# A stable multiplicative dynamical low-rank discretization for the linear Boltzmann-BGK equation

Lena Baumann[a], Lukas Einkemmer[b], Christian Klingenberg[a], Jonas Kusch[c]

[a] *University of Wuerzburg, Department of Mathematics, Wuerzburg, Germany, lena.baumann@uni-wuerzburg.de (Lena Baumann), christian.klingenberg@uni-wuerzburg.de (Christian Klingenberg)*

[b] *University of Innsbruck, Numerical Analysis and Scientific Computing, Innsbruck, Austria, lukas.einkemmer@uibk.ac.at*

[c] *Norwegian University of Life Sciences, Scientific Computing, Ås, Norway, jonas.kusch@nmbu.no*

**Abstract**

The numerical method of dynamical low-rank approximation (DLRA) has recently been applied to various kinetic equations showing a significant reduction of the computational effort. In this paper, we apply this concept to the linear Boltzmann-Bhatnagar-Gross-Krook (Boltzmann-BGK) equation which due its high dimensionality is challenging to solve. Inspired by the special structure of the non-linear Boltzmann-BGK problem, we consider a multiplicative splitting of the distribution function. We propose a rank-adaptive DLRA scheme making use of the basis update & Galerkin integrator and combine it with an additional basis augmentation to ensure numerical stability, for which an analytical proof is given and a classical hyperbolic Courant–Friedrichs–Lewy (CFL) condition is derived. This allows for a further acceleration of computational times and a better understanding of the underlying problem in finding a suitable discretization of the system. Numerical results of a series of different test examples confirm the accuracy and efficiency of the proposed method compared to the numerical solution of the full system.

*Keywords:* linear Boltzmann equation, BGK relaxation model, dynamical low-rank approximation, multiplicative splitting, numerical stability, rank adaptivity

## 1. Introduction

Numerically solving kinetic equations usually requires immense computational and memory efforts due to the high-dimensional phase space containing all possible states of the system. The state of a kinetic system is described by a distribution function $f$ which can be interpreted as the corresponding particle density in phase space. Instead of solving one high-dimensional equation, the concept of dynamical low-rank approximation (DLRA) [23] allows us to split the problem into three lower dimensional subequations leading to an appropriate approximation of the solution. In particular, in a one-dimensional setting we approximate the distribution function $f(t, x, v)$, with $t \in \mathbb{R}^+$ denoting the time, $x \in D \subset \mathbb{R}$ the spatial, and $v \in \mathbb{R}$ the velocity variable, by

$$f(t, x, v) \approx \sum_{i,j=1}^{r} X_i(t, x) S_{ij}(t) V_j(t, v),$$

and evolve the corresponding low-rank factors in three substeps further in time. The sets $\{X_i : i = 1, .., r\}$ and $\{V_j : j = 1, .., r\}$ contain the orthonormal basis functions in space and in velocity, respectively, and $r$ is called the rank of this approximation. DLRA has recently gained increasing interest and has been studied in various fields including radiation transport [2, 30, 13, 31, 34], radiation therapy [27], plasma physics [16, 18, 15, 19], chemical kinetics [32, 17] and Boltzmann type transport problems [12, 14, 11, 22]. The core idea of this method is to project the solution to a manifold of low-rank functions of the above form and

constrain the solution dynamics there. Different time integrators which are able to ensure this behaviour and are robust to the presence of small singular values are available. Frequently used integrators for kinetic problems are the *projector-splitting* [29] as well as the *(rank-adaptive) basis update & Galerkin* (BUG) [10, 8] and the *parallel* integrator [9]. For the rank-adaptive BUG and the parallel integrator extensions to schemes with proven second-order robust error bounds have been derived in [7, 25].

For a large number of collisions, the solution $f$ of the Boltzmann-BGK equation stays close to the Maxwellian equilibrium distribution $M$ which in general is not a low-rank function. Inspired by [14, 24], we use the multiplicative splitting $f = Mg$, for which in [14] it has been shown that $g$ is a low-rank function even if this if not the case for $f$. Hence, we derive an evolution equation for $g$ and apply the low-rank approach to this part of the distribution function. Difficulties may arise in the discretization as it is per se not clear how to treat spatial derivatives.

In this paper we propose a stable dynamical low-rank discretization for the linear Boltzmann-BGK equation. The main features of this paper are:

- *A multiplicative splitting of the distribution function:* As the Maxwellian equilibrium distribution $M$ is generally not a low-rank function, we consider the multiplicative splitting $f = Mg$ and apply the low-rank ansatz to the remaining function $g$. It can be considered as a deviation from the equilibrium and is shown to be of low rank [14].

- *A stable numerical scheme for linear Boltzmann-BGK with rigorous mathematical proofs:* We show that a stable discretization has to be derived carefully and compare it with an intuitive discretization that fails to guarantee numerical stability. We give a rigorous analytical proof of stability and derive a classic hyperbolic CFL condition. This enables us to choose an optimal time step size of $\Delta t = \text{CFL} \cdot \Delta x$ with CFL denoting the CFL number, leading to a reduction of the computational effort.

- *A rank-adaptive integrator:* For the low-rank scheme we use the rank-adaptive BUG integrator from [8], leading to a basis augmentation in both the $K$- and $L$-step of the low-rank algorithm. Compared to the projector-splitting integrator used in [12, 14, 11], this allows us to determine the rank adaptively in each step avoiding the a priori choice of a certain fixed rank.

- *A series of numerical experiments validating the derived properties:* We give a number of numerical examples that validate the derived stability while showing a significant reduction of computational and memory requirements of the low-rank scheme compared to the full order method.

The paper is structured as follows: After the introduction in Section 1, we provide background information on the linear Boltzmann-BGK equation, explain the considered multiplicative structure, and derive two possible systems of equations in Section 2. Both systems are equivalent in the continuous setting. In Section 3, we discretize in velocity and in space, before subsequently time is discretized giving two different fully discretized schemes. It is then shown in Section 4 that a naive discretization can lead to a numerical scheme that is not von Neumann stable whereas a more careful treatment guarantees numerical stability. Section 5 gives a brief introduction to the concept of DLRA and applies this method such that a numerically stable low-rank scheme is obtained. Numerical experiments in both 1D and 2D in Section 6 confirm the derived results. Section 7 gives a brief conclusion and outlook.

## 2. Linear Boltzmann-BGK

The Boltzmann equation is a fundamental model in kinetic theory describing a gas that is not in thermodynamic equilibrium [6, 33]. In its full formulation it makes use of the so called "Stosszahlansatz" leading to a collision operator for which the solution of the Boltzmann equation is demanding. To overcome this, the *BGK model* [4], named after Bhatnagar, Gross and Krook, can be considered. It simplifies the collision term while maintaining the key properties of the equation. In a one-dimensional setting it reads

$$\partial_t f(t, x, v) + v \partial_x f(t, x, v) = \sigma \left( M[f](t, x, v) - f(t, x, v) \right), \tag{1a}$$

where $f(t, x, v)$ denotes the distribution function depending on the time $t \in \mathbb{R}^+$, the spatial variable $x \in D \subset \mathbb{R}$ and the velocity variable $v \in \mathbb{R}$. The constant $\sigma$ describes the collisionality of the particles and $M[f]$ stands for the Maxwellian equilibrium distribution. It depends on the density $\rho(t, x) = \int_{\mathbb{R}} f(t, x, v) \mathrm{d}v$ for which we obtain an evolution equation by integrating (1a) with respect to $v$. This gives

$$\partial_t \rho(t, x) = -\partial_x \int v f(t, x, v) \mathrm{d}v. \tag{1b}$$

In [14] it has been shown that using the multiplicative decomposition

$$f(t, x, v) = M[f](t, x, v) g(t, x, v) \tag{2}$$

is advantageous as $g$ is low-rank even if this is not the case for the Maxwellian (which is not true for the classic additive micro-macro decomposition). In order to reduce computational and memory costs a dynamical low-rank approach has then been applied in [14] to treat the resulting evolution equations for $g$.

In this work, we consider an isothermal Maxwellian without drift, i.e.

$$M[f](t, x, v) = \frac{\rho(t, x)}{\sqrt{2\pi}} \exp\left(-v^2/2\right).$$

This results in a linear model, which we call the *linear Boltzmann-BGK equation*. It has been extensively studied in the PDE community (see, e.g., [20, 5, 1]) as well as from a numerical point of view [3]. In this paper, we provide a stability analysis of this model in the context of dynamical low-rank simulation. The stability analysis for the simplified problem provides insight into the numerical schemes that have been used in the literature [14] for the Boltzmann–BGK equation. In particular, it explains why such schemes need to take relatively small time step sizes even though the collision operator is treated implicitly.

We insert the multiplicative approach (2) into (1a) and (1b) and obtain

$$\partial_t g(t, x, v) = -v \partial_x g(t, x, v) + \sigma \left(1 - g(t, x, v)\right) - \frac{g(t, x, v)}{\rho(t, x)} \partial_t \rho(t, x) - v \frac{g(t, x, v)}{\rho(t, x)} \partial_x \rho(t, x), \tag{3a}$$

$$\partial_t \rho(t, x) = -\frac{1}{\sqrt{2\pi}} \partial_x \int \rho(t, x) g(t, x, v) v e^{-v^2/2} \mathrm{d}v. \tag{3b}$$

This set of equations is called the *advection form* of the multiplicative system. It corresponds to the way the equations are treated in [14]. We can rewrite equation (3a) into a *conservative form*, leading to the system

$$\partial_t g(t, x, v) = -\frac{v}{\rho(t, x)} \partial_x \left(\rho(t, x) g(t, x, v)\right) + \sigma \left(1 - g(t, x, v)\right) - \frac{g(t, x, v)}{\rho(t, x)} \partial_t \rho(t, x), \tag{4a}$$

$$\partial_t \rho(t, x) = -\frac{1}{\sqrt{2\pi}} \partial_x \int \rho(t, x) g(t, x, v) v e^{-v^2/2} \mathrm{d}v. \tag{4b}$$

Note that for both systems we omit initial and boundary conditions for now. It is a challenging task to construct a suitable numerical scheme as it is per se not clear how to treat the spatial derivative in the transport part of the first equations and which of both systems to prefer. Further, the potentially stiff collision term requires an implicit time discretization. In addition, the consideration of higher dimensionality occurring in practical applications leads to prohibitive numerical costs. To overcome this last problem, we make use of the numerical reduced order method of dynamical low-rank approximation after having derived a stable discretization.

## 3. Discretization of the $Mg$ system

In this section we give a full discretization of both versions (3) and (4) of the $Mg$ system. We start with discretizing equations (3) and (4) in velocity and space before a time discretization is presented in the next subsection.

*3.1. Discretization in velocity and in space*

For the discretization in the velocity space we use a nodal approach and prescribe a certain number of grid points $N_v \in \mathbb{N}$. Due to the special structure of (3b) and (4b) we use a Gauss-Hermite quadrature providing the quadrature nodes $v_1, ..., v_{N_v}$ and weights $\omega_1, ..., \omega_{N_v}$ enabling us to approximate integrals as

$$\int_{\mathbb{R}} e^{-v^2} g(t, x, v) \mathrm{d}v \approx \sum_{k=1}^{N_v} \omega_k g(t, x, v_k).$$

For the discretization of the spatial domain $D \subset \mathbb{R}$ we take $N_x \in \mathbb{N}$ grid points and choose a grid $x_1, ..., x_{N_x}$ with equidistant spacing $\Delta x = \frac{1}{N_x}$. We approximate $x$-dependent quantities by

$$\rho_j(t) \approx \rho(t, x_j) \qquad \text{and} \qquad g_{jk}(t) \approx g(t, x_j, v_k).$$

Spatial derivatives $\partial_x$ are approximated by the tridiagonal stencil matrices $\mathbf{D}^x \in \mathbb{R}^{N_x \times N_x}$ corresponding to a first-order central differencing scheme. Further, a tridiagonal second-order central differencing stabilization matrix $\mathbf{D}^{xx} \in \mathbb{R}^{N_x \times N_x}$ approximating $\partial_{xx}$ is added. Their entries are defined as

$$D^x_{j,j\pm1} = \frac{\pm1}{2\Delta x} \,, \qquad D^{xx}_{j,j} = -\frac{2}{(\Delta x)^2} \,, \quad D^{xx}_{j,j\pm1} = \frac{1}{(\Delta x)^2} \,,$$

whereas all other entries are set to zero. Note that from now on we assume periodic boundary conditions. For this reason we set

$$D^x_{1,2} = D^x_{1,N_x} = 0 \,, \qquad D^x_{N_x,N_x-1} = D^x_{N_x,1} = 0 \,,$$
$$D^{xx}_{1,1} = 0 \,, \; D^{xx}_{1,2} = 0 \,, \qquad D^{xx}_{N_x,N_x-1} = 0 \,, D^{xx}_{N_x,N_x} = 0.$$

Similar to the proof in [2], one can show that the stencil matrices $\mathbf{D}^x$ and $\mathbf{D}^{xx}$ fulfill the following properties:

**Lemma 1** (Summation by parts). *Let $y, z \in \mathbb{R}^{N_x}$ with indices $i, j = 1, ..., N_x$. Then it holds*

$$\sum_{i,j=1}^{N_x} y_j D^x_{ji} z_i = -\sum_{i,j=1}^{N_x} z_j D^x_{ji} y_i \,, \qquad \sum_{i,j=1}^{N_x} z_j D^x_{ji} z_i = 0 \,, \qquad \sum_{i,j=1}^{N_x} y_j D^{xx}_{ji} z_i = \sum_{i,j=1}^{N_x} z_j D^{xx}_{ji} y_i.$$

*Moreover, let $\mathbf{D}^+ \in \mathbb{R}^{N_x \times N_x}$ be defined as*

$$D^+_{j,j} = \frac{-1}{\Delta x} \,, \qquad D^+_{j,j+1} = \frac{1}{\Delta x} \,.$$

*Then, $\sum_{i,j=1}^{N_x} z_j D^{xx}_{ji} z_i = -\sum_{j=1}^{N_x} \left( \sum_{i=1}^{N_x} D^+_{ji} z_i \right)^2$.*

We insert the proposed velocity and space discretization into the advection form (3) and add a stabilizing second-order term for $\rho \partial_x g$. This corresponds to the method used in [14] for the non-linear isothermal Boltzmann-BGK equation and leads to the semi-discrete time-continuous system

$$\dot{g}_{jk}(t) = -\sum_{i=1}^{N_x} D^x_{ji} g_{ik}(t) v_k + \frac{\Delta x}{2} \sum_{i=1}^{N_x} D^{xx}_{ji} g_{ik}(t)|v_k| + \sigma\left(1 - g_{jk}(t)\right) - \frac{g_{jk}(t)}{\rho_j(t)} \dot{\rho}_j(t) - \sum_{i=1}^{N_x} \frac{g_{jk}(t)}{\rho_j(t)} D^x_{ji} \rho_i(t) v_k,$$
(5a)

$$\dot{\rho}_j(t) = -\frac{1}{\sqrt{2\pi}} \sum_{i=1}^{N_x} \sum_{k=1}^{N_v} D^x_{ji} \rho_i(t) g_{ik}(t) v_k \omega_k e^{v_k^2/2} + \frac{\Delta x}{2\sqrt{2\pi}} \sum_{i=1}^{N_x} \sum_{k=1}^{N_v} D^{xx}_{ji} \rho_i(t) g_{ik}(t)|v_k|\omega_k e^{v_k^2/2}.$$
(5b)

4

For the conservative form (4) the second-order stabilization term is applied to $\partial_x(\rho g)$. We obtain the semi-discrete system

$$\dot{g}_{jk}(t) = -\sum_{i=1}^{N_x} \frac{1}{\rho_j(t)} D_{ji}^x \rho_i(t) g_{ik}(t) v_k + \frac{\Delta x}{2} \sum_{i=1}^{N_x} \frac{1}{\rho_j(t)} D_{ji}^{xx} \rho_i(t) g_{ik}(t) |v_k| + \sigma\left(1 - g_{jk}(t)\right) - \frac{g_{jk}(t)}{\rho_j(t)} \dot{\rho}_j(t), \quad (6a)$$

$$\dot{\rho}_j(t) = -\frac{1}{\sqrt{2\pi}} \sum_{i=1}^{N_x} \sum_{k=1}^{N_v} D_{ji}^x \rho_i(t) g_{ik}(t) v_k \omega_k e^{v_k^2/2} + \frac{\Delta x}{2\sqrt{2\pi}} \sum_{i=1}^{N_x} \sum_{k=1}^{N_v} D_{ji}^{xx} \rho_i(t) g_{ik}(t) |v_k| \omega_k e^{v_k^2/2}. \quad (6b)$$

### 3.2. Time discretization

The time discretization of both systems has to be derived carefully to ensure numerical stability. We start with the advection form (5) and perform an explicit Euler step for the transport part in (5a) as well as in (5b). The potentially stiff collision term is treated implicitly. For approximating the time derivative $\partial_t \rho$ the corresponding difference quotient is used. We obtain the fully discrete scheme

$$g_{jk}^{n+1} = g_{jk}^n - \Delta t \sum_{i=1}^{N_x} D_{ji}^x g_{ik}^n v_k + \Delta t \frac{\Delta x}{2} \sum_{i=1}^{N_x} D_{ji}^{xx} g_{ik}^n |v_k|$$

$$+ \sigma \Delta t \left(1 - g_{jk}^{n+1}\right) - \Delta t \frac{g_{jk}^{n+1}}{\rho_j^n} \frac{\rho_j^{n+1} - \rho_j^n}{\Delta t} - \Delta t \frac{g_{jk}^n}{\rho_j^n} \sum_{i=1}^{N_x} D_{ji}^x \rho_i^n v_k, \quad (7a)$$

$$\rho_j^{n+1} = \rho_j^n - \Delta t \frac{1}{\sqrt{2\pi}} \sum_{i=1}^{N_x} \sum_{k=1}^{N_v} D_{ji}^x \rho_i^n g_{ik}^n v_k \omega_k e^{v_k^2/2} + \Delta t \frac{\Delta x}{2\sqrt{2\pi}} \sum_{i=1}^{N_x} \sum_{k=1}^{N_v} D_{ji}^{xx} \rho_i^n g_{ik}^n |v_k| \omega_k e^{v_k^2/2}. \quad (7b)$$

For the conservative form (6) we again perform an explicit Euler step for the transport part in (6a) as well as in (6b). The collision term is treated implicitly and a factor $\frac{\rho^{n+1}}{\rho^n}$ coming from the analysis is added. As before, the time derivative $\partial_t \rho$ is approximated by its difference quotient. This leads to the fully discretized equations

$$g_{jk}^{n+1} = g_{jk}^n - \Delta t \sum_{i=1}^{N_x} \frac{1}{\rho_j^n} D_{ji}^x \rho_i^n g_{ik}^n v_k + \Delta t \frac{\Delta x}{2} \sum_{i=1}^{N_x} \frac{1}{\rho_j^n} D_{ji}^{xx} \rho_i^n g_{ik}^n |v_k|$$

$$+ \sigma \Delta t \frac{\rho_j^{n+1}}{\rho_j^n} \left(1 - g_{jk}^{n+1}\right) - \Delta t \frac{g_{jk}^{n+1}}{\rho_j^n} \frac{\rho_j^{n+1} - \rho_j^n}{\Delta t}, \quad (8a)$$

$$\rho_j^{n+1} = \rho_j^n - \Delta t \frac{1}{\sqrt{2\pi}} \sum_{i=1}^{N_x} \sum_{k=1}^{N_v} D_{ji}^x \rho_i^n g_{ik}^n v_k \omega_k e^{v_k^2/2} + \Delta t \frac{\Delta x}{2\sqrt{2\pi}} \sum_{i=1}^{N_x} \sum_{k=1}^{N_v} D_{ji}^{xx} \rho_i^n g_{ik}^n |v_k| \omega_k e^{v_k^2/2}. \quad (8b)$$

Note that the discretizations for $\rho$ given in (7b) and (8b) are exactly the same. The main differences between the naive discretization of the advection form (7) and the proposed scheme (8) are the stabilization of $\partial_x(\rho g)$ in (8a), opposed to a stabilization of $\rho \partial_x g$ as done in (7a), and the additional factor $\frac{\rho^{n+1}}{\rho^n}$ in the collision term of (8a).

In the fully discrete setting, we make use of the following notations.

**Definition 1** (Fully discrete solution and Maxwellian). *The full solution $f$ of the linear Boltzmann-BGK equation in the fully discrete setting at time $t_n$ is given by* $\mathbf{f}^n = (f_{jk}^n) \in \mathbb{R}^{N_x \times N_v}$ *with entries*

$$f_{jk}^n = \frac{1}{\sqrt{2\pi}} \rho_j^n g_{jk}^n e^{-v_k^2/2}.$$

*For the fully discrete Maxwellian* $\mathbf{M}^n = (M_{jk}^n) \in \mathbb{R}^{N_x \times N_v}$ *at time $t_n$, we have* $M_{jk}^n = \frac{1}{\sqrt{2\pi}} \rho_j^n e^{-v_k^2/2}$.

## 4. Numerical stability

Although the derivation of the equations in (7) and (8) is similar, both systems differ drastically in terms of numerical stability. In this section, both fully discretized schemes presented are compared.

### 4.1. Naive discretization

We begin with the naive discretization given in (7) that is comparable to the one chosen in [14] in the sense that the advection form of the multiplicative splitting is used. In [14], numerical experiments are given but no explicit stability analysis is conducted. In the following, we give an example which shows that numerical stability in the sense of von Neumann can not be guaranteed.

**Theorem 1.** *There exist initial values* $\mathbf{g}^n = \left(g_{jk}^n\right) \in \mathbb{R}^{N_x \times N_v}$ *and* $\boldsymbol{\rho}^n = \left(\rho_j^n\right) \in \mathbb{R}^{N_x}$ *such that the numerical scheme proposed in* (7) *for* $\sigma = 0$ *is not von Neumann stable.*

*Proof.* Let us assume a solution $g_{jk}^n$ that is constant in space and velocity, e.g. $g_{jk}^n \equiv 1$. For this solution the terms containing $\mathbf{D}^x \mathbf{g}^n$ and $\mathbf{D}^{xx} \mathbf{g}^n$ are zero. Let us further assume that there is no collisionality, i.e. $\sigma = 0$. We insert this information into (7a) and get

$$g_{jk}^{n+1} = 1 - \Delta t \frac{1}{\rho_j^n} \frac{\rho_j^{n+1} - \rho_j^n}{\Delta t} g_{jk}^{n+1} - \Delta t \frac{1}{\rho_j^n} \sum_{i=1}^{N_x} \left(D_{ji}^x \rho_i^n\right) v_k.$$

After rearranging, we have that

$$\rho_j^{n+1} g_{jk}^{n+1} = \rho_j^n - \Delta t \sum_{i=1}^{N_x} \left(D_{ji}^x \rho_i^n\right) v_k.$$

Multiplication with $\frac{1}{\sqrt{2\pi}} e^{-v_k^2/2}$ then leads to

$$f_{jk}^{n+1} = f_{jk}^n - \Delta t \sum_{i=1}^{N_x} D_{ji}^x f_{ik}^n v_k. \tag{9}$$

This corresponds to a discretization of $\partial_t f + v \partial_x f = 0$ with explicit Euler in time and centered finite differences in space for which it is well-known that it is not von Neumann stable [21, 28]. $\square$

Indeed, one can show that the discretization given in (9) is not von Neumann stable, but stable for relatively small time step sizes [28]. This matches our numerical insights from [14], where the space discretization is comparable to (7) and small time step sizes are required.

### 4.2. Stable discretization

Having seen that for a certain choice of the initial values the system of equations (7) is not von Neumann stable, we now consider equations (8) in terms of numerical stability. We observe that the advection terms are treated explicitly, whereas the collision term is treated implicitly, leading to a removal of the potential stiffness caused by a large number of collisions. We seek a rigorous proof of stability under a classic hyperbolic CFL condition that will be derived in the following norm.

**Definition 2** (Stability norm). *For* $\mathbf{f}^n = (f_{jk}^n) \in \mathbb{R}^{N_x \times N_v}$, *the* $\mathscr{H}$-*norm shall be defined as*

$$\|\mathbf{f}^n\|_{\mathscr{H}}^2 = \sqrt{2\pi} \sum_{j=1}^{N_x} \sum_{k=1}^{N_v} \left(f_{jk}^n\right)^2 \omega_k e^{3v_k^2/2}.$$

*This corresponds to a Frobenius norm* $\|\cdot\|_F$ *with weights* $\sqrt{2\pi}\omega_k e^{3v_k^2/2}$.

The choice of this norm is inspired by the analysis in [1], where hypocoercivity for the linear Boltzmann-BGK equation is shown. Different from there, we use a fully discrete analogue to the considered weighted $L^2$-norm that also takes the Gauss-Hermite quadrature into account. Note that the factor $\sqrt{2\pi}$ does not affect the stability but is added for consistency in the sense that for $g_{jk}^n = 1$ the condition $\frac{1}{\sqrt{2\pi}} \sum_{k=1}^{N_v} g_{jk}^n \omega_k e^{v_k^2/2} = 1$ holds. This is the discrete counterpart of $\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} g e^{-v^2/2} \mathrm{d}v = 1$ which is equivalent to $\int_{\mathbb{R}} f \mathrm{d}v = \rho$ that in a discrete formulation can be written as $\sum_{k=1}^{N_v} f_{jk}^n \omega_k e^{v_k^2} = \rho_j^n$. This relation shall be preserved by the numerical scheme.

**Lemma 2.** *Let us assume that the initial condition for $g$ satisfies $\frac{1}{\sqrt{2\pi}} \sum_{k=1}^{N_v} g_{jk}^0 \omega_k e^{v_k^2/2} = 1$. Then, for all $n \in \mathbb{N}$ we have*

$$\sum_{k=1}^{N_v} f_{jk}^n \omega_k e^{v_k^2} = \rho_j^n.$$

*Proof.* We start from equation (8a), bring the terms containing $g_{jk}^{n+1}$ to the left-hand side and multiply it with $\rho_j^{n+1}$. This gives

$$(1 + \sigma\Delta t) \rho_j^{n+1} g_{jk}^{n+1} = \rho_j^n g_{jk}^n - \Delta t \sum_{i=1}^{N_x} D_{ji}^x \rho_i^n g_{ik}^n v_k + \Delta t \frac{\Delta x}{2} \sum_{i=1}^{N_x} D_{ji}^{xx} \rho_i^n g_{ik}^n |v_k| + \sigma\Delta t \rho_j^{n+1}.$$

Multiplication with $\frac{1}{\sqrt{2\pi}} \omega_k e^{v_k^2/2}$ and summation over $k$ leads to

$$(1 + \sigma\Delta t) \rho_j^{n+1} \frac{1}{\sqrt{2\pi}} \sum_{k=1}^{N_v} g_{jk}^{n+1} \omega_k e^{v_k^2/2} = \rho_j^n \frac{1}{\sqrt{2\pi}} \sum_{k=1}^{N_v} g_{jk}^n \omega_k e^{v_k^2/2} - \Delta t \frac{1}{\sqrt{2\pi}} \sum_{i=1}^{N_x} \sum_{k=1}^{N_v} D_{ji}^x \rho_i^n g_{ik}^n v_k \omega_k e^{v_k^2/2}$$

$$+ \Delta t \frac{\Delta x}{2\sqrt{2\pi}} \sum_{i=1}^{N_x} \sum_{k=1}^{N_v} D_{ji}^{xx} \rho_i^n g_{ik}^n |v_k| \omega_k e^{v_k^2/2}$$

$$+ \sigma\Delta t \rho_j^{n+1} \frac{1}{\sqrt{2\pi}} \sum_{k=1}^{N_v} \omega_k e^{v_k^2/2}.$$

We insert the induction assumption as well as $\frac{1}{\sqrt{2\pi}} \sum_{k=1}^{N_v} \omega_k e^{v_k^2/2} = 1$. Then, together with (8b), we get

$$(1 + \sigma\Delta t) \rho_j^{n+1} \frac{1}{\sqrt{2\pi}} \sum_{k=1}^{N_v} g_{jk}^{n+1} \omega_k e^{v_k^2/2} = (1 + \sigma\Delta t) \rho_j^{n+1}.$$

Cancelling with $(1 + \sigma\Delta t) \rho_j^{n+1}$ gives the desired equality, and completes the proof. $\qquad\square$

Also, the following inequality is indispensable to show numerical stability of the above system.

**Lemma 3.** *Under the time step restriction $\max_k(|v_k|)\Delta t \leq \Delta x$ it holds*

$$\Delta t \left\| \mathbf{D}^x \mathbf{f}^n \operatorname{diag}(v_k) - \frac{\Delta x}{2} \mathbf{D}^{xx} \mathbf{f}^n \operatorname{diag}(|v_k|) \right\|_{\mathscr{H}}^2 - \Delta x \left\| \mathbf{D}^+ \mathbf{f}^n \operatorname{diag}\left(|v_k|^{1/2}\right) \right\|_{\mathscr{H}}^2 \leq 0.$$

*Proof.* We want to apply a Fourier analysis in $x$ that allows us to write the stencil matrices in diagonal form. As in [2, 26], we define the matrix $\mathbf{E} \in \mathbb{C}^{N_x \times N_x}$ as

$$E_{j\alpha} = \sqrt{\Delta x} \exp(i\alpha\pi x_j), \quad j, \alpha = 1, ..., N_x,$$

7

where $i \in \mathbb{C}$ denotes the imaginary unit. It is orthonormal, i.e. $\mathbf{E}\mathbf{E}^H = \mathbf{E}^H\mathbf{E} = \mathbf{I}$, where the superscript $H$ stands for the complex transpose, and diagonalizes the stencil matrices

$$\mathbf{D}^\gamma \mathbf{E} = \mathbf{E}\mathbf{\Lambda}^\gamma, \qquad \text{with } \gamma \in \{x, xx, +\},$$

with $\mathbf{\Lambda}^\gamma \in \mathbb{C}^{N_x \times N_x}$ being the diagonal matrices with entries

$$\lambda_{\alpha\alpha}^x = \frac{1}{2\Delta x}(e^{i\alpha\pi\Delta x} - e^{-i\alpha\pi\Delta x}) = \frac{i}{\Delta x}\sin(\nu_\alpha),$$

$$\lambda_{\alpha\alpha}^{xx} = \frac{1}{(\Delta x)^2}\left(e^{i\alpha\pi\Delta x} - 2 + e^{-i\alpha\pi\Delta x}\right) = \frac{2}{(\Delta x)^2}\left(\cos(\nu_\alpha) - 1\right),$$

$$\lambda_{\alpha\alpha}^+ = \frac{1}{\Delta x}\left(e^{i\alpha\pi\Delta x} - 1\right) = \frac{1}{\Delta x}\left(\cos(\nu_\alpha) + i\sin(\nu_\alpha) - 1\right),$$

and $\nu_\alpha := \alpha\pi\Delta x$. Let us denote $\widehat{\mathbf{f}}^n = \left(\widehat{f}_{\alpha k}^n\right) \in \mathbb{C}^{N_x \times N_v}$ with entries $\widehat{f}_{\alpha k}^n = \sum_{j=1}^{N_x} E_{\alpha j} f_{jk}^n$. With Parseval's identity we obtain

$$\Delta t \left\|\mathbf{D}^x\mathbf{f}^n \operatorname{diag}(v_k) - \frac{\Delta x}{2}\mathbf{D}^{xx}\mathbf{f}^n \operatorname{diag}(|v_k|)\right\|_{\mathscr{H}}^2 - \Delta x \left\|\mathbf{D}^+\mathbf{f}^n \operatorname{diag}\left(|v_k|^{1/2}\right)\right\|_{\mathscr{H}}^2$$

$$= \Delta t \left\|\mathbf{D}^x\mathbf{f}^n \operatorname{diag}\left(v_k\omega_k^{1/2}e^{3v_k^2/4}\right) - \frac{\Delta x}{2}\mathbf{D}^{xx}\mathbf{f}^n \operatorname{diag}\left(|v_k|\omega_k^{1/2}e^{3v_k^2/4}\right)\right\|_F^2$$

$$- \Delta x \left\|\mathbf{D}^+\mathbf{f}^n \operatorname{diag}\left(|v_k|^{1/2}\omega_k^{1/2}e^{3v_k^2/4}\right)\right\|_F^2$$

$$\overset{\text{Parseval}}{=} \Delta t \left\|\mathbf{\Lambda}^x\widehat{\mathbf{f}}^n \operatorname{diag}\left(v_k\omega_k^{1/2}e^{3v_k^2/4}\right) - \frac{\Delta x}{2}\mathbf{\Lambda}^{xx}\widehat{\mathbf{f}}^n \operatorname{diag}\left(|v_k|\omega_k^{1/2}e^{3v_k^2/4}\right)\right\|_F^2$$

$$- \Delta x \left\|\mathbf{\Lambda}^+\widehat{\mathbf{f}}^n \operatorname{diag}\left(|v_k|^{1/2}\omega_k^{1/2}e^{3v_k^2/4}\right)\right\|_F^2$$

$$= 2\sum_{\alpha=1}^{N_x}\sum_{k=1}^{N_v}\left(\Delta t\frac{|v_k|^2}{(\Delta x)^2}(1 - \cos(\nu_\alpha)) - \frac{|v_k|}{\Delta x}(1 - \cos(\nu_\alpha))\right)\omega_k e^{3v_k^2/2}\left|\widehat{f}_{\alpha k}^n\right|^2.$$

A sufficient condition to ensure negativity is that

$$\Delta t\frac{|v_k|^2}{(\Delta x)^2}(1 - \cos(\nu_\alpha)) \le \frac{|v_k|}{\Delta x}(1 - \cos(\nu_\alpha))$$

must hold for each index $k$. This leads to the time step restriction $\max_k(|v_k|)\Delta t \le \Delta x$, which proves the lemma. $\qquad\square$

We can now show numerical stability of the proposed system.

**Theorem 2.** *Under the time step restriction $\max_k(|v_k|)\Delta t \le \Delta x$, the fully discrete system (8) is numerically stable in the $\mathscr{H}$-norm, i.e.*

$$\|\mathbf{f}^{n+1}\|_{\mathscr{H}}^2 \le \|\mathbf{f}^n\|_{\mathscr{H}}^2.$$

*Proof.* We multiply (8a) with $\rho_j^{n+1}\rho_j^n g_{jk}^{n+1}$ and bring the last term of the equation from the right-hand to the left-hand side. This gives

$$\left(\rho_j^{n+1}g_{jk}^{n+1}\right)^2 = \rho_j^n g_{jk}^n\rho_j^{n+1}g_{jk}^{n+1} - \Delta t\sum_{i=1}^{N_x}\rho_j^{n+1}g_{jk}^{n+1}D_{ji}^x\rho_i^n g_{ik}^n v_k + \Delta t\frac{\Delta x}{2}\sum_{i=1}^{N_x}\rho_j^{n+1}g_{jk}^{n+1}D_{ji}^{xx}\rho_i^n g_{ik}^n|v_k|$$

$$+ \sigma\Delta t\rho_j^{n+1}g_{jk}^{n+1}\left(\rho_j^{n+1} - \rho_j^{n+1}g_{jk}^{n+1}\right).$$

Multiplication with $2 \left( \frac{1}{\sqrt{2\pi}} e^{-v_k^2/2} \right)^2$ then leads to

$$2 \left( f_{jk}^{n+1} \right)^2 = 2 f_{jk}^n f_{jk}^{n+1} - 2\Delta t \sum_{i=1}^{N_x} f_{jk}^{n+1} D_{ji}^x f_{ik}^n v_k + \Delta t \Delta x \sum_{i=1}^{N_x} f_{jk}^{n+1} D_{ji}^{xx} f_{ik}^n |v_k| + 2\sigma \Delta t f_{jk}^{n+1} \left( M_{jk}^{n+1} - f_{jk}^{n+1} \right).$$

Note that it holds

$$2 f_{jk}^n f_{jk}^{n+1} = \left( f_{jk}^{n+1} \right)^2 + \left( f_{jk}^n \right)^2 - \left( f_{jk}^{n+1} - f_{jk}^n \right)^2.$$

We insert this relation and obtain

$$\left( f_{jk}^{n+1} \right)^2 = \left( f_{jk}^n \right)^2 - \left( f_{jk}^{n+1} - f_{jk}^n \right)^2 - 2\Delta t \sum_{i=1}^{N_x} f_{jk}^{n+1} D_{ji}^x f_{ik}^n v_k + \Delta t \Delta x \sum_{i=1}^{N_x} f_{jk}^{n+1} D_{ji}^{xx} f_{ik}^n |v_k|$$
$$+ 2\sigma \Delta t f_{jk}^{n+1} \left( M_{jk}^{n+1} - f_{jk}^{n+1} \right).$$

In the next step, we multiply with $\sqrt{2\pi} \omega_k e^{3v_k^2/2}$ and sum over $j$ and $k$. This gives

$$\|\mathbf{f}^{n+1}\|_{\mathcal{H}}^2 = \|\mathbf{f}^n\|_{\mathcal{H}}^2 - \sqrt{2\pi} \sum_{j=1}^{N_x} \sum_{k=1}^{N_v} \left( f_{jk}^{n+1} - f_{jk}^n \right)^2 \omega_k e^{3v_k^2/2}$$
$$- 2\sqrt{2\pi} \Delta t \sum_{i,j=1}^{N_x} \sum_{k=1}^{N_v} f_{jk}^{n+1} D_{ji}^x f_{ik}^n v_k \omega_k e^{3v_k^2/2} + \sqrt{2\pi} \Delta t \Delta x \sum_{i,j=1}^{N_x} \sum_{k=1}^{N_v} f_{jk}^{n+1} D_{ji}^{xx} f_{ik}^n |v_k| \omega_k e^{3v_k^2/2}$$
$$+ 2\sqrt{2\pi} \sigma \Delta t \sum_{j=1}^{N_x} \sum_{k=1}^{N_v} f_{jk}^{n+1} \left( M_{jk}^{n+1} - f_{jk}^{n+1} \right) \omega_k e^{3v_k^2/2}.$$

From Lemma 2 we have that $\sum_{k=1}^{N_v} f_{jk}^{n+1} \omega_k e^{v_k^2} = \rho_j^{n+1}$. Hence, we can derive that the following term $2\sqrt{2\pi} \sigma \Delta t \sum_{j=1}^{N_x} \sum_{k=1}^{N_v} M_{jk}^{n+1} \left( M_{jk}^{n+1} - f_{jk}^{n+1} \right) \omega_k e^{3v_k^2/2}$ is equal to zero. Lemma 2 gives that also the term $2\sqrt{2\pi} \Delta t \sum_{i,j=1}^{N_x} \sum_{k=1}^{N_v} f_{jk}^n D_{ji}^x f_{ik}^n v_k \omega_k e^{3v_k^2/2}$ is equal to zero. We subtract both and add an additional zero by adding and subtracting the second-order term $\sqrt{2\pi} \Delta t \Delta x \sum_{i,j=1}^{N_x} \sum_{k=1}^{N_v} f_{jk}^n D_{ji}^{xx} f_{ik}^n |v_k| \omega_k e^{3v_k^2/2}$. This leads to

$$\|\mathbf{f}^{n+1}\|_{\mathcal{H}}^2 = \|\mathbf{f}^n\|_{\mathcal{H}}^2 - \sqrt{2\pi} \sum_{j=1}^{N_x} \sum_{k=1}^{N_v} \left( f_{jk}^{n+1} - f_{jk}^n \right)^2 \omega_k e^{3v_k^2/2}$$

$$- 2\sqrt{2\pi} \Delta t \sum_{i,j=1}^{N_x} \sum_{k=1}^{N_v} \left( f_{jk}^{n+1} - f_{jk}^n \right) D_{ji}^x f_{ik}^n v_k \omega_k e^{3v_k^2/2} \tag{I}$$

$$+ \sqrt{2\pi} \Delta t \Delta x \sum_{i,j=1}^{N_x} \sum_{k=1}^{N_v} \left( f_{jk}^{n+1} - f_{jk}^n \right) D_{ji}^{xx} f_{ik}^n |v_k| \omega_k e^{3v_k^2/2} \tag{II}$$

$$+ \sqrt{2\pi} \Delta t \Delta x \sum_{i,j=1}^{N_x} \sum_{k=1}^{N_v} f_{jk}^n D_{ji}^{xx} f_{ik}^n |v_k| \omega_k e^{3v_k^2/2} \tag{III}$$

$$- 2\sqrt{2\pi} \sigma \Delta t \sum_{j=1}^{N_x} \sum_{k=1}^{N_v} \left( f_{jk}^{n+1} - M_{jk}^{n+1} \right)^2 \omega_k e^{3v_k^2/2}.$$

Now, we consider the parts (I), (II) and (III) separately. Let us start with (I) and (II) and apply Young's

inequality which states that for $a, b \in \mathbb{R}$ it holds that $a \cdot b \leq \frac{a^2}{2} + \frac{b^2}{2}$. This gives for (I) +(II)

$$-2\sqrt{2\pi}\Delta t \sum_{i,j=1}^{N_x}\sum_{k=1}^{N_v}\left(f_{jk}^{n+1}-f_{jk}^n\right)D_{ji}^x f_{ik}^n v_k\omega_k e^{3v_k^2/2} + \sqrt{2\pi}\Delta t\Delta x \sum_{i,j=1}^{N_x}\sum_{k=1}^{N_v}\left(f_{jk}^{n+1}-f_{jk}^n\right)D_{ji}^{xx} f_{ik}^n |v_k|\omega_k e^{3v_k^2/2}$$

$$=\sum_{j=1}^{N_x}\sum_{k=1}^{N_v}\left(-\sqrt{2}\left(2\pi\right)^{1/4}\left(f_{jk}^{n+1}-f_{jk}^n\right)\omega_k^{1/2}e^{3v_k^2/4}\right)\left(\sqrt{2}\left(2\pi\right)^{1/4}\Delta t\sum_{i=1}^{N_x}\left(D_{ji}^x f_{ik}^n v_k - \frac{\Delta x}{2}D_{ji}^{xx}f_{ik}^n|v_k|\right)\omega_k^{1/2}e^{3v_k^2/4}\right)$$

$$\leq \sqrt{2\pi}\sum_{j=1}^{N_x}\sum_{k=1}^{N_v}\left(f_{jk}^{n+1}-f_{jk}^n\right)^2\omega_k e^{3v_k^2/2} + \sqrt{2\pi}\left(\Delta t\right)^2\sum_{j=1}^{N_x}\sum_{k=1}^{N_v}\left(\sum_{i=1}^{N_x}\left(D_{ji}^x f_{ik}^n v_k - \frac{\Delta x}{2}D_{ji}^{xx}f_{ik}^n|v_k|\right)\right)^2\omega_k e^{3v_k^2/2}.$$

For (III) we get with the properties of the stencil matrices given in Lemma 1 that

$$\sqrt{2\pi}\Delta t\Delta x\sum_{i,j=1}^{N_x}\sum_{k=1}^{N_v}f_{jk}^n D_{ji}^{xx}f_{ik}^n|v_k|\omega_k e^{3v_k^2/2} = -\sqrt{2\pi}\Delta t\Delta x\sum_{j=1}^{N_x}\sum_{k=1}^{N_v}\left(\sum_{i=1}^{N_x}D_{ji}^+ f_{ik}^n|v_k|^{1/2}\right)^2\omega_k e^{3v_k^2/2}.$$

Inserting these relations then gives

$$\|\mathbf{f}^{n+1}\|_{\mathscr{H}}^2 \leq \|\mathbf{f}^n\|_{\mathscr{H}}^2 + \sqrt{2\pi}\left(\Delta t\right)^2\sum_{i,j=1}^{N_x}\sum_{k=1}^{N_v}\left(D_{ji}^x f_{ik}^n v_k - \frac{\Delta x 2}{D}{}^{xx}_{ji}f_{ik}^n|v_k|\right)^2\omega_k e^{3v_k^2/2}$$

$$- \sqrt{2\pi}\Delta t\Delta x\sum_{j=1}^{N_x}\sum_{k=1}^{N_v}\left(\sum_{i=1}^{N_x}D_{ji}^+ f_{ik}^n|v_k|^{1/2}\right)^2\omega_k e^{3v_k^2/2}$$

$$- 2\sqrt{2\pi}\sigma\Delta t\sum_{j=1}^{N_x}\sum_{k=1}^{N_v}\left(f_{jk}^{n+1}-M_{jk}^{n+1}\right)^2\omega_k e^{3v_k^2/2}.$$

With Lemma 3 we can conclude that under the CFL condition $\max_k(|v_k|)\Delta t \leq \Delta x$ it holds $\|\mathbf{f}^{n+1}\|_{\mathscr{H}}^2 \leq \|\mathbf{f}^n\|_{\mathscr{H}}^2$. Hence, with this time step the proposed fully discrete system (8) is numerically stable in the $\mathscr{H}$-norm. □

## 5. Dynamical low-rank approximation for the stable $Mg$ system

In practical applications, the implementation of the full system given in (8) may lead to prohibitive numerical costs, especially when computing in higher dimensional settings. To reduce computational and memory demands, we apply dynamical low-rank approximation to the distribution function $g$.

### 5.1. Background on dynamical low-rank approximation

The concept of dynamical low-rank approximation has been introduced in a semi-discrete time-dependent matrix setting [23]. Let us consider $\mathbf{g} \in \mathbb{R}^{N_x \times N_v}$ being the solution of the matrix differential equation

$$\dot{\mathbf{g}}(t) = \mathbf{F}\left(t, \mathbf{g}(t)\right),$$

where $\mathbf{F} : \mathbb{R}^{N_x \times N_v} \to \mathbb{R}^{N_x \times N_v}$ denotes the right-hand side of the equation. We then seek for an approximation of $\mathbf{g}$ in the following form

$$\mathbf{g}_r(t) = \mathbf{X}(t)\mathbf{S}(t)\mathbf{V}(t)^\top, \tag{10}$$

with $\mathbf{X} \in \mathbb{R}^{N_x \times r}$ and $\mathbf{V} \in \mathbb{R}^{N_v \times r}$ denoting the orthonormal spatial and orthonormal velocity basis, respectively. The slim matrix $\mathbf{S} \in \mathbb{R}^{r \times r}$ is called the coefficient or coupling matrix and determines the rank $r$ of

the approximation. The set of all matrices of the above form (10) constitute the low-rank manifold $\mathcal{M}_r$. Its corresponding tangent space at $\mathbf{g}_r(t)$ shall be denoted by $\mathcal{T}_{\mathbf{g}_r(t)}\mathcal{M}_r$. We now look for $\mathbf{g}_r(t) \in \mathcal{M}_r$ such that at all times $t$ the minimization problem

$$\min_{\dot{\mathbf{g}}_r(t) \in \mathcal{T}_{\mathbf{g}_r(t)}\mathcal{M}_r} \| \dot{\mathbf{g}}_r(t) - \mathbf{F}\left(t, \mathbf{g}_r(t)\right) \|_F$$

is fulfilled. Following [23], this minimization constraint is equivalent to determining $\dot{\mathbf{g}}_r(t) \in \mathcal{T}_{\mathbf{g}_r(t)}$ by an orthogonal projection onto the tangent space such that

$$\dot{\mathbf{g}}_r(t) = \mathbf{P}(\mathbf{g}_r(t))\mathbf{F}(\mathbf{g}_r(t)), \tag{11}$$

where the orthogonal projector $\mathbf{P}$ onto $\mathcal{T}_g\mathcal{M}_r$ applied to an arbitrary quantity $\mathbf{G}$ can explicitly be given as

$$\mathbf{P}(\mathbf{g}_r(t))\mathbf{G} = \mathbf{X}\mathbf{X}^\top\mathbf{G} - \mathbf{X}\mathbf{X}^\top\mathbf{G}\mathbf{V}\mathbf{V}^\top + \mathbf{G}\mathbf{V}\mathbf{V}^\top.$$

Different robust time integrators to solve (11) exist. Frequently used integrators are the projector-splitting integrator [29], the basis update & Galerkin integrator [10] as well as its rank-adaptive extension [8] and the parallel integrator [9]. In this paper, we make use of the rank-adaptive BUG integrator whose concept shall be explained in the following.

In the first two steps, the rank-adaptive BUG integrator updates and augments the bases $\mathbf{X}$ and $\mathbf{V}$ in parallel such that their rank increases from $r$ to $2r$, for the spatial and velocity basis respectively. We denote the augmented quantities of rank $2r$ with hats. Then, for the augmented bases a Galerkin step is performed before in a last step a new rank $r_1 \leq 2r$ is determined using a truncation with prescribed tolerance. In detail, the rank-adaptive BUG integrator performs the following steps in order to update the matrix $\mathbf{g}_r^n = \mathbf{X}^n\mathbf{S}^n\mathbf{V}^{n,\top}$ at time $t_n$ to $\mathbf{g}_r^{n+1} = \mathbf{X}^{n+1}\mathbf{S}^{n+1}\mathbf{V}^{n+1,\top}$ at time $t_{n+1} = t_n + \Delta t$:

**K-Step**: Let us fix the velocity basis $\mathbf{V}^n$ at time $t_n$ and introduce the notation $\mathbf{K}(t) = \mathbf{X}(t)\mathbf{S}(t)$. The spatial basis $\mathbf{X}^n$ is updated and augmented by first solving the PDE

$$\dot{\mathbf{K}}(t) = \mathbf{F}\left(t, \mathbf{K}(t)\mathbf{V}^{n,\top}\right)\mathbf{V}^n, \quad \mathbf{K}(t_n) = \mathbf{X}^n\mathbf{S}^n,$$

and then determining $\widehat{\mathbf{X}}^{n+1} \in \mathbb{R}^{N_x \times 2r}$ as an orthonormal basis of $[\mathbf{K}(t_{n+1}), \mathbf{X}^n] \in \mathbb{R}^{N_x \times 2r}$, e.g. by QR-decomposition. Then, we compute $\widehat{\mathbf{M}} = \widehat{\mathbf{X}}^{n+1,\top}\mathbf{X}^n \in \mathbb{R}^{2r \times r}$.

**L-Step**: Let us fix the spatial basis $\mathbf{X}^n$ at time $t_n$ and introduce the notation $\mathbf{L}(t) = \mathbf{V}(t)\mathbf{S}(t)^\top$. The velocity basis $\mathbf{V}^n$ is updated and augmented by first solving the PDE

$$\dot{\mathbf{L}}(t) = \mathbf{F}\left(t, \mathbf{X}^n\mathbf{L}(t)^\top\right)^\top\mathbf{X}^n, \quad \mathbf{L}(t_n) = \mathbf{V}^n\mathbf{S}^{n,\top},$$

and then determining $\widehat{\mathbf{V}}^{n+1} \in \mathbb{R}^{N_v \times 2r}$ as an orthonormal basis of $[\mathbf{L}(t_{n+1}), \mathbf{V}^n] \in \mathbb{R}^{N_v \times 2r}$, e.g. by QR-decomposition. Then, we compute $\widehat{\mathbf{N}} = \widehat{\mathbf{V}}^{n+1,\top}\mathbf{V}^n \in \mathbb{R}^{2r \times r}$.

**S-step**: Update the coupling matrix from $\mathbf{S}^n \in \mathbf{R}^{r \times r}$ to $\widehat{\mathbf{S}}^{n+1} \in \mathbb{R}^{2r \times 2r}$ by solving the ODE

$$\dot{\widehat{\mathbf{S}}}(t) = \widehat{\mathbf{X}}^{n+1,\top}\mathbf{F}\left(t, \widehat{\mathbf{X}}^{n+1}\widehat{\mathbf{S}}(t)\widehat{\mathbf{V}}^{n+1,\top}\right)\widehat{\mathbf{V}}^{n+1}, \quad \widehat{\mathbf{S}}(t_n) = \widehat{\mathbf{M}}\mathbf{S}^n\widehat{\mathbf{N}}^\top.$$

**Truncation**: Compute the singular value decomposition of $\widehat{\mathbf{S}}^{n+1} = \widehat{\mathbf{P}}\boldsymbol{\Sigma}\widehat{\mathbf{Q}}^\top$ with $\boldsymbol{\Sigma} = \text{diag}(\sigma_j)$. For a prescribed tolerance parameter $\vartheta$ the new rank $r_1 \leq 2r$ is chosen such that

$$\left( \sum_{j=r_1+1}^{2r} \sigma_j^2 \right)^{1/2} \leq \vartheta.$$

11

Let now $\mathbf{S}^{n+1} \in \mathbb{R}^{r_1 \times r_1}$ contain the $r_1$ largest singular values and $\mathbf{P}^{n+1} \in \mathbb{R}^{2r \times r_1}$ and $\mathbf{Q}^{n+1} \in \mathbb{R}^{2r \times r_1}$ contain the first $r_1$ columns of $\widehat{\mathbf{P}}$ and $\widehat{\mathbf{Q}}$, respectively. Then, the time-updated spatial basis can be determined as $\mathbf{X}^{n+1} = \widehat{\mathbf{X}}^{n+1}\mathbf{P}^{n+1} \in \mathbb{R}^{N_x \times r_1}$ and the time-updated velocity basis as $\mathbf{V}^{n+1} = \widehat{\mathbf{V}}^{n+1}\mathbf{Q}^{n+1} \in \mathbb{R}^{N_v \times r_1}$.

Altogether, the time-updated approximation of the distribution function after one time step is given by $\mathbf{g}_r^{n+1} = \mathbf{X}^{n+1}\mathbf{S}^{n+1}\mathbf{V}^{n+1,\top}$.

## 5.2. DLRA algorithm for multiplicative linear Boltzmann-BGK

In this section, we apply DLRA to the stable and fully discretized system (8). Bringing all term containing $g_{jk}^{n+1}$ to the left-hand side of (8a) and multiplying it with $\frac{\rho_j^n}{\rho_j^{n+1}}$, equations (8) can be equivalently written as

$$g_{jk}^{n+1}(1 + \sigma\Delta t) = \frac{\rho_j^n}{\rho_j^{n+1}}g_{jk}^n - \Delta t \sum_{i=1}^{N_x} \frac{1}{\rho_j^{n+1}} D_{ji}^x (\rho_i^n g_{ik}^n) v_k + \Delta t \frac{\Delta x}{2} \sum_{i=1}^{N_x} \frac{1}{\rho_j^{n+1}} D_{ji}^{xx} (\rho_i^n g_{ik}^n) |v_k| + \sigma\Delta t, \quad (12a)$$

$$\rho_j^{n+1} = \rho_j^n - \Delta t \frac{1}{\sqrt{2\pi}} \sum_{i=1}^{N_x} \sum_{k=1}^{N_v} D_{ji}^x \rho_i^n g_{ik}^n v_k \omega_k e^{v_k^2/2} + \Delta t \frac{\Delta x}{2\sqrt{2\pi}} \sum_{i=1}^{N_x} \sum_{k=1}^{N_v} D_{ji}^{xx} \rho_i^n g_{ik}^n |v_k| \omega_k e^{v_k^2/2}. \quad (12b)$$

We propose a low-rank implementation that uses the rank-adaptive BUG integrator [8] introduced in the previous subsection for equation (12a) together with an additional basis augmentation and a suitable truncation strategy. For this scheme we show numerical stability. Note that in the following, for simplicity, we write $\mathbf{g} = (g_{jk})$ instead of $\mathbf{g}_r$.

Starting from (12), we apply the rank-adaptive BUG integrator leading to a splitting of (12a) into three substeps, the $K$-, $L$-, and $S$-step. In the $K$- as well as in the $L$-step we perform an additional basis augmentation ensuring certain quantities to be contained in the basis at all times. The augmented bases are then used to determine the $S$-step. Note that the scattering term $(1 + \sigma\Delta t)$ is only applied in the $S$-step as it does not affect the span of the basis functions derived in the $K$- and $L$-step. We obtain the following DLRA scheme:

Substituting $g_{jk}^n = \sum_{m,\ell=1}^r X_{jm}^n S_{m\ell}^n V_{k\ell}^n$ into the update equation (12b) yields

$$\rho_j^{n+1} = \rho_j^n - \Delta t \frac{1}{\sqrt{2\pi}} \sum_{i=1}^{N_x} D_{ji}^x \rho_i^n \sum_{m,\ell=1}^r X_{im}^n S_{m\ell}^n \sum_{k=1}^{N_v} V_{k\ell}^n v_k \omega_k e^{v_k^2/2}$$
$$+ \Delta t \frac{\Delta x}{2\sqrt{2\pi}} \sum_{i=1}^{N_x} D_{ji}^{xx} \rho_i^n \sum_{m,\ell=1}^r X_{im}^n S_{m\ell}^n \sum_{k=1}^{N_v} V_{k\ell}^n |v_k| \omega_k e^{v_k^2/2}. \quad (13a)$$

For the $K$-step we introduce the notation $K_{j\ell}^n = \sum_{m=1}^r X_{jm}^n S_{m\ell}^n$ and solve

$$K_{jp}^{n+1} = \frac{\rho_j^n}{\rho_j^{n+1}} K_{jp}^n - \Delta t \frac{1}{\rho_j^{n+1}} \sum_{i=1}^{N_x} D_{ji}^x \rho_i^n \sum_{\ell=1}^r K_{i\ell}^n \sum_{k=1}^{N_v} V_{k\ell}^n v_k V_{kp}^n$$
$$+ \Delta t \frac{\Delta x}{2} \frac{1}{\rho_j^{n+1}} \sum_{i=1}^{N_x} D_{ji}^{xx} \rho_i^n \sum_{\ell=1}^r K_{i\ell}^n \sum_{k=1}^{N_v} V_{k\ell}^n |v_k| V_{kp}^n + \sigma\Delta t \sum_{k=1}^{N_v} V_{kp}^n. \quad (13b)$$

This gives the updated matrix $\mathbf{K}^{n+1} = (K_{jp}^{n+1})$ to which together with the old basis $\mathbf{X}^n = (X_{jm}^n)$ a QR-decomposition is applied, giving the new augmented basis $\widehat{\mathbf{X}}^{n+1} = (\widehat{X}_{jm}^{n+1})$.

In addition, we augment this basis according to

$$\widehat{\widehat{\mathbf{X}}}^{n+1} = \mathrm{qr}\left([\widehat{\mathbf{X}}^{n+1}, (\boldsymbol{\rho}^{n+1})^2 \widehat{\mathbf{X}}^{n+1}]\right) \quad (13c)$$

leading to a new augmented basis $\widehat{\widehat{\mathbf{X}}}^{n+1} = (\widehat{\widehat{X}}^{n+1}_{jm})$ of rank $4r$, and we compute $\widehat{\widehat{\mathbf{M}}} = \widehat{\widehat{\mathbf{X}}}^{n+1,\top}\mathbf{X}^n$. Note that all quantities of rank $2r$ are denoted with one hat and all quantities of rank $4r$ with double hats.

For the $L$-step we write $L^n_{mk} = \sum_{\ell=1}^r S^n_{\ell m} V^n_{\ell k}$ and solve

$$
L^{n+1}_{pk} = \sum_{m=1}^r L^n_{mk} \sum_{j=1}^{N_x} X^n_{mj} \frac{\rho^n_j}{\rho^{n+1}_j} X^n_{pj} - \Delta t \sum_{m=1}^r v_k L^n_{mk} \sum_{i=1}^{N_x} X^n_{mi} \rho^n_i \sum_{j=1}^{N_x} D^x_{ij} \frac{1}{\rho^{n+1}_j} X^n_{pj}
$$
$$
+ \Delta t \frac{\Delta x}{2} \sum_{m=1}^r |v_k| L^n_{mk} \sum_{i=1}^{N_x} X^n_{mi} \rho^n_i \sum_{j=1}^{N_x} D^{xx}_{ij} \frac{1}{\rho^{n+1}_j} X^n_{pj} + \sigma \Delta t \sum_{j=1}^{N_x} X^n_{pj}. \tag{13d}
$$

This gives the updated matrix $\mathbf{L}^{n+1} = (L^{n+1}_{pk})$ to which together with the old basis $\mathbf{V}^n = (V^n_{\ell k})$ a QR-decomposition is applied, giving the new augmented basis $\widehat{\mathbf{V}}^{n+1} = (\widehat{V}^{n+1}_{\ell k})$.

In addition, we augment this basis according to

$$
\widehat{\widehat{\mathbf{V}}}^{n+1} = \mathrm{qr}\left([\widehat{\mathbf{V}}^{n+1}, \omega e^{\mathbf{v}^2/2} \widehat{\mathbf{V}}^{n+1}]\right) \tag{13e}
$$

leading to a new augmented basis $\widehat{\widehat{\mathbf{V}}}^{n+1} = (\widehat{\widehat{V}}^{n+1}_{\ell k})$ of rank $4r$, and we compute $\widehat{\widehat{\mathbf{N}}} = \widehat{\widehat{\mathbf{V}}}^{n+1,\top}\mathbf{V}^n$.

For the $S$-step we denote $\widetilde{S}^n_{m\ell} = \sum_{j,k=1}^r \widehat{\widehat{M}}_{mj} S^n_{jk} \widehat{\widehat{N}}_{\ell k}$ and insert the expressions $g^n_{jk} = \sum_{m,\ell=1}^{4r} \widehat{\widehat{X}}^{n+1}_{jm} \widetilde{S}^n_{m\ell} \widehat{\widehat{V}}^{n+1}_{k\ell}$ and $g^{n+1}_{jk} = \sum_{m,\ell=1}^{4r} \widehat{\widehat{X}}^{n+1}_{jm} \widehat{\widehat{S}}^{n+1}_{m\ell} \widehat{\widehat{V}}^{n+1}_{k\ell}$ into (12a). We multiply with $\widehat{\widehat{X}}^{n+1}_{jq} \widehat{\widehat{V}}^{n+1}_{kp}$ and sum over $j$ and $k$. This leads to

$$
\widehat{\widehat{S}}^{n+1}_{qp} = \frac{1}{1+\sigma\Delta t}\left( \sum_{j=1}^{N_x} \widehat{\widehat{X}}^{n+1}_{jq} \frac{\rho^n_j}{\rho^{n+1}_j} \sum_{m,\ell=1}^{4r} \widehat{\widehat{X}}^{n+1}_{jm} \widetilde{S}^n_{m\ell} \sum_{k=1}^{N_v} \widehat{\widehat{V}}^{n+1}_{k\ell} \widehat{\widehat{V}}^{n+1}_{kp} \right.
$$
$$
- \Delta t \sum_{j=1}^{N_x} \widehat{\widehat{X}}^{n+1}_{jq} \frac{1}{\rho^{n+1}_j} \sum_{i=1}^{N_x} D^x_{ji} \rho^n_i \sum_{m,\ell=1}^{4r} \widehat{\widehat{X}}^{n+1}_{im} \widetilde{S}^n_{m\ell} \sum_{k=1}^{N_v} \widehat{\widehat{V}}^{n+1}_{k\ell} v_k \widehat{\widehat{V}}^{n+1}_{kp}
$$
$$
+ \Delta t \frac{\Delta x}{2} \sum_{j=1}^{N_x} \widehat{\widehat{X}}^{n+1}_{jq} \frac{1}{\rho^{n+1}_j} \sum_{i=1}^{N_x} D^{xx}_{ji} \rho^n_i \sum_{m,\ell=1}^{4r} \widehat{\widehat{X}}^{n+1}_{im} \widetilde{S}^n_{m\ell} \sum_{k=1}^{N_v} \widehat{\widehat{V}}^{n+1}_{k\ell} |v_k| \widehat{\widehat{V}}^{n+1}_{kp}
$$
$$
\left. + \sigma\Delta t \sum_{j=1}^{N_x} \widehat{\widehat{X}}^{n+1}_{jq} \sum_{k=1}^{N_v} \widehat{\widehat{V}}^{n+1}_{kp} \right). \tag{13f}
$$

The last step consists in truncating the augmented quantities $\widehat{\widehat{\mathbf{X}}}^{n+1}$, $\widehat{\widehat{\mathbf{V}}}^{n+1}$ and $\widehat{\widehat{\mathbf{S}}}^{n+1}$ from rank $4r$ to a new rank $r_1$. We use a modification of the truncation strategy described in Section 5.1 that ensures that $\frac{1}{\sqrt{2\pi}} \sum_{\ell,m=1}^r \sum_{k=1}^{N_v} X^n_{jm} S^n_{m\ell} V^n_{k\ell} \omega_k e^{v_k^2/2} = 1$ stays valid in each time step and works as follows:

Let us denote $\mathbf{Z} \in \mathbb{R}^{N_v}$ being the vector with entries $Z_k = \frac{1}{\sqrt{2\pi}} \omega_k e^{v_k^2/2}$ and let $\mathbf{z} = \frac{\mathbf{Z}}{\|\mathbf{Z}\|_E}$, where $\|\cdot\|_E$ stands for the Euclidean norm. We then want to have

$$
\mathbf{1} = \widehat{\widehat{\mathbf{X}}}^{n+1} \widehat{\widehat{\mathbf{S}}}^{n+1} \widehat{\widehat{\mathbf{V}}}^{n+1,\top} \mathbf{Z} = \left( \widehat{\widehat{\mathbf{X}}}^{n+1} \widehat{\widehat{\mathbf{S}}}^{n+1} \widehat{\widehat{\mathbf{V}}}^{n+1,\top} \mathbf{z}\mathbf{z}^\top + \widehat{\widehat{\mathbf{X}}}^{n+1} \widehat{\widehat{\mathbf{S}}}^{n+1} \widehat{\widehat{\mathbf{V}}}^{n+1,\top} \left(\mathbf{I} - \mathbf{z}\mathbf{z}^\top\right) \right) \mathbf{Z} =: (\mathbf{H}_1 + \mathbf{H}_2)\mathbf{Z},
$$

with $\mathbf{H}_1 = \widehat{\widehat{\mathbf{X}}}^{n+1} \widehat{\widehat{\mathbf{S}}}^{n+1} \widehat{\widehat{\mathbf{V}}}^{n+1,\top} \mathbf{z}\mathbf{z}^\top$, $\mathbf{H}_2 = \widehat{\widehat{\mathbf{X}}}^{n+1} \widehat{\widehat{\mathbf{S}}}^{n+1} \widehat{\widehat{\mathbf{V}}}^{n+1,\top} \left(\mathbf{I} - \mathbf{z}\mathbf{z}^\top\right)$, $\mathbf{I} \in \mathbb{R}^{N_v \times N_v}$ denoting the identity matrix and $\mathbf{1} \in \mathbb{R}^{N_x}$ the vector containing ones at each entry. $\mathbf{H}_1$ is a matrix of rank 1. We determine its low-rank factors by a singular value decomposition such that $\mathbf{X}^{\mathbf{H}_1} \mathbf{S}^{\mathbf{H}_1} \mathbf{V}^{\mathbf{H}_1,\top} = \mathrm{svd}\left(\widehat{\widehat{\mathbf{S}}}^{n+1} \widehat{\widehat{\mathbf{V}}}^{n+1,\top} \mathbf{z}\mathbf{z}^\top\right)$ with $\mathbf{X}^{\mathbf{H}_1} \in \mathbb{R}^{4r}$, $\mathbf{S}^{\mathbf{H}_1} \in \mathbb{R}$, and $\mathbf{V}^{\mathbf{H}_1} \in \mathbb{R}^{N_v}$. For $\mathbf{H}_2$, it holds that $\mathbf{H}_2\mathbf{Z} = 0$. We apply the truncation strategy

from Section 5.1 to $\mathbf{H}_2$ and obtain $\mathbf{X}^*, \mathbf{S}^*$, and $\mathbf{V}^*$, which shall be of rank $\widetilde{r}_1$. Finally, we combine both parts by performing a QR-decomposition of

$$\mathbf{X}^{n+1}\mathbf{R}^1 = \left[\widehat{\widehat{\mathbf{X}}}^{n+1}\mathbf{X}^{\mathbf{H}_1}, \mathbf{X}^*\right], \quad \text{and} \quad \mathbf{V}^{n+1}\mathbf{R}^2 = \left[\mathbf{V}^{\mathbf{H}_1}, \mathbf{V}^*\right],$$

and setting

$$\mathbf{S}^{n+1} = \mathbf{R}^1 \begin{bmatrix} \mathbf{S}^{\mathbf{H}_1} & 0 \\ 0 & \mathbf{S}^* \end{bmatrix} \mathbf{R}^{2,\top}.$$

The new rank $r_1$ is then given by $r_1 = \widetilde{r}_1 + 1$. With the time-updated low-rank factors $\mathbf{X}^{n+1}$, $\mathbf{V}^{n+1}$ and $\mathbf{S}^{n+1}$, the updated low-rank approximation of the solution is $g_{jk}^{n+1} = \sum_{m,\ell=1}^{r_1} X_{jm}^{n+1} S_{m\ell}^{n+1} V_{k\ell}^{n+1}$. The steps of the proposed DLRA scheme are visualized in Figure 1.



Figure 1: Flowchart of the (simplified) stable DLRA scheme (13).

From the low-rank approximation of $g$ we can regain the full solution $f$ as follows.

**Definition 3** (Low-rank approximation of the full solution). *The DLRA approximation of the full solution $f$ of the linear Boltzmann-BGK equation in the fully discrete setting at time $t_n$ is given by $\mathbf{f}^n = (f_{jk}^n) \in \mathbb{R}^{N_x \times N_v}$ with entries*

$$f_{jk}^n = \frac{1}{\sqrt{2\pi}}\rho_j^n \sum_{m,\ell=1}^r X_{jm}^n S_{m\ell}^n V_{k\ell}^n e^{-v_k^2/2}.$$

14

*5.3. Stability of the proposed low-rank scheme*

We can then show that algorithm (13) is numerically stable.

**Theorem 3.** *Under the time step restriction* $\max_k(|v_k|)\Delta t \leq \Delta x$, *the fully discrete DLRA scheme (13) is numerically stable in the $\mathscr{H}$-norm, i.e*

$$\|\mathbf{f}^{n+1}\|_{\mathscr{H}}^2 \leq \|\mathbf{f}^n\|_{\mathscr{H}}^2.$$

*Proof.* We start from the $S$-step in (13f), multiply it with $\widehat{\widehat{X}}_{\alpha q}^{n+1}$ and $\widehat{\widehat{V}}_{\beta p}^{n+1}$ and sum over $q$ and $p$. For simplicity of notation, we introduce the projections $P_{\alpha j}^X = \sum_{q=1}^{4r} \widehat{\widehat{X}}_{\alpha q}^{n+1} \widehat{\widehat{X}}_{jq}^{n+1}$ and $P_{k\beta}^V = \sum_{p=1}^{4r} \widehat{\widehat{V}}_{kp}^* \widehat{\widehat{V}}_{\beta p}^*$. This gives

$$g_{\alpha\beta}^{n+1} = \frac{1}{1+\sigma\Delta t}\left( \sum_{j=1}^{N_x}\sum_{k=1}^{N_v} \frac{\rho_j^n}{\rho_j^{n+1}} g_{jk}^n P_{\alpha j}^X P_{k\beta}^V - \Delta t \sum_{i,j=1}^{N_x}\sum_{k=1}^{N_v} \frac{1}{\rho_j^{n+1}} D_{ji}^x \rho_i^n g_{ik}^n v_k P_{\alpha j}^X P_{k\beta}^V \right.$$

$$\left. + \Delta t \frac{\Delta x}{2} \sum_{i,j=1}^{N_x}\sum_{k=1}^{N_v} \frac{1}{\rho_j^{n+1}} D_{ji}^{xx} \rho_i^n g_{ik}^n |v_k| P_{\alpha j}^X P_{k\beta}^V + \sigma\Delta t \sum_{j=1}^{N_x}\sum_{k=1}^{N_v} P_{\alpha j}^X P_{k\beta}^V \right).$$

Multiplication with $\frac{2}{\sqrt{2\pi}}\left(\rho_\alpha^{n+1}\right)^2 g_{\alpha\beta}^{n+1} \omega_\beta e^{v_\beta^2/2}(1+\sigma\Delta t)$ and summation over $\alpha$ and $\beta$ leads to

$$\frac{2}{\sqrt{2\pi}} \sum_{\alpha=1}^{N_x}\sum_{\beta=1}^{N_v} \left(\rho_\alpha^{n+1} g_{\alpha\beta}^{n+1}\right)^2 \omega_\beta e^{v_\beta^2/2}\left(1+\sigma\Delta t\right)$$

$$= \frac{2}{\sqrt{2\pi}} \sum_{j=1}^{N_x}\sum_{k=1}^{N_v} \frac{\rho_j^n}{\rho_j^{n+1}} g_{jk}^n \sum_{\alpha=1}^{N_x} P_{\alpha j}^X g_{\alpha\beta}^{n+1} \left(\rho_\alpha^{n+1}\right)^2 \sum_{\beta=1}^{N_v} P_{k\beta}^V \omega_\beta e^{v_\beta^2/2}$$

$$- \frac{2\Delta t}{\sqrt{2\pi}} \sum_{i,j=1}^{N_x}\sum_{k=1}^{N_v} \frac{1}{\rho_j^{n+1}} D_{ji}^x \rho_i^n g_{ik}^n v_k \sum_{\alpha=1}^{N_x} P_{\alpha j}^X g_{\alpha\beta}^{n+1} \left(\rho_\alpha^{n+1}\right)^2 \sum_{\beta=1}^{N_v} P_{k\beta}^V \omega_\beta e^{v_\beta^2/2}$$

$$+ \Delta t \frac{\Delta x\sqrt{2\pi}}{\sum} \sum_{i,j=1}^{N_x}\sum_{k=1}^{N_v} \frac{1}{\rho_j^{n+1}} D_{ji}^{xx} \rho_i^n g_{ik}^n |v_k| \sum_{i,j,\alpha=1}^{N_x} P_{\alpha j}^X g_{\alpha\beta}^{n+1} \left(\rho_\alpha^{n+1}\right)^2 \sum_{\beta=1}^{N_v} P_{k\beta}^V \omega_\beta e^{v_\beta^2/2}$$

$$+ \frac{2\sigma\Delta t}{\sqrt{2\pi}} \sum_{j,\alpha=1}^{N_x}\sum_{k=1}^{N_v} P_{\alpha j}^X g_{\alpha\beta}^{n+1} \left(\rho_\alpha^{n+1}\right)^2 \sum_{\beta=1}^{N_v} P_{k\beta}^V \omega_\beta e^{v_\beta^2/2}.$$

We now use the fact that we have augmented the bases in (13c) and (13e) such that

$$\sum_{\alpha=1}^{N_x} P_{\alpha j}^X g_{\alpha\beta}^{n+1}(\rho_\alpha^{n+1})^2 = g_{j\beta}^{n+1}(\rho_j^{n+1})^2 \quad \text{and} \quad \sum_{\beta=1}^{N_v} P_{k\beta}^V g_{j\beta}^{n+1} \omega_\beta e^{v_\beta^2/2} = g_{jk}^{n+1} \omega_k e^{v_k^2/2}$$

holds. We insert these relations and, to be consistent in notation, change the summation indices on the

left-hand side from $\alpha$ to $j$ and $\beta$ to $k$ giving

$$\frac{2}{\sqrt{2\pi}}\sum_{j=1}^{N_x}\sum_{k=1}^{N_v}\left(\rho_j^{n+1}g_{jk}^{n+1}\right)^2\omega_k e^{v_k^2/2}\left(1+\sigma\Delta t\right) = \frac{2}{\sqrt{2\pi}}\sum_{j=1}^{N_x}\sum_{k=1}^{N_v}\rho_j^n g_{jk}^n\rho_j^{n+1}g_{jk}^{n+1}\omega_k e^{v_k^2/2}$$

$$-\frac{2\Delta t}{\sqrt{2\pi}}\sum_{i,j=1}^{N_x}\sum_{k=1}^{N_v}\rho_j^{n+1}g_{jk}^{n+1}D_{ji}^x\rho_i^n g_{ik}^n v_k\omega_k e^{v_k^2/2}$$

$$+\Delta t\frac{\Delta x}{\sqrt{2\pi}}\sum_{i,j=1}^{N_x}\sum_{k=1}^{N_v}\rho_j^{n+1}g_{jk}^{n+1}D_{ji}^{xx}\rho_i^n g_{ik}^n|v_k|\omega_k e^{v_k^2/2}$$

$$+\frac{2\sigma\Delta t}{\sqrt{2\pi}}\sum_{j=1}^{N_x}\sum_{k=1}^{N_v}\left(\rho_j^{n+1}\right)^2 g_{jk}^{n+1}\omega_k e^{v_k^2/2}.$$

Analogously to the proof of Theorem 2 and, as the truncation step is designed not to alter these expressions, we can conclude that the proposed DLRA scheme is numerically stable under the time step restriction $\max_k(|v_k|)\Delta t \leq \Delta x$. □

## 6. Numerical results

To validate the theoretical considerations and the proposed DLRA scheme, numerical results for different test examples in 1D as well as in 2D are given in this section.

### 6.1. 1D Plane source

We start with a one-dimensional analogue to the plane source problem, which is a common test case for radiative transfer [2, 8, 26, 31], and compare the solution of the full equations (8) to the solution obtained by the DLRA scheme given in (13). The spatial domain shall be set to $D = [-10, 10]$. As initial conditions we choose the density $\rho$ to be a cutoff Gaussian

$$\rho(t=0,x) = \max\left(10^{-4}, \frac{1}{\sqrt{2\pi\sigma_{\mathrm{IC}}^2}}\exp\left(-\frac{x^2}{2\sigma_{\mathrm{IC}}^2}\right)\right),$$

where $\sigma_{\mathrm{IC}} = 0.3$ denotes a constant deviation, and the function $g$ to be constant in space and velocity, i.e. $g(t=0,x,v) = 1$. We consider relatively large collisionality by choosing $\sigma = 10$. Computations are started with an initial rank of $r = 20$. As computational parameters we use $N_x = 1000$ grid points in the spatial as well as $N_v = 500$ grid points in the velocity domain. Due to this choice, we obtain $\max_k(|v_k|) \approx 31.05$, which shall be adjusted to the next larger integer value such that the time step size is determined by $\Delta t = \mathrm{CFL}\cdot\frac{\Delta x}{32}$ with $\mathrm{CFL} = 0.99$, according to the corresponding CFL condition. Practical implementations now show that the basis augmentations to rank $4r$ in in (13c) and (13e) needed for the theoretical proof of the numerical stability may not be necessary for numerical examples and that the standard basis augmentations to rank $2r$ provide similar solutions while being significantly faster. For this reason, we propose to leave out the basis augmentations (13c) and (13e) in practical applications. In this case, all quantities with double hats related to rank $4r$ reduce to quantities of rank $2r$ with one single hat. The simplified scheme with rank $2r$ is also visualized (in brackets) in the flowchart of Figure 1. In Figure 2, we now compare the results for the full solution $f(t,x,v)$, computed with the full solver (Mg (reference)), the DLRA scheme with rank $2r$ (Mg DLRA) and the basis augmented DLRA scheme with rank $4r$ (Mg DLRA BasisAug) at different times up to $t_{\mathrm{End}} = 8$. We observe that the reduced as well as the augmented DLRA algorithm capture the main characteristics of the full reference Mg system. This is also true for the computational results for the density $\rho(t,x)$, displayed in Figure 3. Figure 4 shows the evolution of the rank, which for a chosen tolerance parameter of $\vartheta = 10^{-5}\|\boldsymbol{\Sigma}\|_2$ increases up to $r = 76$ before it significantly reduces over time. Note that the evolution of the rank for the reduced as well as for the basis augmented algorithm show good agreement

Figure 2: Numerical results for the solution $f(t, x, v)$ of the 1D plane source analogue at time $t = 0$ (first column), $t = 2$ (second column), $t = 4$ (third column), and $t = 6$ (fourth column), computed with the full solver (Mg (reference)) (first row), the reduced DLRA scheme (Mg DLRA) (second row) and the basis augmented DLRA scheme (Mg DLRA BasisAug) (third row).

as the new rank is displayed after the corresponding truncation step. Further, the behavior of the norm $\|\mathbf{f}\|_{\mathscr{H}}^2$ is depicted. As expected, it decreases smoothly over time for all considered systems. Additionally, we display the quantities $\kappa^+ := \max_j \left( \frac{1}{\sqrt{2\pi}} \sum_{k=1}^{N_v} g_{jk} \omega_k e^{v_k^2/2} \right)$ and $\kappa^- := \min_j \left( \frac{1}{\sqrt{2\pi}} \sum_{k=1}^{N_v} g_{jk} \omega_k e^{v_k^2/2} \right)$. It is essential that they are equal to 1, which for the low-rank schemes is ensured by the adjusted truncation step. It can be observed that this property is fulfilled up to order $\mathcal{O}\left(10^{-11}\right)$.

### 6.2. 2D Plane source

In higher dimensions, the computational advantages of DLRA schemes are enhanced. For this reason, we give some two-dimensional test examples starting with the two-dimensional version of the plane source problem considered in the previous section. The corresponding two-dimensional set of equations becomes

$$\partial_t g(t, \mathbf{x}, \mathbf{v}) = -\frac{\mathbf{v}}{\rho(t, \mathbf{x})} \cdot \nabla_{\mathbf{x}} \left( \rho(t, \mathbf{x}) g(t, \mathbf{x}, \mathbf{v}) \right) + \sigma \left( 1 - g(t, \mathbf{x}, \mathbf{v}) \right) - \frac{g(t, \mathbf{x}, \mathbf{v})}{\rho(t, \mathbf{x})} \partial_t \rho(t, \mathbf{x}),$$

$$\partial_t \rho(t, \mathbf{x}) = -\frac{1}{2\pi} \nabla_{\mathbf{x}} \cdot \int \rho(t, \mathbf{x}) g(t, \mathbf{x}, \mathbf{v}) \mathbf{v} e^{-|\mathbf{v}|^2/2} \mathrm{d}\mathbf{v},$$

where $\mathbf{x} = (x_1, x_2) \in D \subset \mathbb{R}^2$ and $\mathbf{v} = (v_1, v_2) \in \mathbb{R}^2$. We compare the solution of this system with the solution of the DLRA scheme for which the extension to the two-dimensional setting is straightforward. For this test example we choose the spatial domain $D = [-3, 3] \times [-3, 3]$ and prescribe the initial condition for

17

Figure 3: Numerical results for the density $\rho(t,x)$ of the 1D plane source analogue at time $t = 0$, $t = 2$, $t = 4$, and $t = 6$, computed with the full solver (Mg (reference)), the reduced DLRA scheme (Mg DLRA) and the basis augmented DLRA scheme (Mg DLRA BasisAug).



Figure 4: Left: Evolution of the rank in time for the 1D plane source analogue for the reduced DLRA scheme (Mg DLRA) and the basis augmented DLRA scheme (Mg DLRA BasisAug). Middle: Evolution of the $\mathscr{H}$-norm in time for the full solver (Mg (reference)), the reduced DLRA scheme (Mg DLRA) and the basis augmented DLRA scheme (Mg DLRA BasisAug). Right: Evolution of $\kappa^{\pm}$ in time for the full solver (Mg (reference)), the reduced DLRA scheme (Mg DLRA) and the basis augmented DLRA scheme (Mg DLRA BasisAug). The red line has the constant value 1. The deviations of the DLRA schemes from 1 are of order $\mathcal{O}\left(10^{-11}\right)$.

the density by

$$\rho(t=0,\mathbf{x}) = \frac{1}{4\pi} \max\left(10^{-1}, \frac{10^2}{4\pi\sigma_{\mathrm{IC}}^2} \exp\left(-\frac{|\mathbf{x}|^2}{4\sigma_{\mathrm{IC}}^2}\right)\right),$$

with a constant deviation of $\sigma_{\mathrm{IC}} = 0.3$. The function $g$ shall be set to $g(t=0,\mathbf{x},\mathbf{v}) = 1$ and the collision coefficient to $\sigma = 100$. We allow an initial rank of $r = 20$. Computations are performed on a spatial grid with $N_{x_1} = 128$ grid points in $x_1$ and $N_{x_2} = 128$ grid points in $x_2$. For the velocity grid we choose $N_{v_1} = 32$ grid

points in $v_1$ and $N_{v_2} = 32$ grid points in $v_2$. Due to this choice, we obtain $\max_k (|\mathbf{v}_k|) \approx 10.08$, which shall be adjusted to the next larger integer value such that the time step size is determined by $\Delta t = \text{CFL} \cdot \frac{\Delta x}{11}$ with a CFL number of CFL $= 0.7$. Figure 5 compares the density $\rho(t, x)$ at different times up to $t_{\text{End}} = 3.0$, computed with the full solver (Mg (reference)) and the reduced DLRA scheme with rank $2r$ (Mg DLRA). Note that we refrain from computations with the basis augmented $4r$ scheme as in two space and velocity dimensions this would lead to extremely increased computational costs. We observe that the solution of the DLRA scheme matches the solution of the full system. To determine the evolution of the rank, we use a tolerance pa-



Figure 5: Numerical results for the density $\rho(t, \mathbf{x})$ of the 2D plane source analogue at time $t = 0$ (first column), $t = 1$ (second column), $t = 2$ (third column, and $t = 3$ (fourth column), computed with the full solver (Mg (reference)) (first row) and the reduced DLRA scheme (Mg DLRA) (second row).

rameter of $\vartheta = 10^{-5} \|\mathbf{\Sigma}\|_2$. In Figure 6, we observe an increasing up to $r = 73$ before it decreases continuously over time. Further, the norm $\|\mathbf{f}\|_{\mathscr{H}}^2$ is displayed. It decreases smoothly over time for all considered systems. In addition, we plot the quantities $\kappa^+ := \max_j \left( \frac{1}{2\pi} \sum_{k=1}^{N_{v_1}} \sum_{\ell=1}^{N_{v_2}} g(t, \mathbf{x}_j, v_k^1, v_\ell^2) \omega_k^1 \omega_\ell^2 e^{\left(v_k^1\right)^2/2} e^{\left(v_k^2\right)^2/2} \right)$ and $\kappa^- := \min_j \left( \frac{1}{2\pi} \sum_{k=1}^{N_{v_1}} \sum_{\ell=1}^{N_{v_2}} g(t, \mathbf{x}_j, v_k^1, v_\ell^2) \omega_k^1 \omega_\ell^2 e^{\left(v_k^1\right)^2/2} e^{\left(v_k^2\right)^2/2} \right)$, which are equal to 1 up to order $\mathcal{O}\left(10^{-3}\right)$. For this setup, the running time of the DLRA scheme compared to the full solver is clearly faster. It reduces



Figure 6: Left: Evolution of the rank in time for the 2D plane source analogue for the reduced DLRA scheme (Mg DLRA). Middle: Evolution of the $\mathscr{H}$-norm in time for the full solver (Mg (reference)) and the reduced DLRA scheme (Mg DLRA). Right: Evolution of $\kappa^\pm$ in time for the full solver (Mg (reference)) and the reduced DLRA scheme (Mg DLRA).

by a factor of approximately 2 from 3954 seconds to 1926 seconds, confirming the computational advantages of the DLRA scheme.

## 6.3. 2D Beam

As a second two-dimensional test example we consider a beam in the spatial domain $D = [-5, 5] \times [-5, 5]$ starting at $(0,0)$ in the middle of the spatial plane and moving to the bottom left. The initial values are given by

$$\rho(t = 0, \mathbf{x}) = \frac{1}{4\pi} \max\left(10^{-1}, \frac{10^2}{4\pi\sigma_{\text{IC}}^2} \exp\left(-\frac{|\mathbf{x}|^2}{4\sigma_{\text{IC}}^2}\right)\right),$$

$$g(t = 0, \mathbf{x}, \mathbf{v}) = C\frac{10^2}{4\pi\sigma_{\text{IC}}^2} \exp\left(-\frac{|\mathbf{v} - \mathbf{v}_{\text{beam}}|^2}{4\sigma_{\text{IC}}^2}\right),$$

where $\sigma_{\text{IC}} = 0.01$ and $C$ is a normalization constant such that $\frac{1}{2\pi}\int g(t = 0, \mathbf{x}, \mathbf{v})e^{-|\mathbf{v}|^2/2}\mathrm{d}\mathbf{v} = \mathbf{1}$ holds. The beam velocity $\mathbf{v}_{\text{beam}}$ is set to $\mathbf{v}_{\text{beam}} = \begin{pmatrix} -1 \\ -1 \end{pmatrix}$ and the collisionality to $\sigma = 1.5$. The initial rank is given as $r = 20$. Computations are performed on a spatial grid with $N_{x_1} = 128$ grid points in $x_1$ and $N_{x_2} = 128$ grid points in $x_2$. For the velocity grid we choose $N_{v_1} = 32$ grid points in $v_1$ and $N_{v_2} = 32$ grid points in $v_2$. As before, this leads to a time step size of $\Delta t = \text{CFL} \cdot \frac{\Delta x}{11}$ with a CFL number of CFL $= 0.7$. Figure 7 compares the density $\rho(t, x)$ at different times up to $t_{\text{End}} = 3.0$, computed with the full solver (Mg (reference)) and the reduced DLRA scheme with rank $2r$ (Mg DLRA). At all displayed time steps the DLRA solution resembles the solution of the full system. In Figure 8, the evolution of the rank in time is shown. We use a tolerance param-



Figure 7: Numerical results for the density $\rho(t, \mathbf{x})$ of the 2D beam test problem at time $t = 0$ (first column), $t = 1$ (second column), $t = 2$ (third column, and $t = 3$ (fourth column), computed with the full solver (Mg reference)) (first row) and the reduced DLRA scheme (Mg DLRA) (second row).

eter of $\vartheta = 10^{-4}\|\boldsymbol{\Sigma}\|_2$ and allow a maximal rank of 200. Due to the choice of $\sigma$ the solution of the problem is not low-rank. For this reason, we observe an increasing of the rank up to the maximal allowed value. Also, the norm $\|\mathbf{f}\|_{\mathcal{H}}^2$ is depicted over time. It decreases continuously, matching our theoretical considerations. Further, it can be observed that the quantities $\kappa^+ := \max_j \left(\frac{1}{2\pi} \sum_{k=1}^{N_{v_1}} \sum_{\ell=1}^{N_{v_2}} g(t, \mathbf{x}_j, v_k^1, v_\ell^2)\omega_k^1\omega_\ell^2 e^{\left(v_k^1\right)^2/2} e^{\left(v_k^2\right)^2/2}\right)$ and $\kappa^- := \min_j \left(\frac{1}{2\pi} \sum_{k=1}^{N_{v_1}} \sum_{\ell=1}^{N_{v_2}} g(t, \mathbf{x}_j, v_k^1, v_\ell^2)\omega_k^1\omega_\ell^2 e^{\left(v_k^1\right)^2/2} e^{\left(v_k^2\right)^2/2}\right)$ are equal to 1 up to order $\mathcal{O}\left(10^{-12}\right)$. Due to the high rank, the running time of the DLRA scheme compared with the full solver is slightly increased. This example illustrates the relation between the choice of $\sigma$ and the low-rank structure of the solution. It is expected that for larger values of $\sigma$ the solution becomes low-rank and hence the computational benefits of the DLRA scheme are enhanced.

Figure 8: Left: Evolution of the rank in time for the 2D beam test problem for the reduced DLRA scheme (Mg DLRA). The rank increases up to the maximal allowed value of $r = 200$. Middle: Evolution of the $\mathscr{H}$-norm in time for the full solver (Mg (reference)) and the reduced DLRA scheme (Mg DLRA). Right: Evolution of $\kappa^{\pm}$ in time for the full solver (Mg (reference)) and the reduced DLRA scheme (Mg DLRA). The red line has the constant value 1. The deviation of the DLRA scheme from 1 is of order $\mathcal{O}\left(10^{-12}\right)$.

## 7. Conclusion and outlook

We have derived a multiplicative DLRA discretization for the linear Boltzmann-BGK problem that in contrast to another presented naive discretization is numerically stable. To show this, we have conducted a stability analysis leading to a concrete hyperbolic CFL condition. In addition, numerical examples in 1D and 2D confirm the stability, accuracy and efficiency of the proposed DLRA scheme.

The insights gained from this article can be helpful for future work as the employed multiplicative splitting is attached to the investigation of more complicated equations, e.g. the non-linear Boltzmann-BGK equation treated in [14]. However, a direct transition of knowledge is hardly possible as for the non-linear case most of the theoretical concepts applied here are not available, making the analysis much more difficult to consider. Further, in the non-linear case a discretization of the conservative form of the equations, i.e. by not splitting up the term $\partial_x (Mg)$, is possible but cannot be efficiently implemented as the Maxwellian $M$ is generally not low-rank. Even though, we propose to reconsider the chosen discretization in [14] in terms of stabilization, which may allow for a larger time step size leading to an even more efficient DLRA algorithm.

### Acknowledgements

# References

[1] F. Achleitner, A. Arnold, and E. A. Carlen. On linear hypocoercive BGK models. In P. Gonçalves and A. J. Soares, editors, *From Particle Systems to Partial Differential Equations III. Springer Proceedings in Mathematics & Statistics*, volume 162, pages 1–37, Cham, 2016. Springer.

[2] L. Baumann, L. Einkemmer, C. Klingenberg, and J. Kusch. Energy stable and conservative dynamical low-rank approximation for the Su-Olson problem. *SIAM Journal on Scientific Computing*, 46(2):B137–B158, 2024.

[3] M. Bessemoulin-Chatard, M. Herda, and T. Rey. Hypocoercivity and diffusion limit of a finite volume scheme for linear kinetic equations. *Mathematics of Computation*, 89(323):1093–1133, 2020.

[4] E. P. Bhatnagar, T. Gross, and M. Krook. A model for collision processes in gases. I. Small amplitude processes in charged and neutral one-component systems. *Physical Review*, 94(3):511–525, 1954.

[5] J. A. Cañizo, C. Cao, J. Evans, and H. Yoldaş. Hypocoercivity of linear kinetic equations via Harris's Theorem. *Kinetic and Related Models*, 13(1):97–128, 2020.

[6] C. Cercignani. *The Boltzmann Equation and Its Applications*. Springer, New York, 1988.

[7] G. Ceruti, L. Einkemmer, J. Kusch, and C. Lubich. A robust second-order low-rank BUG integrator based on the midpoint rule. *arXiv preprint arXiv:2402.08607*, 2024.

[8] G. Ceruti, J. Kusch, and C. Lubich. A rank-adaptive robust integrator for dynamical low-rank approximation. *BIT Numerical Mathematics*, 62:1149–1174, 2022.

[9] G. Ceruti, J. Kusch, and C. Lubich. A parallel rank-adaptive integrator for dynamical low-rank approximation. *SIAM Journal on Scientific Computing*, 46(3):B205–B228, 2024.

[10] G. Ceruti and C. Lubich. An unconventional robust integrator for dynamical low-rank approximation. *BIT Numerical Mathematics*, 62(1):23–44, 2022.

[11] Z. Ding, L. Einkemmer, and Q. Li. Dynamical low-rank integrator for the linear Boltzmann equation: error analysis in the diffusion limit. *SIAM Journal on Numerical Analysis*, 59(4):2254–2285, 2021.

[12] L. Einkemmer. A low-rank algorithm for weakly compressible flow. *SIAM Journal on Scientific Computing*, 41(5):A2795–A2814, 2019.

[13] L. Einkemmer, J. Hu, and J. Kusch. Asymptotic-preserving and energy stable dynamical low-rank approximation. *SIAM Journal on Numerical Analysis*, 62(1):73–92, 2024.

[14] L. Einkemmer, J. Hu, and L. Ying. An efficient dynamical low-rank algorithm for the Boltzmann-BGK equation close to the compressible viscous flow regime. *SIAM Journal on Scientific Computing*, 43(5):B1057–B1080, 2021.

[15] L. Einkemmer and I. Joseph. A mass, momentum, and energy conservative dynamical low-rank scheme for the Vlasov equation. *Journal of Computational Physics*, 443:110493, 2021.

[16] L. Einkemmer and C. Lubich. A low-rank projector-splitting integrator for the Vlasov-Poisson equation. *SIAM Journal on Scientific Computing*, 40(5):B1330–B1360, 2018.

[17] L. Einkemmer, J. Mangott, and M. Prugger. A low-rank complexity reduction algorithm for the high-dimensional kinetic chemical master equation. *Journal of Computational Physics*, 503:112827, 2024.

[18] L. Einkemmer, A. Ostermann, and C. Piazzola. A low-rank projector-splitting integrator for the Vlasov–Maxwell equations with divergence correction. *Journal of Computational Physics*, 403:109063, 2020.

[19] L. Einkemmer, A. Ostermann, and C. Scalone. A robust and conservative dynamical low-rank algorithm. *Journal of Computational Physics*, 484:112060, 2023.

[20] J. Evans. Hypocoercivity in Phi-entropy for the linear relaxation Boltzmann equation on the torus. *SIAM Journal on Mathematical Analysis*, 53(2):1357–1378, 2021.

[21] E. Hairer and G. Wanner. *Solving Ordinary Differential Equations II: Stiff and Differential-Algebraic Problems*. Springer Berlin, Heidelberg, 1996.

[22] J. Hu and Y. Wang. An adaptive dynamical low rank method for the nonlinear Boltzmann equation. *Journal of Scientific Computing*, 92:75, 2022.

[23] O. Koch and C. Lubich. Dynamical low-rank approximation. *SIAM Journal on Matrix Analysis and Applications*, 29(2):434–454, 2007.

[24] K. Kormann and E. Sonnendrücker. Sparse grids for the Vlasov–Poisson equation. In J. Garcke and D. Pflüger, editors, *Sparse Grids and Applications – Stuttgart 2014. Lecture Notes in Computational Science and Engineering*, volume 109, pages 163–190, Cham, 2016. Springer.

[25] J. Kusch. Second-order robust parallel integrators for dynamical low-rank approximation. *arXiv preprint arXiv:2403.02834*, 2024.

[26] J. Kusch, L. Einkemmer, and G. Ceruti. On the stability of robust dynamical low-rank approximations for hyperbolic problems. *SIAM Journal on Scientific Computing*, 45(1):A1–A24, 2023.

[27] J. Kusch and P. Stammer. A robust collision source method for rank adaptive dynamical low-rank approximation in radiation therapy. *ESAIM: M2AN*, 57(2):865–891, 2023.

[28] R. J. LeVeque. *Finite Volume Methods for Hyperbolic Problems*. Cambridge University Press, Cambridge, 2002.

[29] C. Lubich and I. V. Oseledets. A projector-splitting integrator for dynamical low-rank approximation. *BIT Numerical Mathematics*, 54(1):171–188, 2014.

[30] C. Patwardhan, M. Frank, and J. Kusch. Asymptotic-preserving and energy stable dynamical low-rank approximation for thermal radiative transfer equations. *arXiv preprint arXiv:2402.16746*, 2024.

[31] Z. Peng, R. G. McClarren, and M. Frank. A low-rank method for two-dimensional time-dependent radiation transport calculations. *Journal of Computational Physics*, 421:109735, 2020.

[32] M. Prugger, L. Einkemmer, and C. F. Lopez. A dynamical low-rank approach to solve the chemical master equation for

biological reaction networks. *Journal of Computational Physics*, 489:112250, 2023.

[33] H. Struchtrup. *Macroscopic Transport Equations for Rarefied Gas Flows. Approximation Methods in Kinetic Theory.* Springer, Berlin, Heidelberg, 2005.

[34] P. Yin, E. Endeve, C. D. Hauck, and S. R. Schnake. A semi-implicit dynamical low-rank discontinuous Galerkin method for space homogeneous kinetic equations. Part I: emission and absorption. *arXiv preprint arXiv:2308.05914*, 2023.