

# Reduction and Inflation of Linear Models with an Application to Moment Closures of the Linearized Boltzmann Equation

C. David Levermore

Department of Mathematics *and*  
Institute for Physical Science and Technology  
University of Maryland, College Park  
[lvrmr@math.umd.edu](mailto:lvrmr@math.umd.edu)

presented 2 April 2019 during the KI-Net Workshop  
on *Dimension Reduction in Physical and Data Sciences*,  
held 1-3 April 2019 at the *Information Initiative at Duke*,  
Gross Hall, Duke University, Durham, NC

## Outline

1. Introduction
2. Model Reduction
3. Model Inflation
4. Dissipative Structure
5. Moment Closures for the Linearized Boltzmann Equation
6. Conclusion

## 1. Introduction

We discuss the following.

- Model reduction builds a smaller model from a larger one, but there are limitations on how this can be done.
- Model inflation builds a larger model from a smaller one within the framework of a large family of models. (This is a learning algorithm.)
- Dissipative structure can and should be preserved by model reduction.
- These ideas can be used to build a large family of well-posed moment closures for the linearized Boltzmann equation.

## 2. Model Reduction

Consider the linear initial-value problem

$$\frac{dU}{dt} + AU = F(t), \quad U(0) = U^{\text{in}}, \quad (1)$$

where  $U(t) \in \mathbb{R}^N$ ,  $A \in \mathbb{R}^{N \times N}$ ,  $F(t) \in \mathbb{R}^N$ , and  $U^{\text{in}} \in \mathbb{R}^N$ . We will assume that  $A$  is nonnegative definite in the sense that

$$V^{\text{T}}AV \geq 0, \quad \text{for every } V \in \mathbb{R}^N.$$

Suppose that

$$U = \begin{pmatrix} U_1 \\ U_2 \end{pmatrix}, \quad A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}, \quad F = \begin{pmatrix} F_1 \\ F_2 \end{pmatrix}, \quad U^{\text{in}} = \begin{pmatrix} U_1^{\text{in}} \\ U_2^{\text{in}} \end{pmatrix},$$

where  $U_1(t) \in \mathbb{R}^n$ ,  $A_{11} \in \mathbb{R}^{n \times n}$ ,  $F_1(t) \in \mathbb{R}^n$ , and  $U_1^{\text{in}} \in \mathbb{R}^n$  for some  $n < N$ . Our goal will be to find a reduced system that describes the evolution of  $U_1(t)$ .

Then (1) becomes the system

$$\begin{aligned} \frac{d}{dt} \begin{pmatrix} U_1 \\ U_2 \end{pmatrix} + \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} \begin{pmatrix} U_1 \\ U_2 \end{pmatrix} &= \begin{pmatrix} F_1(t) \\ F_2(t) \end{pmatrix}, \\ \begin{pmatrix} U_1(0) \\ U_2(0) \end{pmatrix} &= \begin{pmatrix} U_1^{\text{in}} \\ U_2^{\text{in}} \end{pmatrix}. \end{aligned} \tag{2}$$

By making different balances, we can derive different reduced models that approximate  $U_1$  by the solution  $u$  of an initial-value problem in the form

$$m \frac{du}{dt} + au = f(t), \quad u(0) = u^{\text{in}}, \tag{3}$$

where  $u(t) \in \mathbb{R}^n$ ,  $m \in \mathbb{R}^{n \times n}$ ,  $a \in \mathbb{R}^{n \times n}$ ,  $f(t) \in \mathbb{R}^n$ , and  $u^{\text{in}} \in \mathbb{R}^n$ . We will see that there are limitations to reduced models of this form.

The *Galerkin approximation* sets  $U_2 = 0$  in the first equation of system (2). This approximation is

$$U_1 = u \quad \text{and} \quad U_2 = 0,$$

where  $u$  satisfies the reduced initial-value problem (3) with

$$\begin{aligned} m &= I, \\ a &= A_{11}, \\ f(t) &= F_1(t), \\ u^{\text{in}} &= U_1^{\text{in}}. \end{aligned}$$

This approximation is dissipative because  $A \geq 0$  implies that

$$a = A_{11} \geq 0.$$

However, when  $F_1$  and  $F_2$  are constants then its stationary solution is generally *not* correct because  $F_2$  is missing from the approximation.

Better approximations can be obtained by also considering the second equation of system (2). By formally solving this equation for  $U_2(t)$  we find that

$$U_2(t) = e^{-tA_{22}}U_2^{\text{in}} + \int_0^t e^{-(t-t')A_{22}}(F_2(t') - A_{21}U_1(t')) dt'.$$

Upon placing this result into the first equation in (2) we obtain

$$\begin{aligned} \frac{dU_1}{dt} + A_{11}U_1 - \int_0^t A_{12}e^{-(t-t')A_{22}}A_{21}U_1(t') dt' \\ = F_1(t) - \int_0^t A_{12}e^{-(t-t')A_{22}}F_2(t') dt' - A_{12}e^{-tA_{22}}U_2^{\text{in}}, \end{aligned} \quad (4)$$

$$U_1(0) = U_1^{\text{in}}.$$

This is a good starting point for deriving better reduced models. We can avoid initial-layer asymptotics by assuming that  $U_2^{\text{in}} = 0$ .

For example, if we assume that

- the eigenvalues of  $A_{22}^{-1}$  each have a positive real part ,
- $A_{22}^{-1}$  is small compared to the time scale  $\tau$  over which  $F_1$  and  $F_2$  vary ,

then we may use the *Laplace approximation* (Taylor expanding  $U_1(t')$  and  $F_2(t')$  about  $t$  and neglecting all terms containing  $e^{-tA_{22}}$ ) to obtain

$$\int_0^t e^{-(t-t')A_{22}} A_{21} U_1(t') dt' = A_{22}^{-1} A_{21} U_1 - A_{22}^{-2} A_{21} \frac{dU_1}{dt} + A_{22}^{-3} A_{21} \frac{d^2 U_1}{dt^2} + \dots ,$$

$$\int_0^t e^{-(t-t')A_{22}} F_2(t') dt' = A_{22}^{-1} F_2(t) - A_{22}^{-2} \frac{dF_2}{dt}(t) + A_{22}^{-3} \frac{d^2 F_2}{dt^2}(t) + \dots .$$
(5)

These expansions are uniformly asymptotic in  $\frac{1}{\tau} A_{22}^{-1}$ .



When  $F_1$  and  $F_2$  are comparable the *relaxation approximation* is

$$U_1 = u \quad \text{and} \quad U_2 = -A_{22}^{-1} A_{21} u,$$

where  $u$  satisfies the reduced initial-value problem (3) with

$$\begin{aligned} m &= I, \\ a &= A_{11} - A_{12} A_{22}^{-1} A_{21}, \\ f(t) &= F_1(t), \\ u^{\text{in}} &= U_1^{\text{in}}. \end{aligned}$$

This approximation is dissipative because  $A \geq 0$  implies that

$$a = A_{11} - A_{12} A_{22}^{-1} A_{21} \geq 0.$$

However, when  $F_1$  and  $F_2$  are constants then its stationary solution is generally *not* correct because  $F_2$  is missing from the approximation.

When  $F_2$  is larger than  $F_1$  the *relaxation approximation* becomes

$$U_1 = u \quad \text{and} \quad U_2 = A_{22}^{-1} F_2(t) - A_{22}^{-1} A_{21} u,$$

where  $u$  satisfies the reduced initial-value problem (3) with

$$\begin{aligned} m &= I, \\ a &= A_{11} - A_{12} A_{22}^{-1} A_{21}, \\ f(t) &= F_1(t) - A_{12} A_{22}^{-1} F_2(t), \\ u^{\text{in}} &= U_1^{\text{in}}. \end{aligned}$$

This approximation is also dissipative because  $A \geq 0$  implies that

$$a = A_{11} - A_{12} A_{22}^{-1} A_{21} \geq 0.$$

When  $F_1$  and  $F_2$  are constants then its stationary solution is correct. This relaxation approximation is identical to the first when  $F_2 = 0$ .

When  $F_1$  and  $F_2$  are comparable the *quasi-relaxation approximation* is

$$U_1 = u \quad \text{and} \quad U_2 = A_{22}^{-1}F_2(t) - A_{22}^{-1}A_{21}u + A_{22}^{-2}A_{21}\frac{du}{dt},$$

where  $u$  satisfies the reduced initial-value problem (3) with

$$\begin{aligned} m &= I + A_{12}A_{22}^{-2}A_{21}, \\ a &= A_{11} - A_{12}A_{22}^{-1}A_{21}, \\ f(t) &= F_1(t) - A_{12}A_{22}^{-1}F_2(t), \\ u^{\text{in}} &= U_1^{\text{in}}. \end{aligned}$$

This approximation is generally *not* dissipative because  $m$  can behave badly. It will be dissipative if  $A^2 \geq 0$  as well as  $A \geq 0$ . For example, if  $A^T = A$  then we have  $A_{22}^T = A_{22} > 0$ ,  $A_{21}^T = A_{12}$ , and

$$m = m^T = I + A_{21}^T A_{22}^{-2} A_{21} > 0.$$

When  $F_1$  and  $F_2$  are constants then its stationary solution is correct.

When  $F_2$  is larger than  $F_1$  the *quasi-relaxation approximation* becomes

$$U_1 = u \quad \text{and} \quad U_2 = A_{22}^{-1}F_2(t) - A_{22}^{-1}A_{21}u + A_{22}^{-2}A_{21}\frac{du}{dt},$$

where  $u$  satisfies the reduced initial-value problem (3) with

$$\begin{aligned} m &= I + A_{12}A_{22}^{-2}A_{21}, \\ a &= A_{11} - A_{12}A_{22}^{-1}A_{21}, \\ f(t) &= F_1(t) - A_{12}A_{22}^{-1}F_2(t) + A_{12}A_{22}^{-2}\frac{dF_2}{dt}(t), \\ u^{\text{in}} &= U_1^{\text{in}}. \end{aligned}$$

This approximation also is generally *not* dissipative because  $m$  can behave badly. When  $F_1$  and  $F_2$  are constants then its stationary solution is correct. This quasi-relaxation approximation is identical to the first when  $\frac{d}{dt}F_2 = 0$ .

**Remark.** Because the third- and higher-order terms in the first asymptotic expansion of (5) contain second-order and higher-order derivatives of  $U_1$ , *these five are the only such temporally-local, first-order reductions.*

**Remark.** The Laplace transform approach is to approximate the 11 block of  $(sI - A)^{-1}$  by  $(sm - a)^{-1}m$ . The results are the same.

**Remark.** We can relax the assumption  $U_2^{\text{in}} = 0$  with some initial layer asymptotics. The result does not modify  $u^{\text{in}}$ .

**Remark.** We can approximate  $A_{22}^{-1}$  in the relaxation and quasi-relaxation approximations. Each such approximation gives rise to a new reduced model based upon the original approximation.

### **3. Model Inflation**

Model inflation builds a larger model from a smaller one. Here we present an approach that has four ingredients.

1. Imbed your current model in a large family of larger models.
2. Identify a few objective functions that you wish to predict.
3. Use adjoint sensitivity analysis to compute the infinitesimal response of your current objective functions to every parameter in the family.
4. Enlarge your current model to capture the parameters to which your objective functions are most sensitive and repeat.

We illustrate the approach in the nonlinear setting in which the family of larger models for the unknown vector  $U$  of dimension  $M$  has the form

$$F(U, P) = 0,$$

where the parameter vector  $P$  has dimension  $N \gg M$ . We assume that for every  $P$  there exists a unique solution  $U(P)$  of this equation. We also assume that when  $P = P_o$  this family reduces to our current model.

Let the objective function  $G(U, P)$  have values in dimension  $K \ll M$ . (Usually  $K = 1$  or some other small integer.)

We are interested in the sensitivity of the response  $R(P) = G(U(P), P)$ . Specifically, we set  $U_o = U(P_o)$  and want to compute

$$\partial_P R(P_o) = G_U(U_o, P_o)U_P(P_o) + G_P(U_o, P_o).$$

Here we know  $G_U(U_o, P_o)$  and  $G_P(U_o, P_o)$ , but not  $U_P(P_o)$ .

By differentiating the family of models we see that  $U_P(P_o)$  satisfies the linearized model

$$0 = F_U(U_o, P_o)U_P(P_o) + F_P(U_o, P_o).$$

This approach generally requires inverting the matrix  $F_U(U_o, P_o)$ .

Rather, we find the row-vector  $J_o$  that solves the adjoint problem

$$0 = J_o F_U(U_o, P_o) + G_U(U_o, P_o).$$

Then the sensitivity becomes

$$\begin{aligned}\partial_P R(P_o) &= -J_o F_U(U_o, P_o)U_P(P_o) + G_P(U_o, P_o) \\ &= J_o F_P(U_o, P_o) + G_P(U_o, P_o).\end{aligned}$$

Because we know  $F_P(U_o, P_o)$  and  $G_P(U_o, P_o)$ , this is easy to compute once  $J_o$  is computed.



The point here is that computing  $J_o$  only requires solving  $K$  linear systems. The cost of doing this will be roughly  $K$  times the cost of solving for  $U_o$ . Because  $K \ll M$  this will be much less than the cost of solving the linearized model for  $U_P(P_o)$ , which will generally be  $M$  times the cost of solving for  $U_o$ .

Let  $\Delta P_o$  be a diagonal matrix of uncertainties associated with  $P_o$ . If an entry of  $P_o$  is the mean of some data then the corresponding entry of  $\Delta P_o$  might be the standard deviation. If an entry of  $P_o$  is zero in order to turn off or decouple some physics then the corresponding entry of  $\Delta P_o$  might be your best guess at the expected value of that parameter.

Now let  $\Delta R_o$  be the  $K \times N$  matrix whose entries are the absolute values of the entries of  $\partial_P R(P_o) \Delta P_o$ . The entries of this matrix are your best guesses of the uncertainties in the objective functions.

The idea is now to use the matrix  $\Delta R_o$  to learn which parameters that are zero in the current model should be turned on so as to enlarge the model. One way to do this is to simply choose those parameters corresponding to the largest entry in each row of  $\Delta R_o$ . We can also consider some convex combinations of the entries in each row of  $\Delta R_o$  that respect LTE or some other balance. We then enlarge the current model and repeat the process.

The model inflation stops when none of the largest uncertainties are due to parameters that are zero. When this happens you can not make better predictions by enlarging the model.

**A model should be as simple as possible, but no simpler!**

The family of models  $F(U, P) = 0$  might be a system of either ordinary or partial differential equations, while  $G(U, P)$  might be some temporal or spatial-temporal averages of a solution of this system.

## 4. Dissipative Structure

Let  $G$  be a symmetric, positive definite  $N \times N$  real matrix:

$$G^T = G > 0.$$

Let  $J$  and  $K$  be  $N \times N$  real matrices such that  $J$  is skew-adjoint and  $K$  is self-adjoint, nonnegative definite with respect to  $G$ . This means that

$$GJ + J^T G = 0, \quad GK = K^T G \geq 0. \quad (6)$$

We consider the linear dynamical system

$$\frac{dU}{dt} + JU + KU = 0. \quad (7)$$

This is (1) with  $A = J + K$  and  $F(t) = 0$ . By (6) solutions of (7) satisfy

$$\frac{1}{2} \frac{d}{dt} (U^T G U) + U^T G K U = 0. \quad (8)$$

This shows that the scalar product associated with  $G$  is dissipated.

We will show how this dissipative structure is preserved by the reduced models presented earlier. We will work in a more general setting than we did earlier.

Let  $n < N$  and  $S$  be an  $N \times n$  real matrix of rank  $n$ . Then  $g = S^T G S$  is a symmetric, positive definite  $n \times n$  real matrix. This means that

$$g^T = g > 0.$$

Let  $R = g^{-1} S^T G$ . This means that  $S = G^{-1} R^T g$  and that  $RS = I$ . The matrix  $P = SR$  satisfies

$$P^2 = P, \quad GP = P^T G. \quad (9)$$

We conclude that  $P$  is the orthogonal projection onto the range of  $S$  with respect to the scalar product associated with  $G$ . Then  $\tilde{P} = I - P$  is the orthogonal projection onto the subspace of  $\mathbb{R}^N$  that is orthogonal to the range of  $S$  with respect to the scalar product associated with  $G$ .

Let  $u = RU$ . Then  $PU = Su$  and the orthogonal decomposition of  $U$  is  $U = Su + \tilde{U}$  where  $\tilde{U} = \tilde{P}U$ . If  $U$  satisfies system (7) then  $u$  satisfies

$$\frac{du}{dt} + RJSu + RKSu + RJ\tilde{U} + RK\tilde{U} = 0. \quad (10)$$

We would like to find closed systems for  $u$  that approximate the dynamics of (10).

**Remark.** If the range of  $S$  is an invariant subspace of  $J + K$  then it can be shown that

$$RJ\tilde{U} + RK\tilde{U} = 0.$$

In that case (10) is the closed system

$$\frac{du}{dt} + RJSu + RKSu = 0.$$

However, in general  $S$  will not be an invariant subspace of  $J + K$ , in which case (10) will not be a closed system for  $u$ .

The *Galerkin approximation* is obtained by setting  $\tilde{U} = 0$  in (10). This can be recast as

$$\frac{du}{dt} + ju + ku = 0, \quad (11)$$

where

$$j = RJS, \quad k = RKS. \quad (12)$$

Then  $j$  and  $k$  are  $n \times n$  real matrices such that  $j$  is skew-adjoint and  $k$  is self-adjoint, nonnegative definite with respect to the scalar product associated with  $g$ . This means that

$$gj + j^T g = 0, \quad gk = k^T g \geq 0. \quad (13)$$

These structural relations imply that solutions of (11) satisfy

$$\frac{1}{2} \frac{d}{dt} (u^T g u) + u^T g k u = 0. \quad (14)$$

This shows that the scalar product associated with  $g$  is dissipated.

Better approximations can be constructed by examining the dynamics of the so-called deviation  $\tilde{U}$ . The so-called deviation equation is

$$\frac{d\tilde{U}}{dt} + \tilde{J}\tilde{U} + \tilde{K}\tilde{U} = -\tilde{P}JSu - \tilde{P}KSu, \quad (15)$$

where

$$\tilde{J} = \tilde{P}J\tilde{P}, \quad \tilde{K} = \tilde{P}K\tilde{P}.$$

It follows from (6) and (9) that

$$G\tilde{J} + \tilde{J}^\top G = 0, \quad G\tilde{K} = \tilde{K}^\top G \geq 0. \quad (16)$$

We will build different approximations of  $\tilde{U}$  by balancing different terms in the deviation equation (15).

**Remark.** We can recover the special case presented earlier by setting

$$G = I, \quad S = \begin{pmatrix} I_n \\ 0 \end{pmatrix}, \quad R = \begin{pmatrix} I_n & 0 \end{pmatrix}, \quad P = \begin{pmatrix} I_n & 0 \\ 0 & 0 \end{pmatrix}, \quad g = I_n.$$

where  $I$  is the  $N \times N$  identity and  $I_n$  is the  $n \times n$  identity. Then

$$u = U_1, \quad \tilde{U} = \begin{pmatrix} 0 \\ U_2 \end{pmatrix},$$

and

$$\begin{aligned} RJSu + RKSu &= A_{11}U_1, & RJ\tilde{U} + RK\tilde{U} &= A_{12}U_2, \\ \tilde{P}JSu - \tilde{P}KSu &= \begin{pmatrix} 0 \\ A_{21}U_1 \end{pmatrix}, & \tilde{J}\tilde{U} + \tilde{K}\tilde{U} &= \begin{pmatrix} 0 \\ A_{22}U_2 \end{pmatrix}. \end{aligned}$$

Therefore equations (10) and (15) respectively play the roles of the first and second equations in system (2). Because of the decomposition  $A = J + K$  of  $A$  into skew-adjoint and self-adjoint matrices, there is now an additional balance to consider in the deviation equation (15).



First consider the new *semi-relaxation balance*

$$\widetilde{K}\widetilde{U} = -\widetilde{P}JSu - \widetilde{P}KSu.$$

Assuming that  $\widetilde{K}$  is invertible on the range of  $\widetilde{P}$ , we obtain

$$\widetilde{U} = -\widetilde{K}^{-1}JSu - \widetilde{K}^{-1}KSu.$$

When this *semi-relaxation approximation* is placed into (10) we see that  $u$  satisfies (11) with  $j$  and  $k$  given by

$$\begin{aligned} j &= RJS - RJ\widetilde{K}^{-1}KS - RK\widetilde{K}^{-1}JS, \\ k &= RKS - RK\widetilde{K}^{-1}KS - RJ\widetilde{K}^{-1}JS. \end{aligned} \tag{17}$$

These  $j$  and  $k$  satisfy the structural relations (13).

Next consider the full *relaxation balance*

$$(\widetilde{K} + \widetilde{J})\widetilde{U} = -\widetilde{P}JSu - \widetilde{P}KSu.$$

If  $\widetilde{K}$  is invertible on the range of  $\widetilde{P}$  then so is  $\widetilde{K} + \widetilde{J}$ .

In that case we obtain the *relaxation approximation*

$$\tilde{U} = -(\tilde{K} + \tilde{J})^{-1}JSu - (\tilde{K} + \tilde{J})^{-1}KSu. \quad (18)$$

When this approximation is placed into (10) we see that  $u$  satisfies (11) with  $j$  and  $k$  given by

$$\begin{aligned} j &= RJS - RJ(\tilde{K} - \tilde{J})^{-1}K(\tilde{K} + \tilde{J})^{-1}KS \\ &\quad - RK(\tilde{K} - \tilde{J})^{-1}K(\tilde{K} + \tilde{J})^{-1}JS \\ &\quad + RK(\tilde{K} - \tilde{J})^{-1}J(\tilde{K} + \tilde{J})^{-1}KS \\ &\quad + RJ(\tilde{K} - \tilde{J})^{-1}J(\tilde{K} + \tilde{J})^{-1}JS, \\ k &= RKS - RK(\tilde{K} - \tilde{J})^{-1}K(\tilde{K} + \tilde{J})^{-1}KS \\ &\quad - RJ(\tilde{K} - \tilde{J})^{-1}K(\tilde{K} + \tilde{J})^{-1}JS \\ &\quad + RJ(\tilde{K} - \tilde{J})^{-1}J(\tilde{K} + \tilde{J})^{-1}KS \\ &\quad + RK(\tilde{K} - \tilde{J})^{-1}J(\tilde{K} + \tilde{J})^{-1}JS. \end{aligned} \quad (19)$$

These  $j$  and  $k$  satisfy the structural relations (13).

The semi-relaxation approximation replaces  $(\widetilde{K} + \widetilde{J})^{-1}$  with  $\widetilde{K}^{-1}$  in the relaxation approximation (18). In other words, it assumes that  $\widetilde{K}^{-1}\widetilde{J}$  is small when multiplying those vectors to which it is applied.

This smallness assumption can be used for each nonnegative integer  $\ell$  to construct an approximation that formally lie between the semi-relaxation and relaxation approximations by using the Neuman approximation

$$(\widetilde{K} + \widetilde{J})^{-1} \approx \sum_{i=0}^{\ell} \left(-\widetilde{K}^{-1}\widetilde{J}\right)^i \widetilde{K}^{-1}.$$

However it can be shown that the  $j$  and  $k$  resulting from this approximation will generally satisfy the structural relations (13) if and only if  $\ell = 4m$  or  $\ell = 4m + 1$  for some nonnegative integer  $m$ , i.e. if and only if

$$\ell \in \{0, 1, 4, 5, 8, 9, 12, 13, \dots\}.$$

The choice  $\ell = 0$  yields the semi-relaxation approximation (17).

The choice  $\ell = 1$  yields the so-called *first-correction approximation*

$$\begin{aligned}
 j &= RJS - RJ\tilde{K}^{-1}KS - RK\tilde{K}^{-1}JS \\
 &\quad + RK\tilde{K}^{-1}J\tilde{K}^{-1}KS + RJ\tilde{K}^{-1}J\tilde{K}^{-1}JS, \\
 k &= RKS - RK\tilde{K}^{-1}KS - RJ\tilde{K}^{-1}JS \\
 &\quad + RJ\tilde{K}^{-1}J\tilde{K}^{-1}KS + RK\tilde{K}^{-1}J\tilde{K}^{-1}JS.
 \end{aligned} \tag{20}$$

The terms in the first line of each equation are those of the semi-relaxation approximation (17). Those in the second line are the correction terms. This approximation can be obtained directly from the relaxation approximation (19) by replacing  $(\tilde{K} - \tilde{J})^{-1}$  and  $(\tilde{K} + \tilde{J})^{-1}$  with  $\tilde{K}^{-1}$ .

**Remark.** In practice there is little reason to go beyond the first-correction approximation due to the resulting increased complexity. It is usually easier to simply increase the dimension  $n$  of the reduced model. So we stop here.

## 5. Moment Closures for the Linearized Boltzmann Equation

The Boltzmann equation linearized about a global Maxwellian  $M(v)$  is

$$\partial_t g + v \cdot \nabla_x g + \mathcal{L}g = 0. \quad (21)$$

Here  $g(t, x, v)$  is the relative kinetic density and  $\mathcal{L}$  is the linearized collision operator, which is given by

$$\mathcal{L}g = \iint_{\mathbb{S}^{D-1} \times \mathbb{R}^D} (g + g_* - g' - g'_*) b(|v - v_*|, n \cdot \omega) d\omega M(v_*) dv_*, \quad (22)$$

where  $b(|v - v_*|, n \cdot \omega)$  is the collision kernel,  $n = \frac{v - v_*}{|v - v_*|}$ , and

$$g_* = g(t, x, v_*), \quad g' = g(t, x, v'), \quad g'_* = g(t, x, v'_*),$$

with  $v' = v - \omega \omega \cdot (v - v_*)$  and  $v'_* = v_* + \omega \omega \cdot (v - v_*)$ . To avoid boundary conditions we consider  $x \in \mathbb{T}^D$ , the  $D$ -dimensional torus. Without loss of generality we can take  $M(v) = (2\pi)^{-\frac{D}{2}} \exp(-\frac{1}{2}|v|^2)$ .

- The operator  $\mathcal{L}$  acts only on the  $v$  variable.

- It is self-adjoint over  $L^2(Mdv)$  in the sense that

$$\int_{\mathbb{R}^D} g \mathcal{L} h M(v) dv = \int_{\mathbb{R}^D} h \mathcal{L} g M(v) dv .$$

- It is nonnegative definite over  $L^2(Mdv)$  in the sense that

$$\int_{\mathbb{R}^D} g \mathcal{L} g M(v) dv \geq 0 .$$

- Its null space in  $L^2(Mdv)$  is

$$\text{Null}(\mathcal{L}) = \left\{ \alpha + v \cdot \beta + |v|^2 \gamma : (\alpha, \beta, \gamma) \in \mathbb{R} \times \mathbb{R}^D \times \mathbb{R} \right\} .$$

Now introduce the notation

$$\langle h \rangle = \int_{\mathbb{R}^D} h M(v) dv .$$

The foregoing properties imply that solutions of the linearized Boltzmann equation (21) satisfy the local conservation laws

$$\begin{aligned} \partial_t \langle g \rangle + \nabla_x \cdot \langle v g \rangle &= 0 , \\ \partial_t \langle v g \rangle + \nabla_x \cdot \langle v \otimes^2 g \rangle &= 0 , \\ \partial_t \langle \frac{1}{2} |v|^2 g \rangle + \nabla_x \cdot \langle v \frac{1}{2} |v|^2 g \rangle &= 0 . \end{aligned}$$

These are the local conservation laws of mass, momentum, and energy.

Solutions of the linearized Boltzmann equation (21) also satisfy the local dissipation law

$$\partial_t \langle g^2 \rangle + \nabla_x \cdot \langle v g^2 \rangle + \langle g \mathcal{L} g \rangle = 0 .$$

This reflects the Boltzmann entropy local dissipation law for the linearized Boltzmann equation (21).

The natural scalar product for the linearized Boltzmann equation (21) is

$$(g | h) = \int_{\mathbb{T}^D} \langle g h \rangle dx = \iint_{\mathbb{R}^D \times \mathbb{T}^D} g h M(v) dv dx .$$

With respect to this scalar product

1.  $v \cdot \nabla_x$  is skew-adjoint,
2.  $\mathcal{L}$  is self-adjoint and nonnegative definite.

Hence, the linearized Boltzmann equation (21) has a dissipative structure analogous to that of the linear dynamical system (7) presented in the last section with the identifications

$$U \longleftrightarrow g, \quad J \longleftrightarrow v \cdot \nabla_x, \quad K \longleftrightarrow \mathcal{L} .$$

Approximations of the linearized Boltzmann equation can be constructed analogous to the Galerkin, semi-relaxation, and first-correction reduced models of the last section.



Let  $\mathbf{m}$  be a vector of linearly independent tensor powers of  $v$  whose span includes  $\{1, v, |v|^2\}$  and that respects rotational symmetry. For example, we can choose

$$\mathbf{m} = \begin{pmatrix} 1 \\ v \\ |v|^2 \end{pmatrix}, \quad \mathbf{m} = \begin{pmatrix} 1 \\ v \\ v^{\otimes 2} \end{pmatrix}, \quad \mathbf{m} = \begin{pmatrix} 1 \\ v \\ v^{\otimes 2} \\ |v|^2 v \end{pmatrix}, \quad \mathbf{m} = \begin{pmatrix} 1 \\ v \\ v^{\otimes 2} \\ v^{\otimes 3} \\ |v|^2 v^{\otimes 2} \end{pmatrix}.$$

The associated moments of  $g$  are

$$\rho(t, x) = \langle \mathbf{m} g(t, x, v) \rangle = \int_{\mathbb{R}^D} \mathbf{m} g(t, x, v) M(v) dv.$$

In  $\mathbb{R}^3$  the dimensions of  $\rho$  for the above choices of  $\mathbf{m}$  are 5, 10, 13, and 26 respectively. The first choice corresponds to just the fluid moments.

Taking moments of (21) gives

$$\partial_t \langle \mathbf{m} g \rangle + \nabla_x \cdot \langle v \mathbf{m} g \rangle + \langle \mathbf{m} \mathcal{L} g \rangle = 0. \quad (23)$$

We decompose  $g$  as

$$g = \mathbf{m}^\top \alpha + \tilde{g}, \quad \text{where} \quad \langle \mathbf{m} \tilde{g} \rangle = 0. \quad (24)$$

This is an orthogonal decomposition of  $g$  into its moment component  $\mathbf{m}^\top \alpha$  and a deviation  $\tilde{g}$ .

The vector  $\alpha$  is related to the moments  $\rho$  by

$$\rho(t, x) = \langle \mathbf{m} \mathbf{m}^\top \rangle \alpha(t, x).$$

The matrix  $\langle \mathbf{m} \mathbf{m}^\top \rangle$  is positive definite, so we have

$$\alpha(t, x) = \langle \mathbf{m} \mathbf{m}^\top \rangle^{-1} \rho(t, x).$$

Setting decomposition (24) into (23) gives

$$\begin{aligned} \langle \mathbf{m} \mathbf{m}^T \rangle \partial_t \alpha + \langle \mathbf{m} \mathbf{m}^T v \rangle \cdot \nabla_x \alpha + \langle \mathbf{m} \mathcal{L} \mathbf{m} \rangle \alpha \\ + \langle \mathbf{m} \mathcal{L} \tilde{g} \rangle + \nabla_x \cdot \langle v \mathbf{m} \tilde{g} \rangle = 0. \end{aligned} \quad (25)$$

This is not a closed system for  $\alpha$  because of the terms that involve  $\tilde{g}$ .

The goal of a *moment closure* is to express the terms in (25) that involve  $\tilde{g}$  in terms of  $\alpha$  so that the resulting system approximately governs  $\alpha$ , and thereby approximately governs  $\rho$ . For any choice of  $\mathbf{m}$  we will make the identifications

$$\begin{aligned} u &\longleftrightarrow \rho, & Su &\longleftrightarrow \mathbf{m}^T \alpha, & \tilde{U} &\longleftrightarrow \tilde{g}, \\ RU &\longleftrightarrow \langle \mathbf{m} g \rangle, & PV &\longleftrightarrow \mathbf{m}^T \langle \mathbf{m} \mathbf{m}^T \rangle^{-1} \langle \mathbf{m} h \rangle, \end{aligned}$$

and derive moment closures from the Galerkin, semi-relaxation, and first-correction approximations presented in the last section.

By setting  $\tilde{g} = 0$  in (25) we obtain the *Galerkin closure*

$$\langle \mathbf{m} \mathbf{m}^T \rangle \partial_t \alpha + \langle \mathbf{m} \mathbf{m}^T v \rangle \cdot \nabla_x \alpha + \langle \mathbf{m} \mathcal{L} \mathbf{m}^T \rangle \alpha = 0. \quad (26)$$

Its solutions satisfy the local dissipation law

$$\partial_t \left[ \frac{1}{2} \alpha^T \langle \mathbf{m} \mathbf{m}^T \rangle \alpha \right] + \nabla_x \cdot \left[ \frac{1}{2} \alpha^T \langle v \mathbf{m} \mathbf{m}^T \rangle \alpha \right] + \alpha^T \langle \mathbf{m} \mathcal{L} \mathbf{m}^T \rangle \alpha = 0.$$

The Navier-Stokes approximation of this closure will usually have incorrect values for its transport coefficients (its viscosity and thermal conductivity).

When  $\mathbf{m}$  corresponds to the fluid moments then  $\mathcal{L} \mathbf{m} = 0$  and the Galerkin closure (26) reduces to

$$\langle \mathbf{m} \mathbf{m}^T \rangle \partial_t \alpha + \langle \mathbf{m} \mathbf{m}^T v \rangle \cdot \nabla_x \alpha = 0.$$

This is exactly the linearized Euler system of gas dynamics. In particular, its transport coefficients are zero, which is incorrect. However, every Galerkin closure (26) recovers this correct Euler system in fluid regimes, while most recover an incorrect Navier-Stokes system.

The deviation  $\tilde{g}$  satisfies the so-called deviation equation

$$\partial_t \tilde{g} + \tilde{\mathcal{P}} v \cdot \nabla_x \tilde{g} + \tilde{\mathcal{L}} \tilde{g} = -\tilde{\mathcal{P}} \mathcal{L} \mathbf{m}^\top \alpha - \tilde{\mathcal{P}} \mathbf{m}^\top v \cdot \nabla_x \alpha. \quad (27)$$

Here  $\tilde{\mathcal{P}} = \mathcal{I} - \mathcal{P}$  where  $\mathcal{P}$  is the orthogonal projection

$$\mathcal{P} = \mathbf{m}^\top \langle \mathbf{m} \mathbf{m}^\top \rangle^{-1} \mathbf{m},$$

and  $\tilde{\mathcal{L}} = \tilde{\mathcal{P}} \mathcal{L} \tilde{\mathcal{P}}$ . The *semi-relaxation balance* gives

$$\tilde{g} = -\tilde{\mathcal{L}}^{-1} \mathcal{L} \mathbf{m}^\top \alpha - \tilde{\mathcal{L}}^{-1} \mathbf{m}^\top v \cdot \nabla_x \alpha.$$

where  $\tilde{\mathcal{L}}^{-1}$  is the pseudo-inverse of  $\tilde{\mathcal{L}}$ . Setting this into (25) gives the *semi-relaxation closure*

$$\begin{aligned} & \langle \mathbf{m} \mathbf{m}^\top \rangle \partial_t \alpha + \langle \mathbf{m} \mathbf{m}^\top v \rangle \cdot \nabla_x \alpha + \langle \mathbf{m} \mathcal{L} \mathbf{m}^\top \rangle \alpha \\ & - \langle \mathbf{m} \mathcal{L} \tilde{\mathcal{L}}^{-1} \mathcal{L} \mathbf{m}^\top \rangle \alpha - \nabla_x \cdot \left[ \langle v \mathbf{m} \tilde{\mathcal{L}}^{-1} \mathcal{L} \mathbf{m}^\top \rangle \alpha \right] \\ & - \langle \mathbf{m} \mathcal{L} \tilde{\mathcal{L}}^{-1} \mathbf{m}^\top v \rangle \cdot \nabla_x \alpha - \nabla_x \cdot \left[ \langle v \mathbf{m} \tilde{\mathcal{L}}^{-1} \mathbf{m}^\top v \rangle \cdot \nabla_x \alpha \right] = 0. \end{aligned} \quad (28)$$

It can be shown that

$$\langle \mathbf{m} \mathcal{L} \mathbf{m}^\top \rangle - \langle \mathbf{m} \mathcal{L} \tilde{\mathcal{L}}^{-1} \mathcal{L} \mathbf{m}^\top \rangle = \langle \mathbf{m} \mathbf{m}^\top \rangle \langle \mathbf{m} \mathcal{L}^{-1} \mathbf{m}^\top \rangle^{-1} \langle \mathbf{m} \mathbf{m}^\top \rangle,$$

where

- $\mathcal{L}^{-1}$  is the pseudo-inverse of the operator  $\mathcal{L}$ , and
- $\langle \mathbf{m} \mathcal{L}^{-1} \mathbf{m}^\top \rangle^{-1}$  is the pseudo-inverse of the matrix  $\langle \mathbf{m} \mathcal{L}^{-1} \mathbf{m}^\top \rangle$ .

This fact insures that the correct transport coefficients will arise in the Navier-Stokes approximation of this closure. Because

$$\langle \mathbf{m} \mathbf{m}^\top \rangle \langle \mathbf{m} \mathcal{L}^{-1} \mathbf{m}^\top \rangle^{-1} \langle \mathbf{m} \mathbf{m}^\top \rangle \geq 0,$$

it also shows that solutions of the semi-relaxation closure satisfy the local dissipation law

$$\begin{aligned} & \partial_t \left[ \frac{1}{2} \boldsymbol{\alpha}^\top \langle \mathbf{m} \mathbf{m}^\top \rangle \boldsymbol{\alpha} \right] + \nabla_x \cdot \left[ \frac{1}{2} \boldsymbol{\alpha}^\top \langle v \mathbf{m} \mathbf{m}^\top \rangle \boldsymbol{\alpha} \right] + \boldsymbol{\alpha}^\top \langle \mathbf{m} \mathcal{L} \mathbf{m}^\top \rangle \boldsymbol{\alpha} \\ & \quad - \boldsymbol{\alpha}^\top \langle \mathbf{m} \mathcal{L} \tilde{\mathcal{L}}^{-1} \mathcal{L} \mathbf{m}^\top \rangle \boldsymbol{\alpha} - \nabla_x \cdot \left[ \boldsymbol{\alpha}^\top \langle v \mathbf{m} \tilde{\mathcal{L}}^{-1} \mathcal{L} \mathbf{m}^\top \rangle \boldsymbol{\alpha} \right] \\ & - \nabla_x \cdot \left[ \boldsymbol{\alpha}^\top \langle v \mathbf{m} \tilde{\mathcal{L}}^{-1} \mathbf{m}^\top v \rangle \cdot \nabla_x \boldsymbol{\alpha} \right] + \nabla_x \boldsymbol{\alpha}^\top \cdot \langle v \mathbf{m} \tilde{\mathcal{L}}^{-1} \mathbf{m}^\top v \rangle \cdot \nabla_x \boldsymbol{\alpha} = 0. \end{aligned}$$

When  $\mathbf{m}$  corresponds to the fluid moments then  $\mathcal{L}\mathbf{m} = 0$  and  $\tilde{\mathcal{L}} = \mathcal{L}$ , so that the semi-relaxation closure (28) reduces to

$$\langle \mathbf{m} \mathbf{m}^\top \rangle \partial_t \alpha + \langle \mathbf{m} \mathbf{m}^\top v \rangle \cdot \nabla_x \alpha = \nabla_x \cdot \left[ \langle v \mathbf{m} \mathcal{L}^{-1} \mathbf{m}^\top v \rangle \cdot \nabla_x \alpha \right].$$

This is exactly the linearized Navier-Stokes system of gas dynamics. In particular, its transport coefficients are correct. Moreover, every semi-relaxation closure (28) recovers this Navier-Stokes system in fluid regimes.

The *relaxation balance* gives

$$\tilde{g} + \tilde{\mathcal{L}}^{-1} v \cdot \nabla_x \tilde{g} = -\tilde{\mathcal{L}}^{-1} \mathcal{L} \mathbf{m}^\top \alpha - \tilde{\mathcal{L}}^{-1} \mathbf{m}^\top v \cdot \nabla_x \alpha.$$

If we assume that  $\tilde{\mathcal{L}}^{-1} v \cdot \nabla_x \tilde{g}$  is smaller than  $\tilde{g}$  then we can approximate this  $\tilde{g}$  by the *first-correction* to the semi-relaxation balance, which gives

$$\begin{aligned} \tilde{g} = & -\tilde{\mathcal{L}}^{-1} \mathcal{L} \mathbf{m}^\top \alpha + \tilde{\mathcal{L}}^{-1} v \cdot \nabla_x \tilde{\mathcal{L}}^{-1} \mathcal{L} \mathbf{m}^\top \alpha \\ & - \tilde{\mathcal{L}}^{-1} \mathbf{m}^\top v \cdot \nabla_x \alpha + \tilde{\mathcal{L}}^{-1} v \cdot \nabla_x \tilde{\mathcal{L}}^{-1} \mathbf{m}^\top v \cdot \nabla_x \alpha. \end{aligned}$$

Setting this into (25) gives the *first-correction closure*

$$\begin{aligned}
& \langle \mathbf{m} \mathbf{m}^\top \rangle \partial_t \alpha + \langle \mathbf{m} \mathbf{m}^\top v \rangle \cdot \nabla_x \alpha + \langle \mathbf{m} \mathcal{L} \mathbf{m}^\top \rangle \alpha \\
& - \langle \mathbf{m} \mathcal{L} \tilde{\mathcal{L}}^{-1} \mathcal{L} \mathbf{m}^\top \rangle \alpha - \nabla_x \cdot \left[ \langle v \mathbf{m} \tilde{\mathcal{L}}^{-1} \mathcal{L} \mathbf{m}^\top \rangle \alpha \right] \\
& - \langle \mathbf{m} \mathcal{L} \tilde{\mathcal{L}}^{-1} \mathbf{m}^\top v \rangle \cdot \nabla_x \alpha - \nabla_x \cdot \left[ \langle v \mathbf{m} \tilde{\mathcal{L}}^{-1} \mathbf{m}^\top v \rangle \cdot \nabla_x \alpha \right] \\
& \quad + \langle \mathbf{m} \mathcal{L} \tilde{\mathcal{L}}^{-1} v \cdot \nabla_x \tilde{\mathcal{L}}^{-1} \mathcal{L} \mathbf{m}^\top \rangle \alpha \\
& \quad + \nabla_x \cdot \left[ \langle v \mathbf{m} \tilde{\mathcal{L}}^{-1} v \cdot \nabla_x \tilde{\mathcal{L}}^{-1} \mathcal{L} \mathbf{m}^\top \rangle \alpha \right] \\
& \quad + \langle \mathbf{m} \mathcal{L} \tilde{\mathcal{L}}^{-1} v \cdot \nabla_x \tilde{\mathcal{L}}^{-1} \mathbf{m}^\top v \rangle \cdot \nabla_x \alpha \\
& \quad + \nabla_x \cdot \left[ \langle v \mathbf{m} \tilde{\mathcal{L}}^{-1} v \cdot \nabla_x \tilde{\mathcal{L}}^{-1} \mathbf{m}^\top v \rangle \cdot \nabla_x \alpha \right] = 0.
\end{aligned} \tag{29}$$

The last four terms above do not appear in the semi-relaxation closure (28). We will not give the local dissipation law associated with (29) here. Remarkably, the four new terms in (29) add new flux terms to the local dissipation law, but do not change its dissipation terms!



When  $\mathbf{m}$  corresponds to the fluid moments then  $\mathcal{L}\mathbf{m} = 0$  and  $\tilde{\mathcal{L}} = \mathcal{L}$ , so that the first-correction closure (29) reduces to

$$\begin{aligned} & \langle \mathbf{m} \mathbf{m}^\top \rangle \partial_t \alpha + \langle \mathbf{m} \mathbf{m}^\top v \rangle \cdot \nabla_x \alpha \\ &= \nabla_x \cdot \left[ \langle v \mathbf{m} \mathcal{L}^{-1} \mathbf{m}^\top v \rangle \cdot \nabla_x \alpha \right] \\ & \quad - \nabla_x \cdot \left[ \langle v \mathbf{m} \mathcal{L}^{-1} v \cdot \nabla_x \mathcal{L}^{-1} \mathbf{m}^\top v \rangle \cdot \nabla_x \alpha \right]. \end{aligned}$$

This is the linearization of a well-posed Burnett system of gas dynamics. In particular, its transport coefficients are correct. Moreover, every first-correction closure (29) recovers this Burnett system in fluid regimes. This Burnett system is not recovered in fluid regimes by essentially all of the semi-relaxation closures (28).

**Remark.** Corrections to the semi-relaxation closure beyond the first quickly become too complicated to be useful. The full relaxation closure is spatially nonlocal, so is also too complicated to be useful. Therefore we stop here.

## 6. Conclusion

We have applied principles of model reduction that preserve dissipative structure to construct a family of well-posed moment closures for the linearized Boltzmann equation (21), the members of which are specified by

- the choice of  $m$ ,
- the choice of the Galerkin, semi-relaxation, or first-correction closure.

Model inflation can be used to select a closure from this family to model a problem that has a few objectives.

**Remark.** What we have called the Galerkin closure can be viewed as a Petrov-Galerkin closure because its test functions are polynomials in  $v$  while its trial functions for the particle density  $f$  are the Maxwellian  $M(v)$  times polynomials in  $v$ . However, this simple relationship between the test and trial functions makes the moniker Galerkin appropriate.