \in \mathcal{R}_{\mathbf{m}}^+$ . Then there exists  $\Psi_{\boldsymbol{\rho}} \in \mathcal{C}_c^\infty(\mathbb{R}^3)$ ,  $\Psi_{\boldsymbol{\rho}} \geq 0$  s.t.  $\int \Psi_{\boldsymbol{\rho}} \mathbf{m} d\mathbf{v} = \boldsymbol{\rho}$ . Thus the set  $D^+(\boldsymbol{\rho})$  defined as

$$D^+(\boldsymbol{\rho}) = \{g \geq 0, \int g \mathbf{m} = \boldsymbol{\rho}, \mathcal{H}(g) \leq \mathcal{H}(\Psi_{\boldsymbol{\rho}})\}$$

is non empty and convex. Moreover,  $\inf \mathcal{H}$  exists on  $D^+(\boldsymbol{\rho})$  but is not necessarily attained by a function in  $D^+(\boldsymbol{\rho})$ . Then  $\mathcal{R}_{\mathbf{m}}^+ \subset \text{dom}(h_{\mathbf{m}})$ . If  $\boldsymbol{\rho} \notin \mathcal{R}_{\mathbf{m}}^+$ , then there

are no nonnegative function that realizes  $\boldsymbol{\rho}$ . Hence,  $h_{\mathbf{m}}(\boldsymbol{\rho}) = +\infty$  by definition of  $h_{\mathbf{m}}$ . So  $\text{dom}(h_{\mathbf{m}}) = \mathcal{R}_{\mathbf{m}}^+$ .

- It is very clear that any ball of  $\mathbb{R}^q$  which contains  $\mathbf{0}$  contains moment  $\boldsymbol{\rho} \notin \mathcal{R}_{\mathbf{m}}^+$ . Then any subset of  $\mathcal{R}_{\mathbf{m}}^+$  which contains  $\mathbf{0}$  is not open. On the other hand the set  $\mathcal{R}_{\mathbf{m}}^{+*}$  is open. Then it is obviously the biggest open set (in sense of inclusion) which is included in  $\mathcal{R}_{\mathbf{m}}^+$ . Then we have straightforwards:  $\text{int}(\text{dom}(h_{\mathbf{m}})) = \text{int}(\mathcal{R}_{\mathbf{m}}^+) = \mathcal{R}_{\mathbf{m}}^{+*}$ .
- Now we need to prove that the function  $h_{\mathbf{m}} : \mathbb{R}^d \mapsto \mathbb{R} \cup \{+\infty\}$  is convex. Let  $\rho_1, \rho_2 \in \mathcal{R}_{\mathbf{m}}^{+*}$ . Then  $(\forall \varepsilon > 0), \exists g_1, g_2 \in \mathbb{L}^1(\mathbf{m})$  with  $\int g_i \mathbf{m} d\mathbf{v} = \rho_i$  such that  $h(\rho_i) > \mathcal{H}(g_i) - \varepsilon$  for  $i \in \{1; 2\}$ . Thus,

$$(\forall \lambda \in [0, 1]) \lambda h_{\mathbf{m}}(\rho_1) + (1 - \lambda)h_{\mathbf{m}}(\rho_2) > \lambda \mathcal{H}(g_1) + (1 - \lambda)\mathcal{H}(g_2) - \varepsilon.$$

$\mathcal{H}$  being strictly convex it comes that

$$(\forall \lambda \in [0, 1]) \lambda h_{\mathbf{m}}(\rho_1) + (1 - \lambda)h_{\mathbf{m}}(\rho_2) > \mathcal{H}(\lambda g_1 + (1 - \lambda)g_2) - \varepsilon.$$

By definition of  $h_{\mathbf{m}}$  it holds that

$$(\forall \lambda \in [0, 1]) \mathcal{H}(\lambda g_1 + (1 - \lambda)g_2) \geq h_{\mathbf{m}}(\lambda \rho_1 + (1 - \lambda)\rho_2)$$

and the convexity of  $h_{\mathbf{m}}$  follows. □

*Proof.* (Theorem 16)

1. Let  $\boldsymbol{\rho} \in \mathbb{R}^q$  and recall that  $D(\boldsymbol{\rho}) = \{g \in \mathbb{L}^1(\mathbf{m}), \int \mathbf{m} g = \boldsymbol{\rho}\}$  (which is never the empty set). There is by definition:

$$h_{\mathbf{m}}(\boldsymbol{\rho}) = \inf_{g \in D(\boldsymbol{\rho})} \int \phi\left(\frac{g}{\mathcal{M}}\right) \mathcal{M}$$

Noting that for any  $g \in D(\boldsymbol{\rho})$  there holds  $\boldsymbol{\rho} = \int \mathbf{a} g$  we get

$$\boldsymbol{\alpha} \cdot \boldsymbol{\rho} + \sup_{g \in D(\boldsymbol{\rho})} \left[ - \int \phi\left(\frac{g}{\mathcal{M}}\right) \mathcal{M} \right] = \sup_{g \in D(\boldsymbol{\rho})} \left( \int \left[ \mathbf{m} \cdot \boldsymbol{\alpha} g - \phi\left(\frac{g}{\mathcal{M}}\right) \mathcal{M} \right] \right)$$

As  $\boldsymbol{\alpha} \cdot \boldsymbol{\rho} - h_{\mathbf{m}}(\boldsymbol{\rho}) = h_{\mathbf{m}}^*(\boldsymbol{\alpha})$  we have

$$h_{\mathbf{m}}^*(\boldsymbol{\alpha}) = \sup_{\boldsymbol{\rho} \in \mathbb{R}^q} \sup_{g \in D(\boldsymbol{\rho})} \left( \int \left[ \mathbf{m} \cdot \boldsymbol{\alpha} g - \phi\left(\frac{g}{\mathcal{M}}\right) \mathcal{M} \right] \right).$$

The next step consists to show that  $\sup_{\boldsymbol{\rho} \in \mathbb{R}^q}$  and  $\sup_{g \in D(\boldsymbol{\rho})}$  can be permuted in the previous formula.

Let  $\boldsymbol{\rho}$  be fixed. It is clear that  $\forall \boldsymbol{\rho} \in \mathbb{R}^q$

$$\sup_{g \in D(\boldsymbol{\rho})} \left( \int \left[ \mathbf{m} \cdot \boldsymbol{\alpha} g - \phi \left( \frac{g}{\mathcal{M}} \right) \mathcal{M} \right] \right) \leq \sup_{g \in \mathbb{L}^1} \left( \int \left[ \mathbf{m} \cdot \boldsymbol{\alpha} g - \phi \left( \frac{g}{\mathcal{M}} \right) \mathcal{M} \right] \right)$$

then obviously we have:

$$\sup_{\boldsymbol{\rho} \in \mathbb{R}^q} \sup_{g \in D(\boldsymbol{\rho})} \left( \int \left[ \mathbf{m} \cdot \boldsymbol{\alpha} g - \phi \left( \frac{g}{\mathcal{M}} \right) \mathcal{M} \right] \right) \leq \sup_{g \in \mathbb{L}^1} \left( \int \left[ \mathbf{m} \cdot \boldsymbol{\alpha} g - \phi \left( \frac{g}{\mathcal{M}} \right) \mathcal{M} \right] \right)$$

On the other hand, let  $g \in \mathbb{L}^1(\mathbf{m})$  and note by  $\boldsymbol{\rho}(g) = \int \mathbf{m} g \in \mathbb{R}^q$ . Then

$$\int \left[ \mathbf{m} \cdot \boldsymbol{\alpha} g - \phi \left( \frac{g}{\mathcal{M}} \right) \mathcal{M} \right] \leq \sup_{\psi \in D(\boldsymbol{\rho}(g))} \int \left[ \mathbf{m} \cdot \boldsymbol{\alpha} \psi - \phi \left( \frac{\psi}{\mathcal{M}} \right) \mathcal{M} \right].$$

We have then

$$\sup_{g \in \mathbb{L}^1(\mathbf{a})} \int \left[ \mathbf{m} \cdot \boldsymbol{\alpha} g - \phi \left( \frac{g}{\mathcal{M}} \right) \mathcal{M} \right] \leq \sup_{g \in \mathbb{L}^1(\mathbf{a})} \sup_{\psi \in D(\boldsymbol{\rho}(g))} \int \left[ \mathbf{m} \cdot \boldsymbol{\alpha} \psi - \phi \left( \frac{\psi}{\mathcal{M}} \right) \mathcal{M} \right].$$

But there is  $\mathbb{R}^q = \{ \int \mathbf{m} g, g \in \mathbb{L}^1 \}$ . Then

$$\sup_{g \in \mathbb{L}^1(\mathbf{a})} \sup_{\psi \in D(\boldsymbol{\rho}(g))} = \sup_{\boldsymbol{\rho} \in \mathbb{R}^q} \sup_{\psi \in D(\boldsymbol{\rho})}$$

and finally we have

$$\sup_{g \in \mathbb{L}^1(\mathbf{m})} \int \left[ \mathbf{m} \cdot \boldsymbol{\alpha} g - \phi \left( \frac{g}{\mathcal{M}} \right) \mathcal{M} \right] \leq \sup_{\boldsymbol{\rho} \in \mathbb{R}^q} \sup_{\psi \in D(\boldsymbol{\rho})} \int \left[ \mathbf{m} \cdot \boldsymbol{\alpha} \psi - \phi \left( \frac{\psi}{\mathcal{M}} \right) \mathcal{M} \right]$$

and we get formula (67).

2. Pick  $\boldsymbol{\alpha} \in \mathbb{R}^3$  and consider the polynomial  $\pi_{\boldsymbol{\alpha}} := \boldsymbol{\alpha} \cdot \mathbf{m}$  and  $G_{\boldsymbol{\alpha}} = \phi^{*'}(\pi_{\boldsymbol{\alpha}}) \mathcal{M}$ . From the characterization of sub-differential of  $\phi^*$  at (real) point  $G_{\boldsymbol{\alpha}}/\mathcal{M}$ , and taking into account that  $\phi$  is semi-lower continuous, we have

$$\phi \left( \frac{G_{\boldsymbol{\alpha}}}{\mathcal{M}} \right) + \phi^* (\pi_{\boldsymbol{\alpha}}) = \pi_{\boldsymbol{\alpha}} \frac{G_{\boldsymbol{\alpha}}}{\mathcal{M}}$$

Using Young inequality for  $\phi^*$  we have straightforwards

$$\phi^* (2\pi_{\boldsymbol{\alpha}}) - \phi^* (\pi_{\boldsymbol{\alpha}}) \geq \pi_{\boldsymbol{\alpha}} \phi^{*'} (\pi_{\boldsymbol{\alpha}}) = \pi_{\boldsymbol{\alpha}} \frac{G_{\boldsymbol{\alpha}}}{\mathcal{M}}.$$

Recall that we have for any convex function  $\phi : \mathbb{R} \mapsto \overline{\mathbb{R}}$

$$-\phi^* (0) = \inf_{y \in \mathbb{R}} \phi (y)$$

Then there holds immediately

$$-\phi^*(0) \leq \phi\left(\frac{G_{\alpha}}{\mathcal{M}}\right) = \pi_{\alpha} \frac{G_{\alpha}}{\mathcal{M}} - \phi^*(\pi_{\alpha}) \leq \phi^*(2\pi_{\alpha}) - 2\phi^*(\pi_{\alpha})$$

Since  $\pi_{\alpha}$  is a polynomial function, and since by the assumption 3 any  $\phi^*(\pi)$  is in  $\mathbb{L}^1(\mathcal{M}d\mathbf{v})$ , then the former inequalities prove that  $\phi\left(\frac{G_{\alpha}}{\mathcal{M}}\right) \in \mathbb{L}^1(\mathcal{M}d\mathbf{v})$ . Using then

$$\phi\left(\frac{G_{\alpha}}{\mathcal{M}}\right) + \phi^*(\pi_{\alpha}) = \pi_{\alpha} \frac{G_{\alpha}}{\mathcal{M}}$$

we deduce finally that  $\frac{G_{\alpha}}{\mathcal{M}}\pi_{\alpha}$  is also in  $\mathbb{L}^1(\mathcal{M}d\mathbf{v})$ . By Young inequality we have for any  $g$

$$\frac{g}{\mathcal{M}}\pi_{\alpha} - \phi\left(\frac{g}{\mathcal{M}}\right) \leq \phi^*(\pi_{\alpha}).$$

Multiplying by  $\mathcal{M}$  and integrating (any term can be computed) there holds:

$$\int \left[ \pi_{\alpha} g - \phi\left(\frac{g}{\mathcal{M}}\right) \mathcal{M} \right] \leq \int \phi^*(\pi_{\alpha}) \mathcal{M}.$$

By having the infimum:

$$\sup_{g \in \mathbb{L}^1(\mathbf{a})} \int \left[ \pi_{\alpha} g - \phi\left(\frac{g}{\mathcal{M}}\right) \mathcal{M} \right] \leq \int \phi^*(\pi_{\alpha}) \mathcal{M}$$

which gives us finally:

$$h_{\mathbf{m}}^*(\alpha) \leq \int \phi^*(\pi_{\alpha}) \mathcal{M} = \int \phi^*(\alpha \cdot \mathbf{m}) \mathcal{M}.$$

On the other hand, since we have

$$\phi^*(\pi_{\alpha}) \mathcal{M} = \left( \phi\left(\frac{G_{\alpha}}{\mathcal{M}}\right) - \frac{G_{\alpha}}{\mathcal{M}}\pi_{\alpha} \right) \mathcal{M}.$$

Any term can be integrated and by having integration we have:

$$\int \phi^*(\pi_{\alpha}) \mathcal{M} = \int \left( \phi\left(\frac{G_{\alpha}}{\mathcal{M}}\right) - \frac{G_{\alpha}}{\mathcal{M}}\pi_{\alpha} \right) \mathcal{M} \leq \sup_{g \in \mathbb{L}^1(\mathbf{a})} \int \left[ \pi_{\alpha} g - \phi\left(\frac{g}{\mathcal{M}}\right) \mathcal{M} \right]$$

Then we have exactly (68)

3. Finally, consider for  $\epsilon \in (0, 1]$  the function

$$f_{\epsilon}(\beta) = \left( \frac{\phi^*(\pi_{\alpha} + \epsilon\pi_{\beta}) - \phi^*(\pi_{\alpha})}{\epsilon} - \phi^{*\prime}(\pi_{\alpha})\pi_{\beta} \right) \mathcal{M}.$$

Young inequality gives  $\epsilon\phi^{*'}(\pi_\alpha)\pi_\beta\mathcal{M} \leq [\phi^*(\pi_\alpha + \epsilon\pi_\beta) - \phi^*(\pi_\alpha)]\mathcal{M}$ . So  $f_\epsilon(\beta) \geq 0$ . On the other hand, by using  $\pi_\alpha + \epsilon\pi_\beta = \epsilon(\pi_\alpha + \pi_\beta) + (1 - \epsilon)\pi_\alpha$ , the convexity of  $\phi^*$  gives

$$\phi^*(\pi_\alpha + \epsilon\pi_\beta) \leq \epsilon\phi^*(\pi_\alpha + \pi_\beta) + (1 - \epsilon)\phi^*(\pi_\alpha).$$

Hence we get

$$\frac{\phi^*(\pi_\alpha + \epsilon\pi_\beta) - \phi^*(\pi_\alpha)}{\epsilon} \leq \phi^*(\pi_\alpha + \pi_\beta) - \phi^*(\pi_\alpha).$$

So, after multiplication by  $\mathcal{M}$

$$f_\epsilon(\beta) \leq [\phi^*(\pi_\alpha + \pi_\beta) - \phi^*(\pi_\alpha) - \phi^{*'}(\pi_\alpha)\pi_\beta]\mathcal{M} = f_1$$

that is  $0 \leq f_\epsilon(\beta) \leq f_1(\beta)$ . Since  $\lim_{\epsilon \rightarrow 0} f_\epsilon(\beta) = 0$  pointwise, using Lebesgues dominated convergence theorem there holds:

$$\forall \beta, \lim_{\epsilon \rightarrow 0^+} \frac{h^*(\alpha + \epsilon\beta) - h^*(\alpha)}{\epsilon} = \int \phi^*(\pi_\alpha)\pi_\beta\mathcal{M} = \beta \cdot \left[ \int \phi^*(\alpha \cdot \mathbf{a})\mathbf{a}\mathcal{M} \right].$$

The convex function  $\alpha \in \mathbb{R}^q \mapsto h^*(\alpha)$  has partial derivatives in any direction at any point. Then, using a classical result of convex analysis, it is  $C^1$  smooth and we have (69). □

*Proof.* (Theorem 14) Let us finally prove Theorem 14 step by step.

1. We first prove that  $h_{\mathbf{m}}$  is closed in its domain.  $h_{\mathbf{m}}$  being convex, this amounts to prove that  $h_{\mathbf{m}}^{**} = h_{\mathbf{m}}$  in  $\mathcal{R}_{\mathbf{m}}^+$ . The delicate point is to prove this relation at  $\mathbf{0}$ . Let us observe the following: there is  $h_{\mathbf{m}}(\mathbf{0}) = 0$ . This is just because the only non negative function which is able to realize  $\mathbf{0}$  is  $g = 0$ . Since  $\phi(0) = 0$  then we deduce immediately that  $h_{\mathbf{m}}(\mathbf{0}) = 0$ . Let us now compute  $h_{\mathbf{m}}^{**}(\mathbf{0}) := \sup_{\alpha \in \mathbb{R}^q} (\mathbf{0} \cdot \alpha - h_{\mathbf{m}}^*(\alpha)) = -\inf_{\alpha \in \mathbb{R}^q} h_{\mathbf{m}}^*(\alpha) \leq 0$ . Remark then that  $h_{\mathbf{m}}^{**}(\mathbf{0}) \leq 0$  from the expression given of  $h^*$  (theorem 16). Recall from Lemma 2 that  $\phi^*(y) \geq 0$  and for  $y \leq p_0$ ,  $\phi^*(y) = 0$ . Chose now  $\alpha_0 = (y_0, \mathbf{0}, \dots, \mathbf{0})$ . There is  $\phi^*(\alpha_0 \cdot \mathbf{m}(\mathbf{v})) = 0$ . So  $h_{\mathbf{m}}^*(\alpha_0) = 0$  and finally  $h_{\mathbf{m}}^{**}(\mathbf{0}) = 0 = h_{\mathbf{m}}(\mathbf{0})$ .  
Finally,  $h_{\mathbf{m}}$  being convex,  $h_{\mathbf{m}}$  is continuous on  $\text{int}(\text{dom}(h_{\mathbf{m}}))$ . As a consequence there is  $h_{\mathbf{m}}^{**} = h_{\mathbf{m}}$  in  $\text{dom}(h_{\mathbf{m}}) = \mathcal{R}_{\mathbf{m}}^+$ .
2. Let  $\rho \in \text{int}(\text{dom}(h_{\mathbf{m}}))$ .  $h_{\mathbf{m}}$  being continuous at this point, there is  $\partial h_{\mathbf{m}}(\rho) \neq \emptyset$  [39]. Pick some  $\alpha \in \partial h(\rho)$ .  $h_{\mathbf{m}}$  being proper convex and closed at  $\rho$ , there is  $\rho \in \partial h_{\mathbf{m}}(\alpha)$ . But  $h_{\mathbf{m}}^*$  is  $C^1$  in  $\mathbb{R}^q$  so  $\partial h_{\mathbf{m}}(\alpha) = \nabla h_{\mathbf{m}}^*(\alpha)$  and  $\rho = \nabla h_{\mathbf{m}}^*(\alpha)$ .
3. Let  $\rho \in \text{int}(\text{dom}(h_{\mathbf{m}}))$  and  $\alpha \in \partial h_{\mathbf{m}}(\rho)$ . We prove that the function  $G = \mathcal{M}\phi^{*'}(\alpha \cdot \mathbf{m}(\mathbf{v}))$  is the unique solution to the primal problem. We firstly have

$$\rho = \nabla h_{\mathbf{m}}^*(\alpha) = \int \phi^{*'}(\alpha \cdot \mathbf{m}(\mathbf{v}))\mathbf{m}(\mathbf{v})\mathcal{M}d\mathbf{v}.$$



Next  $\phi^*$  is a  $C^1$  convex function in  $\mathbb{R}$  and thus

$$(\forall y \in \mathbb{R}), \phi^*(y) + \phi^{**}((\phi^*)'(y)) = y(\phi^*)'(y).$$

$\phi$  is also convex, proper and semi lower continuous. So  $\phi^{**} = \phi$  and

$$(\forall y \in \mathbb{R}), \phi^*(y) + \phi((\phi^*)'(y)) = y(\phi^*)'(y).$$

Put  $y = \boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v})$  in the above equation, multiply by  $\mathcal{M}$  and integrate w.r.t  $v$  gives

$$h_{\mathbf{m}}^*(\boldsymbol{\alpha}) + \int \mathcal{M}\phi((\phi^*)'(\boldsymbol{\alpha} \cdot \mathbf{m})) = \int \mathcal{M}(\boldsymbol{\alpha} \cdot \mathbf{m})(\phi^*)'(\boldsymbol{\alpha} \cdot \mathbf{m}) = \boldsymbol{\alpha} \cdot \boldsymbol{\rho}.$$

One then deduces from the subdifferential equation (64) that

$$h_{\mathbf{m}}(\boldsymbol{\rho}) = \int \mathcal{M}\phi((\phi^*)'(\boldsymbol{\alpha} \cdot \mathbf{m})) = \mathcal{H}(G).$$

Recall that  $\mathcal{H}$  is strictly convex and thus  $G$  is the unique solution to the primal problem. From the form  $G = \mathcal{M}\phi^{*'}(\boldsymbol{\alpha} \cdot \mathbf{m}(\mathbf{v}))$  which is necessarily strictly positive on a set of non-zero measure - that is for those velocities  $\mathbf{v}$  for which  $\boldsymbol{\alpha} \cdot \mathbf{m}(v) > p_0$  -  $\boldsymbol{\alpha}$  is found to be unique. This in turn proves that the subdifferential of  $h$  at interior point of  $\mathcal{R}_{\mathbf{m}}^+$  is reduced to one point. As a consequence  $\nabla h^*$  is a bijection from  $C^\circ$  to  $\mathcal{R}_{\mathbf{m}}^+$  where  $C^\circ$  is defined in (66). And there is  $\nabla h^*(\boldsymbol{\alpha}) = 0$  in the complementary set of  $C^\circ$ .

4. Let us finally prove that  $h_{\mathbf{m}}$  is strictly convex on  $\mathcal{R}_{\mathbf{m}}^+$ . Let  $\boldsymbol{\rho}_1 \neq \boldsymbol{\rho}_2 \in \mathcal{R}_{\mathbf{m}}^+$ . Consider  $\boldsymbol{\rho}(t) = (1-t)\boldsymbol{\rho}_1 + t\boldsymbol{\rho}_2$ ,  $t \in [0, 1]$  and  $\boldsymbol{\alpha}(t)$  s.t.  $\nabla h_{\mathbf{m}}^*(\boldsymbol{\alpha}(t)) = \boldsymbol{\rho}(t)$ . In particular, we have  $\boldsymbol{\rho}(0) = \boldsymbol{\rho}_1$  and  $\boldsymbol{\rho}(1) = \boldsymbol{\rho}_2$ . Define  $f_1$  and  $f_2$  by  $f_1 = \mathcal{M}(\phi^*)'(\boldsymbol{\alpha}(0) \cdot \mathbf{m})$  and  $f_2 = \mathcal{M}(\phi^*)'(\boldsymbol{\alpha}(1) \cdot \mathbf{m})$ . They satisfy the relation

$$h_{\mathbf{m}}(\boldsymbol{\rho}_1) = \mathcal{H}(f_1), \quad h_{\mathbf{m}}(\boldsymbol{\rho}_2) = \mathcal{H}(f_2). \quad (83)$$

Moreover,  $(1-t)f_1 + tf_2$  is a nonnegative function which moment is  $\boldsymbol{\rho}(t)$ . Then

$$(\forall t \in ]0, 1[) \quad h_{\mathbf{m}}(\boldsymbol{\rho}(t)) \leq \mathcal{H}((1-t)f_1 + tf_2).$$

$\mathcal{H}$  being strictly convex, we get from (83) for any  $t \in ]0, 1[$ ,

$$h_{\mathbf{m}}(\boldsymbol{\rho}(t)) < (1-t)\mathcal{H}(f_1) + t\mathcal{H}(f_2) = (1-t)h_{\mathbf{m}}(\boldsymbol{\rho}_1) + th_{\mathbf{m}}(\boldsymbol{\rho}_2).$$

□

### 6.3 Proof of section 5

*Proof.* (Proposition 17). We may first consider the rotation around the mean velocity  $u$  since they play an important role to obtain the right hydrodynamic limit. So we

define  $\tau = \tau_{\mathbf{u}^{-1}\theta\mathbf{u}} = \tau_{\mathbf{u}^{-1}}\tau_\theta\tau_{\mathbf{u}}$  and let us prove the result for this  $\tau$ . By definition of  $G$ ,  $G(f) = \mathcal{M}(\phi^*)'(\boldsymbol{\alpha} \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}))$ , where  $\boldsymbol{\alpha}$  is the polar variable of  $L(\boldsymbol{\rho}_f)$ . In other words,  $\nabla h^*(\boldsymbol{\alpha}) = L(\boldsymbol{\rho}_f)$ . By definition of  $\mathcal{M}$ ,  $\mathcal{M}$  remains unchanged with the transformation  $\tau$ . So  $\tau(G(f)) = \mathcal{M}(\phi^*)'(\boldsymbol{\alpha} \cdot \mathbf{a}(\tau(\mathbf{v}) - \mathbf{u}))$ , with

$$\mathbf{a}(\tau(\mathbf{v}) - \mathbf{u}) = (1, \theta(\mathbf{v} - \mathbf{u}), \frac{(\mathbf{v} - \mathbf{u})^2}{2} - \frac{3}{2}k_B T, \theta(\mathbf{v} - \mathbf{u}) \otimes \theta(\mathbf{v} - \mathbf{u}) - \frac{1}{3}(\mathbf{v} - \mathbf{u})^2 I_d, \theta(\mathbf{v} - \mathbf{u})\left(\frac{\mathbf{v} - \mathbf{u}}{2} - \frac{5}{2}\right)).$$

This last vector can be written as

$$\mathbf{a}(\tau(\mathbf{v}) - \mathbf{u}) = (1, \theta(\mathbf{v} - \mathbf{u}), a_2(\mathbf{v} - \mathbf{u}), \theta\mathbb{A}(\mathbf{v} - \mathbf{u})\theta^t, \theta\mathbf{b}(\mathbf{v} - \mathbf{u})) \quad (84)$$

A simple computation leads to  $\boldsymbol{\alpha} : \mathbf{a}(\tau(\mathbf{v}) - \mathbf{u}) = \Theta(\boldsymbol{\alpha}) : \mathbf{a}(\mathbf{v} - \mathbf{u})$ , with

$$\Theta(\boldsymbol{\alpha}) = (\alpha_0, \theta^t\alpha_1, \alpha_2, \theta^t\alpha_3\theta, \theta^t\alpha_4). \quad (85)$$

And thus  $\tau(G(f)) = \mathcal{M}(\phi^*)'(\Theta(\boldsymbol{\alpha}) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}))$ . Remark that  $\tau G(f)$  has the form of the solution of Theorem 14 for some moment. Let us compute this moment.  $\boldsymbol{\rho}_{\tau G(f)}$  is defined by

$$\begin{aligned} \boldsymbol{\rho}_{\tau G(f)} &= \int_{\mathbb{R}^3} \tau G(f)(\mathbf{v}) \mathbf{a}(\mathbf{v} - \mathbf{u}) d\mathbf{v} \\ &= \int_{\mathbb{R}^3} G(\tau\mathbf{v}) \mathbf{a}(\mathbf{v} - \mathbf{u}) d\mathbf{v}. \end{aligned}$$

By using the change of variable  $\mathbf{w} = \tau(\mathbf{v}) = \theta(\mathbf{v} - \mathbf{u}) + \mathbf{u}$ ,

$$\boldsymbol{\rho}_{\tau G(f)} = \int_{\mathbb{R}^3} G(\mathbf{w}) \mathbf{a}(\theta^t(\mathbf{w} - \mathbf{u})) d\mathbf{w}.$$

But as  $\Theta(\mathbf{a}(\mathbf{v} - \mathbf{u})) = \mathbf{a}(\theta^t(\mathbf{v} - \mathbf{u}))$ . Then

$$\boldsymbol{\rho}_{\tau G(f)} = \Theta \left( \int_{\mathbb{R}^3} G(\mathbf{w}) \mathbf{a}(\mathbf{w} - \mathbf{u}) d\mathbf{w} \right) = \Theta(\boldsymbol{\rho}_G).$$

Likewise with the same computations gives  $\boldsymbol{\rho}_{\tau(f)} = \Theta(\boldsymbol{\rho}_f)$ . Thus  $G(\tau(f))$  is the solution of Theorem 14 for  $\boldsymbol{\rho} = L(\boldsymbol{\rho}_f)$ . Now remark for the definition of  $L$  and  $\Theta$  that  $\Theta L(\boldsymbol{\rho}_f) = L(\Theta(\boldsymbol{\rho}_f))$ . As a conclusion  $\tau G(f) = G(\tau f)$ .

It remains to prove the result for the translations i.e.  $\tau_{\mathbf{z}} G(f) = G(\tau_{\mathbf{z}} f)$ . There is  $G = \mathcal{M}(\phi^*)'(\boldsymbol{\alpha} \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}))$ . Moreover

$$\tau_{\mathbf{z}} G(f) = \frac{n}{(2\pi T)^{\frac{3}{2}}} \exp\left(-\frac{(\mathbf{v} - \mathbf{u} - \mathbf{z})^2}{2T}\right) \phi^*(\boldsymbol{\alpha} \cdot \mathbf{a}((\mathbf{v} - \mathbf{u} - \mathbf{z})))$$

$G(\tau_{\mathbf{z}}f)$  is the solution of the minimization problem when changing the Grad basis into the framework related to  $\tau_{\mathbf{z}}f$  which is precisely moving at velocity  $\mathbf{u} + \mathbf{z}$ .

$$\int_{\mathbb{R}^3} \tau_{\mathbf{z}}f(\mathbf{1}, \mathbf{v}, \mathbf{v}^2) d\mathbf{v} = (n, n(\mathbf{u} + \mathbf{z}), \frac{1}{2}n(\mathbf{u} + \mathbf{z})^2 + \frac{3}{2}nT).$$

In the corresponding Grad basis  $\tau_{\mathbf{z}}f$  has the same macroscopic value  $\boldsymbol{\rho}_f$  as  $f$  in  $\mathbf{a}(\mathbf{v} - \mathbf{u} - \mathbf{z})$ . Namely

$$\boldsymbol{\rho}_f = \int_{\mathbb{R}^3} f(\mathbf{v})\mathbf{a}(\mathbf{v} - \mathbf{u})d\mathbf{v} = \int_{\mathbb{R}^3} \tau_{\mathbf{z}}f(\mathbf{v})\mathbf{a}(\mathbf{v} - \mathbf{u})d\mathbf{v}$$

So

$$L(\boldsymbol{\rho}_f) = \int_{\mathbb{R}^3} \mathcal{M}(\mathbf{v} - \mathbf{z})(\phi^*)'(\boldsymbol{\alpha} \cdot \mathbf{a}((\mathbf{v} - \mathbf{u} - \mathbf{z})\mathbf{a}(\mathbf{v} - \mathbf{u} - \mathbf{z}))d\mathbf{v}.$$

This means that the relation between  $\boldsymbol{\alpha}$  and  $L(\boldsymbol{\rho}_f)$  is valid whatever is framework. The minimisation problem remains unchanged by changing both  $\mathcal{M}$  and  $f$  in  $\tau_{\mathbf{z}}\mathcal{M}$  and  $\tau_{\mathbf{z}}f$ .  $\square$

*Proof.* (Theorem 18). Remark that  $\mathcal{M}$  is the unique minimizer of  $\mathcal{H}$  just under the constraints of conservation laws. Indeed there is  $\mathcal{M} = \mathcal{M} \times 1 = \mathcal{M}\phi^{*'}(\alpha_0)$  for some  $\alpha_0 \in \mathbb{R}$  since  $\phi^{*'}$  is a bijection from  $\mathbb{R}^+$  to  $\mathbb{R}^+$ . Thus  $\mathcal{M}$  has the form of the solution of the primal problem given by theorem 14 when  $\mathbf{m}(\mathbf{v}) = \{\mathbf{1}, \mathbf{v}, \mathbf{v}^2\}$  and the constraint is  $\int G\mathbf{m}d\mathbf{v} = (n, 0, 0)$ . Adding more constraints (w.r.t.  $\mathbb{A}(\mathbf{v} - \mathbf{u})$  and  $\mathbf{b}(\mathbf{v} - \mathbf{u})$ ) prove that  $\mathcal{H}(\mathcal{M}) = h(\boldsymbol{\rho}_{eq})$  since again  $\mathcal{M}$  has the shape of the solution and satisfies the constraints. We have

$$\begin{aligned} (\forall \boldsymbol{\rho} \in \mathcal{R}_{\mathbf{a}}^+), \text{ with } \boldsymbol{\rho} = (n, 0, 0, 0, 0), h_{\mathbf{a}}(\boldsymbol{\rho}_{eq}) \leq h_{\mathbf{a}}(\boldsymbol{\rho}) \\ \text{with equality } h_{\mathbf{a}}(\boldsymbol{\rho}_{eq}) = h_{\mathbf{a}}(\boldsymbol{\rho}) \text{ iff } \boldsymbol{\rho}_{eq} = \boldsymbol{\rho} \end{aligned} \quad (86)$$

thanks to the strict convexity of  $h$  in  $\mathcal{R}_{\mathbf{a}}^+$  (Theorem 14). In other words,  $\mathcal{M}$  is the unique minimizer of  $\mathcal{H}$  of all functions in  $\mathbb{L}^1(\mathbf{a})$  having the same mass, momentum and energy as  $f$ .

We define the function  $F$  which satisfies

$$\int F\mathbf{a}(\mathbf{v} - \mathbf{u})d\mathbf{v} = \boldsymbol{\rho}_f \quad \text{and} \quad \mathcal{H}(F) = h_{\mathbf{a}}(\boldsymbol{\rho}_f). \quad (87)$$

$F$  is unique thanks to Theorem 14 and reads  $F = \mathcal{M}(\phi^*)'(\boldsymbol{\alpha}_F \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}))$ . Consider as in the proof of Lemma 17  $\tau f$  with  $\tau = \tau_{-\mathbf{u}}\tau_{-I_d}\tau_{\mathbf{u}}$ . Then

$$\Theta(\boldsymbol{\rho}_f) = \int_{\mathbb{R}^3} \tau f(\mathbf{v})\mathbf{a}(\mathbf{v} - \mathbf{u})d\mathbf{v} = (n, \mathbf{0}, 0, n\Pi, -n\mathbf{b}), \quad h(\tilde{\boldsymbol{\rho}}) = h(\boldsymbol{\rho}).$$

There is  $\tau\mathcal{M} = \mathcal{M}$  in such a way that  $\mathcal{H}(\tau f) = \mathcal{H}(f)$  and  $h_{\mathbf{a}}(\Theta\boldsymbol{\rho}) = h_{\mathbf{a}}(\boldsymbol{\rho})$ , since  $\tau f$  is solution to the minimisation problem by changing  $\boldsymbol{\rho}$  and  $\Theta\boldsymbol{\rho}$ .

From the strict convexity of  $h_{\mathbf{a}}$ , one finds

$$h_{\mathbf{a}}(n, \mathbf{0}, 0, n\Pi, \mathbf{0}) = h_{\mathbf{a}}\left(\frac{1}{2}\boldsymbol{\rho} + \frac{1}{2}\Theta(\boldsymbol{\rho})\right) \leq \frac{1}{2}h_{\mathbf{a}}(\boldsymbol{\rho}) + \frac{1}{2}h_{\mathbf{a}}(\Theta(\boldsymbol{\rho})) = h_{\mathbf{a}}(\boldsymbol{\rho}). \quad (88)$$

with equality only if  $\mathbf{b} = \mathbf{0}$ .

Let then  $\lambda_a, \lambda_b \in [0, 1]$ . Without a loss of generality, assume that  $\lambda_b \leq \lambda_a$  and take  $\lambda \in [0, 1]$  such that  $\lambda_b = \lambda\lambda_a$  (if not, one sets  $\lambda_a = \lambda\lambda_b$ ). We have

$$(n, \mathbf{0}, 0, n\Pi, \lambda\mathbf{b}) = \lambda(n, \mathbf{0}, 0, n\Pi, \mathbf{b}) + (1 - \lambda)(n, \mathbf{0}, 0, n\Pi, \mathbf{0}),$$

and as a consequence  $h_{\mathbf{a}}(n, \mathbf{0}, 0, n\Pi, \lambda\mathbf{b}) \leq h_{\mathbf{a}}(n, \mathbf{0}, 0, n\Pi, \mathbf{b})$  where we have used (88). Again, the equality holds only if  $\mathbf{b} = \mathbf{0}$ . Likewise, we have

$$(\rho_1, \mathbf{0}, 0, \lambda_{\mathbb{A}}n\Pi, \lambda_b\mathbf{b}) = (1 - \lambda_{\mathbb{A}})(n, \mathbf{0}, 0, \mathbb{O}, \mathbf{0}) + \lambda_{\mathbb{A}}(n, \mathbf{0}, 0, n\Pi, \lambda\mathbf{b}),$$

and finally (using former inequalities)

$$h_{\mathbf{a}}(n, \mathbf{0}, 0, \lambda_{\mathbb{A}}n\Pi, \lambda_b\mathbf{b}) \leq h_{\mathbf{a}}(n, \mathbf{0}, 0, n\Pi, \mathbf{b}). \quad (89)$$

with equality only if  $\Pi = \mathbb{O}$  and  $\mathbf{b} = \mathbf{0}$  or  $\lambda_{\mathbb{A}} = \lambda_b = 1$ . There is

$$I := \left\langle K(f) \phi' \left( \frac{f}{\mathcal{M}} \right) \right\rangle = \int \nu(G - f) \phi' \left( \frac{f}{\mathcal{M}} \right) d\mathbf{v}.$$

From the expression of  $G$  there holds

$$I = \int \left[ \phi^{*'}(\boldsymbol{\alpha}(\mathbf{L}(\boldsymbol{\rho}_f)) \cdot \mathbf{a}) - \frac{f}{\mathcal{M}} \right] \phi' \left( \frac{f}{\mathcal{M}} \right) \mathcal{M} d\mathbf{v}.$$

Use Young inequality for  $\phi$ :  $\phi(y) - \phi(x) \geq \phi'(x)(y - x)$ , with

$$x = \frac{f}{\mathcal{M}}, \quad y = \phi^{*'}(\boldsymbol{\alpha}(\mathbf{L}(\boldsymbol{\rho}_f)) \cdot \mathbf{a}),$$

multiply by  $\mathcal{M}$  and integrate over  $\mathbb{R}^3$  gives

$$\begin{aligned} I &\leq \int \nu \left( \phi(\phi^{*'}(\boldsymbol{\alpha}(\mathbf{L}(\boldsymbol{\rho}_f)) \cdot \mathbf{a})) \mathcal{M} - \int \phi \left( \frac{f}{\mathcal{M}} \right) \mathcal{M} \right) = \nu(\mathcal{H}(G) - \mathcal{H}(f)) \\ &= \nu(h_{\mathbf{a}}(\mathbf{L}(\boldsymbol{\rho}_f)) - \mathcal{H}(f)). \end{aligned}$$

Hence  $I \leq \nu(h_{\mathbf{a}}(\mathbf{L}(\boldsymbol{\rho}_f)) - \mathcal{H}(f))$ . This reads also

$$I \leq \nu(h_{\mathbf{a}}(\mathbf{L}(\boldsymbol{\rho}_f)) - h_{\mathbf{a}}(\boldsymbol{\rho}_f) + h_{\mathbf{a}}(\boldsymbol{\rho}_f) - \mathcal{H}(f)).$$

By definition of the entropy  $h_{\mathbf{a}}(\boldsymbol{\rho})$  there is  $h_{\mathbf{a}}(\boldsymbol{\rho}_f) - \mathcal{H}(f) \leq 0$ . But on the other hand we have by computing the moments on  $\mathbf{a}(\mathbf{v} - \mathbf{u})$

$$\boldsymbol{\rho}_f = (n, \mathbf{0}, 0, n\mathbb{A}, n\mathbf{b}), \mathbf{L}(\boldsymbol{\rho}_f) = (n, \mathbf{0}, 0, \lambda_{\mathbb{A}}n\mathbb{A}, \lambda_{\mathbf{b}}n\mathbf{b}), \lambda_{\mathbb{A}}, \lambda_{\mathbf{b}} \in [0, 1).$$

As  $h_{\mathbf{a}}(\mathbf{L}(\boldsymbol{\rho}_f)) \leq h_{\mathbf{a}}(\boldsymbol{\rho}_f)$ , the entropy theorem is proved.

From (89) and (87)  $I = 0$  iff  $\mathbf{L}(\boldsymbol{\rho}_f) = \boldsymbol{\rho}_f$  and  $f = F$ . But  $\mathbf{L}(\boldsymbol{\rho}_f) = \boldsymbol{\rho}_f$  iff  $\boldsymbol{\rho}_f = \boldsymbol{\rho}_{eq}$  that is  $f = \mathcal{M}$ .  $\square$

*Proof.* (Proposition 19). Remark that thanks to Galilean invariance in the Grad space, the choice of the basis functions for defining the constraints does not change the result of the minimization problem. This means that we may write

$$G[f](\mathbf{v}) = \mathcal{M}(\tilde{\boldsymbol{\mu}}(f) \cdot \mathbf{a}(\mathbf{v})) \phi^{*'}(\tilde{\boldsymbol{\alpha}}(f) \cdot \mathbf{a}(\mathbf{v})).$$

But it is more convenient to write those Lagrange multipliers when the basis is  $\mathbf{a}(\mathbf{v} - \mathbf{u})$ . That is

$$G[f](\mathbf{v}) = \exp(\boldsymbol{\mu}(f) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u})) \phi^{*'}(\boldsymbol{\alpha}(f) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u})).$$

Then there is

$$\exp(\boldsymbol{\mu}(\mathcal{M}) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u})) \phi^{*'}(\boldsymbol{\alpha}(\mathcal{M}) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u})) = \mathcal{M},$$

since  $\mathcal{M}$  is the unique solution to the variational problem with the moment constraints

$$\int f \mathbf{m}(\mathbf{v}) d\mathbf{v} = (n, \mathbf{0}, 0, 0, 0) = \boldsymbol{\rho}_{\mathcal{M}}.$$

In the above equation  $\phi^{*'}(\boldsymbol{\alpha}(\mathcal{M}) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u})) = 1$  implies  $\boldsymbol{\alpha}(\mathcal{M}) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}) = \alpha$  for some constant  $\alpha$  because  $\phi^{*'}$  is a bijection from  $[p_0, +\infty)$  into  $\mathbb{R}^+$ . So finally there holds also

$$\phi^{*''}(\boldsymbol{\alpha}(\mathcal{M}) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u})) = c,$$

for some constant  $c$  if  $\phi^*$  is twice differentiable.

Note that the properties of  $G[f]$  implies that  $G[\mathcal{M}] = \mathcal{M}$  so  $K(\mathcal{M}) = 0$ . Let us consider  $f = \mathcal{M}(1 + \epsilon g)$ . Differentiating formally the function  $\boldsymbol{\alpha}$  and  $\boldsymbol{\mu}$  there holds :

$$\boldsymbol{\alpha}(f) = \boldsymbol{\alpha}(\mathcal{M}) + \epsilon d\boldsymbol{\alpha}_{\mathcal{M}}(g) + O(\epsilon^2), \quad \boldsymbol{\mu}(f) = \boldsymbol{\mu}(\mathcal{M}) + \epsilon d\boldsymbol{\mu}_{\mathcal{M}}(g) + O(\epsilon^2)$$

Then we compute for  $f = \mathcal{M}(1 + \epsilon g)$  the following approximation:

$$\begin{aligned} \exp(\boldsymbol{\mu}(f) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u})) &= \exp(\boldsymbol{\mu}(\mathcal{M}) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u})) (1 + \epsilon d\boldsymbol{\mu}_{\mathcal{M}}(g) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}) + o(\epsilon)) \\ &= \mathcal{M} (1 + \epsilon d\boldsymbol{\mu}_{\mathcal{M}}(g) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}) + o(\epsilon)) \end{aligned}$$

$$\begin{aligned}\phi^{*'}(\boldsymbol{\alpha}(f) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u})) &= \phi^{*'}(\boldsymbol{\alpha}(\mathcal{M}) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}) + o(\epsilon)) \\ &\quad + \epsilon d\boldsymbol{\alpha}_{\mathcal{M}}(g) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}) \phi^{*''}(\boldsymbol{\alpha}(\mathcal{M}) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}) + o(\epsilon)).\end{aligned}$$

so that

$$\begin{aligned}G[f](\mathbf{v}) &= \mathcal{M}(\mathbf{v}) (1 + \epsilon(d\boldsymbol{\mu}_{\mathcal{M}}(g) + c d\boldsymbol{\alpha}_{\mathcal{M}}(g)) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}) + o(\epsilon)) \\ &= \mathcal{M}(\mathbf{v}) (1 + \epsilon\Lambda_{\mathcal{M}}(g) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}) + o(\epsilon))\end{aligned}$$

Finally, there is

$$K(\mathcal{M}(1 + \epsilon g)) = \nu [\mathcal{M} + \epsilon\Lambda_{\mathcal{M}}(g) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}) \mathcal{M}(\mathbf{v}) - \mathcal{M}(1 + \epsilon g) + o(\epsilon)]$$

and by definition of  $\mathcal{L}$  (6)

$$\mathcal{L}(g) = \nu (\Lambda_{\mathcal{M}}(g) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}) - g). \quad (90)$$

It is convenient to write this relaxation equation with  $a_i$  being a scalar function of  $\mathbf{v} - \mathbf{u}$  rather than tensors. We need now to focus on  $\Lambda_{\mathcal{M}}(g)$  in order to identify the expression of  $\mathcal{L}_{\mathcal{M}}$ . To do that we use the prescribe condition on moment:

$$\int G[\mathcal{M}(1 + \epsilon g)] \mathbf{a}_i(\mathbf{v} - \mathbf{u}) = \left(1 - \frac{\nu_i}{\nu}\right) \int \mathcal{M}(1 + \epsilon g) \mathbf{a}_i(\mathbf{v} - \mathbf{u})$$

So with  $K(f) = \nu(G[f] - f)$  we conclude that:

$$\int K(\mathcal{M}(1 + \epsilon g)) \mathbf{a}_i(\mathbf{v} - \mathbf{u}) = -\epsilon \nu_i \int g \mathcal{M}(\mathbf{v}) \mathbf{a}_i(\mathbf{v} - \mathbf{u})$$

using the linear approximation for  $K$  we have

$$\epsilon \int \nu (\Lambda_{\mathcal{M}}(g) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u}) - g) \mathbf{a}_i(\mathbf{v} - \mathbf{u}) \mathcal{M} = -\epsilon \nu_i \int g \mathcal{M}(\mathbf{v}) \mathbf{a}_i(\mathbf{v} - \mathbf{u}).$$

That is finally for any component  $i$  we have:

$$\int (\Lambda_{\mathcal{M}}(g) \cdot \mathbf{a}(\mathbf{v} - \mathbf{u})) \mathbf{a}_i(\mathbf{v} - \mathbf{u}) \mathcal{M} = \left(1 - \frac{\nu_i}{\nu}\right) \int g \mathcal{M}(\mathbf{v}) \mathbf{a}_i(\mathbf{v} - \mathbf{u})$$

by expanding the dot product, we have formally:

$$\sum_j \int (\Lambda_{\mathcal{M}}^j(g) \cdot \mathbf{a}_j(\mathbf{v} - \mathbf{u})) \mathbf{a}_i(\mathbf{v} - \mathbf{u}) \mathcal{M} = \left(1 - \frac{\nu_i}{\nu}\right) \int g \mathcal{M}(\mathbf{v}) \mathbf{a}_i(\mathbf{v} - \mathbf{u}).$$

So, by using orthogonality relations

$$\Lambda_{\mathcal{M}}^i \|\mathbf{a}_i\|^2 = \left(1 - \frac{\nu_i}{\nu}\right) \int_{\mathbb{R}^3} g \mathcal{M} \mathbf{a}_i(\mathbf{v} - \mathbf{u}).$$

Then, according to formula (90) we conclude that

$$\mathcal{L}g = \nu \left( \sum_i \left( 1 - \frac{\nu_i}{\nu} \right) \mathcal{P}_{\mathbf{a}_i} g - g \right)$$

where  $\nu_i = 0$  if  $\mathbf{a}_i \in \mathbb{K}$ . From this, it is easy to see that  $\mathcal{L}$  is self adjoint, Fredholm with  $\text{Ker}(\mathcal{L}) = \mathbb{K}$ .

We end this proof by showing that

$$\forall f, \left[ \int K(f) \phi = 0 \right] \Leftrightarrow \phi \in \mathbb{K}$$

The right implication is just a consequence of (13). Assume now that

$$\forall f, \int K(f) \phi(\mathbf{v}) d\mathbf{v} = 0$$

Expanding  $K(\mathcal{M}(1 + \varepsilon g))$ , we have for any  $g \in L^2(\mathcal{M})$

$$\forall g, \int \mathcal{M} \mathcal{L}(g) \phi(\mathbf{v}) d\mathbf{v} = 0$$

Since  $\mathcal{L}$  is self adjoint, we have

$$\forall g, \int \mathcal{M} g \mathcal{L}(\phi) d\mathbf{v} = 0$$

which proves that  $\phi \in \mathbb{K}$ . □

## 7 Conclusion

In the present article, we have proposed a methodology to construct relaxation operators. The derivation is performed in three steps. We first consider the projection of the inverse linearized Boltzmann operator  $\mathcal{L}_B^{-1}$  on a polynomial space of finite dimension. We then state relaxation equations on the moments of the probability distribution  $f$  basing on its diagonalization. The model must satisfy those equations together with the conservation laws. From this one derives linear relations between the moments of  $f$  and the target function  $G$  to be found. The later is then found by solving a variational problem. Different mathematical problems related to this construction have been addressed. We have firstly revisited a theorem by Junk [30] relating realizable moments (i.e moments of nonnegative integrable functions) to nonnegative polynomials. From this we have derived necessary conditions for the realizability of the moments of  $G$  and proved that it allows to specify the admissible relaxation equations on the Grad thirteen moments. The variational problem has been studied in detail by using different functional to be minimized under moment constraints. We have reestablished a theorem of Csiszar [20] on the existence of solution to such minimization problems

by using convex analysis and exactly derived the shape of the solution by duality. In the last part of the article, we have proposed different models from this construction and analyzed their well-posedness. In particular, when relaxations occur on the Grad thirteen moments, the model satisfies almost all properties of the original Boltzmann equation: nonnegativity of the solution, conservation laws, H theorem, Galilean invariance and the right hydrodynamic limit up to Navier-Stokes level. However, the control of the entropy defined by the  $\phi$ -divergence is only local. In the general case, those properties are also preserved but the control of the entropy is not yet proved. Finally, the present approach encompasses the derivation of many known models and for some of them their generalization.

There are many perspectives and questions related to this work. In principle the new model based only on Grad thirteen moments should not bring more than the ESBGK or Shakhov models. It remains however to compare them. We also intend to study the generalization of those relaxation operators beyond the Grad case. In such cases, the present method does not require the effective computation of the relaxation operator if moment methods such as in [2, 34] are used. The computation of the approximate inverse linearized Boltzmann operator is of the same order of complexity than that of the transport coefficients for multicomponent fluids for which there exists plenty efficient methods. One may then in a first time compare from a numerical point of view the solution of this general model to that of the linearized Boltzmann equation and in a second time compare it to that of the known relaxation models and to the Boltzmann equation itself. Also, some study related to existence of solutions to the generalized Shakhov model can be addressed the framework of Bae and Yun [5].

**Data Availability.** No data was used for the research described in the article.

**Conflict of Interest.** The authors declare that they have no conflicts of interest.

**Acknowledgements.** The authors would like to thank J.J. Alibert from the laboratory IMATH (university of Toulon) for helpful discussions in the approach of the variational problem.

## References

- [1] Abdelmalik, M., Cai, M., Pichard, T.: On the Renormalization Maps for the phi-Divergence Moment Closures Applied in Radiative Transfer, *J.Comput. Theor. Transport* **52**, 6, 399-428 (2023)
- [2] Abdelmalik, M. R. A., van Brummelen, H.: Moment closure approximations of the Boltzmann equation based on  $\phi$ -divergences., *J. Stat. Phys.* **164**, 77-104 (2016)
- [3] Akhiezer, N.I.: *The classical moment problem*, Amsterdam (1965)
- [4] Andries, P., Le Tallec, P., Perlat, J.P., Perthame, B.: The Gaussian-BGK Model of Boltzmann Equation with Small Prandtl Number, *Eur. J. Mech. B Fluids.* **19**, 813-830 (2000)



- [5] Bae, G., Yun, S.B.: The Shakhov model near global Maxwellian, *Nonlinear analysis real world application* **52**, 6, 399-428 (2023)
- [6] Benoist, O.: Writing positive polynomials as sums of (few) squares, *EMS. Newsletter*, 8-13 (2017)
- [7] Bhatnagar, P. L. and Gross, E. P. and Krook, M.: A Model for Collision Processes in Gases. I. Small Amplitude Processes in Charged and Neutral One-Component Systems, *Physical Review* **94**, 3, 511-525 (1954)
- [8] Borwein, J.M., Lewis, S.: Duality relationships for entropy-like minimization problems, *SIAM J. Control Optim.* **29**, 2, 325-338 (1991)
- [9] Bouchut, F., Perthame, B.: A BGK model for small Prandtl number in the Navier-Stokes approximation, *J. Statist. Phys.* **71**, 191-207 (1993)
- [10] Bourgat, J.F., Desvillettes, L., Le Tallec, P., Perthame, B.: Microreversible collisions for polyatomic gases and Boltzmann's theorem, *European J. Mech. B Fluids* **13**, 2, 237-254 (1994)
- [11] Brull, S., Pavan, V., Schneider, J.: Derivation of a BGK model for mixture, *European Journal of Mechanics - B/Fluids* **33**, 74-86 (2012)
- [12] Brull, S., Schneider, J.: A new approach of the Ellipsoidal Statistical Model *Cont. Mech. Thermodyn.* **20**, 63-74 (2008)
- [13] Brull, S., Schneider, J.: On the Ellipsoidal Statistical Model for polyatomic gases, *Cont. Mech. Thermodyn.* **20**, 489-508 (2009)
- [14] Brull, S.: An Ellipsoidal Statistical Model for gas mixtures, *Comm. Math. Sci.* **13**, 1-13 (2015)
- [15] Brull, S.: An Ellipsoidal Statistical Model for a monoatomic and polyatomic gas mixture, *Comm. Math. Sci.* **19**, 2177-2194 (2021)
- [16] Cercignani C.: *The Boltzmann Equation and its Application*, Springer Verlag (1992)
- [17] Chapman, S., Cowling, T.G.: *The Mathematical Theory of non-uniform Gases*, Third edition, Cambridge University Press (1990)
- [18] Csiszar, I.: A class of measures of informativity of observation channels, *Period. Math. Hung.* **2**, 191-213 (1972)
- [19] Csiszar, I.: I divergence geometry of probability distribution function and minimization problems Sanov property, *Ann. Probab.* **3**, 146-158 (1975)

- [20] Csiszar, I.: Generalized projections for nonnegative functions, *Acta. Math. Hung.* **68**, 161-185 (1995)
- [21] Curto, R., Fialkow, L.: Recursiveness, positivity, and truncated moment problem, *Houton J. Math.* **17**, 4, 603-635 (1991)
- [22] Curto, R., Fialkow, L.: A duality proof of Tchakalff's theorem, *J. Math. Anal. Appl.* **262**, 519-532 (2002)
- [23] Curto, R., Fialkow, L.: Truncated K-moment problems in several variables, *J. Functional Analysis* **54**, 1, 189-226 (2005)
- [24] Curto, R., Fialkow, L.: An analogue of the Riesz-Haviland theorem for the truncated moment problem, *J. Functional Analysis* **255**, 2709-2731 (2008)
- [25] Fialkow, L.: The truncated K-moment problem: a Survey., *Theta Ser. Adv. Math.* **18**, 25-51 (2016)
- [26] Hauck, C., Levermore, C.D., Tits, A.L.: Convex duality and entropy-based moment closures: characterizing degenerate densities, *SIAM J. Control Optimal* **47**, 4, 1977-2015 (2008)
- [27] Haviland, E.K.: On the momentum problem for distributions in more than one dimension ii, *Amer. J. Math.* **57**, 3, 562-568 (1935)
- [28] Haviland, E.K.: On the momentum problem for distributions in more than one dimension, *Amer. J. Math.* **58**, 1, 164-168 (1936)
- [29] Holway, L. H.: New statistical models for kinetic theory: methods of construction, *Phys. Fluids.* **9**, 1658-1673 (1966)
- [30] Junk M.: Maximum entropy for reduced moment problems., *Math. Models and Methods in Appl. Sci.* **10**, 7, 1001-1025 (2000)
- [31] Junk, M., Unterreiter, A.: A maximum entropy moment systems and Galilean invariance., *Contin. Mech. Thermodyn.* **14**, 6, 563-576 (2002)
- [32] Krasnoselskii, M.A., Rutikii, Ja. B., *Convex Functions and Orlicz Spaces* (in Russian), Gos. Izd. Fiz.-Mat. Lit. (Moscow, 1958).
- [33] Lasserre, J.B.: *Moments, Positive Polynomials and Their Applications*, Imperial College Press, London (2009)
- [34] Levermore, C. D.: Moment closure hierarchies for kinetic theories., *J. Stat. Phys.* **83**, 1021-1065 (1996)
- [35] Mieussens L., Struchtrup H.: Numerical comparison of BGK-models with proper Prandtl number, *Phys. Fluids* **16**, 8, 2797-2813 (2004)

- [36] Pavan, V.: General entropic approximations for canonical systems described by kinetic equations, *J. Stat. Phys* **142**, 4, 792-827 (2011)
- [37] Pichard, T.: A moment closure based on a projection on the boundary of the realizability domain: 1D case, *Kin. Rel. Models* **13**, 6, 1243-1280 (2020)
- [38] Pichard, T.: A moment closure based on a projection on the boundary of the realizability domain: Extension and analysis, *Kin. Rel. Models* **15**, 5, 793-822 (2022)
- [39] Rockafellar, R.T.: *Convex analysis*, Princeton University Press (1970)
- [40] Schneider, J.: Entropic approximation in kinetic theory, *Math. Model. Numer. Anal.* **38**, 3, 541-561 (2004)
- [41] Schneider, J.: A well-posed simulation model for multicomponent reacting gases, *Comm. Math. Sci.* **13**, 5, 1075-1103 (2015)
- [42] Shakhov, E. M.: Generalization of the Krook kinetic relaxation equation, *Fluid Dynamics.* **3**, 95-96 (1968)
- [43] Struchtrup, H.: The BGK model with velocity dependant collision frequency, *Cont. Mech. Thermodyn.* **9**, 23-31 (2004)
- [44] Struchtrup, H., Zheng, Y.: Ellipsoidal Statistical Bhatnagar-Gross-Krook model with velocity-dependent collision frequency, *Phys. Fluids.* **12** (2005)