# CONVERGENCE ANALYSIS OF GRAD'S HERMITE EXPANSION FOR LINEAR KINETIC EQUATIONS[*]

NEERAJ SARNA[†], JAN GIESSELMANN[‡], AND MANUEL TORRILHON[†]

**Abstract.** In [*Commun. Pure Appl. Math.*, 2 (1949), pp. 331–407], Grad proposed a Hermite series expansion for approximating solutions to kinetic equations that have an unbounded velocity space. However, for initial boundary value problems, poorly imposed boundary conditions lead to instabilities in Grad's Hermite expansion, which could result in nonconverging solutions. For linear kinetic equations, a method for posing stable boundary conditions was recently proposed for (formally) arbitrary order Hermite approximations. In the present work, we study $L^2$-convergence of these stable Hermite approximations and prove explicit convergence rates under suitable regularity assumptions on the exact solution. We confirm the presented convergence rates through numerical experiments involving the linearized BGK equation of rarefied gas dynamics.

**Introduction.** Evolution of charged or neutral particles (under certain conditions of interaction) can be modeled by linear kinetic equations. The explicit form of these kinetic equations depends on the physical system they model and many of these forms have been extensively studied in the past; see [11, 12, 14, 28]. Broadly speaking, different forms of kinetic equations have mainly three differentiating factors: the space of possible velocities of particles, i.e., the so-called velocity space; the external or the internal forces that act on the particles; and the collision operator that models the interaction between different particles. In the present work, we are concerned with linear kinetic equations that have the whole $\mathbb{R}^d$ ($1 \leq d \leq 3$) as their velocity space, have no external force acting on the particles, and have a collision operator that is bounded and negative semidefinite on $L^2(\mathbb{R}^d)$. Such kinetic equations usually arise from the kinetic gas theory after the linearization of the nonlinear Boltzmann or the BGK equation [4].

Mostly, an exact solution to a kinetic equation is not known and one seeks an approximation through a temporal, spatial, and velocity space discretization. In the present work, we analyze a Galerkin-type velocity space approximation where we approximate the solution's velocity dependence in a finite-dimensional space [13, 20]. Our finite-dimensional space is the span of a finite number of Grad's tensorial Hermite polynomials, which results in the so-called Grad's moment approximation [14].

[†]Center for Computational Engineering and Department of Mathematics, RWTH, Aachen University, Aachen, 52062, Germany (sarna@mathcces.rwth-aachen.de, mt@mathcces.rwth-aachen.de).

[‡]Department of Mathematics, Technical University of Darmstadt, Darmstadt, 64293, Germany (giesselmann@mathematik.tu-darmstadt.de).

We consider initial boundary value problems (IBVPs) and equip the Hermite approximation with boundary conditions that lead to its $L^2$-stability [21].

The convergence behavior of moment approximations, particularly for IBVPs, is not very well understood. Lack of understanding originates from expecting a monotonic (and test case independent) decrease in the error as the number of moments is increased but such a decrease is usually not observed in practice [26]. It is known that convergence of Galerkin methods is solution's regularity dependent, which is in turn test case dependent. Therefore, one possible way to understand the test case dependent convergence of moment approximations is to reformulate them as Galerkin methods [9, 10, 23]. We use such a reformulation for the Grad's moment approximation to prove that it converges (in the $L^2$-sense) to the kinetic equation's solution.

Reformulation of a moment approximation as a Galerkin method allows us to use the following (standard) steps for convergence analysis. First, we define a projection onto the Hermite approximation space and use it to split the approximation error into two parts: (i) one part containing the error in the expansion coefficients (or the moments) and (ii) the other part containing the projection error. Second, we bound the error in the expansion coefficients in terms of the projection error. To develop this bound, we exploit the $L^2$-stability property of the Hermite approximation, which is possible by defining the projection such that it satisfies the same boundary conditions as those satisfied by the moment approximation. We complete our analysis by proving that the projection error converges to zero.

It is worth noting that the orthogonal projection onto the approximation space does not satisfy the same boundary condition as the numerical solution and, thus, the $L^2$-stability results are not available. Indeed, from a technical perspective, defining a suitable projection operator is a key contribution of this work.

In previous works [20, 23], for kinetic equations with an unbounded velocity space, authors have analyzed convergence of Galerkin methods that use a grid in the velocity space. Although easier to implement, such methods fail to preserve the Galilean and the rotational invariance of kinetic equations. In contrast, Grad's tensorial Hermite polynomials cannot be mapped to a velocity space grid but they do preserve especially rotational invariance of kinetic equations. This allows for an approximation that is physically more sound. To the best of our knowledge, the present work is the first step toward analyzing the convergence of a rotational invariant Galerkin method for IBVPs involving kinetic equations with an unbounded velocity domain.

Other approximation schemes that lead to a rotational invariant approximation (for both bounded and unbounded velocity spaces) use spherical harmonics instead of Grad's Hermite polynomials; see [2, 5, 10]. Preliminary analysis shows that our framework is extendable to such approximations. Indeed, using our current framework one can even analyze the convergence of a general rotational invariant Galerkin scheme for a general rotational invariant kinetic equation considered in [1]. Moreover, our framework has an extension to linear approximations of the nonlinear Boltzmann equation [13]. We leave an extension of our framework to other linear kinetic equations as a part of our future work.

A summary of the article's structure is as follows: the first section discusses the kinetic equation and its Grad's moment approximation; the second section discusses the projection operator and contains the main convergence result; the fourth section discusses an example of the linear kinetic equation that arises from the kinetic gas theory; and the fifth section contains our numerical experiment.

**1. Linear kinetic equation.** With $f : (0, T) \times \Omega \times \mathbb{R}^d \to \mathbb{R}$ we represent the solution to our kinetic equation where $\Omega$ is the physical space, $(0, T)$ is a bounded temporal domain and $\mathbb{R}^d$ is the velocity space. For simplicity, we focus most of

our discussion on the case for which the spatial domain is the open half-space $\Omega :=$ $\mathbb{R}^- \times \mathbb{R}^{d-1}$ ($1 \leq d \leq 3$). In subsection 2.2 we discuss how our framework can be extended to general $C^2$ spatial domains. With $V := (0,T) \times \Omega$ we represent the space-time domain and with $D := V \times \mathbb{R}^d$ we represent our space-time-velocity domain. With $\nabla_{t,x} := (\partial_t, \partial_{x_1}, \ldots, \partial_{x_d})$ we denote the gradient operator along the space-time domain and using it we define the following operator:

$$(1) \qquad \begin{aligned} \mathcal{L} &:= \partial_t + \sum_{i=1}^{d} \xi_i \partial_{x_i} - Q, \ \xi \in \mathbb{R}^d, \\ &= (1, \xi) \cdot \nabla_{t,x} - Q, \end{aligned}$$

where $Q : L^2(\mathbb{R}^d) \to L^2(\mathbb{R}^d)$ is the collision operator. The second form of the above operator will be helpful in understanding the regularity of a strong solution of an IBVP involving $\mathcal{L}$. We restrict our analysis to the case for which the operator $Q$ satisfies the conditions enlisted below. Later, in section 3, we give examples of collision operators that satisfy the assumption below.

*Assumption* 1. We assume that $Q : L^2(\mathbb{R}^d) \to L^2(\mathbb{R}^d)$ is (i) linear, (ii) bounded, (iii) negative semidefinite, and (iv) self-adjoint.

We consider $\mathcal{L}$ as a mapping from $H_{\mathcal{L}}$ to $L^2(D)$ where $H_{\mathcal{L}}$ is the graph space of $\mathcal{L}$ and is defined as

$$(2) \qquad H_{\mathcal{L}} := \{v \in L^2(D) \ : \ \mathcal{L}v \in L^2(D)\} \quad \text{where} \quad \|f\|_{H_{\mathcal{L}}}^2 := \|f\|_{L^2(D)}^2 + \|\mathcal{L}f\|_{L^2(D)}^2.$$

For IBVPs involving the operator $\mathcal{L}$, we need to define trace operators over $H_{\mathcal{L}}$. To define these trace operators, we first define the following boundaries of the set $D = (0,T) \times \Omega \times \mathbb{R}^d$:

$$\Sigma^{\pm} := (0,T) \times \partial\Omega_{\xi}^{\pm}, \quad V^{\pm} := \{T^{\pm}\} \times \Omega \times \mathbb{R}^d, \quad \partial D := \Sigma^+ \cup \Sigma^- \cup V^+ \cup V^-,$$

where we set $T^+ = T$ and $T^- = 0$. Moreover, $\partial\Omega_{\xi}^{\pm}$ is a result of splitting $\partial\Omega \times \mathbb{R}^d$ into two nonoverlapping parts and is defined as $\partial\Omega_{\xi}^{\pm} := \partial\Omega \times \mathbb{R}^{\pm} \times \mathbb{R}^{d-1}$. Thus $\partial\Omega_{\xi}^+$ and $\partial\Omega_{\xi}^-$ are sets containing points in $\partial\Omega \times \mathbb{R}^d$ corresponding to outgoing and incoming velocities, respectively. Using these boundary sets, in the following we define the relevant trace operators. A detailed derivation of these operators can be found in [28].

DEFINITION 1.1. *Traces of functions in $H_{\mathcal{L}}$ are well-defined in $L^2(\partial D, |\xi_1|)$, i.e., in the $L^2$ space of functions over $\partial D$ with the Lebesgue measure weighted with $|\xi_1|$. We denote the trace operator by*

$$\gamma_D : H_{\mathcal{L}} \to L^2(\partial D, |\xi_1|).$$

*To restrict $\gamma_D$ to $\Sigma^{\pm}$ and $\Sigma = \Sigma^+ \cup \Sigma^-$, we define $\gamma^{\pm} f = \gamma_D f|_{\Sigma^{\pm}}$ and $\gamma f = \gamma_D f|_{\Sigma}$. Similarly, we interpret $f(T^{\pm})$ as $f(T^{\pm}) = \gamma_D f|_{V^{\pm}}$.*

Using the above trace operators, we give the following IBVP:

$$(3) \qquad \mathcal{L}f = 0 \ \text{ in } \ D, \quad f(0) = f_I \ \text{ on } \ V^-, \quad \gamma^- f = f_{in} \ \text{ on } \ \Sigma^-,$$

where $f_I \in L^2(\Omega \times \mathbb{R}^d)$ and $f_{in} \in L^2(\Sigma^-; |\xi_1|) \cap L^2(\mathbb{R}^- \times \mathbb{R}^{d-1}; H^{1/2}(\partial\Omega \times (0,T)))$ are some suitable initial and boundary data, respectively. Here $H^{\frac{1}{2}}$ denotes a standard fractional Sobolev space. The reason behind assuming $f_I$ to be in $L^2(\Omega \times \mathbb{R}^d)$

and $f_{in}$ to be in $L^2(\Sigma^-, |\xi_1|)$ is clear from the definition of trace operators, whereas the assumption that $f_{in} \in L^2(\mathbb{R}^- \times \mathbb{R}^{d-1}; H^{1/2}(\partial\Omega \times (0,T)))$ will be made clear in Assumption 2.

We stick to strong solutions of the above IBVP and we define them as follows [28].

DEFINITION 1.2. *Let $f \in H_{\mathcal{L}}$ where $H_{\mathcal{L}}$ is as given in (2). Then, $f$ is a strong solution to the linear kinetic equation if it satisfies*

$$\langle v, \mathcal{L}f \rangle_{L^2(D)} = 0 \quad \forall \quad v \in L^2(D), \quad \gamma^- f = f_{in}, \quad f(0) = f_I.$$

It has been shown in [28] that the IBVP (3) has a unique strong solution and for our convergence analysis, we will make additional regularity assumptions on this strong solution. We start with defining the notion of moments.

**1.1. Moments and Hermite polynomials.** We define tensorial Hermite polynomials with the help of the multi-index $\beta^{(i)}$ as

$$(4) \qquad \psi_{\beta^{(i)}}(\xi) := \prod_{p=1}^{d} He_{\beta_p^{(i)}}(\xi_p), \quad \beta^{(i)} := \left(\beta_1^{(i)}, \ldots, \beta_d^{(i)}\right),$$

where the Hermite polynomials $(He_k)$ enjoy the property of orthogonality and recursion

(5a)

$$\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} He_i(\xi) He_j(\xi) \exp\left(-\frac{\xi^2}{2}\right) d\xi = \delta_{ij} \quad \Rightarrow \quad \int_{\mathbb{R}^d} \psi_{\beta^{(k)}} \psi_{\beta^{(l)}} f_0 d\xi = \prod_{p=1}^{d} \delta_{\beta_p^{(k)} \beta_p^{(l)}},$$

$$(5b) \qquad \sqrt{i+1} He_{i+1}(\xi) + \sqrt{i} He_{i-1}(\xi) = \xi He_i(\xi).$$

Above, $f_0$ is a Gaussian weight given as

$$(6) \qquad f_0(\xi) := \exp\left(-\xi \cdot \xi/2\right) / \sqrt[d]{2\pi}.$$

The quantity $\|\beta^{(i)}\|_{l^1}$ is the so-called degree of the basis function $\psi_{\beta^{(i)}}$. Below we define the $\|\beta^{(i)}\|_{l^1}$th order moment of a function in $L^2(\mathbb{R}^d)$.

DEFINITION 1.3. *Let $n(m)$ represent the total number of tensorial Hermite polynomials (i.e., $\psi_{\beta^{(i)}}(\xi)$) of degree $m$ and let $\psi_m(\xi) \in \mathbb{R}^{n(m)}$ represent a vector containing all such basis functions. Using $\psi_m(\xi)$, we define $\lambda_m : L^2(\mathbb{R}^d) \to \mathbb{R}^{n(m)}$ as $\lambda_m(r) = \int_{\mathbb{R}^d} \sqrt{f_0} \psi_m(\xi) r(\xi) d\xi, \forall r \in L^2(\mathbb{R}^d)$. Thus, $\lambda_m(r)$ represents a vector containing all the $m$th order moments of $r$. To collect all the moments of $r$ which are of order less than or equal to $M$ ($m \le M$), we additionally define*

$$\Psi_M(\xi) = (\psi_0(\xi)', \psi_1(\xi)', \ldots, \psi_M(\xi)')', \quad \Lambda_M(r) = (\lambda_0(r)', \lambda_1(r)', \ldots, \lambda_M(r)')',$$

*where $\Psi_M(\xi) \in \mathbb{R}^{\Xi^M}$ and $\Lambda_M : L^2(\mathbb{R}^d) \to \mathbb{R}^{\Xi^M}$ with $\Xi^M = \sum_{m=0}^{M} n(m)$ being the total number of moments. Above and in all of our following discussion, a prime ( $'$ ) over a vector will represent its transpose.*

**1.2. Regularity assumptions.** For further discussion we recall that $V = \Omega \times (0,T)$ and $D = V \times \mathbb{R}^d$. With $C^k([0,T]; X)$ we denote a $k$-times continuously differential function of time with values in some Hilbert space $X$. We equip $C^k([0,T]; X)$ with the norm $\|g\|_{C^k([0,T];X)} = \max_{j \le k} \|\partial_t^j g\|_{C^0([0,T];X)}$ where $\|g\|_{C^0([0,T];X)} = \max_{t \in [0,T]} \|g(t)\|_X$.

To capture velocity space regularity of solutions, we make use of the Hermite–Sobolev space $W_H^k(\mathbb{R}^d)$ which is the image of $L^2(\mathbb{R}^d)$ under the inverse of the Hermite Laplacian operator $(\Delta_H)^k = (-2\Delta + \frac{1}{2}\xi \cdot \xi)^k$; see [25] for details. One can show that a tensorial Hermite polynomial $(\psi_{\beta^{(m)}})$ is an eigenfunction of $\Delta_H$ with an eigenvalue of $(2m + d)$ and therefore, one can define norm of functions in $L^2(\Omega; W_H^k(\mathbb{R}^d))$ as

$$\|f\|_{L^2(\Omega; W_H^k(\mathbb{R}^d))} := \left( \sum_{m=0}^{\infty} (2m + d)^{2k} \|\lambda_m(f(t, ., .))\|_{L^2(\Omega; \mathbb{R}^{n(m)})}^2 \right)^{1/2}.$$

For further discussion we assume that the solution to our IBVP, along with its derivatives, lies in $C^0([0, T]; L^2(\Omega; W_H^k(\mathbb{R}^d)))$ for some $k$. We summarize this assumption in the following.

*Assumption* 2. Let $f$ be a strong solution to the kinetic equation (3). We assume that there exist numbers $k^{e/o} \geq 0$, $k_t^{e/o} \geq 0$, and $k_x^{e/o} \geq 0$ such that

$$f^{e/o} \in C^0([0, T]; L^2(\Omega; W_H^{k^{e/o}}(\mathbb{R}^d))), \ (\partial_t f)^{e/o} \in C^0([0, T]; L^2(\Omega; W_H^{k_t^{e/o}}(\mathbb{R}^d))),$$

$$(\partial_{x_i} f)^{e/o} \in C^0([0, T]; L^2(\Omega; W_H^{k_x^{e/o}}(\mathbb{R}^d))) \quad \forall \, i \in \{1, \ldots, d\}.$$

Above, $(.)^e$ and $(.)^o$ denote the even and odd parts (of the various quantities) defined with respect to $\xi_1$, i.e.,

$$f^o(\xi_1, \xi_2, \xi_3) = \frac{1}{2} \left( f(\xi_1, \xi_2, \xi_3) - f(-\xi_1, \xi_2, \xi_3) \right),$$

$$f^e(\xi_1, \xi_2, \xi_3) = \frac{1}{2} \left( f(\xi_1, \xi_2, \xi_3) + f(-\xi_1, \xi_2, \xi_3) \right).$$

Note that for simplicity we have assumed the same degree of regularity for all spatial derivatives. Extending the forthcoming results to cases where different spatial derivatives have different degrees of regularity is straightforward.

To understand the relation between a standard Sobolev space and the Hermite–Sobolev space, we recall the following result [25] (see Theorem 2.1):

$$W_H^k(\mathbb{R}^d) \subseteq H^{2k}(\mathbb{R}^d) \subseteq L^2(\mathbb{R}^d) \quad \forall \, k \geq 0,$$

where $H^k(\mathbb{R}^d)$ represents a standard Sobolev space and the last inclusion results from its definition. The above relation and the assumption in Assumption 2 trivially imply that the space-time gradient of $f$ (i.e., $\nabla_{t,x} f$) is in $L^2(D; \mathbb{R}^{d+1})$, which further leads to

(7) $$f \in L^2(\mathbb{R}^d; H^1(\Omega)) \cap H_{\mathcal{L}}.$$

Later, during the convergence analysis error terms will appear along the boundary $(\partial\Omega \times (0, T))$ involving the moments of the traces of $f$, i.e., $\lambda_m(\gamma f)$, and due to Assumption 2 these error terms are well-defined. Indeed, $\lambda_m(\gamma f)$ is an element of $H^{\frac{1}{2}}(\partial\Omega \times (0, T); \mathbb{R}^{n(m)})$. Note that for strong solutions, the moments of the traces are not necessarily well-defined. The fact that $\gamma f \in L^2(\mathbb{R}^d; H^{\frac{1}{2}}(\partial\Omega \times (0, T)))$ is required by our analysis is the reason why we assume the boundary data ($f_{in}$ in (3)) to be in $L^2(\Sigma^-; |\xi_1|) \cap L^2(\mathbb{R}^- \times \mathbb{R}^{d-1}; H^{1/2}(\partial\Omega \times (0, T)))$, since for compatibility we want $\gamma^- f = f_{in}$ on $\Sigma^-$.

### 1.3. Moment approximation.

**Even and odd basis functions.** To formulate boundary conditions for our moment approximation (discussed next), we first need the notion of even and odd moments.

DEFINITION 1.4. *Let $n_o(m)$ and $n_e(m)$ denote the total number of tensorial Hermite polynomials in $\psi_m(\xi)$ which are odd and even, with respect to $\xi_1$, respectively. Similarly, let $\psi_m^o(\xi) \in \mathbb{R}^{n_o(m)}$ and $\psi_m^e(\xi) \in \mathbb{R}^{n_e(m)}$ represent vectors containing those basis functions out of $\psi_m(\xi)$ which are odd and even, with respect to $\xi_1$, respectively. Then, we define $\lambda_m^o : L^2(\mathbb{R}^d) \to \mathbb{R}^{n_o(m)}$ and $\lambda_m^e : L^2(\mathbb{R}^d) \to \mathbb{R}^{n_e(m)}$ as $\lambda_m^o(r) = \left\langle \psi_m^o \sqrt{f_0}, r \right\rangle_{L^2(\mathbb{R}^d)}$ and $\lambda_m^e(r) = \left\langle \psi_m^e \sqrt{f_0}, r \right\rangle_{L^2(\mathbb{R}^d)}$, where $r \in L^2(\mathbb{R}^d)$. To collect all the odd and even moments of $r$ which have a degree less than or equal to $M$ ($m \leq M$), we define*

$$\Psi_M^o(\xi) = (\psi_1^o(\xi)', \psi_2^o(\xi)', \dots \psi_M^o(\xi)')', \quad \Psi_M^e(\xi) = (\psi_0^e(\xi)', \psi_1^e(\xi)', \dots \psi_M^e(\xi)')',$$
$$\Lambda_M^o(r) = (\lambda_1^o(r)', \lambda_2^o(r)', \dots \lambda_M^o(r)')', \quad \Lambda_M^e(r) = (\lambda_0^e(r)', \lambda_1^e(r)', \dots \lambda_M^e(r)')',$$

*where $\Lambda_M^o : L^2(\mathbb{R}^d) \to \mathbb{R}^{\Xi_o^M}$, $\Lambda_M^e : L^2(\mathbb{R}^d) \to \mathbb{R}^{\Xi_e^M}$, $\Psi_M^o(\xi) \in \mathbb{R}^{\Xi_o^M}$, and $\Psi_M^e(\xi) \in \mathbb{R}^{\Xi_e^M}$. We represent the total number of odd and even moments of degree less than or equal to $M$ through $\Xi_o^M = \sum_{i=1}^M n_o(i)$ and $\Xi_e^M = \sum_{i=0}^M n_e(i)$, respectively.*

Expressions for boundary conditions become compact if we define the following matrices.

DEFINITION 1.5. *We define*

$$A_\psi^{(p,r)} = \left\langle \Psi_p^o \xi_1 \sqrt{f_0}, (\psi_r^e)' \sqrt{f_0} \right\rangle_{L^2(\mathbb{R}^d)}, \quad A_\Psi^{(p,q)} = \left( A_\psi^{(p,1)}, A_\psi^{(p,2)}, \dots, A_\psi^{(p,q)} \right).$$

*We interpret $\langle \Psi_p^o \xi_1 \sqrt{f_0}, (\psi_r^e)' \sqrt{f_0} \rangle_{L^2(\mathbb{R}^d)}$ as a matrix whose elements contain $L^2(\mathbb{R}^d)$ inner product between different elements of vectors $\Psi_p^o \sqrt{f_0}$ and $\xi_1 \psi_r^e \sqrt{f_0}$. Therefore, $A_\psi^{(p,r)}$ is a matrix with real entries of dimension $\Xi_o^p \times n_e(r)$. Moreover by definition, $A_\psi^{(p,r)}$ are the different groups of columns of $A_\Psi^{(p,q)}$ for $r \in \{1, \dots, q\}$.*

Recall that both $\Psi_q^e(\xi)$ and $\psi_q^e(\xi)$ are vectors but $\Psi_q^e(\xi)$ contains all those basis functions that have a degree less than or equal to $q$, whereas $\psi_q^e(\xi)$ contains basis functions of degree equal to $q$. Similar to the above matrices, we define the following matrices, which also contain the inner products between Hermite polynomials but on a half velocity space.

DEFINITION 1.6. *We define*

$$B_\psi^{(p,r)} = 2 \left\langle \Psi_p^o \sqrt{f_0}, (\psi_r^e)' \sqrt{f_0} \right\rangle_{L^2(\mathbb{R}^+ \times \mathbb{R}^{d-1})}, \quad B_\Psi^{(p,q)} = \left( B_\psi^{(p,1)}, B_\psi^{(p,2)}, \dots, B_\psi^{(p,q)} \right),$$

*where $B_\psi^{(p,r)} \in \mathbb{R}^{\Xi_o^p \times n_e(r)}$. Similar to $A_\psi^{(p,r)}$ defined above, $B_\psi^{(p,r)} \in \mathbb{R}^{\Xi_o^p \times n_e(r)}$ are the different groups of columns of $B_\Psi^{(p,q)}$ for $r \in \{1, \dots, q\}$.*

**Test and trial space.** To approximate the strong solution (see Theorem 1.2) to our kinetic equation (3), we use a Petrov–Galerkin type approach where we approximate the velocity dependence in the test space (i.e., $L^2(D)$) and in the solution space (i.e., $H_{\mathcal{L}}$) through a finite Hermite series expansion (4). Indeed, for our

Petrov–Galerkin approach, we choose the following test $(X_M)$ and the solution space $(H_M)$:

$$(8) \quad \left(L^2(\mathbb{R}^d; H^1(V)) \cap H_{\mathcal{L}}\right) \supset H_M := \{\alpha \cdot \Psi_M \sqrt{f_0} \ : \ \alpha \in H^1(V; \mathbb{R}^{\Xi^M})\},$$
$$L^2(D) \supset X_M := \{\alpha \cdot \Psi_M \sqrt{f_0} \ : \ \alpha \in L^2(V; \mathbb{R}^{\Xi^M})\},$$

where $\Psi_M$ is a vector containing all the Hermite polynomials up to a degree $M$; see Theorem 1.3. Since $\alpha \in H^1(V; \mathbb{R}^{\Xi^M})$, trivially, $H_M$ is a subset of $L^2(\mathbb{R}^d; H^1(V))$, which means that our Galerkin method is conforming. However, the fact that $H_M \subset H_{\mathcal{L}}$ is not obvious and we prove it in the following result.

LEMMA 1.7. *Let $H_M$ be as defined in* (8); *then,* $H_M \subset H_{\mathcal{L}}$.

*Proof.* Let $f \in H_M$. To prove our claim we need to show that $\mathcal{L}f \in L^2(D)$ for which we only need to show that $\xi \cdot \nabla_x f \in L^2(D)$; the definition of $H_M$ and boundedness of $Q$ on $L^2(\mathbb{R}^d)$ already implies that $\partial_t f \in L^2(D)$ and $Q(f) \in L^2(D)$. We show that $\xi \cdot \nabla_x f \in L^2(D)$ by proving that $\xi_i \partial_{x_i} f \in L^2(D) \ \forall \ i \in \{1, \ldots, d\}$. For brevity we consider $i = 1$; for other values of $i$ the result follows analogously. Computing $\|\xi_1 \partial_{x_1} f\|^2_{L^2(D)}$ by expressing $f$ as $f = \alpha \cdot \Psi_M \sqrt{f_0}$, we find

$$\|\xi_1 \partial_{x_1} f\|^2_{L^2(D)} = \|(\partial_{x_1}\alpha)' A \partial_{x_1}\alpha\|_{L^2(V)} \leq C \|\partial_{x_1}\alpha\|^2_{L^2(V;\mathbb{R}^{\Xi^M})} < \infty,$$

where $A = \langle \Psi_M \sqrt{f_0}, \xi_1^2 \Psi_M \sqrt{f_0}\rangle_{L^2(\mathbb{R}^d)}$. Above, the first inequality is a result of each entry of $A$ being bounded and the last inequality is a result of $\alpha \in H^1(V; \mathbb{R}^{\Xi^M})$. $\quad\square$

*Remark* 1. Note that for the BGK and the Boltzmann collision operator (given in section 3), $\sqrt{f_0}$ is the global equilibrium. Therefore, for both of these operators, an approximation in $H_M$ (given in (8)) is equivalent to expanding around the global equilibrium. This ensures that there exists a finite $M$ such that

$$(9) \quad \ker(Q) \subseteq \mathrm{span}\{\psi_{\beta^{(i)}}\sqrt{f_0}\}_{\|\beta^{(i)}\|_{l^1}=1,\ldots,M}.$$

The equilibrium state of the kinetic equation belongs to $\ker(Q)$ and the above conditions allow one to compute the same numerically. Note that for the linearized Boltzmann and the BGK operator, the above condition holds for $M = 2$ [4].

Collision operators of practical relevance known to us have $\sqrt{f_0}$ (or $f_0$ depending on the scaling) as their global equilibrium. If the global equilibrium is different from $f_0$, say, $\hat{f}_0$, then an expansion around $\hat{f}_0$ results in an approximation space different from $H_M$. If this approximation space has basis functions that satisfy the property of recursion (5b), orthogonality (5a), totality in $L^2(\mathbb{R}^d)$, even/odd parity (given in Theorem 1.4), etc., then we expect to have results similar to what we propose here. Considering a different approximation space is out of the scope of the present work.

**Variational formulation.** To develop our Galerkin approximation, in the definition of the strong solution (given in Theorem 1.2), we restrict the test space and the trial space to $X_M$ and $H_M$, respectively. This provides

*Find $f_M \in H_M$ such that*
$$(10a) \quad \langle v, \mathcal{L}f_M\rangle_{L^2(D)} = 0 \quad \forall \, v \in X_M, \quad \Lambda_M(f_M(0)) = \Lambda_M(f_I) \text{ on } \Omega,$$
$$(10b) \quad \Lambda_M^o(\gamma f_M) = R^{(M)} A_\Psi^{(M,M)} \Lambda_M^e(\gamma f_M) + \mathcal{G}(f_{in}) \text{ on } (0,T) \times \partial\Omega,$$

where $R^{(M)} \in \mathbb{R}^{\Xi_o^M \times \mathbb{R}^{\Xi_o^M}}$ is an s.p.d. matrix given as [22]

$$(11) \qquad R^{(M)} = B_\Psi^{(M,M-1)} \left( A_\Psi^{(M,M-1)} \right)^{-1}.$$

Invertibility of the matrix $A_\Psi^{(M,M-1)}$ follows from the recursion relation (5b) and is discussed in detail in Appendix B. Moreover, $\mathcal{G} : L^2(\mathbb{R}^- \times \mathbb{R}^{d-1}) \to \mathbb{R}^{\Xi_o^M}$ is defined as $\mathcal{G}(f_{in}) := \langle \Psi_M^o, f_{in} \rangle_{L^2(\mathbb{R}^- \times \mathbb{R}^{d-1})}$. Thus, $\mathcal{G}(f_{in})$ is a vector containing all the half-space odd moments of $f_{in}$. The variational form in (10a) and its initial condition follow trivially from the definition of a strong solution given in Theorem 1.2. However, the derivation of boundary conditions (10b) is more involved and one can find details of this derivation in [19, 21, 22]. For brevity, we refrain from discussing these details here.

The Galerkin formulation (10a) is $L^2$-stable and its stability results from the specific form of the boundary conditions given in (10b). Since stability will be crucial for developing error bounds, we present a brief derivation of the stability estimate. We choose $v$ as $f_M$ in (10a), consider (for simplicity) $f_{in} = 0$, use the negative semi-definiteness of $Q$, and perform integration-by-parts on the space-time derivatives to find

$$(12) \quad \begin{aligned} &\|f_M(T)\|_{L^2(\Omega \times \mathbb{R}^d)}^2 - \|f_M(0)\|_{L^2(\Omega \times \mathbb{R}^d)}^2 \\ &\leq -2 \left\langle \Lambda_M^o(\gamma f_M), A_\Psi^{(M,M)} \Lambda_M^e(\gamma f_M) \right\rangle_{L^2((0,T) \times \partial\Omega; \mathbb{R}^{\Xi_o^M})} \\ &= -2 \left\langle A_\Psi^{(M,M)} \Lambda_M^e(\gamma f_M), R^{(M)} A_\Psi^{(M,M)} \Lambda_M^e(\gamma f_M) \right\rangle_{L^2((0,T) \times \partial\Omega; \mathbb{R}^{\Xi_o^M})} \\ &\leq 0, \end{aligned}$$

where the last inequality is a result of $R^{(M)}$ being s.p.d. and all the boundary integrals are well-defined because $\Lambda_M(\gamma f_M) \in L^2(V; \mathbb{R}^{\Xi^M})$, which is a result of our definition of $H_M$ given in (8). Moreover, the integral on the boundary involving $A_\Psi^{(M,M)}$ results from the following, which results from the orthogonality of even and odd Hermite polynomials:

$$(13) \quad \begin{aligned} \int_{\mathbb{R}^d} \xi_1 (\gamma f_M)^2 d\xi &= 2 \int_{\mathbb{R}^d} \xi_1 (\gamma f_M)^o (\gamma f_M)^e d\xi \\ &= 2 \int_{\mathbb{R}^d} \left( \Lambda_M^o(\gamma f_M) \cdot \Psi_M^o(\xi)\sqrt{f_0} \right) \xi_1 \left( \Psi_M^e(\xi) \cdot \Lambda_M^e(\gamma f_M)\sqrt{f_0} \right) d\xi \\ &= 2 \left\langle \Lambda_M^o(\gamma f_M), A_\Psi^{(M,M)} \Lambda_M^o(\gamma f_M) \right\rangle_{\mathbb{R}^{\Xi_o^M}}. \end{aligned}$$

*Remark* 2. The variational form in (10a) is the same as the one that leads to the Grad's moment equations [14]. However, through (10a), we only recover the so-called *full moment* approximations [3, 26].

*Remark* 3. Grad [14] prescribes boundary conditions through $\Lambda_M^o(\gamma f_M) = B_\Psi^{(M,M)} \Lambda_M^e(\gamma f_M) + \mathcal{G}(f_{in})$ but they lead to $L^2$-instabilities [19, 21]. To see the difference between Grad's boundary conditions and those which lead to stability (10b), we use the expression for $R^{(M)}$ from (11) and subtract the boundary matrix in (10b) with the Grad's boundary matrix to find

$$(14) \qquad R^{(M)} A_\Psi^{(M,M)} - B_\Psi^{(M,M)} = \left( 0, \left[ R^{(M)} A_\psi^{(M,M)} - B_\psi^{(M,M)} \right] \right).$$

The above relation implies that the two boundary conditions differ only in terms of the highest order even moments of $f_M$, i.e., through $\lambda_M^e(f_M(t,x,.))$. This difference will show up in the convergence analysis and will influence the convergence order of our moment approximation.

*Remark* 4. In [10], the authors consider an IBVP for the radiative transport equation and develop an $L^2$-stable moment approximation for the same. Comparing our approach to that proposed in [10] is ongoing research and we hope to cater to it in the future. The framework proposed in [10] considers a bounded velocity domain, which does not have a radial direction. Therefore, the first step is to extend this framework to an unbounded velocity domain, and then to compare it to ours.

**2. Convergence analysis.** We outline the forthcoming convergence analysis in the following steps:

(i) *Define a projection operator:* We define a projection operator $\hat{\Pi}_M : L^2(\mathbb{R}^d; H^1(V)) \to H_M$ (with $H_M$ as defined in (8)) such that the trace of the projection satisfies the same type of boundary conditions as those satisfied by the moment approximation (10b). Such a projection operator helps us exploit the stability of the moment approximation (12) during error analysis.

(ii) *Decompose the error:* We decompose the moment approximation error into two parts:

$$(15) \qquad E_M = f - f_M = \underbrace{f - \hat{\Pi}_M f}_{P_M} + \underbrace{\hat{\Pi}_M f - f_M}_{e_M}.$$

Above, $e_M$ is the error in moments (or the expansion coefficients) and $P_M$ is the projection error.

(iii) *Bound for the projection error:* We derive a bound for $\|P_M\|_{L^2(D)}$ in terms of the moments of the solution, and using our regularity assumption (see Assumption 2) we show that $\|P_M\|_{L^2(D)} \to 0$ as $M \to \infty$.

(iv) *Bound for the error in moments:* Using stability of our moment approximation (12), we bound $\|e_M\|_{L^2(D)}$ in terms of $\|\mathcal{L}P_M\|_{L^2(D)}$, where $\mathcal{L}$ is the projection operator. We complete the analysis by showing that $\|\mathcal{L}P_M\|_{L^2(D)} \to 0$ as $M \to \infty$.

**2.1. The projection operator.** We sketch our formulation of the projection operator $\hat{\Pi}_M : L^2(\mathbb{R}^d; H^1(V)) \to H_M$. Let $r \in L^2(\mathbb{R}^d; H^1(V))$. We represent the projection $\hat{\Pi}_M r$ generically through $\hat{\Pi}_M r = (\hat{\Lambda}_M^o(r) \cdot \Psi_M^o + \hat{\Lambda}_M^e(r) \cdot \Psi_M^e(r))\sqrt{f_0}$, where $\hat{\Lambda}_M^o$ and $\hat{\Lambda}_M^e$ are linear functionals defined over $L^2(\mathbb{R}^d)$. For now assume that $\hat{\Pi}_M r \in H_M$ and that the trace of the projection (i.e., $\gamma\hat{\Pi}_M r$) is such that $\gamma(\hat{\Pi}_M r) = (\hat{\Lambda}_M^o(\gamma r) \cdot \Psi_M^o + \hat{\Lambda}_M^e(\gamma r) \cdot \Psi_M^e)\sqrt{f_0}$. Once we define $\hat{\Lambda}_M^o$ and $\hat{\Lambda}_M^e$, it will be trivial that both of these assumptions are satisfied. As mentioned earlier, we want $\gamma(\hat{\Pi}_M r)$ to satisfy the moment approximation's boundary conditions (10b). Since these boundary conditions have no restriction over the even moments, we choose $\hat{\Lambda}_M^e(r)$ to be the same as the even moments of $r$, i.e., we choose $\hat{\Lambda}_M^e(r) = \Lambda_M^e(r)$. However, coefficients of the odd basis functions are constrained by the moment approximation's boundary conditions (10b) and thus we choose them as $\hat{\Lambda}_M^o(r) = R^{(M)}A_\Psi^{(M,M)}\Lambda_M^e(r) + \mathcal{G}(r)$. Such a choice of $\hat{\Lambda}_M^o(r)$ ensures that, provided the inflow part of $r$ coincides with $f_{in}$, we have $\hat{\Lambda}_M^o(\gamma r) = R^{(M)}A_\Psi^{(M,M)}\Lambda_M^e(\gamma r) + \mathcal{G}(f_{in})$ along the boundary, i.e., the projection satisfies the boundary conditions of the moment approximation (10b). In the following, we summarize our projection operator and, for convenience, we also define the orthogonal projection operator.

DEFINITION 2.1. *We define* $\hat{\Pi}_M : L^2(\mathbb{R}^d; H^1(V)) \to H_M$ *as*

$$r(\cdot) \mapsto \left( \hat{\Lambda}_M^o(r) \cdot \Psi_M^o(\cdot) + \Lambda_M^e(r) \cdot \Psi_M^e(\cdot) \right) \sqrt{f_0(\cdot)} \quad with$$

$$\hat{\Lambda}_M^o(r) := R^{(M)} A_\Psi^{(M,M)} \Lambda_M^e(r) + \mathcal{G}(r).$$

*Similarly, with* $X_M$ *as given in* (8), *we define the orthogonal projection operator* $\Pi_M : L^2(D) \to X_M$ *as*

$$(\Pi_M r)(\xi) = (\Lambda_M^o(r) \cdot \Psi_M^o(\xi) + \Lambda_M^e(r) \cdot \Psi_M^e(\xi)) \sqrt{f_0(\xi)}, \quad r \in L^2(D).$$

*Remark* 5. In (10a), we prescribe the initial conditions using the orthogonal projection operator, but there is no unique way of doing so. Our convergence analysis covers all projection or interpolation operators which introduce errors that decay at least as fast as the moment approximation error $(E_M)$. Upcoming convergence analysis will clarify the fact that both $\hat{\Pi}_M$ and $\Pi_M$ satisfy these criteria. Therefore, for simplification, we prescribe the initial conditions through $f_M(0) = \hat{\Pi}_M f_I$, which ensures that $e_M(0) = 0$. Note that implementing $\hat{\Pi}_M$ is cumbersome and therefore for implementation, one might want to prescribe initial conditions using $\Pi_M$ or some other (easier to implement) interpolation.

*Remark* 6. Due to our definition of the projection operator $\hat{\Pi}_M$, the projection error $P_M$ (defined in (15)) is not orthogonal to the approximation space $H_M$. This is in contrast to the analysis in [12, 23], where the use of an orthogonal projection operator leads to a $P_M$ that is orthogonal to the approximation space.

**2.2. Extension to spatial domains with $C^2$ boundaries.** Velocity perpendicular to our spatial domain's boundary is $\xi_1$ and we have defined the projection operator $(\hat{\Pi}_M)$ with respect to this velocity; this is implicit in the definition of the operators $\mathcal{G}$ and $A_\Psi^{(M,M)}$. Since for the half-space $(\Omega = \mathbb{R}^- \times \mathbb{R}^{d-1})$ the boundary normal is the same at every boundary point, the definition of the projection operator remains the same for all boundary points. However, for a spatial domain other than the half-space, the normal along the boundary varies, which results in different boundary points having different projection operators. We briefly discuss a methodology to construct the projection operators for a $C^2$-domain, which can have a normal that varies along the boundary.

Let $\Omega \subset \mathbb{R}^d$ be a domain with a $C^2$ boundary. Then, for every point $x_0 \in \partial\Omega$ we can define a line which passes through $x_0$ and points toward the interior of the domain in the direction opposite to the normal at $x_0$ $(n(x_0))$: $L_{x_0} := \{x \in \Omega : x - x_0 = \alpha n(x_0), \alpha \in \mathbb{R}^-\}$. Since the boundary is $C^2$, there exists some $\delta > 0$ such that $\Omega_\delta := \{x \in \Omega : \operatorname{dist}(x, \partial\Omega) \geq \delta\}$ has the property that no two lines $L_{x_0}$ and $L_{x_1}$, for any $x_0, x_1 \in \partial\Omega$, intersect within $\Omega_\delta^c$.

Inside $\Omega_\delta$ we use the orthogonal projection $\Pi_M$, whereas outside of $\Omega_\delta$ we proceed as follows. For every $x \in \Omega_\delta^c$ (by definition of $\Omega_\delta$) there exists a unique $x_0$ such that $x \in L_{x_0}$. Let $\hat{\Pi}_M^{x_0}$ denote the projection operator accounting for the boundary conditions at $x_0$. Then at $x$ we define the projection operator to be the linear combination of the projection operator which satisfies the boundary conditions, $\hat{\Pi}_M^{x_0}$, and the orthogonal projection operator $\Pi_M$

$$\hat{\Pi}_M^x := \left( 1 - \frac{|x - x_0|}{\delta} \right) \hat{\Pi}_M^{x_0} + \frac{|x - x_0|}{\delta} \Pi_M.$$

In this way, $x \mapsto \hat{\Pi}_M^x(f_M(., x, .))$ satisfies the desired boundary conditions and is $C^1$.

*Remark* 7. We emphasize that the projection operator defined in Theorem 2.1 is an analytical tool defined such that the projection satisfies the same boundary conditions as those satisfied by the moment approximation. It is nowhere needed for computing the moment approximation. This is also clear from the variational formulation given in (10a), where we set to zero the orthogonal projection of the evolution equation onto the approximation space.

**2.3. Main result.** In the following, we summarize our main convergence result.

THEOREM 2.2. *We can bound the error in the moment approximation, $E_M = f - f_M$, as*

(16)
$$\|E_M(T)\|_{L^2(\Omega\times\mathbb{R}^d)} \le \|f(T) - \hat{\Pi}_M f(T)\|_{L^2(\Omega\times\mathbb{R}^d)} + T\left(A_1(T) + \|Q\|A_2(T) + A_3(T)\right),$$

*where*

$$A_1(T) = \Bigg(\Theta^{(M)}\|\lambda_M^e(\partial_t f)\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n_e(M)}))}$$

(17a)
$$+ \sqrt{2}\sum_{\beta\in\{e,o\}}\frac{1}{(2(M+1)+d)^{k_t^\beta}}\|(\partial_t f)^o\|_{C^0([0,T];L^2(\Omega;W_H^{k_t^\beta}(\mathbb{R}^d)))}\Bigg),$$

$$A_2(T) = \Bigg(\Theta^{(M)}\|\lambda_M^e(f)\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n_e(M)}))}$$

(17b)
$$+ \sqrt{2}\sum_{\beta\in\{e,o\}}\frac{1}{(2(M+1)+d)^{k^\beta}}\|f^\beta\|_{C^0([0,T];L^2(\Omega;W_H^{k^\beta}(\mathbb{R}^d)))}\Bigg),$$

$$A_3(T) = \sum_{i=1}^d\Big(\Theta^{(M)}\|A_\Psi^{(M,M)}\|_2\|\lambda_M^e(\partial_{x_i}f)\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n_e(M)}))}$$

$$+ \sqrt{(M+1)}\|\lambda_{M+1}(\partial_{x_i}f)\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n(M+1)}))}\Big)$$

(17c)
$$+ \frac{\|A_\Psi^{(M,M)}\|_2}{(2(M+1)+d)^{k_x^e}}\sum_{i=1}^d\|(\partial_{x_i}f)^e\|_{C^0([0,T];L^2(\Omega;W_H^{k_x^e}(\mathbb{R}^d)))},$$

(17d)  $$\Theta^{(M)} = \|R^{(M)}A_\psi^{(M,M)} - B_\psi^{(M,M)}\|_2.$$

*As $M\to\infty$, we have the convergence rate*

(18)  $$\|E_M(T)\|_{L^2(\Omega\times\mathbb{R}^d)} \le \frac{C}{M^\omega},\quad \omega = \min\left\{k^{e/o} - \frac{1}{2}, k_t^{e/o} - \frac{1}{2}, k_x^e - 1, k_x^o - \frac{1}{2}\right\}.$$

The motivation behind decomposing the right-hand side into the different $A_i$'s is that each of these terms vanishes in different physical settings. The term $A_1$ vanishes for steady state problems, i.e., for $\partial_t f = 0$, the term $A_2$ vanishes in the absence of collisions, and the term $A_3$ vanishes under spatial homogeneity, i.e., for $\partial_{x_i}f = 0$.

An alternative way to understand the right-hand side of the error bound given in Theorem 2.2 is to identify the following four different types of errors:

(i) *Projection error:* This is the first term appearing on the right side of the error bound in (16) and is the $P_M$ defined in (15).

(ii) *Closure error:* This is the second term appearing in $A_3(T)$, (17c), and involves the $M + 1$th order moment of $\partial_{x_i}f$. The term accounts for the influence of

the flux of the $M + 1$th order moment which was dropped out during the moment approximation.

(iii) *Boundary stabilization error:* These are all the terms involving $\Theta^{(M)}$ and are all the first terms appearing in (17a)–(17c). These terms are a result of the difference between the boundary conditions proposed by Grad [14] and those given in (10b) which lead to a stable moment approximation; Remark 3 explains the difference between the two boundary conditions. Since the two boundary conditions only differ in the coefficients of the highest order even moment (see (14)), this error depends only upon this highest order even moment.

(iv) *Boundary truncation error:* These are all the terms which are not included in the above definitions. They are a result of ignoring contributions from all those even (and odd) moments which have an order greater than $M$ and do not appear in the boundary conditions for the moment approximation (10b).

We prove Theorem 2.2 in the next few pages.

**2.4. Error equation.** To derive a bound for the moment approximation error (i.e., for $\|E_M(T)\|_{L^2(\Omega \times \mathbb{R}^d)}$) we first derive a bound for the error in the expansion coefficients (i.e., for $\|e_M(T)\|_{L^2(\Omega \times \mathbb{R}^d)}$) and then use triangle's inequality to arrive at a bound for $\|E_M(T)\|_{L^2(\Omega \times \mathbb{R}^d)}$; see (15) for definition of $E_M$ and $e_M$. In the following discussion we suppress dependencies on $x$ and $\xi$, for brevity.

We start with adding and subtracting $\mathcal{L}(\hat{\Pi}_M f)$ in the definition of a strong solution given in Theorem 1.2. For all $v \in X_M$, and $\forall\, t \in (0, T)$, considering the integral over $\Omega \times \mathbb{R}^d$ provides

$$\left\langle v(t), \mathcal{L}(\hat{\Pi}_M f(t)) \right\rangle_{L^2(\Omega \times \mathbb{R}^d)} = \left\langle v(t), \mathcal{L}(\hat{\Pi}_M f(t) - f(t)) \right\rangle_{L^2(\Omega \times \mathbb{R}^d)}$$
$$= \left\langle v(t), \Pi_M \mathcal{L}(\hat{\Pi}_M f(t) - f(t)) \right\rangle_{L^2(\Omega \times \mathbb{R}^d)},$$

where $X_M \subset L^2(D)$ is as defined in (8). For the last equality we have used the trivial relation $\langle v(t), w(t) \rangle_{L^2(\Omega \times \mathbb{R}^d)} = \langle v(t), \Pi_M w(t) \rangle_{L^2(\Omega \times \mathbb{R}^d)}$ $\forall (v, w) \in X_M \times L^2(D)$. Subtracting the above relation from our moment approximation (10a), and using the linearity of $\mathcal{L}$, we find

(19)
$$\langle v(t), \mathcal{L}(e_M(t)) \rangle_{L^2(\Omega \times \mathbb{R}^d)} = \left\langle v(t), \Pi_M \mathcal{L}(f(t) - \hat{\Pi}_M f(t)) \right\rangle_{L^2(\Omega \times \mathbb{R}^d)} \quad \forall v \in X_M, \forall t \in (0, T),$$

where $e_M$ is as given in (15). To derive a bound for $e_M$, we want to use the stability of our moment approximation (12). We do so by choosing $v(t) = e_M(t)$ in the above expression and by performing integration-by-parts on the spatial derivatives, which provides

(20) $\quad \langle e_M(t), \partial_t e_M(t) \rangle_{L^2(\Omega \times \mathbb{R}^d)} - \langle e_M(t), Q e_M(t) \rangle_{L^2(\Omega \times \mathbb{R}^d)}$

$$\leq \left\langle e_M(t), \Pi_M \mathcal{L}(f(t) - \hat{\Pi}_M f(t)) \right\rangle_{L^2(\Omega \times \mathbb{R}^d)} - \underbrace{\oint_{\partial \Omega} \int_{\mathbb{R}^d} \xi_1 (\gamma e_M(t))^2 d\xi ds}_{\geq 0}.$$

Later (section 3) we present physically relevant examples where the nondimensionalization of the kinetic equation results in the so-called Knudsen number, the inverse of which scales the collision operator. Depending on whether or not we are interested in

the low Knudsen number regime, we can proceed with the above bound in different ways. Here we consider a Knudsen number that is large enough and postpone the discussion of small Knudsen numbers to subsection 2.7. Since $Q$ is negative semidefinite, using the Cauchy–Schwarz inequality in the above bound provides

$$(21) \quad \langle e_M(t), \partial_t e_M(t) \rangle_{L^2(\Omega \times \mathbb{R}^d)} \leq \|e_M(t)\|_{L^2(\Omega \times \mathbb{R}^d)} \|\Pi_M \mathcal{L}(f(t) - \hat{\Pi}_M f(t))\|_{L^2(\Omega \times \mathbb{R}^d)}.$$

The integral over the boundary is positive because the trace of the projection (i.e., $\gamma \hat{\Pi}_M f$) satisfies the same boundary conditions as those satisfied by our moment approximation (10b). To see this more clearly, consider the following relation which results from the even-odd decoupling (13) and the moment equation's boundary conditions:

$$\oint_{\partial \Omega} \int_{\mathbb{R}^d} \xi_1 (\gamma e_M(t))^2 d\xi ds = \oint_{\partial \Omega} (\Lambda_M^o(\gamma e_M(t)))' A_\Psi^{(M,M)} \Lambda_M^e(\gamma e_M(t)) ds$$

$$= \oint_{\partial \Omega} (\Lambda_M^e(\gamma e_M(t)))' \left( A_\Psi^{(M,M)} \right)' R^{(M)} A_\Psi^{(M,M)} \Lambda_M^e(\gamma e_M(t)) ds \geq 0.$$

The last inequality is a result of $R^{(M)}$ being s.p.d. Using the fact that

$$\langle e_M(t), \partial_t e_M(t) \rangle_{L^2(\Omega \times \mathbb{R}^d)} = \|e_M(t)\|_{L^2(\Omega \times \mathbb{R}^d)} \partial_t \|e_M(t)\|_{L^2(\Omega \times \mathbb{R}^d)}$$

in (21), dividing throughout by $\|e_M(t)\|_{L^2(\Omega \times \mathbb{R}^d)}$ (the result is trivial for $e_M = 0$) and integrating over time provides the following bound:

$$(22) \quad \begin{aligned} \|e_M(T)\|_{L^2(\Omega \times \mathbb{R}^d)} &\leq \int_0^T \|\Pi_M \mathcal{L}(f(t) - \hat{\Pi}_M f(t))\|_{L^2(\Omega \times \mathbb{R}^d)} dt \\ &\leq T \|\Pi_M \mathcal{L}(f(t) - \hat{\Pi}_M f(t))\|_{C^0([0,T];L^2(\Omega \times \mathbb{R}^d))}. \end{aligned}$$

Above, our choice of the initial conditions (see Remark 5) results in $e_M(0) = 0$. To spell out the above term on the right, we use the definition of $\mathcal{L}$ from (1), the boundedness assumption on $Q$, and the triangle's inequality to find

(23)
$$\begin{aligned} \|\Pi_M \mathcal{L}(f(t) - \hat{\Pi}_M f(t))\|_{L^2(\Omega \times \mathbb{R}^d)} &\leq \|\partial_t f(t) - \hat{\Pi}_M \partial_t f(t)\|_{L^2(\Omega \times \mathbb{R}^d)} \\ &\quad + \|Q\| \|f(t) - \hat{\Pi}_M f(t)\|_{L^2(\Omega \times \mathbb{R}^d)} \\ &\quad + \sum_{i=1}^d \left\| \Pi_M \left( \xi_i \left( \partial_{x_i} f(t) - \hat{\Pi}_M \partial_{x_i} f(t) \right) \right) \right\|_{L^2(\Omega \times \mathbb{R}^d)}. \end{aligned}$$

We can further simplify $\|\Pi_M(\xi_i(\partial_{x_i} f(t) - \hat{\Pi}_M \partial_{x_i} f(t)))\|_{L^2(\Omega \times \mathbb{R}^d)}$ by adding and subtracting $\Pi_M \xi_i \Pi_M \partial_{x_i} f(t)$. Then, the triangle's inequality provides

$$(24) \quad \begin{aligned} &\left\| \Pi_M \left( \xi_i \left( \partial_{x_i} f(t) - \hat{\Pi}_M \partial_{x_i} f(t) \right) \right) \right\|_{L^2(\Omega \times \mathbb{R}^d)} \\ &\leq \left( \left\| \Pi_M \left( \xi_i \left( \Pi_M \partial_{x_i} f(t) - \hat{\Pi}_M \partial_{x_i} f(t) \right) \right) \right\|_{L^2(\Omega \times \mathbb{R}^d)} \right. \\ &\quad \left. + \left\| \Pi_M \left( \xi_i \left( \partial_{x_i} f(t) - \Pi_M \partial_{x_i} f(t) \right) \right) \right\|_{L^2(\Omega \times \mathbb{R}^d)} \right). \end{aligned}$$

To simplify the first term on the right we use [23, p. 80]

$$(25) \quad \begin{aligned} &\left\| \Pi_M \left( \xi_i \left( \Pi_M \partial_{x_i} f(t) - \hat{\Pi}_M \partial_{x_i} f(t) \right) \right) \right\|_{L^2(\Omega \times \mathbb{R}^d)} \\ &\leq \left\| A_\Psi^{(M,M)} \right\|_2 \left\| \left( \Pi_M \partial_{x_i} f(t) - \hat{\Pi}_M \partial_{x_i} f(t) \right) \right\|_{L^2(\Omega \times \mathbb{R}^d)}. \end{aligned}$$

Moreover, to simplify the second term on the right in (24) we use the orthogonality and the recursion of Hermite polynomials to find

(26)
$$\begin{aligned}
&\|\Pi_M \left( \xi_i \left( \partial_{x_i} f(t) - \Pi_M \partial_{x_i} f(t) \right) \right)\|_{L^2(\Omega \times \mathbb{R}^d)} \\
&= \left\| \Pi_M \left( \xi_i \left( \lambda_{M+1}(\partial_{x_i} f(t)) \cdot \psi_{M+1} \right) \sqrt{f_0} \right) \right\|_{L^2(\Omega \times \mathbb{R}^d)} \\
&\leq \sqrt{(M+1)} \|\lambda_{M+1}(\partial_{x_i} f(t))\|_{L^2(\Omega; \mathbb{R}^{n(M+1)})}.
\end{aligned}$$

Substituting (24)–(26) into (23) and substituting the resulting expression into the bound for $e_M$, we find the following bound for $\|E_M(T)\|_{L^2(\Omega \times \mathbb{R}^d)}$:

(27)
$$\begin{aligned}
\|E_M(T)\|_{L^2(\Omega \times \mathbb{R}^d)} &\leq \|f(T) - \hat{\Pi}_M f(T)\|_{L^2(\Omega \times \mathbb{R}^d)} + \|e_M(T)\|_{L^2(\Omega \times \mathbb{R}^d)} \\
&\leq \|f(T) - \hat{\Pi}_M f(T)\|_{L^2(\Omega \times \mathbb{R}^d)} + T\left( \tilde{A}_1(T) + \|Q\| \tilde{A}_2(T) + \tilde{A}_3(T) \right),
\end{aligned}$$

with

$$\begin{aligned}
\tilde{A}_1(T) &:= \|\partial_t f - \hat{\Pi}_M \partial_t f\|_{C^0([0,T]; L^2(\Omega \times \mathbb{R}^d))}, \\
\tilde{A}_2(T) &:= \|f - \hat{\Pi}_M f\|_{C^0([0,T]; L^2(\Omega \times \mathbb{R}^d))},
\end{aligned}$$

(28)
$$\begin{aligned}
\tilde{A}_3(T) &:= \sqrt{(M+1)} \sum_{i=1}^d \|\lambda_{M+1}(\partial_{x_i} f)\|_{C^0([0,T]; L^2(\Omega; \mathbb{R}^{n(M+1)}))} \\
&\quad + \|A_\Psi^{(M,M)}\|_2 \sum_{i=1}^d \|\Pi_M \partial_{x_i} f - \hat{\Pi}_M \partial_{x_i} f\|_{C^0([0,T]; L^2(\Omega \times \mathbb{R}^d))}.
\end{aligned}$$

The above expression is a bound for the moment approximation error in terms of the *closure error* and the *projection error* of different quantities. Rate of convergence for the *closure error* will trivially follow from the velocity space regularity assumption made in Assumption 2. Therefore, to complete our proof of Theorem 2.2 we develop a bound for the norm of $A_\Psi^{(M,M)}$ and a bound for the *projection error*. In particular, Theorem 2.5 will show

(29)
$$\tilde{A}_i(T) \leq A_i(T) \quad \text{for } i = 1, 2, 3,$$

where $A_i(T)$ are as defined in Theorem 2.2.

**2.5. Projection error.** The following result shows that we can express the odd moments of any $r \in L^2(\mathbb{R}^d)$ in terms of its even moments and the function $\mathcal{G}$ defined in (10b). The result will allow us to quantify the projection error in terms of the odd and the even moments of degree higher than $M$ which were left out while defining the projection operator $\hat{\Pi}_M$.

LEMMA 2.3. *For every* $r \in L^2(\mathbb{R}^d)$, *it holds that*

(30)
$$\left\langle \Psi_M^o \sqrt{f_0}, r^o \right\rangle_{L^2(\mathbb{R}^d)} = 2 \left\langle \Psi_M^o \sqrt{f_0}, r^e \right\rangle_{L^2(\mathbb{R}^+ \times \mathbb{R}^{d-1})} + \mathcal{G}(r),$$

*or equivalently* $\Lambda_M^o(r) = \lim_{q \to \infty} B_\Psi^{(M,q)} \Lambda_q^e(r) + \mathcal{G}(r)$ *where* $r^o$ *and* $r^e$ *are the odd and even parts of* $r$, *with respect to* $\xi_1$, *respectively, and* $\mathcal{G}$ *is as given in* (10b). *We interpret* $\lim_{q \to \infty} B_\Psi^{(M,q)} \Lambda_q^e(r)$ *as* $\lim_{q \to \infty} (B_\Psi^{(M,q)} \Lambda_q^e(r))$ *where* $B_\Psi^{(M,q)}$ *is as given in Theorem 1.6 and the limit is well-defined* $\forall r \in L^2(\mathbb{R}^d)$.

*Proof.* See Appendix A.                                                   □

In the following result, we collect all the relevant bounds on different matrices and operators. We will use these bounds to formulate the convergence rate of the *projection error*.

LEMMA 2.4.
(i) *For* $\lim_{q\to\infty} B_\Psi^{(M,q)}$ *it holds that* $\|\lim_{q\to\infty} B_\Psi^{(M,q)}\| \leq 1$ *where* $\lim_{q\to\infty} B_\Psi^{(M,q)}$ *is as given in Theorem* 2.3.
(ii) *For* $A_\Psi^{(M,M)}$ *and* $A_\Psi^{(M,M-1)}$ *it holds that* $\|(A_\Psi^{(M,M-1)})^{-1} A_\psi^{(M,M)}\|_2 \leq C\sqrt{M}$ *and* $\|A_\Psi^{(M,M)}\|_2 \leq C\sqrt{M}$.

*Proof.* See Appendix C. $\square$

Using the above results, in the following we develop a convergence rate and an error bound for the projection error.

LEMMA 2.5. *Let* $r^{e/o} \in C^0([0,T]; L^2(\Omega; W_H^{k^{e/o}}(\mathbb{R}^d)))$; *then we can bound* $\|\hat{\Pi}_M r(t) - r(t)\|_{L^2(\Omega\times\mathbb{R}^d)}^2$ *as*

$$\|\hat{\Pi}_M r(t) - r(t)\|_{L^2(\Omega\times\mathbb{R}^d)}^2 \leq (\Theta^{(M)})^2 \|\lambda_M^e(r(t))\|_{L^2(\Omega;\mathbb{R}^{n_e(M)})}^2$$
$$+ 2 \sum_{\beta\in\{e,o\}} \frac{1}{(2(M+1)+d)^{2k^\beta}} \|r^\beta(t)\|_{L^2(\Omega;W_H^{k^\beta}(\mathbb{R}^d))}^2,$$

*where* $\Theta^{(M)} = \|R^{(M)} A_\psi^{(M,M)} - B_\psi^{(M,M)}\|_2$ *and dependency on* $x$ *and* $\xi$ *is hidden for brevity. Similarly, we can bound the difference between the orthogonal projection and the projection that satisfies the boundary conditions as*

$$\|\hat{\Pi}_M r(t) - \Pi_M r(t)\|_{L^2(\Omega\times\mathbb{R}^d)}^2 \leq (\Theta^{(M)})^2 \|\lambda_M^e(r(t))\|_{L^2(\Omega;\mathbb{R}^{n_e(M)})}^2$$
$$+ \frac{1}{(2(M+1)+d)^{2k^e}} \|r^e(t)\|_{L^2(\Omega;W_H^{k^e}(\mathbb{R}^d))}^2.$$

*As* $M\to\infty$, *we have the convergence rate*

$$\|\hat{\Pi}_M r - r\|_{C^0([0,T];L^2(\Omega\times\mathbb{R}^d))} \leq CM^{-\tilde{\omega}}, \quad \|\hat{\Pi}_M r - \Pi_M r\|_{C^0([0,T];L^2(\Omega\times\mathbb{R}^d))} \leq CM^{-(k^e-\frac{1}{2})},$$

*where* $\tilde{\omega} = \min\left\{k^o - \frac{1}{2}, k^e - \frac{1}{2}\right\}$.

*Proof.* We express $r$ in terms of tensorial Hermite polynomials and use Theorem 2.3 to find

$$r = \sum_{m=0}^M (\lambda_m^o(r) \cdot \psi_m^o(\xi) + \lambda_m^e(r) \cdot \psi_m^e(\xi)) \sqrt{f_0}, \text{ with } \Lambda_M^o(r) = \lim_{q\to\infty} B_\Psi^{(M,q)} \Lambda_q^e(r) + \mathcal{G}(r),$$

where $\Lambda_M^o = (\lambda_1^o(r)', \ldots, \lambda_M^o(r)')$ and $\Lambda_M^e = (\lambda_0^e(r)', \ldots, \lambda_M^e(r)')$. Moreover, the definition of $\hat{\Pi}_M r$ (see Theorem 2.1) provides

$$\hat{\Pi}_M r = \sum_{m=0}^M \left(\hat{\Lambda}_m^o(r) \cdot \Psi_m^o(\xi) + \Lambda_m^e(r) \cdot \Psi_m^e(\xi)\right) \sqrt{f_0}$$
$$\text{with } \hat{\Lambda}_M^o(r) = R^{(M)} A_\Psi^{(M,M)} \Lambda_M^e(r) + \mathcal{G}(r),$$

where $\hat{\Lambda}_M^o = (\hat{\lambda}_1^o(r)', \ldots, \hat{\lambda}_M^o(r)')$. Subtracting $r$ from $\hat{\Pi}_M r$, using $\lim_{q\to\infty} B_\Psi^{(M,q)} \Lambda_q^e(r) = \sum_{q=0}^\infty B_\psi^{(M,q)} \lambda_q^e(r)$ and the simplified expression for $R^{(M)} A_\Psi^{(M,M)} - B_\Psi^{(M,M)}$ from (14), we find

$$\hat{\Pi}_M r - r = \left( \left( R^{(M)} A_\psi^{(M,M)} - B_\psi^{(M,M)} \right) \lambda_M^e(r) \right) \cdot \psi_M^o(\xi) \sqrt{f_0}$$

(31)
$$- \sum_{q=M+1}^{\infty} \left( B_\psi^{(M,q)} \lambda_q^e(r) \right) \cdot \psi_M^o(\xi) \sqrt{f_0}$$

$$- \sum_{q=M+1}^{\infty} \left( \lambda_q^e(r) \cdot \psi_q^e(\xi) + \lambda_q^o(r) \cdot \psi_q^o(\xi) \right) \sqrt{f_0},$$

where $B_\psi^{(M,M)}$ is as defined in Theorem 1.6. The matrices $B_\psi^{(M,q)}$ and the operator $\lim_{q\to\infty} B_\psi^{(M,q)}$ appearing above can be looked upon as restrictions of the operator $\lim_{q\to\infty} B_\Psi^{(M,q)}$ given in Theorem 2.4; thus all of their norms can be bounded by one. This provides

(32)
$$\|\hat{\Pi}_M r(t) - r(t)\|_{L^2(\Omega\times\mathbb{R}^d)}^2 \le \left( \Theta^{(M)} \right)^2 \|\lambda_M^e(r(t))\|_{L^2(\Omega;\mathbb{R}^{n_e(M)})}^2$$

$$+ 2 \sum_{\beta\in\{e,o\}} \sum_{q=M+1}^{\infty} \|\lambda_q^\beta(r(t))\|_{L^2(\Omega;\mathbb{R}^{n_\beta(q)})}^2$$

$$\le \left( \Theta^{(M)} \right)^2 \|\lambda_M^e(r(t))\|_{L^2(\Omega;\mathbb{R}^{n_e(M)})}^2$$

$$+ 2 \sum_{\beta\in\{e,o\}} \sum_{q=M+1}^{\infty} \frac{(2q+d)^{2k^\beta}}{(2(M+1)+d)^{2k^\beta}} \|\lambda_q^\beta(r(t))\|_{L^2(\Omega;\mathbb{R}^{n_\beta(q)})}^2$$

$$\le \left( \Theta^{(M)} \right)^2 \|\lambda_M^e(r(t))\|_{L^2(\Omega;\mathbb{R}^{n_e(M)})}^2$$

$$+ 2 \sum_{\beta\in\{e,o\}} \frac{1}{(2(M+1)+d)^{2k^\beta}} \|r^\beta(t)\|_{L^2(\Omega;W_H^{k^\beta}(\mathbb{R}^d))}^2,$$

where for the last inequality we use the definition

$$\|r^e(t)\|_{L^2(\Omega;W_H^{k^e}(\mathbb{R}^d))}^2 = \sum_{q=0}^{\infty} (2q+d)^{2k^e} \|\lambda_q^e(r(t))\|_{L^2(\Omega;\mathbb{R}^{n_o(q)})}^2.$$

The above relation proves the bound for $\|\hat{\Pi}_M r - r\|_{L^2(\Omega\times\mathbb{R}^d)}$. To prove the convergence rate we use the last inequality in (32). The convergence rate of terms involving $\|r^{e/o}(t)\|_{L^2(\Omega;W_H^{k^{e/o}}(\mathbb{R}^d))}$ follows trivially, and to obtain a convergence rate for the term involving $\Theta^{(M)}$ we use the definition of $R^{(M)}$ to find

$$\left( \Theta^{(M)} \right)^2 \|\lambda_M^e(r)\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n_e(M)}))}^2$$

$$= \|R^{(M)} A_\psi^{(M,M)} - B_\psi^{(M,M)}\|_2^2 \|\lambda_M^e(r)\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n_e(M)}))}^2$$

(33)
$$\le \left( \left\| \left( A_\Psi^{(M,M-1)} \right)^{-1} A_\psi^{(M,M)} \right\|_2 + \|B_\psi^{(M,M)}\|_2 \right)^2 \|\lambda_M^e(r)\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n_e(M)}))}^2$$

$$\le \frac{C}{M^{2k^e-1}}.$$

The last inequality in the above relation follows from the matrix norm bound given in Theorem 2.4 and from the following estimate:

$$
\begin{aligned}
\|\lambda_M^e(r(t))\|_{L^2(\Omega;\mathbb{R}^{n_e(M)})}^2 &\leq \sum_{m=M}^{\infty} \|\lambda_m^e(r(t))\|_{L^2(\Omega;\mathbb{R}^{n_e(M)})}^2 \\
&\leq \sum_{m=M}^{\infty} \left(\frac{2m+d}{2M+d}\right)^{2k^e} \|\lambda_m^e(r(t))\|_{L^2(\Omega;\mathbb{R}^{n_e(M)})}^2 \\
&\leq \frac{1}{(2M+d)^{2k^e}} \|r(t)\|_{L^2(\Omega;W_H^{k^e}(\mathbb{R}^d))}^2 .
\end{aligned}
$$
(34)

In a similar way, we prove the bound and the convergence rate for $\|\Pi_M r - \hat{\Pi}_M r\|_{C^0([0,T];L^2(\Omega\times\mathbb{R}^d)}$. Using the definition of $\Pi_M$ and $\hat{\Pi}_M$ from Theorem 2.1 we find

$$
\begin{aligned}
\hat{\Pi}_M r - \Pi_M r = \left(\left(R^{(M)}A_\psi^{(M,M)} - B_\psi^{(M,M)}\right)\lambda_M^e(r)\right) \\
\cdot \psi_M^o \sqrt{f_0} - \sum_{q=M+1}^{\infty}\left(B_\psi^{(M,q)}\lambda_q^e(r)\right)\cdot\psi_M^o(\xi)\sqrt{f_0},
\end{aligned}
$$

which implies

$$
\begin{aligned}
\|\hat{\Pi}_M r(t) - \Pi_M r(t)\|_{L^2(\Omega\times\mathbb{R}^d)}^2 &\leq \left(\Theta^{(M)}\right)^2 \|\lambda_M^e(r(t))\|_{L^2(\Omega;\mathbb{R}^{n_e(M)})}^2 \\
&+ \sum_{q=M+1}^{\infty}\|\lambda_q^e(r(t))\|_{L^2(\Omega;\mathbb{R}^{n_e(q)})}^2 .
\end{aligned}
$$

The above inequality is the same as the first inequality in (32) but without any contribution from the odd moments of degree higher than $M$. Therefore, we get the bound for $\|\hat{\Pi}_M r - \Pi_M r\|_{L^2(\Omega\times\mathbb{R}^d)}^2$ and its corresponding convergence rate from (32) and (33) by removing contribution from the odd moments of order higher than $M$. □

Using the result from Theorem 2.5 in the upper bound for $E_M$ (27) proves the error bound given in Theorem 2.2. To arrive at the convergence rate given in Theorem 2.2, first we split the bound for the *closure error* in Theorem 2.2 as

$$
\begin{aligned}
\sqrt{(M+1)}\|\lambda_{M+1}(\partial_{x_i}f)\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n(M+1)}))} \\
\leq \sqrt{(M+1)}\left(\|\lambda_{M+1}^o(\partial_{x_i}f)\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n_o(M+1)}))}\right. \\
\left.+\|\lambda_{M+1}^e(\partial_{x_i}f)\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n_e(M+1)}))}\right),
\end{aligned}
$$
(35)

which results from acknowledging that $\lambda_{M+1}(\partial_{x_i}f) = \left(\lambda_{M+1}^o(\partial_{x_i}f)', \lambda_{M+1}^e(\partial_{x_i}f)'\right)$. The bound for the individual moments of $r \in L^2(\Omega;W_H^k(\mathbb{R}^d))$ in terms of $\|r\|_{L^2(\Omega;W_H^k(\mathbb{R}^d))}$ (see (34)) implies that, with respect to $M$, the *closure error* decays as $\mathcal{O}(\min\{k_x^e - \frac{1}{2}, k_x^o - \frac{1}{2}\})$. The convergence rate for all the other terms in the error bound for $E_M$ follows from the fact that $\|A_\Psi^{(M,M)}\|_2 \leq C\sqrt{M}$ and from the convergence rate of the projection error.

**2.6. Sharper estimate.** As already noted in [12], a bound for the individual moments of $r \in L^2(\Omega;W_H^k(\mathbb{R}^d))$ in terms of $\|r\|_{L^2(\Omega;W_H^k(\mathbb{R}^d))}$ is pessimistic; see the relation in (34). Therefore, one can make the error bound in Theorem 2.2 sharper

by additionally assuming that the individual moments decay at a certain rate. The following result provides such a sharpened error bound, which is useful during numerical experiments because solutions to most numerical experiments have moments that decay at a certain rate [12, 26].

THEOREM 2.6. *In addition to Assumption* 2, *assume that*

$$(36) \quad \|\lambda_m^\beta(f)\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n_\beta}))} < \frac{C}{m^{k^\beta + \frac{1}{2}}}, \quad \|\lambda_m^\beta(\partial_t f)\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n_\beta}))} < \frac{C}{m^{k_t^\beta + \frac{1}{2}}},$$

$$(37) \quad \|\lambda_m^\beta(\partial_{x_i} f)\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n_\beta}))} < \frac{C}{m^{k_x^\beta + \frac{1}{2}}} \quad \forall\, i \in \{1, \ldots, d\},$$

*where* $\beta \in \{e, o\}$. *Then, we can sharpen the convergence rate presented in Theorem 2.2 to*

$$(38) \qquad \omega_{\mathrm{shp}} = \min\left\{ k^{e/o}, k_t^{e/o}, k_x^{e/o} - \frac{1}{2} \right\}.$$

*Proof.* The result trivially follows from the above analysis by using the assumed moment decay rate (36) instead of the pessimistic bound in (34). □

*Remark* 8. Note that the Hermite–Sobolev index in $W_H^k(\mathbb{R}^d)$, i.e., $k$, does not provide a decay rate for individual moments. However, if moments decay at a certain rate, i.e., if $\|\lambda_m(r)\|_{L^2(\Omega;\mathbb{R}^{n(m)})} \le \frac{C}{m^s}$, then $r \in L^2(\Omega; W_H^k(\mathbb{R}^d))$ for $k < s - \frac{1}{2}$. A detailed discussion can be found on p. *12* of [12].

**2.7. Uniform in Knudsen number estimate.** Here we are interested in the small Knudsen number regime and, in particular, we assume $\|Q\| > 0$. For convenience we define the seminorm

$$(39) \qquad |f|_Q := -\langle f, Q(f) \rangle_{L^2(\Omega \times \mathbb{R}^d)},$$

which is well-defined because of Assumption 1. We show that by treating the bound in (20) differently, we get a bound for $\|e_M(t)\|_{L^2(\Omega \times \mathbb{R}^d)}$ that scales with $\sqrt{\|Q\|}$, which (for small Knudsen numbers) is better than the scaling of $\|Q\|$ considered in Theorem 2.2. Moreover, we derive a uniform-in-Knudsen-number bound for the part of the error that is orthogonal to the null-space of $Q$. Precisely, for any function $f$ the seminorm $|f|_Q$ scales with $\mathrm{Kn}^{-1}$ by definition and we derive a linear-in-$\mathrm{Kn}^{-1}$-number bound for $|e_M|_Q$. Recall that the Knudsen number results from the nondimensionalization of the kinetic equation and is explicitly given below in (51).

From (20) we can infer

$$(40) \qquad \frac{d}{dt}\|e_M(t)\|_{L^2(\Omega \times \mathbb{R}^d)}^2 + |e_M(t)|_Q^2$$

$$\le (\bar{A}_1(t) + \bar{A}_3(t))\|e_M(t)\|_{L^2(\Omega \times \mathbb{R}^d)} + \|(-Q)^{\frac{1}{2}}\|\bar{A}_2(t)|e_M(t)|_Q$$

with

$$\bar{A}_1(t) := \|\Pi_M \partial_t f(t) - \hat{\Pi}_M \partial_t f(t)\|_{L^2(\Omega \times \mathbb{R}^d)},$$
$$\bar{A}_2(t) := \|f(t) - \hat{\Pi}_M f(t)\|_{L^2(\Omega \times \mathbb{R}^d)},$$
$$\bar{A}_3(t) := \sum_i \|\Pi_M(\xi_i(\partial_{x_i} f(t) - \hat{\Pi}_M \partial_{x_i} f(t)))\|_{L^2(\Omega \times \mathbb{R}^d)},$$

where we have used that $Q$ is self-adjoint and negative semidefinite, so that $-Q$ admits a square root. The discussion in (23)–(26) and Theorem 2.5 shows that $\forall\, t \in [0,T]$ and $i \in \{1,2,3\}$, we have

$$\tag{41} \bar{A}_i(t) \leq \tilde{A}_i(T) \leq A_i(T),$$

such that we infer that

$$\tag{42}
\begin{aligned}
\frac{d}{dt}\|e_M(t)\|^2_{L^2(\Omega\times\mathbb{R}^d)} &+ \frac{1}{2}|e_M(t)|^2_Q \\
&\leq (A_1(T) + \|Q\|^{\frac{1}{2}}A_2(T) + A_3(T))\|e_M(t)\|_{L^2(\Omega\times\mathbb{R}^d)} + \|Q\|A_2(T)^2 \\
&\leq \sqrt{2\left((A_1(T)+\|Q\|^{\frac{1}{2}}A_2(T)+A_3(T))^2\|e_M(t)\|^2_{L^2(\Omega\times\mathbb{R}^d)}+\|Q\|^2A_2(T)^4\right)}.
\end{aligned}
$$

Thus, $\forall\, t \in [0,T]$, $\|e_M(t)\|^2_{L^2(\Omega\times\mathbb{R}^d)}$ is bounded by $z(t)$ where $z$ solves

$$\tag{43} \frac{d}{dt}z(t) = \sqrt{2\left((A_1(T)+\|Q\|^{\frac{1}{2}}A_2(T)+A_3(T))^2 z(t) + \|Q\|^2A_2(T)^4\right)}$$

with $z(0) = \|e_M(0)\|^2_{L^2(\Omega\times\mathbb{R}^d)} = 0$. The solution $z$ satisfies

$$\tag{44}
\begin{aligned}
\sqrt{(A_1(T)+\|Q\|^{\frac{1}{2}}A_2(T)+A_3(T))^2 z(t) + \|Q\|^2A_2(T)^4} \\
= \frac{1}{\sqrt{2}}(A_1(T)+\|Q\|^{\frac{1}{2}}A_2(T)+A_3(T))^2 t + \|Q\|A_2(T)^2.
\end{aligned}
$$

The above relation provides

$$\tag{45}
\begin{aligned}
(A_1(T)+\|Q\|^{\frac{1}{2}}A_2(T)+A_3(T))^2 z(t) \\
\leq (A_1(T)+\|Q\|^{\frac{1}{2}}A_2(T)+A_3(T))^4 t^2 + \|Q\|^2A_2(T)^4,
\end{aligned}
$$

which results in

$$\tag{46}
\sup_{t\in[0,T]}\|e_M(t)\|^2_{L^2(\Omega\times\mathbb{R}^d)} \leq z(T) \leq (A_1(T)+\|Q\|^{\frac{1}{2}}A_2(T)+A_3(T))^2 T^2 + \|Q\|A_2(T)^2,
$$

and

$$\tag{47}
\begin{aligned}
\sup_{t\in[0,T]}\|e_M(t)\|_{L^2(\Omega\times\mathbb{R}^d)} \leq \sqrt{z(T)} &\leq (A_1(T)+\|Q\|^{\frac{1}{2}}A_2(T)+A_3(T))T \\
&+ \|Q\|^{\frac{1}{2}}A_2(T) =: B(T).
\end{aligned}
$$

It is worthwhile to note that the decay of $B(T)$ with respect to $M$ is the same as the decay of the bound derived in Theorem 2.2. Moreover, both the above bound and the bound in Theorem 2.2 are linear in time. However, while the bound in Theorem 2.2 scaled (for small Knudsen numbers) with $\|Q\|$, the bound in (47) scales with $\|Q\|^{\frac{1}{2}}$. In order to obtain a uniform-in-Knudsen-number bound for $|e_M(t)|_Q$, we return to (40) and integrate on $[0,T]$. This leads to the following.

THEOREM 2.7.

$$(48) \quad \int_0^T \frac{1}{2}|e_M(t)|_Q^2 dt \leq \int_0^T \left( (A_1(T) + A_3(T)) \|e_M(t)\|_{L^2(\Omega \times \mathbb{R}^d)} + \|Q\| A_2(T)^2 \right) dt,$$

$$\leq T \cdot \left( (A_1(T) + A_3(T))B(T) + \|Q\| A_2(T)^2 \right),$$

where $|\cdot|_Q$ is as defined in (39), $A_1, A_2,$ and $A_3$ are as defined in (17a)–(17c), and $B$ is as defined in (47).

We note the following for the above result:
1. The right-hand side in (48) is a bound for the square of the error and it decays with twice the rate of the right-hand side in Theorem 2.2.
2. Both sides of (48) scale with $\|Q\|$, i.e., it provides a uniform-in-Knudsen-number bound. It must be noted that $|e_M(t)|_Q$ is a seminorm and it does not quantify the part of $e_M(t)$ that is in the null-space of $Q$.

### 2.8. Discussion.

**Improved boundary conditions.** The model for the matrix $R^{(M)}$ (see (11)) is not unique and can be altered to enhance the accuracy of a moment approximation. For example, in [19] the authors did such an alteration for the R-13 moment equations using a data-driven approach. However, due to the absence of an error bound they did not analyze the correlation between the matrix $R^{(M)}$ and the R-13 moment approximation error.

With the error bound of the projection error, we develop some insight into the extent to which the matrix $R^{(M)}$ influences the convergence rate of a moment approximation. Consider the bound for the *projection error* given in Theorem 2.5. We decompose this bound into two parts:

$$\tilde{a} = \sum_{\beta \in \{e,o\}} \frac{1}{(2(M+1)+d)^{2k^\beta}} \|r^\beta\|^2_{L^2(\Omega; W_H^{k^\beta}(\mathbb{R}^d))} \text{ and } a_{\Theta^{(M)}} = (\Theta^{(M)})^2 \|\lambda_M^e(r)\|^2_{L^2(\Omega; \mathbb{R}^{n_e(M)})},$$

where $r^\beta \in L^2(\Omega; W_H^{k^\beta}(\mathbb{R}^d))$ for $\beta \in \{e,o\}$, and for simplicity we consider $k^e = k^o = k$. Clearly, $\tilde{a}$ is independent of $R^{(M)}$, whereas $a_{\Theta^{(M)}}$ is dependent upon $\Theta^{(M)}$, which then depends upon $R^{(M)}$.

Trivially, $\tilde{a}$ is $\mathcal{O}(M^{-k})$, whereas, since $\Theta^{(M)}$ is $\mathcal{O}(\sqrt{M})$, $\tilde{a}_{\Theta^{(M)}}$ is $\mathcal{O}(M^{-(k-\frac{1}{2})})$. Thus if one can improve the model for $R^{(M)}$ such that $\Theta^{(M)}$ decays faster than $\mathcal{O}(\sqrt{M})$, then one can obtain a moment approximation which converges faster than the one presented here. Development of such a $R^{(M)}$ is beyond our present scope and will be discussed in detail elsewhere.

**Suboptimality.** The convergence analysis presented in this paper is suboptimal. What we mean by optimality is twofold. First, optimality means that the difference between the numerical and the exact solution decays with the same rate as the best approximation error of the exact solution. Second, optimality would require that no additional conditions are imposed on the exact solution. For the case at hand, the rate of convergence of the best approximation error is the Hermite–Sobolev index. Our analysis requires additional assumptions in the sense that not only the solution but also its derivatives need to have some Hermite–Sobolev regularity. This is a common feature of the analysis of numerical schemes for hyperbolic problems; see, e.g., [6, 8, 10].

Recalling the convergence rate presented in Theorem 2.2, we find

$$(49) \qquad \omega = \min \left\{ k^{e/o} - \frac{1}{2}, k_t^{e/o} - \frac{1}{2}, k_x^e - \frac{1}{2} - \frac{1}{\underline{2}}, k_x^o - \frac{1}{2} \right\},$$

where $\omega$ is suboptimal with respect to the different Hermite–Sobolev indices, i.e., with respect to the different values of $k$. We elaborate on this particular suboptimality and show (through an example) that it results from the velocity domain in the kinetic equation being unbounded (3). Loss of half an order in all indices is a result of the boundary stabilisation error $(\Theta_M)$, which grows with $\sqrt{M}$. This error gets multiplied by $\|A_\Psi^{(M,M)}\|_2$, which grows with $\sqrt{M}$, and results in a suboptimality of an extra half appearing in the contribution from spatial derivatives; see the terms involving $A_3$ in Theorem 2.2.

Growth in $\|A_\Psi^{(M,M)}\|_2$, which also causes the growth in $\Theta_M$, is a result of the recursion relation of Hermite polynomials (5b), which states that the product of $\xi$ with an $M$th order Hermite polynomial equals a linear combination of an $(M-1)$th and an $(M+1)$th order Hermite polynomial but with factors which grow with $\sqrt{M}$. This growth results in the coefficients of $A_\Psi^{(M,M)}$ growing as $\mathcal{O}(\sqrt{M})$, which leads to a growth in the norm of $A_\Psi^{(M,M)}$. See Appendices B and C for details of the structure of $A_\Psi^{(M,M)}$ and $\Theta_M$, respectively. The use of Hermite polynomials as basis functions (and thus the growth in $\|A_\Psi^{(M,M)}\|_2$) is related to the velocity domain of the kinetic equation (3) being unbounded. For kinetic equations with a bounded velocity space, it might be possible to have basis functions such that $\|A_\Psi^{(M,M)}\|_2$ does not grow with $M$, which would remove the additional suboptimality in the Hermite–Sobolev indices of the spatial derivatives. As an example, consider the radiation transport equation for which the velocity space is a unit sphere and is thus bounded. A moment approximation can, therefore, be developed with the help of spherical harmonics, and contrary to Hermite polynomials, the recursion relation of spherical harmonics is such that $\|A_\Psi^{(M,M)}\|_2 \to 1$ as $M \to \infty$ [2, 10, 12]. Figure 1 shows a comparison between the norm of $A_\Psi^{(M,M)}$ for an $\mathbb{S}^2$ and an $\mathbb{R}^3$ velocity domain. Clearly, as $M$ is increased, for an $\mathbb{S}^2$ velocity space $\|A_\Psi^{(M,M)}\|_2$ approaches its limiting value of one, whereas for an $\mathbb{R}^3$ velocity space $\|A_\Psi^{(M,M)}\|_2$ grows with $\mathcal{O}(\sqrt{M})$. Thus for radiation transport, owing to the boundedness of $\|A_\Psi^{(M,M)}\|_2$ with $M$, we expect that one can entirely remove the second type of suboptimality present in $\omega$, i.e., one can get a convergence
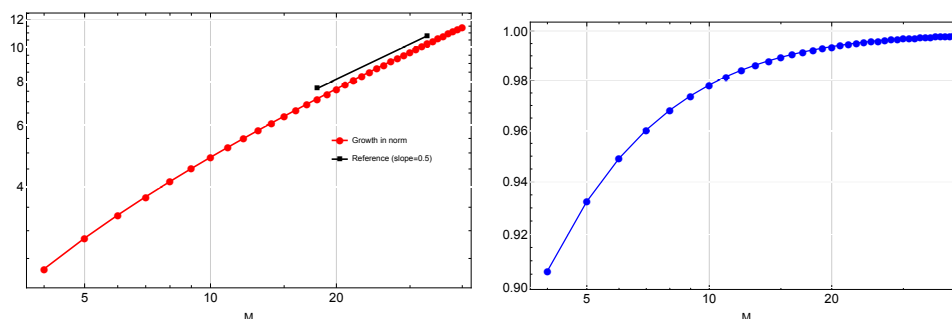


FIG. 1. *Growth in $\|A_\Psi^{(M,M)}\|_2$ with M for* (i) *left*, $\mathbb{R}^3$ *velocity space, and* (ii) *right*, $\mathbb{S}^2$ *velocity space.*

rate which is the same as the Hermite–Sobolev indices. Such a result would be in agreement with the error estimates presented in [10, 12].

**3. Examples: Linearized Boltzmann and BGK equations.** We give examples of kinetic equations which fall into the framework presented above. In particular, we discuss the conditions under which the linearized Boltzmann and the linearized BGK equation fall into our framework.

With $\bar{f} : D \to \mathbb{R}^+$, $(t, x, \xi) \mapsto \bar{f}(t, x, \xi)$, we denote the phase density function of a gas and we normalize $\bar{f}$ such that the density ($\bar{\rho}$), the mean flow velocity ($\bar{v}$), and the temperature in energy units ($\bar{\theta}$) of the gas are given as $\bar{\rho} = \int_{\mathbb{R}^d} \bar{f} d\xi$, $\bar{\rho}\bar{v} = \int_{\mathbb{R}^d} \xi \bar{f} d\xi$, $\bar{\rho}\bar{v} \cdot \bar{v} + d\bar{\rho}\bar{\theta} = \int_{\mathbb{R}^d} \xi \cdot \xi \bar{f} d\xi$. For convenience, we nondimensionalize all quantities with some reference density $\rho_0$, temperature $\theta_0$, and length scale $L$. The evolution of $\bar{f}$ is governed by the nonlinear kinetic equation given as [24]

$$(50) \qquad (1, \xi) \cdot \nabla_{(t,x)} \bar{f} = \frac{1}{\mathrm{Kn}} \bar{Q}(\bar{f}, \bar{f}),$$

where Kn is the so-called Knudsen number which results from nondimensionalisation, and $\bar{Q}$ is a nonlinear collision operator. We consider $\bar{Q}$ to be either the Boltzmann or the BGK collision operator given as

$$\text{Boltzmann Operator: } \bar{Q}_{\mathrm{BE}}(\bar{f}, \bar{f}) = \int_{\mathbb{R}^d} \int_{\mathbb{S}^{d-1}} \mathcal{B}(\xi - \xi_*, \kappa) \left( f(\xi^{'}) f_0(\xi_*^{'}) - f(\xi) f_0(\xi_*) \right) d\kappa d\xi_*;$$

$$\text{BGK Operator: } \bar{Q}_{\mathrm{BGK}}(\bar{f}, \bar{f}) = (\bar{f}_{\mathcal{M}} - \bar{f}).$$

Above, the velocities $\xi_*^{'}$ and $\xi^{'}$ are postcollisional and result from the precollisional velocities $\xi_*$ and $\xi$. The collision kernel ($\mathcal{B}$) depends on the interaction potential between the gas molecules and is nonnegative by physical assumptions. Moreover, $\bar{f}_{\mathcal{M}}$ is a Maxwell–Boltzmann distribution function given as

$$\bar{f}_{\mathcal{M}}(\xi; \bar{\rho}, \bar{v}, \bar{\theta}) = \frac{\bar{\rho}}{\sqrt[d]{2\pi\bar{\theta}}} \exp\left[ -\frac{(\xi - \bar{v}) \cdot (\xi - \bar{v})}{2\bar{\theta}} \right].$$

For low Mach number flows, we assume $\bar{f}$ to be a small perturbation of a ground state $f_0 = \bar{f}_{\mathcal{M}}(\xi; \rho_0, 0, \theta_0)$, i.e., $\bar{f} = f_0 + \epsilon\sqrt{f_0} f$, where $\epsilon$ is some smallness parameter. Substituting the linearization into the nonlinear kinetic equation (50) and considering only $\mathcal{O}(\epsilon)$ terms, we find the evolution equation for $f$

$$(51) \qquad (1, \xi) \cdot \nabla_{(t,x)} f = \frac{1}{\mathrm{Kn}} Q(f),$$

where $Q$ is the linearization of $\bar{Q}_{\mathrm{BE/BGK}}$ about $f_0$ and is given as

$$\text{Linearized Boltzmann Operator: } Q_{\mathrm{BE}}(\bar{f}) = \int_{\mathbb{R}^d} \int_{\mathbb{S}^{d-1}} \mathcal{B}(\xi - \xi_*, \kappa) \sqrt{f_0(\xi_*) f_0(\xi)}$$

$$\left( \frac{f(\xi^{'})}{\sqrt{f_0(\xi^{'})}} + \frac{f(\xi_*^{'})}{\sqrt{f_0(\xi_*^{'})}} - \frac{f(\xi_*)}{\sqrt{f_0(\xi_*)}} - \frac{f(\xi)}{\sqrt{f_0(\xi)}} \right) d\kappa d\xi_*;$$

$$\text{Linearized BGK Operator: } Q_{\mathrm{BGK}}(f) = (f_{\mathcal{M}} - \bar{f}).$$

Above, $f_{\mathcal{M}}\sqrt{f_0}$ is a linearization of $\bar{f}_{\mathcal{M}}$ about $f_0$ and is given as

$$(52) \qquad f_{\mathcal{M}}(\xi; \rho, v, \theta) := \left( \rho + v \cdot \xi + \frac{\theta}{2} (\xi \cdot \xi - 3) \right) \sqrt{f_0(\xi)},$$

where $\rho$, $v$, and $\theta$ are deviations of $\bar{\rho}$, $\bar{v}$, and $\bar{\theta}$ from their respective ground states.

We discuss whether the collision operators $Q_{\mathrm{BE/BGK}}$ satisfy Assumption 1. One can show that both $Q_{\mathrm{BE/BGK}}$ are negative semidefinite and self-adjoint and that $Q_{\mathrm{BGK}}$ is bounded on $L^2(\mathbb{R}^d)$; see [4] for details. Thus $Q_{\mathrm{BGK}}$ satisfies Assumption 1. Below in Remark 9 we summarize the assumptions that make $Q_{\mathrm{BE}}$ a bounded operator, which results in $Q_{\mathrm{BE}}$ satisfying Assumption 1.

As compared to the general kinetic equation (3), our example of the linearized Boltzmann (or the BGK) equation (51) has an additional factor of $1/\mathrm{Kn}$, which scales the collision operator. From the bound on $\|e_M(t)\|_{L^2(\Omega \times \mathbb{R}^d)}$ (in (47)) we find that such a scaling introduces a factor of $1/\sqrt{\mathrm{Kn}}$ in front of the term $\|Q\|^{\frac{1}{2}} A_2(T)$ appearing in the error bound. An asymptotic analysis in terms of the Knudsen number can tell us how the error bound (or equivalently $A_2(T)$) behaves as the Knudsen number is chosen smaller and smaller. The authors in [16] conduct such an analysis for initial value problems. For IBVPs, an asymptotic analysis is available only for the simplified Broadwell equation [17]. We hope to cover the asymptotic study of the error bound in our future work. Although the bound on $\|e_M\|_{L^2(\Omega \times \mathbb{R}^d)}$ is suboptimal in Kn, the bound on $|e_M|_Q$ (given in (48)) is uniform in Kn. However, the seminorm $|e_M|_Q$ only quantifies the part of the error that is orthogonal to the null-space of $Q$, and it is unclear how to get a uniform in Kn bound for the error in the null-space of $Q$.

*Remark* 9. Assume that we can split $Q_{\mathrm{BE}}$ as

$$(53) \quad Q_{\mathrm{BE}}(f)(\xi) = \tilde{Q}(f)(\xi) - v(\xi)f(\xi), \quad v(\xi) = \int_{\mathbb{R}^3} \int_{\mathbb{S}^2} \mathcal{B}(\xi - \xi_*, \kappa)\sqrt{f_0(\xi_*)}d\kappa d\xi_*,$$

where $v(\xi) \geq 0$ is the collision frequency and $\tilde{Q}$ is the remaining integral operator. The explicit form of $\tilde{Q}$ can be found in [7]. We can bound $Q$ on $L^2(\mathbb{R}^d)$ by bounding $\tilde{Q}$ and $v(\xi)$ on $L^2(\mathbb{R}^d)$ and $\mathbb{R}^+$, respectively.

We discuss assumptions that allow for the above splitting of $Q$ and for a bound on $\tilde{Q}$ and $v(\xi)$. Details related to our assumptions can be found in [4, 7, 15]. Assuming an inverse power law potential, we express $\mathcal{B}(\xi - \xi_*, \kappa)$ as

$$\mathcal{B}(\xi - \xi_*, \kappa) = \Psi(|\xi - \xi_*|)b(\cos\theta),$$

$$\Psi(|\xi - \xi_*|) = |\xi - \xi_*|^\gamma, \quad \gamma \in (-3, 1], \quad \cos\theta = \frac{\xi - \xi_*}{|\xi - \xi_*|} \cdot \kappa.$$

Assuming Grad's angular cut-off results in $\theta \mapsto b(\cos\theta) \in L^1([0, \pi])$. This makes $v(\xi)$ well-defined and allows us to split $Q$ as above in (53). The operator $\tilde{Q}$ is bounded on $L^2(\mathbb{R}^d)$ for $\gamma \in (-3, 1]$. Moreover, $|v(\xi)|$ is bounded $\forall \gamma \in (-3, 0]$. Therefore, $Q_{\mathrm{BE}}$ is bounded on $L^2(\mathbb{R}^d)$ for inverse power law potentials with an angular cut-off and $\gamma \in (-3, 0]$.

**4. Numerical results.** Through numerical experiments, we validate the convergence rates presented in the earlier sections by comparing the observed convergence rate with the predicted one. The solution to our numerical experiment has moments that decay at a certain rate and hence we use the sharper estimate presented in Theorem 2.6. With $f_{\mathrm{ref}}$ we denote the reference solution and we set $f_{\mathrm{ref}} = f_{M_{\mathrm{ref}}}$ with $M_{\mathrm{ref}}$ being sufficiently large. To compute the observed convergence rate, which we denote by $\omega_{\mathrm{obs}}$, we first compute the moment approximation error through $E_M(T) = f_{\mathrm{ref}}(T) - f_M(T)$. Then, we compute $\omega_{\mathrm{obs}}$ as the slope of the linear curve that minimizes the $L^2$ distance to the curve $(\log(M), \log(\|E_M(T)\|_{L^2(\Omega \times \mathbb{R}^d)}))$. The predicted convergence rate, which we denote by $\omega_{\mathrm{pre}}$, follows from Theorem 2.6 and is given as

$$\omega_{\text{pre}} = \min\left\{k^{e/o}, k_t^{e/o}, k_x^{e/o} - \frac{1}{2}\right\}.$$

To compute the different values of $k$ we first define the $L^2$ norms of the moments of $f_{\text{ref}}$ and its derivatives

(54)
$$N_m^{(x_i)} := \|\lambda_m(\partial_{x_i} f_{\text{ref}})\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n(m)}))}, \quad N_m^{(t)} := \|\lambda_m(\partial_t f_{\text{ref}})\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n(m)}))},$$
$$N_m := \|\lambda_m(f_{\text{ref}})\|_{C^0([0,T];L^2(\Omega;\mathbb{R}^{n(m)}))}.$$

Let $s^o$ represent the slope of the linear curve that has the minimum $L^2$ distance to the curve $(\log(m), \log(N_m^o))$ with $N_m^o$ being the same as $N_m$ but with a dependency on only the odd moments. We approximate $k^o$, and similarly the other $k$'s, by $k^o \approx s^o - 1/2$. Once values of $k$ are known we can compute $\omega_{\text{pre}}$ using the above expression. To quantify the discrepancy between the observed and the predicted convergence rates, we define

$$\Delta_\omega = \omega_{\text{obs}} - \omega_{\text{pre}}.$$

For simplicity, we stick to a one-dimensional physical and velocity space, i.e., $d = 1$ and $\Omega = (0,1)$. To discretize the one-dimensional physical space we use a discontinuous Galerkin (DG) discretization with first order polynomials and 500 elements. For temporal discretization, we use a fourth order explicit Runge–Kutta scheme. Our DG scheme is based upon a weak boundary implementation that preserves the stability of the moment approximation (12) on a spatially discrete level; see [27] for details. Note that in Theorem 2.2 we assumed $\Omega$ to be the half-plane but we can extend the analysis to $\Omega = (0,1)$ through the following argument. The projection operator ($\hat{\Pi}_M$ in Theorem 2.1) is defined with respect to the boundary conditions at $x = 1$ and a similar projection operator can also be constructed for the boundary conditions at $x = 0$. By taking a linear combination of the projection operation defined with respect to boundary conditions at $x = 0$ and $x = 1$, results analogous to those presented in Theorem 2.2 (and Theorem 2.6) can be obtained for $\Omega = (0,1)$.

As initial data we consider $f_I(x,\xi) = \frac{\rho_I(x)}{\sqrt{2\pi}} \exp(-\frac{\xi^2}{2})$ with $\rho_I(x) := \exp[-(x - 0.5)^2 \times 100]$, which corresponds to a Gaussian density profile with all the higher order moments being zero. As boundary data we consider vacuum at both the ends ($x = 0$ and $x = 1$), i.e., $f_{in} = 0$. As final time we consider $T = 0.3$, and we choose $M_{ref} = 200$.

Figure 2 shows the decay in the $L^2$ norm of the moments defined in (54), and the corresponding Hermite–Sobolev indices are given in Table 1. The moments of the solution and its derivatives have a Hermite–Sobolev index that is close to 1.5, which signifies that the reference solution is sufficiently regular along the velocity space. As expected, the moment approximation error decreases as the value of $M$ is increased; see Figure 3. However, contrary to the previous results [26], the convergence behavior of the approximation error does not show any oscillations.

Table 2 shows the observed and the predicted convergence rate. The observed approximation error converges with an order of 1.16 and is underpredicted by a value of 0.19. For the sake of validation, we also compute the convergence rates with the reference solution obtained through a discrete velocity method (DVM); see [18] for details of a DVM. With DVM as the reference, we obtain $\omega_{\text{pre}} = 0.98$, $\omega_{\text{obs}} = 1.15$, and $\Delta_\omega = \omega_{\text{obs}} - \omega_{\text{pre}} = 0.17$, which is very similar to the results obtained with a moment reference solution Table 2.
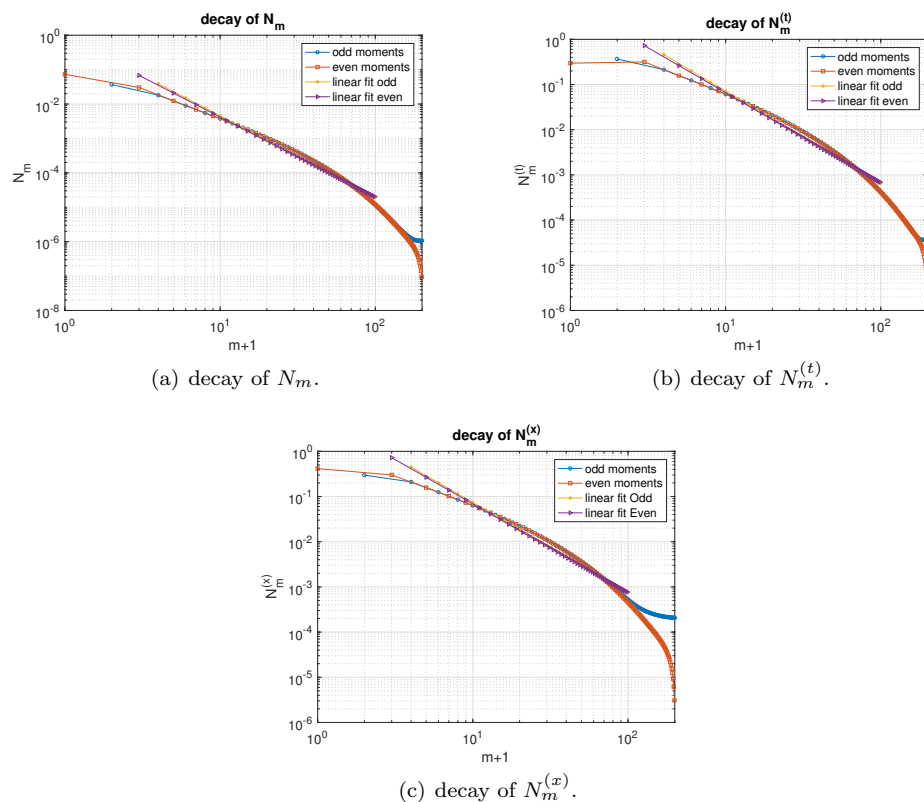
(a) decay of $N_m$.

(b) decay of $N_m^{(t)}$.



(c) decay of $N_m^{(x)}$.

FIG. 2. *Plots depict the decay of the various quantities, defined in* (54), *obtained through a refined moment approximation* ($M = 200$). *All plots are on a log-log scale.*
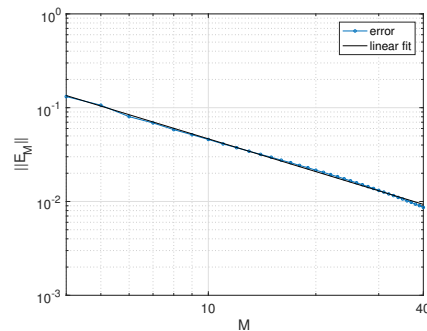
TABLE 1
*Hermite–Sobolev indices corresponding to the time integrated magnitude of moments defined in* (54).

| Quantity | Hermite–Sobolev index (= Decay Rate-0.5) |
|:---:|:---:|
| $N_m$ | 1.8 (= $k^e = k^o$) |
| $N_m^{(t)}$ | 1.45 (= $k_t^e = k_t^o$) |
| $N_m^{(x)}$ | 1.47 (= $k_x^e = k_x^o$) |

*Remark* 10. Authors in [12] observed that moment decay rates computed using $f_{\text{ref}}$ might show some artifacts for higher order moments. To remove these artifacts we follow the methodology proposed in [12], i.e., we compute decay rates from only those values of $N_m$'s whose values computed through $M_{ref}$ and $M_{ref} - 1$ differ by less than 3 percent.

**5. Conclusion.** Using a Galerkin type approach, under certain regularity assumptions on the solution, the global convergence of Grad's Hermite approximation to a linear kinetic equation was proved. The speed of convergence was quantified by proving the convergence rate which, as was expected, depends on the velocity space Sobolev regularity of the solution. The proposed convergence rate was found to be

FIG. 3. *Decay of the approximation error, on a log-log scale, for different values of $M$.*

TABLE 2
*Observed and predicted convergence rates.*

| Values of M | $\omega_{\mathrm{pre}}$ | $\omega_{\mathrm{obs}}$ | $\Delta_\omega = \omega_{\mathrm{obs}} - \omega_{\mathrm{pre}}$ |
|---|---|---|---|
| Odd | 0.97 | 1.16 | 0.19 |
| Even | 0.97 | 1.16 | 0.19 |

suboptimal, in the sense that it is one order lower than the convergence rate of the best approximation in the Galerkin spaces under consideration. Growth in the norm of the Jacobian corresponding to the flux of moment equations was found to be the reason for this suboptimality. For validation of the proven convergence rate, a numerical experiment involving the linearized BGK equation was conducted. For a moderately high Knudsen number ($\mathrm{Kn} = 0.1$), the observed convergence rate matched with the predicted convergence rate with acceptable accuracy.

**Appendix A. Proof of Lemma 2.1.** By splitting the integral over $\xi_1$, we find $\left\langle \Psi_M^o \sqrt{f_0}, r \right\rangle_{L^2(\mathbb{R}^d)} = \left\langle \Psi_M^o \sqrt{f_0}, r \right\rangle_{L^2(\mathbb{R}^+ \times \mathbb{R}^{d-1})} + \frac{1}{2}\mathcal{G}(r)$. Expressing $r$ as $r = r^e + r^o$ and using $\left\langle \Psi_M^o \sqrt{f_0}, r^e \right\rangle_{L^2(\mathbb{R}^d)} = 0$ in the previous expression, we find the desired result. To derive an expression equivalent to (30), we express $r^o$ and $r^e$ as $r^o = \sum_{m=1}^\infty \lambda_m^o(r) \cdot \psi_m^o \sqrt{f_0(\xi)}$ and $r^e = \sum_{m=0}^\infty \lambda_m^e(r) \cdot \psi_m^e \sqrt{f_0(\xi)}$, respectively, and replace these expansion in (30) to find $\Lambda_M^o(r) = \lim_{q\to\infty} B_\Psi^{(M,q)} \Lambda_q^e(r) + \mathcal{G}(r)$.

We consider $\lim_{q\to\infty} B_\Psi^{(M,q)}$ to be an operator defined over $l^2$ in the sense of

$$\left( \lim_{q\to\infty} B_\Psi^{(M,q)} \right) x := \left( \lim_{q\to\infty} B_\Psi^{(M,q)} x \right) \quad \forall\, x \in l^2.$$

We now show that $\lim_{q\to\infty} B_\Psi^{(M,q)}$ is well defined on $l^2$, which is equivalent to showing that the limit $\lim_{q\to\infty} B_\Psi^{(M,q)} x$ is well defined. Let $x \in l^2$ and let $x^q \in \mathbb{R}^q$ be a vector containing the first $q$ elements of $x$. To extend $x^q$ by zeros, we additionally define $\bar{x}^q \in l^2$, which has the same first $q$ elements as $x$ and whose other elements are all zero. From the definition of $B_\Psi^{(M,q)}$ (i.e., Theorem 1.6) we find $B_\Psi^{(M,q)} x^q = 2 \left\langle \Psi_M^o \sqrt{f_0}, g^q \right\rangle_{L^2(\mathbb{R}^+ \times \mathbb{R}^{d-1})}$, where $g^q = (\Psi_q^e \cdot x^q)\sqrt{f_0}$. Trivially, $\bar{x}^q$ converges to $x$ in $l^2$. This implies that $g^q$ converges in $L^2(\mathbb{R}^d)$. Then, by the continuity of the inner product of $L^2(\mathbb{R}^+ \times \mathbb{R}^{d-1})$, we have the convergence of $B_\Psi^{(M,q)} x^q$ in $\mathbb{R}^{\Xi_o^M}$.

**Appendix B. Structure of $A_\Psi^{(MM)}$.** We discuss in detail the structure of $A_\Psi^{(M,M)}$, which will be needed for the proof of Theorem 2.4. From the definition of $A_\Psi^{(M,M)}$ it is clear that it contains blocks of the integral

$$D^{(k,l)} = \left\langle \psi_k^o(\xi)\sqrt{f_0}, \xi_1 \psi_l^e(\xi)' \sqrt{f_0} \right\rangle_{L^2(\mathbb{R}^d)} \text{ and } D^{(M,M+1)} = 0,$$

where the second relation is a result of only considering basis functions up to degree $M$ in our moment approximation (10a). Recursion of the Hermite polynomials (5b) provides $\psi_k^o(\xi)\xi_1 = d^{(k,k-1)}\psi_{k-1}^e(\xi) + d^{(k,k+1)}\hat{\psi}_{k+1}^e$, where $\hat{\psi}_{k+1}^e$ is vector containing the first $n_o(k)$ components of $\psi_{k+1}^e$. Moreover, matrices $d^{(k,k-1)}, d^{(k,k+1)} \in \mathbb{R}^{n_o(k)\times n_o(k)}$ are diagonal matrices containing the square root entries appearing in the recursion relation. Using orthogonality of basis functions, we express $D^{(k,l)}$ as

$$(55) \quad D^{(k,l)} = \begin{cases} d^{(k,k-1)} \int_{\mathbb{R}^d} \psi_{k-1}^e(\xi)\psi_{k-1}^e(\xi)' f_0 d\xi = d^{(k,k-1)}, & l = k-1, \\ d^{(k,k+1)} \int_{\mathbb{R}^d} \hat{\psi}_{k+1}^e \left(\psi_{k+1}^e(\xi)\right)' f_0 d\xi = \begin{pmatrix} d^{(k,k+1)} & 0 \end{pmatrix}, & l = k+1, \\ 0 & \text{else} \end{cases}$$

Note that $D^{(k,k-1)} \in \mathbb{R}^{n_o(k)\times(n_e(k-1))}$, where $n_e(k-1) = n_o(k)$, whereas $D^{(k,k+1)} \in \mathbb{R}^{n_o(k)\times n_e(k+1)}$. Since $n_e(k) = n_o(k+1)$, $A_\Psi^{(M,M)}$ consists of blocks of $D^{(k,k-1)}$ on its main diagonal and blocks of $D^{(k,k+1)}$ on its off-diagonal with no entries below the main diagonal. From the recursion of Hermite polynomials (5b), we conclude

$$(56) \quad d_{ii}^{(k,k-1)} = \sqrt{\left(\beta_k^{(1,o)}\right)_i}, \quad d_{ii}^{(k,k+1)} = \sqrt{\left(\beta_k^{(1,o)}\right)_i + 1}, \quad i \in \{1, \ldots, n_o(k)\},$$

where $\beta_k^{(1,o)}$ is as defined below

DEFINITION B.1. *Let $\beta_k^o \in \mathbb{R}^{n_o(k)\times d}$ be such that each row of $\beta_k^o$ contains the multi-index of the odd basis functions contained in $\psi_k^o(\xi)$. Moreover, let $\beta_k^{(1,o)} \in \mathbb{R}^{n_o(k)}$ represent the first column of $\beta_k^o$.*

Note that all the entries in $\beta_k^{(1,o)}$ are odd. Therefore, all the entries along the diagonal of $d^{(k,k+1)}$ and $d^{(k,k-1)}$ are square roots of even and odd numbers, respectively. It can be shown that the number of times one appears in $\beta_k^{(1,o)}$ is equal to $k+2$. Thus, $d^{(k,k-1)}$ has the structure

$$(57) \qquad d^{(k,k-1)} = \begin{pmatrix} \tilde{d}^{(k,k-1)} & 0 \\ 0 & I^{k+2} \end{pmatrix},$$

where $\tilde{d}^{(k,k-1)} \in \mathbb{R}^{(n_o(k)-(k+2))\times(n_o(k)-(k+2))}$ and $I^{k+2}$ is an identity matrix of size $(k+2) \times (k+2)$. From (55), (56), and (57) we can conclude that

$$(58) \qquad D^{(k,k-1)} = \begin{pmatrix} \tilde{d}^{(k,k-1)} & 0 \\ 0 & I^{k+2} \end{pmatrix}, \quad D^{(k,k+1)} = \begin{pmatrix} d^{(k,k+1)}, & 0 \end{pmatrix}.$$

The matrix $A_\Psi^{(M,M-1)}$, which can be constructed by ignoring the contribution from $D^{(M-1,M)}$ into $A_\Psi^{(M,M)}$, is upper triangular with blocks of $D^{(k,k-1)}$ along its diagonal. Since $D^{(k,k-1)}$ contains square roots of odd numbers along its diagonal, which are all nonzero, the invertibility of $A_\Psi^{(M,M-1)}$ follows.

**Appendix C. Norms of matrices and operators.**    We will need the following result.

LEMMA C.1. *Let* $A \in \mathbb{R}^{n \times n}$, $n \geq 1$, *be given by* $A_{ij} = \sqrt{2i-1}\delta_{ij} + \sqrt{2i}\delta_{(i+1)j}$. *Then the solution* $x \in \mathbb{R}^n$ *to the linear system*

$$(59) \qquad\qquad A_{ij}x_j = \delta_{in}$$

*is such that* $\|x\|_{l^2} = 1$.

*Proof.* For $n = 1$, the result is trivial and so we consider the $n > 1$ case. From the first $n-1$ equations of the linear system (59) it follows $x_i\sqrt{2i-1} + x_{i+1}\sqrt{2i} = 0$, $i \in \{1, 2, \ldots n-1\}$, with which we can express any $x_p$ ($p \geq 2$) in terms of $x_1$ as

$$(60) \qquad x_p = (-1)^{p-1}\prod_{k=1}^{p-1}\sqrt{\frac{2k-1}{2k}}x_1 = (-1)^{p-1}\sqrt{\frac{(2p-3)!!}{(2p-2)!!}}x_1, \quad p \in \{2, \ldots n\}.$$

Thus

$$(61) \qquad\qquad \|x\|_{l^2}^2 = x_1^2\left(1 + \sum_{p=2}^{n}\frac{(2p-3)!!}{(2p-2)!!}\right) = x_1^2\sum_{p=0}^{n-1}\frac{1}{2^p p!}.$$

From the last equation in (59) and using (60) we have $x_n = 1/\sqrt{2n-1}$, which implies $x_1 = (-1)^{n-1}\sqrt{(2n-2)!!/(2n-1)!!}$. Using the expression for $x_1$ in (61), we find

$$\|x\|_{l^2}^2 = \frac{(2n-2)!!}{(2n-1)!!}\sum_{p=0}^{n-1}\frac{1}{2^p p!}.$$

Finally, induction provides $\sum_{p=0}^{n-1}1/(2^p p!) = (2n-1)!!/(2n-2)!!$, which implies $\|x\|_{l^2}^2 = 1$. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\Box$

(i) *Norm of* $\lim_{q\to\infty}B_\Psi^{(M,q)}$: Let $L = \lim_{q\to\infty}B_\Psi^{(M,q)}$ which is well-defined on $l^2$ due to Theorem 2.3. Define $y \in \mathbb{R}^{\Xi_o^M}$ as $y = Lx = 2\langle\Psi_M^o f_0, r\rangle_{K^+}$ where $r = \sum_{m=0}^{\infty}x_m \cdot \psi_m^e f_0$, $x = (x_0', x_1', \ldots, x_k', \ldots)'$ and $x_k \in \mathbb{R}^{n_e(k)}$. Functions $\sqrt{2}\psi_i^e f_0$ are orthonormal under $\langle.,.\rangle_{K^+}$. This implies $\|r\|_{K^+}^2 = \frac{1}{2}\|x\|_{l^2}^2$. Orthogonal projection of $r$ onto $\{\sqrt{2}\psi_m^o f_0\}_{m\leq M}$ can be given as $\mathcal{P}r = \sum_{m=1}^{M}y_m\cdot\psi_m^o f_0$ where $y = (y_1', y_2', \ldots, y_M')'$ and $y_k \in \mathbb{R}^{n_o(k)}$. Therefore, it holds that $\|\mathcal{P}r\|_{K^+} \leq \|r\|_{K^+}$. Since $\|\mathcal{P}r\|_{K^+}^2 = \|y\|_{l^2}^2/2$ and $\|r\|_{K^+}^2 = \|x\|_{l^2}^2/2$, we obtain $\|y\|_{l^2}^2 \leq \|x\|_{l^2}^2$, which provides $\|L\| \leq 1$.

(ii) *Norm of* $A_\Psi^{(M,M)}$: Let $A = A_\Psi^{(M,M)}(A_\Psi^{(M,M)})'$. Since every row of $A_\Psi^{(M,M)}$ contains two entries, one on the main diagonal and one on the off-diagonal (see Appendix B), every row of $A$ will contain a maximum of three entries. Since the maximum magnitude of entries in $A_\Psi^{(M,M)}$ is $\mathcal{O}(\sqrt{M})$, the maximum magnitude of the entries, in $A$, will be $\mathcal{O}(M)$. The Gerschgorin's circle theorem then implies that the maximum eigenvalue of $A$ will be $\mathcal{O}(M)$ which implies $\|A_\Psi^{(M,M)}\|_2 \leq C\sqrt{M}$.

(iii) *Norm of* $\|(A_\Psi^{(M,M-1)})^{-1}A_\psi^{(M,M)}\|_2$: In the coming discussion we will assume $M$ to be even; for $M$ being odd, the proof follows along similar lines and will not be discussed for brevity. From the definition of $A_\psi^{(M,M)}$ it is clear that it only has a contribution from $D^{(M-1,M)} \in \mathbb{R}^{n_o(M-1)\times n_e(M)}$, with

$D^{(M-1,M)}$ as defined in (58). Let $X \in \mathbb{R}^{\Xi_o^M \times n_o(M-1)}$ represent those columns of $(A_\Psi^{(M,M-1)})^{-1}$ which get multiplied with $D^{(M-1,M)}$ appearing in $A_\psi^{(M,M)}$. As a result $\|(A_\Psi^{(M,M-1)})^{-1} A_\psi^{(M,M)}\|_2 = \|XD^{(M-1,M)}\|_2 \leq \|X\|_2 \|D^{(M-1,M)}\|_2$. From (56) it follows that $\|D^{(M-1,M)}\|_2 \leq C\sqrt{M}$. We show that $X$ is unitary, which proves our claim.

Let $x^{(\omega)}$ denote the $\omega$th column of $X$ with $\omega \in \{1,\ldots,n_o(M-1)\}$. We decompose $x^{(\omega)}$ as $x^{(\omega)} = ((x_{n_e(0)}^{(\omega)})', (x_{n_e(1)}^{(\omega)})', \ldots, (x_{n_e(M-1)}^{(\omega)})')$ where $x_{n_e(q)}^{(\omega)} \in \mathbb{R}^{n_e(q)}$. Different values of $x^{(\omega)}$, for different values of $\omega$, can be found by solving the system of equations (which results from $A_\Psi^{(M,M-1)}(A_\Psi^{(M,M-1)})^{-1} = I$)

$$(62) \quad D^{(k,k-1)}x_{n_e(k-1)}^{(\omega)} + D^{(k,k+1)}x_{n_e(k+1)}^{(\omega)} = 0 \quad D^{(M,M-1)}x_{n_e(M-1)}^{(\omega)} = 0,$$

$$(63) \quad D^{(M-1,M-2)}x_{n_e(M-2)}^{(\omega)} = I_\omega^{n_o(M-1)},$$

where $I_\omega^{n_o(M-1)}$ is a diagonal matrix of size $n_o(M-1) \times n_o(M-1)$ such that $(I_\omega^{n_o(M-1)})_{ii} = \delta_{i\omega}$ and $D^{(k,k-1)}$ (and $D^{(k,k+1)}$) are as defined in (58). From (62) we conclude $x_{n_e(M-1)}^{(\omega)} = 0$, which implies $x_{n_e(M-(2q-1))}^{(\omega)} = 0 \ \forall q \in \{1, \ldots \frac{M}{2}\}$. We express the set of remaining equations as

$$(64) \quad \begin{aligned} D^{(k,k-1)}x_{n_e(k-1)}^{(\omega)} + D^{(k,k+1)}x_{n_e(k+1)}^{(\omega)} &= 0 \quad \forall\, k \in \{1,3,\ldots,M-3\}, \\ D^{(M-1,M-2)}x_{n_e(M-2)}^{(\omega)} &= I_\omega^{n_o(M-1)}. \end{aligned}$$

Orthogonality of solutions to (64) is clear from the structure of the linear system itself. Therefore, to prove our claim we need to show that

$$(65) \qquad \|x^{(\omega)}\|_{l^2} = 1 \ \forall\, \omega \in \{1,\ldots n_o(M-1)\},$$

for which we will claim that solving (64) for a given $\omega$ is equivalent to solving a system of the type (59); the result will then follow from Lemma C.1. From the entries of $d^{(k,k-1)}$ and $d^{(k,k+1)}$ defined in (56), it follows that the system in (64) is equivalent to

$$(66) \quad \begin{pmatrix} 1 & \sqrt{2} & 0 & 0 & \cdots & \cdots \\ 0 & \sqrt{3} & \sqrt{4} & 0 & \cdots & \cdots \\ 0 & 0 & \ddots & \ddots & 0 & \cdots \\ 0 & 0 & 0 & \cdots & \sqrt{(\beta_{M-1}^{(1,o)})_j - 2} & \sqrt{(\beta_{M-1}^{(1,o)})_j - 1} \\ 0 & 0 & 0 & \cdots & \cdots & \sqrt{(\beta_{M-1}^{(1,o)})_j} \end{pmatrix}$$
$$\times \begin{pmatrix} \left(x_{n_e(M-2q)}^{(\omega)}\right)_j \\ \left(x_{n_e(M-2(q-1))}^{(\omega)}\right)_j \\ \vdots \\ \left(x_{n_e(M-2)}^{(\omega)}\right)_j \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ \delta_{j,\omega} \end{pmatrix},$$

where $\beta_k^{(1,o)}$ is as defined in Theorem B.1, $q = ((\beta_{M-1}^{(1,o)})_j + 1)/2$, and for every $\omega, j \in \{1,\ldots,n_o(M-1)\}$. For $j = \omega$, the system in (66) is the same as (59) and hence (65) follows.

## REFERENCES

[1] R. BEALS AND V. PROTOPOPESCU, *Abstract time-dependent transport equations*, J. Math. Anal. Appl., 121 (1987), pp. 370–405, https://doi.org/10.1016/0022-247X(87)90252-6.

[2] T. A. BRUNNER AND J. P. HOLLOWAY, *Two-dimensional time dependent Riemann solvers for neutron transport*, J. Comput. Phys., 210 (2005), pp. 386–399, https://doi.org/10.1016/j.jcp. 2005.04.011.

[3] Z. CAI AND R. LI, *Numerical regularized moment method of arbitrary order for Boltzmann-BGK equation*, SIAM J. Sci. Comput., 32 (2010), pp. 2875–2907, https://doi.org/10.1137/100785466.

[4] C. CERCIGNANI, *The Boltzmann Equation and Its Applications*, Appl. Math. Sci. 67, Springer, New York, 1988.

[5] R. CHRISTIAN, *Numerical methods for the semiconductor Boltzmann equation based on spherical harmonics expansions and entropy discretizations*, Transp. Theory Statist. Phys., 31 (2002), pp. 431–452, https://doi.org/10.1081/TT-120015508.

[6] J. DOUGLAS, T. DUPONT, AND M. F. WHEELER, *A quasi-projection analysis of Galerkin methods for parabolic and hyperbolic equations*, Math. Comp., 32 (1978), pp. 345–362, http://www.jstor. org/stable/2006148.

[7] H. B. DRANGE, *The linearized Boltzmann collision operator for cut-off potentials*, SIAM J. Appl. Math., 29 (1975), pp. 665–676, http://www.jstor.org/stable/2100227.

[8] T. DUPONT, *L2-estimates for Galerkin methods for second order hyperbolic equations*, SIAM J. Numer. Anal., 10 (1973), pp. 880–889, http://www.jstor.org/stable/2156322.

[9] H. EGGER AND M. SCHLOTTBOM, *A mixed variational framework for the radiative transfer equation*, Math. Models Methods Appl. Sci., 22 (2012), 1150014, https://doi.org/10.1142/ S021820251150014X.

[10] H. EGGER AND M. SCHLOTTBOM, *A class of Galerkin schemes for time-dependent radiative transfer*, SIAM J. Numer. Anal., 54 (2016), pp. 3577–3599, https://doi.org/10.1137/ 15M1051336.

[11] L. FALK, *Existence of solutions to the stationary linear Boltzmann equation*, Transp. Theory Statist. Phys., 32 (2003), pp. 37–62, https://doi.org/10.1081/TT-120018651.

[12] M. FRANK, C. HAUCK, AND K. KUPPER, *Convergence of filtered spherical harmonic equations for radiation transport*, Commun. Math. Sci., 14 (2016), pp. 1443–1465.

[13] I. M. GAMBA AND S. RJASANOW, *Galerkin-Petrov approach for the Boltzmann equation*, J. Comput. Phys., 366 (2018), pp. 341–365, https://doi.org/10.1016/j.jcp.2018.04.017.

[14] H. GRAD, *On the kinetic theory of rarefied gases*, Comm. Pure Appl. Math., 2 (1949), pp. 331–407, https://doi.org/10.1002/cpa.3160020403.

[15] H. GRAD, *Asymptotic theory of the Boltzmann equation. II*, in Proceedings of the 3rd International Symposium, Palais de l'UNESCO, Paris, 1962, pp. 26–59, https://ci.nii.ac.jp/naid/ 10031083183/en/.

[16] H. GRAD, *Asymptotic theory of the Boltzmann equation*, Phys. Fluids, 6 (1963), pp. 147–181, https://doi.org/10.1063/1.1706716.

[17] J.-G. LIU AND Z. XIN, *Boundary-layer behavior in the fluid-dynamic limit for a nonlinear model Boltzmann equation*, Arch. Ration. Mech. Anal., 135 (1996), pp. 61–105, https://doi. org/10.1007/BF02198435.

[18] L. MIEUSSENS, *Discrete-velocity models and numerical schemes for the Boltzmann-BGK equation in plane and axisymmetric geometries*, J. Comput. Phys., 162 (2000), pp. 429–466, https://doi.org/10.1006/jcph.2000.6548.

[19] A. S. RANA AND H. STRUCHTRUP, *Thermodynamically admissible boundary conditions for the regularized 13 moment equations*, Phys. Fluids, 28 (2016), 027105, https://doi.org/10.1063/1. 4941293.

[20] C. RINGHOFER, C. SCHMEISER, AND A. ZWIRCHMAYR, *Moment methods for the semiconductor Boltzmann equation on bounded position domains*, SIAM J. Numer. Anal., 39 (2001), pp. 1078–1095, https://doi.org/10.1137/S0036142998335984.

[21] N. SARNA AND M. TORRILHON, *Entropy stable Hermite approximation of the linearised Boltzmann equation for inflow and outflow boundaries*, J. Comput. Phys., 369 (2018), pp. 16–44, https://doi.org/10.1016/j.jcp.2018.04.050.

[22] N. SARNA AND M. TORRILHON, *On stable wall boundary conditions for the Hermite discretization of the linearised Boltzmann equation*, J. Stat. Phys., 170 (2018), pp. 101–126, https://doi.org/10.1007/s10955-017-1910-z.

[23] C. SCHMEISER AND A. ZWIRCHMAYR, *Convergence of moment methods for linear kinetic equations*, SIAM J. Numer. Anal., 36 (1998), pp. 74–88, https://doi.org/10.1137/ S0036142996304516.

[24]  H. STRUCHTRUP, *Macroscopic Transport Equations for Rarefied Gas Flows*, Springer, New York, 2010.

[25]  S. THANGAVELU, *On regularity of twisted spherical means and special Hermite expansions*, Proc. Indian Acad. Sci. Math. Sci., 103 (1993), 303, https://doi.org/10.1007/BF02866993.

[26]  M. TORRILHON, *Convergence study of moment approximations for boundary value problems of the Boltzmann-BGK equation*, Commun. Comput. Phys., 18 (2015), pp. 529–557, https://doi.org/10.4208/cicp.061013.160215a.

[27]  M. TORRILHON AND N. SARNA, *Hierarchical Boltzmann simulations and model error estimation*, J. Comput. Phys., 342 (2017), pp. 66–84, https://doi.org/10.1016/j.jcp.2017.04.041.

[28]  S. UKAI, *Solutions of the Boltzmann equation*, in Patterns and Waves, Stud. Math. Appl. 18, Elsevier, Amsterdam, 1986, pp. 37–96, https://doi.org/10.1016/S0168-2024(08)70128-0.