# An energy stable and conservative multiplicative dynamical low-rank discretization for the Su-Olson problem

**Lena Baumann, Lukas Einkemmer, Christian Klingenberg, Jonas Kusch**

**Abstract** Computing numerical solutions of the thermal radiative transfer equations on a finely resolved grid can be costly due to high computational and memory requirements. A numerical reduced order method that has recently been applied to a wide variety of kinetic partial differential equations is the concept of dynamical low-rank approximation (DLRA). In this paper, we consider the thermal radiative transfer equations with Su-Olson closure, leading to a linearized kinetic model. For the conducted theoretical and practical considerations we use a multiplicative splitting of the distribution function that poses additional challenges in finding an energy stable discretization and deriving a hyperbolic Courant-Friedrichs-Lewy (CFL) condition. We propose such an energy stable DLRA scheme that makes use of the augmented basis update & Galerkin integrator. This integrator allows for additional basis augmentations, enabling us to give a mathematically rigorous proof of energy stability and local mass conservation. Numerical examples confirm the derived properties and show the computational advantages of the DLRA scheme compared to a numerical solution of the full system of equations.

Lena Baumann
Department of Mathematics, University of Würzburg, Emil-Fischer-Strasse 40, DE-97074 Würzburg, Germany. E-mail: lena.baumann@uni-wuerzburg.de
Lukas Einkemmer
Numerical Analysis and Scientific Computing, University of Innsbruck, Technikerstrasse 13, A-6020 Innsbruck, Austria. E-mail: lukas.einkemmer@uibk.ac.at
Christian Klingenberg
Department of Mathematics, University of Würzburg, Emil-Fischer-Strasse 40, DE-97074 Würzburg, Germany. E-mail: christian.klingenberg@uni-wuerzburg.de
Jonas Kusch
Scientific Computing, Norwegian University of Life Sciences, Drøbakveien 31, NO-1433 Ås, Norway. E-mail: jonas.kusch@nmbu.no

# 1 Introduction

Thermal radiative transfer problems are a class of kinetic transport equations modeling the movement of particles that interact with a background medium, for instance by scattering or absorption. By this interaction, the background medium can heat up and itself emit new particles, enforcing the exchange of energy between particles and the background material. This process is described by two coupled equations, one for the particle density $f(t,x,\mu)$ and one for the internal energy $e(t,x)$ of the material. The variable $t \in \mathbb{R}^+$ denotes time, $x \in D \subset \mathbb{R}$ stands for the spatial and $\mu \in [-1,1]$ for the directional variable. The numerical solution of this set of equations is challenging as its high dimensionality requires enormous computational and memory costs. To overcome these problems, the concept of dynamical low-rank approximation (DLRA) [26] can be used. It is a numerical reduced order method providing accurate and efficient approximations of the solution of kinetic partial differential equations and has already been applied in various fields of research. Recent work for instance has been published on radiation transport [1,39,11,42,40,41,29, 46], radiation therapy [28], plasma physics [15,18,13,19], chemical kinetics [44,17] and Boltzmann type transport problems [2,9,12,10,24]. The main idea of DLRA consists in approximating the distribution function as

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

where $\{X_i : i = 1,..,r\}$ are the orthonormal basis functions in space and $\{V_j : j = 1,..,r\}$ the orthonormal basis functions in direction. This splitting can be understood as a continuous analogue to the singular value decomposition of a matrix, explaining why $\mathbf{S} = (S_{ij}) \in \mathbb{R}^{r \times r}$ is called the coefficient or coupling matrix and $r$ the rank of this approximation. The time evolution of the low-rank factors is then determined by a projection of the equation onto the tangent space of the low-rank solution manifold. Integrators that ensure that the solution stays on the low-rank manifold while being robust to small singular values (otherwise this may lead to enormous restrictions on the choice of the time step size [25]) are the *projector-splitting* [32], the *(augmented) basis update & Galerkin* (BUG) [7,5], and the *parallel integrator* [6]. Extensions to schemes with proven second-order robust error bounds are available for the augmented BUG [4] as well as for the parallel integrator [27].

A challenge for constructing a stable dynamical low-rank scheme is to derive a suitable discretization of the system. To account for the angular dependence of the solution, we consider a modal representation that makes use of a finite expansion of the particle density in terms of spherical harmonics ($P_N$). This approach is explained in [3] and for instance used in [1,37,34,35,30,47,36,8]. Also the spatial discretization has to be chosen carefully. In [12] it is shown that for an efficient dynamical low-rank scheme for the isothermal Boltzmann-BGK equation it is useful to consider a multiplicative splitting of the distribution function. To transfer knowledge to such systems, we decide on a multiplicative splitting of the particle density $f$ also for the considered problem. In this case, it is per se not clear how to deal with the spatial derivatives. Recent work on this topic has been published for the linear Boltzmann-BGK equation in [2]. For the time discretization the potentially stiff opacity term has to be taken into account, leading to a coupled-implicit scheme similar to the one treated in [1], which again is complicated to solve.

In this paper we propose an energy stable multiplicative dynamical low-rank discretization for a linearized internal energy model called the Su-Olson problem. The main novelties are:

- *A multiplicative splitting of the distribution function:* Based on the insights from [12], we consider a multiplicative splitting of the distribution function. Similar to [2], we show that the spatial discretization has to be carefully chosen in order to obtain an energy stable numerical scheme.
- *An energy stable numerical scheme with rigorous mathematical proofs:* We show that the derived DLRA scheme is energy stable and give rigorous mathematical proofs, enabling us to deduce a classic hyperbolic CFL condition. This allows to compute up to a maximal time step size of $\Delta t = \text{CFL} \cdot \Delta x$, which again enhances the performance of the algorithm.
- *A mass conservative augmented integrator:* We make use of the augmented BUG integrator from [5], leading to a basis augmentation in two substeps of the numerical scheme. As this integrator allows for further modifications, we include additional basis augmentation steps that ensure the exactness of the projection operators needed for the theoretical proof of energy stability as well as the local conservation of mass.
- *A series of test examples:* A series of test examples compares the results of the low-rank discretization with the solution of the full system. This validates the derived properties and shows the accuracy and the efficiency of the proposed DLRA method.

The structure of the paper is as follows: After the introduction in Section 1, we provide background information on the thermal radiative transfer equations in Section 2. We explain the considered multiplicative structure and derive two possible systems that in the continuous setting are equivalent. In Section 3, we discretize both systems in angle, space and time. Section 4 is devoted to the subject of energy stability. We show that the advection form of the multiplicative Su-Olson problem is generally not numerically stable in the sense of von Neumann, whereas for the conservative form a hyperbolic CFL condition is derived, under which an energy estimate can be given. Section 5 first provides an overview of the method of dynamical low-rank approximation. Then, an energy stable DLRA scheme is derived. In addition, mass conservation can be shown, when using a suitable truncation strategy. The numerical results in Section 6 confirm the derived properties, before in Section 7 a brief conclusion and outlook is given.

## 2 Thermal radiative transfer

In a one-dimensional setting the thermal radiative transfer equation with absorbing background material is given by

$$\frac{1}{c}\partial_t f(t,x,\mu) + \mu \partial_x f(t,x,\mu) = \sigma(B(t,x) - f(t,x,\mu)),$$

$$\partial_t e(t,x) = \sigma(\langle f(t,x,\cdot) - B(t,x)\rangle_\mu),$$

where an integration over the directional domain $[-1,1]$ is denoted by $\langle\cdot\rangle_\mu$. The speed of light is denoted by $c$ and the variable $\sigma$ represents the opacity that encodes the rate

at which particles are absorbed by the background medium. The black body radiation $B(T)$ at the material temperature $T$ can be described by the Stefan-Boltzmann law

$$B(T) = acT^4,$$

where $a = \frac{4\sigma_{\mathrm{SB}}}{c}$ is the radiation density constant and $\sigma_{\mathrm{SB}}$ the Stefan-Boltzmann constant. The above set of equations is not closed. To determine a relation between the temperature $T$ and the internal energy $e(T)$ we follow the ideas of Pomraning [43] and Su and Olson [45] and set $e(T) = \alpha B(T)$. From this point on, we call $\alpha B(T)$ the internal energy of the material. Further, we perform a rescaling $\tau = \frac{t}{c}$ and by an abuse of notation write $t$ instead of $\tau$ in the remainder. This leads to the system

$$\partial_t f(t, x, \mu) + \mu \partial_x f(t, x, \mu) = \sigma(B(t, x) - f(t, x, \mu)), \tag{1a}$$

$$\partial_t B(t, x) = \sigma(\langle f(t, x, \cdot) - B(t, x) \rangle_\mu), \tag{1b}$$

where without loss of generality we assume $\alpha = 1$. This system is a closed linearized internal energy model that is analytically solvable and commonly used as a benchmark for numerical examples [38, 36, 37, 35]. In the following, we call equations (1) the *Su-Olson problem*. Note that for the moment we omit initial and boundary conditions.

In [12] it has been shown that for deriving an efficient dynamical low-rank scheme for the Boltzmann-BGK equation it is crucial to consider a multiplicative splitting of the distribution function. This allows one to separate a generally not low-rank Maxwellian from a remaining low-rank function $g$, to which the DLRA scheme is subsequently applied. The considered Su-Olson problem is similar in structure to this equation. To transfer knowledge of the construction of efficient dynamical low-rank schemes from the Su-Olson problem to more general kinetic equations, we have decided on a multiplicative splitting of the distribution function of the form

$$f(t, x, \mu) = B(t, x)g(t, x, \mu), \tag{2}$$

and apply the low-rank ansatz to $g$. For this system, we give a mathematically rigorous proof of energy stability and derive a hyperbolic CFL condition. In this sense, this paper can be understood as an intermediate step from the Su-Olson problem treated in [1] towards more complicated BGK problems with multiplicative splitting as in [12], where the time step size of the proposed algorithm is not determined theoretically by means of analytical considerations but experimentally chosen small enough such that numerical experiments show good agreement. The idea of applying a multiplicative splitting to a linearized model and deriving a concrete CFL condition is also pursued in [2] for the linear Boltzmann-BGK equation.

We insert ansatz (2) into (1a) and (1b) and obtain the set of equations

$$\partial_t g(t, x, \mu) = -\mu \partial_x g(t, x, \mu) + \sigma(1 - g(t, x, \mu)) - \frac{g(t, x, \mu)}{B(t, x)} \partial_t B(t, x) \tag{3a}$$

$$-\mu \frac{g(t, x, \mu)}{B(t, x)} \partial_x B(t, x),$$

$$\partial_t B(t, x) = \sigma B(t, x) \left( \langle g(t, x, \mu) \rangle_\mu - 2 \right), \tag{3b}$$

that is called the *advection form* of the multiplicative system. Using the product rule, it splits up the spatial derivatives for $B$ and $g$ in (3a). This corresponds to the form in

which the multiplicative splitting in [12] is applied to the Boltzmann-BGK equation. Equation (3a) can be equivalently rewritten into a *conservative form,* leaving the spatial derivative of $Bg$ together and leading to the system

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

$$\partial_t B(t, x) = \sigma B(t, x) \left( \langle g(t, x, \mu) \rangle_\mu - 2 \right).$$
(4b)

In later considerations, we are interested in the conservation properties of our numerical scheme. For the multiplicative Su-Olson problem, the mass and the momentum of the system shall be defined as follows.

**Definition 1 (Macroscopic quantities)** The *mass* of the multiplicative Su-Olson problem is defined as

$$\rho(t, x) = \int f(t, x, \mu) \mathrm{d}\mu + B(t, x) = B(t, x) \int g(t, x, \mu) \mathrm{d}\mu + B(t, x).$$

The *momentum* is given as

$$u(t, x) = \int \mu f(t, x, \mu) \mathrm{d}\mu = B(t, x) \int \mu g(t, x, \mu) \mathrm{d}\mu.$$

In particular, the multiplicative Su-Olson problem satisfies the local conservation law

$$\partial_t \rho(t, x) + \partial_x u(t, x) = 0.$$
(5)

Global conservation of mass is then obtained by integrating over the spatial domain [16].

In the following, we discretize both sets of equations to compare them in terms of numerical stability. We derive an energy stable dynamical low-rank scheme and give a concrete hyperbolic CFL condition. Note that in contrast to [1, 12], but similar to [2], we first discretize the equations and then apply the low-rank ansatz here.

## 3 Discretization of the multiplicative system

In this section, we fully discretize the advection form (3) as well as the conservative form (4) of the multiplicative system. We start with the angular and spatial discretization, followed by the time discretization.

### 3.1 Angular discretization

For the angular discretization a modal approach with normalized Legendre polynomials $P_\ell$ is used. They constitute a complete set of orthogonal functions on the interval $[-1, 1]$ that satisfy $\langle P_k, P_\ell \rangle_\mu = \delta_{k\ell}$. As an approximation, a finite expansion

of the distribution function $g$ with $N_\mu$ expansion coefficients, called the moments, is used. It writes

$$g(t,x,\mu) \approx g_{N_\mu}(t,x,\mu) = \sum_{\ell=0}^{N_\mu-1} v_\ell(t,x) P_\ell(\mu).$$

We insert this representation into (3), multiply (3a) with $P_k(\mu)$ and integrate over $\mu$. Further, we introduce the matrix $\mathbf{A} \in \mathbb{R}^{N_\mu \times N_\mu}$ with entries $A_{k\ell} := \langle P_k, \mu P_\ell \rangle_\mu$ and use that $P_0 = \frac{1}{\sqrt{2}}$. This gives

$$\partial_t v_k(t,x) = -\sum_{\ell=0}^{N_\mu-1} \partial_x v_\ell(t,x) A_{k\ell} + \sigma\left(\sqrt{2}\delta_{k0} - v_k(t,x)\right) - \frac{v_k(t,x)}{B(t,x)}\partial_t B(t,x) \quad \text{(6a)}$$

$$-\sum_{\ell=1}^{N_\mu-1} \frac{v_\ell}{B(t,x)}\partial_x B(t,x) A_{k\ell},$$

$$\partial_t B(t,x) = \sigma B(t,x)\left(\sqrt{2}v_0(t,x) - 2\right). \quad \text{(6b)}$$

Analogously, we obtain for system (4) the angularly discretized equations

$$\partial_t v_k(t,x) = -\sum_{\ell=0}^{N_\mu-1} \frac{1}{B(t,x)}\partial_x\left(B(t,x)v_\ell(t,x)\right) A_{k\ell} + \sigma\left(\sqrt{2}\delta_{k0} - v_k(t,x)\right) \quad \text{(7a)}$$

$$-\frac{v_k(t,x)}{B(t,x)}\partial_t B(t,x),$$

$$\partial_t B(t,x) = \sigma B(t,x)\left(\sqrt{2}v_0(t,x) - 2\right). \quad \text{(7b)}$$

Note that the matrix $\mathbf{A}$ is symmetric and diagonalizable in the form $\mathbf{A} = \mathbf{Q}\mathbf{M}\mathbf{Q}^\top$ with $\mathbf{Q}$ orthonormal and $\mathbf{M} = \text{diag}(\sigma_1, ..., \sigma_{N_\mu})$. We then define $|\mathbf{A}| = \mathbf{Q}|\mathbf{M}|\mathbf{Q}^\top$.

### 3.2 Spatial discretization

For the spatial discretization we prescribe a spatial grid with $N_x$ grid points and equidistant spacing $\Delta x = \frac{1}{N_x}$. Spatially dependent quantities are approximated at the grid points $x_j$ for $j = 1, ..., N_x$, and denoted by

$$B_j(t) \approx B(t,x_j), \qquad v_{jk}(t) \approx v_k(t,x_j).$$

First-order spatial derivatives $\partial_x$ are approximated using the tridiagonal stencil matrices $\mathbf{D}^x \in \mathbb{R}^{N_x \times N_x}$ and a second-order stabilization term $\mathbf{D}^{xx} \in \mathbb{R}^{N_x \times N_x}$ approximating $\frac{1}{2}\Delta x \partial_{xx}$ is added. The entries of those matrices are defined as

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

Note that from now on we assume periodic boundary conditions to which we account by setting

$$D^x_{1,N_x} = \frac{-1}{2\Delta x}, \qquad D^x_{N_x,1} = \frac{1}{2\Delta x},$$

$$D^{xx}_{1,N_x} = D^{xx}_{N_x,1} = \frac{1}{2\Delta x}.$$

The stencil matrices $\mathbf{D}^x$ and $\mathbf{D}^{xx}$ then fulfill the following properties:

**Lemma 1** *Let $y, z \in \mathbb{R}^{N_x}$ with indices $i, j = 1, ..., N_x$. It holds*

$$\sum_{i,j=1}^{N_x} y_j D_{ji}^x z_i = - \sum_{i,j=1}^{N_x} z_j D_{ji}^x y_i , \quad \sum_{i,j=1}^{N_x} z_j D_{ji}^x z_i = 0 , \quad \sum_{i,j=1}^{N_x} y_j D_{ji}^{xx} z_i = \sum_{i,j=1}^{N_x} z_j D_{ji}^{xx} y_i.$$

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

$$D_{j,j}^+ = \frac{-1}{\sqrt{2\Delta x}} , \qquad D_{j,j+1}^+ = \frac{1}{\sqrt{2\Delta x}} .$$

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

*Proof* See [1, Lemma 4.2]. □

We insert the proposed discretization into the angularly discretized advection form (6) of the equations and add a second-order stabilization term for $B\partial_x v$, leading to the angularly and spatially discretized set of equations

$$\dot{v}_{jk}(t) = -\sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu - 1} D_{ji}^x v_{i\ell}(t) A_{k\ell} + \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu - 1} D_{ji}^{xx} v_{i\ell}(t) |A|_{k\ell} \tag{8a}$$

$$+ \sigma \left( \sqrt{2}\delta_{k0} - v_{jk}(t) \right) - \frac{v_{jk}(t)}{B_j(t)} \dot{B}_j(t) - \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu - 1} \frac{v_{j\ell}(t)}{B_j(t)} D_{ji}^x B_i(t) A_{k\ell},$$

$$\dot{B}_j(t) = \sigma B_j(t) \left( \sqrt{2} v_{j0}(t) - 2 \right). \tag{8b}$$

Inserting the discretization into the angularly discretized conservative form (7) of the equations and adding a second-order stabilization term to $\partial_x (Bv)$ gives

$$\dot{v}_{jk}(t) = -\sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu - 1} \frac{1}{B_j(t)} D_{ji}^x B_i(t) v_{i\ell}(t) A_{k\ell} + \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu - 1} \frac{1}{B_j(t)} D_{ji}^{xx} B_i(t) v_{i\ell}(t) |A|_{k\ell}$$

$$\tag{9a}$$

$$+ \sigma \left( \sqrt{2}\delta_{k0} - v_{jk}(t) \right) - \frac{v_{jk}(t)}{B_j(t)} \dot{B}_j(t),$$

$$\dot{B}_j(t) = \sigma B_j(t) \left( \sqrt{2} v_{j0}(t) - 2 \right). \tag{9b}$$

Note that due to the different structure of the equations the stabilization term in (8a) is applied to $B\partial_x v$, whereas in (9a) it is added for $\partial_x (Bv)$.

3.3 Time discretization

From [1] we know that constructing an energy stable scheme for the Su-Olson problem is challenging. For the advection form of the equations we start from (8) and apply an explicit Euler step to transport terms. The potentially stiff absorption terms

are treated implicitly and the time derivative $\partial_t B$ is approximated by its difference quotient. We obtain the following fully discrete space-time discretization

$$v_{jk}^1 = v_{jk}^0 - \Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} D_{ji}^x v_{i\ell}^0 A_{k\ell} + \Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} D_{ji}^{xx} v_{i\ell}^0 |A|_{k\ell} \tag{10a}$$

$$+ \sigma \Delta t \left( \sqrt{2} \delta_{k0} - v_{jk}^1 \right) - \Delta t \frac{1}{B_j^0} \frac{B_j^1 - B_j^0}{\Delta t} v_{jk}^1 - \Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \frac{v_{j\ell}^0}{B_j^0} D_{ji}^x B_i^0 A_{k\ell},$$

$$B_j^1 = B_j^0 + \sigma \Delta t B_j^1 \left( sqrt2 v_{j0}^1 - 2 \right), \tag{10b}$$

that describes one time step from time $t_0$ to time $t_1 = t_0 + \Delta t$ and holds for all further time steps equivalently. For the conservative form (9) we again apply an explicit Euler step to the transport parts, treat the absorption terms implicitly and approximate $\partial_t B$ by its difference quotient. In addition, we add a factor $\frac{B^1}{B^0}$ in the absorption term of (9a). This gives the fully discrete scheme

$$v_{jk}^1 = v_{jk}^0 - \Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \frac{1}{B_j^0} D_{ji}^x B_i^0 v_{i\ell}^0 A_{k\ell} + \Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \frac{1}{B_j^0} D_{ji}^{xx} B_i^0 v_{i\ell}^0 |A|_{k\ell} \tag{11a}$$

$$+ \sigma \Delta t \frac{B_j^1}{B_j^0} \left( \sqrt{2} \delta_{k0} - v_{jk}^1 \right) - \Delta t \frac{1}{B_j^0} \frac{B_j^1 - B_j^0}{\Delta t} v_{jk}^1,$$

$$B_j^1 = B_j^0 + \sigma \Delta t B_j^1 \left( \sqrt{2} v_{j0}^1 - 2 \right). \tag{11b}$$

Note that the evolution equations (10b) and (11b) for the internal energy $B$ are the same in both schemes. The main difference of (10a) and (11a) consists in the distinct second-order stabilization terms and the additional factor $\frac{B^1}{B^0}$ in (11a) that we will explain later when showing energy stability.

## 4 Energy stability

The goal of this section is to investigate energy stability of the derived schemes. Note that this section is closely related to the considerations in [1]. We first introduce the following notations.

**Definition 2** In the following we write $u_{jk}^0 := B_j^0 v_{jk}^0$ and $u_{jk}^1 := B_j^1 v_{jk}^1$ at time $t_0$ and $t_1$, respectively. Note that $\mathbf{u}^0 = (u_{jk}^0) \in \mathbb{R}^{N_x \times N_\mu}$ corresponds to $f(t = 0, x, \mu)$ and $\mathbf{v}^0 = (v_{jk}^0) \in \mathbb{R}^{N_x \times N_\mu}$ corresponds to $g(t = 0, x, \mu)$ in (2).

With this notation we can give the definition of the total energy of a fully discretized system.

**Definition 3 (Total energy)** Let $\mathbf{u}^0 \in \mathbb{R}^{N_x \times N_\mu}$ be the fully discretized angular solution of the full Su-Olson problem and $\mathbf{B}^0 = (B_j^0) \in \mathbb{R}^{N_x}$ the internal energy at time $t_0$. Then, the *total energy* at this time is defined as

$$E^0 := \frac{1}{2} \|\mathbf{u}^0\|_F^2 + \frac{1}{2} \|\mathbf{B}^0\|_E^2,$$

where $\| \cdot \|_F$ denotes the Frobenius and $\| \cdot \|_E$ the Euclidean norm. For $t_1 = t_0 + \Delta t$ this definition shall hold analogously.

4.1 Advection form

We start with the advection form (10) of the Su-Olson problem which is comparable to the considered low-rank discretization in [12] for the isothermal Boltzmann-BGK equation in the sense that the term $\partial_x (Bv)$ is split up into the sum of $B\partial_x v$ and $v\partial_x B$. We can show that this scheme is not, in general, von Neumann stable.

**Theorem 1** *There exist initial values $\mathbf{v}^0 \in \mathbb{R}^{N_x \times N_\mu}$ and $\mathbf{B}^0 \in \mathbb{R}^{N_x}$ such that the advection form (10) of the Su-Olson problem for $\sigma = 0$ is not von Neumann stable.*

*Proof* Let us assume a solution $v_{jk}^0$ that is constant in space and direction, e.g. $v_{jk}^0 = 1$. For this solution all spatial derivatives are zero, i.e. the terms containing $\mathbf{D}^x \mathbf{v}^0$ and $\mathbf{D}^{xx} \mathbf{v}^0$ in (10a) drop out. We further assume that for the opacity it holds $\sigma = 0$, i.e. the Su-Olson problem reduces to a simple advection equation. From (10b) we thus obtain that $B_j^1 = B_j^0 = B_j$, i.e. the internal energy is constant in time. We insert these results into (10a) and get

$$v_{jk}^1 = 1 - \Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \frac{1}{B_j} D_{ji}^x B_i A_{k\ell}.$$

Multiplication with $B_j$ then leads to

$$u_{jk}^1 = u_{jk}^0 - \Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} D_{ji}^x u_{i\ell}^0 A_{k\ell}.$$

This is a discretization of $\partial_t u + \mu \partial_x u = 0$ with an explicit Euler step forward in time and a centered finite difference scheme in space. From [23,31] we know that this discretization is not von Neumann stable.

The consideration of this special case shall serve as a motivation to seek for a generally stable numerical discretization as done in the next subsection.

4.2 Conservative form

For the conservative form of the discretization of the Su-Olson problem given in (11), we can derive a hyperbolic CFL condition and show that under this time step restriction the total energy of the system decreases over time. We start with the following lemma.

**Lemma 2** *Under the time step restriction $\Delta t \leq \Delta x$ it holds*

$$\frac{\Delta t}{2} \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \left( \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left( D_{ji}^x u_{jk}^1 A_{k\ell} - D_{ji}^{xx} u_{jk}^1 |A|_{k\ell} \right) \right)^2$$

$$- \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \left( \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} D_{ji}^+ u_{jk}^1 |A|_{k\ell}^{1/2} \right)^2 \leq 0,$$

*Proof* See [1, Lemma 5.2].

With this relation we can now show that the energy of the conservative form (11) of the Su-Olson system dissipates and hence the system is energy stable.

**Theorem 2** *Under the time step restriction $\Delta t \leq \Delta x$ the fully discrete system* (11) *is energy stable, i.e. it holds $E^1 \leq E^0$.*

*Proof* The proof of this theorem is similar to the proof of [1, Theorem 5.3]. We start with equation (11b) and multiply it with $B_j^1$. This gives

$$\left(B_j^1\right)^2 = B_j^0 B_j^1 + \sigma \Delta t \left(B_j^1\right)^2 \left(\sqrt{2} v_{j0}^1 - 2\right).$$

Note that it holds

$$B_j^0 B_j^1 = \frac{1}{2} \left(B_j^1\right)^2 + \frac{1}{2} \left(B_j^0\right)^2 - \frac{1}{2} \left(B_j^1 - B_j^0\right)^2.$$

We insert this relation and sum over $j$, giving

$$\frac{1}{2} \sum_{j=1}^{N_x} \left(B_j^1\right)^2 = \frac{1}{2} \sum_{j=1}^{N_x} \left(B_j^0\right)^2 - \frac{1}{2} \sum_{j=1}^{N_x} \left(B_j^1 - B_j^0\right)^2 + \sigma \Delta t \sum_{j=1}^{N_x} \left(B_j^1\right)^2 \left(\sqrt{2} v_{j0}^1 - 2\right).$$

$$(12)$$

Next, we multiply equation (11a) with $B_j^1 B_j^0 v_{jk}^1$ and sum over $j$ and $k$. This leads to

$$
\begin{aligned}
\sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} B_j^1 B_j^0 \left(v_{jk}^1\right)^2 =& \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} B_j^0 v_{jk}^0 B_j^1 v_{jk}^1 - \Delta t \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} B_j^1 v_{jk}^1 D_{ji}^x B_i^0 v_{i\ell}^0 A_{k\ell} \\
&+ \Delta t \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} B_j^1 v_{jk}^1 D_{ji}^{xx} B_i^0 v_{i\ell}^0 |A|_{k\ell} \\
&+ \sigma \Delta t \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left(B_j^1\right)^2 v_{jk}^1 \left(\sqrt{2} \delta_{k0} - v_{jk}^1\right) \\
&- \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} B_j^1 \left(v_{jk}^1\right)^2 \left(B_j^1 - B_j^0\right)
\end{aligned}
$$

$$(13)$$

Note that for this step the additional factor $\frac{B_1}{B_0}$ in the absorption term of (11a) is crucial. As above, it holds for the first term on the right-hand side that

$$
\begin{aligned}
&\sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} B_j^0 v_{jk}^0 B_j^1 v_{jk}^1 \\
=& \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left(\frac{1}{2} \left(B_j^1 v_{jk}^1\right)^2 + \frac{1}{2} \left(B_j^0 v_{jk}^0\right)^2 - \frac{1}{2} \left(B_j^1 v_{jk}^1 - B_j^0 v_{jk}^0\right)^2\right).
\end{aligned}
$$

We insert the notation $u_{jk}^0 = B_j^0 v_{jk}^0$ and $u_{jk}^1 = B_j^1 v_{jk}^1$, respectively, as well as insert the above relation into (13), bring the last term of (13) to the left-hand side and rearrange the equation. We obtain

$$
\frac{1}{2} \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left(u_{jk}^1\right)^2 = \frac{1}{2} \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left(u_{jk}^0\right)^2 - \frac{1}{2} \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left(u_{jk}^1 - u_{jk}^0\right)^2
$$
$$
- \Delta t \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} u_{jk}^1 D_{ji}^x u_{i\ell}^0 A_{k\ell} + \Delta t \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} u_{jk}^1 D_{ji}^{xx} u_{i\ell}^0 |A|_{k\ell}
$$
$$
+ \sigma \Delta t \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left(B_j^1\right)^2 v_{jk}^1 \left(\sqrt{2}\delta_{k0} - v_{jk}^1\right).
$$

We now add the zero term $\Delta t \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} u_{jk}^1 D_{ji}^x u_{i\ell}^1 A_{k\ell}$ and add and subtract the second-order term $\Delta t \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} u_{jk}^1 D_{ji}^{xx} u_{i\ell}^1 |A|_{k\ell}$ giving

$$
\frac{1}{2} \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left(u_{jk}^1\right)^2 = \frac{1}{2} \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left(u_{jk}^0\right)^2 - \frac{1}{2} \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left(u_{jk}^1 - u_{jk}^0\right)^2
$$
$$
- \Delta t \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} u_{jk}^1 D_{ji}^x \left(u_{i\ell}^0 - u_{i\ell}^1\right) A_{k\ell}
$$
$$
+ \Delta t \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} u_{jk}^1 D_{ji}^{xx} \left(u_{i\ell}^0 - u_{i\ell}^1\right) |A|_{k\ell}
$$
$$
+ \Delta t \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} u_{jk}^1 D_{ji}^{xx} u_{i\ell}^1 |A|_{k\ell}
$$
$$
+ \sigma \Delta t \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left(B_j^1\right)^2 v_{jk}^1 \left(\sqrt{2}\delta_{k0} - v_{jk}^1\right).
$$

In the next step, we apply Young's inequality, which states that for $a, b \in \mathbb{R}$ we have $a \cdot b \le \frac{a^2}{2} + \frac{b^2}{2}$, to the term

$$
- \Delta t \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} u_{jk}^1 D_{ji}^x \left(u_{i\ell}^0 - u_{i\ell}^1\right) A_{k\ell} + \Delta t \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} u_{jk}^1 D_{ji}^{xx} \left(u_{i\ell}^0 - u_{i\ell}^1\right) |A|_{k\ell}
$$
$$
= - \Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \left(u_{i\ell}^0 - u_{i\ell}^1\right) \left( \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left(D_{ji}^x u_{jk}^1 A_{k\ell} - D_{ji}^{xx} u_{jk}^1 |A|_{k\ell}\right) \right)
$$
$$
\le \frac{1}{2} \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \left(u_{i\ell}^0 - u_{i\ell}^1\right)^2
$$
$$
+ \frac{(\Delta t)^2}{2} \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \left( \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left(D_{ji}^x u_{jk}^1 A_{k\ell} - D_{ji}^{xx} u_{jk}^1 |A|_{k\ell}\right) \right)^2.
$$

In addition, we note that with the properties of the stencil matrices from Lemma 1 we can write

$$
\Delta t \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} u_{jk}^1 D_{ji}^{xx} u_{i\ell}^1 |A|_{k\ell} = -\Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \left( \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} D_{ji}^+ u_{jk}^1 |A|_{k\ell}^{1/2} \right)^2 .
$$

We insert both relations and get

$$
\begin{aligned}
\frac{1}{2} \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left( u_{jk}^1 \right)^2 \leq & \frac{1}{2} \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left( u_{jk}^0 \right)^2 \\
& + \frac{(\Delta t)^2}{2} \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \left( \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left( D_{ji}^x u_{jk}^1 A_{k\ell} - D_{ji}^{xx} u_{jk}^1 |A|_{k\ell} \right) \right)^2 \\
& - \Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \left( \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} D_{ji}^+ u_{jk}^1 |A|_{k\ell}^{1/2} \right)^2 \\
& + \sigma \Delta t \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left( B_j^1 \right)^2 v_{jk}^1 \left( \sqrt{2} \delta_{k0} - v_{jk}^1 \right) .
\end{aligned}
$$

With Lemma 2 we have that under the time step restriction $\Delta t \leq \Delta x$ it holds

$$
\frac{1}{2} \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left( u_{jk}^1 \right)^2 \leq \frac{1}{2} \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left( u_{jk}^0 \right)^2 + \sigma \Delta t \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left( B_j^1 \right)^2 v_{jk}^1 \left( \sqrt{2} \delta_{k0} - v_{jk}^1 \right) .
$$

(14)

To obtain an expression for the total energy of the system, we add (14) and (12). This gives

$$
\begin{aligned}
E^1 \leq \ & E^0 - \frac{1}{2} \sum_{j=1}^{N_x} \left( B_j^1 - B_j^0 \right)^2 + \sigma \Delta t \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left( B_j^1 \right)^2 v_{jk}^1 \left( \sqrt{2} \delta_{k0} - v_{jk}^1 \right) \\
& + \sigma \Delta t \sum_{j=1}^{N_x} \left( B_j^1 \right)^2 \left( \sqrt{2} v_{j0}^1 - 2 \right) .
\end{aligned}
$$

The term $-\frac{1}{2} \sum_{j=1}^{N_x} \left( B_j^1 - B_j^0 \right)^2$ is non-positive. The remaining two terms on the right-hand side can be rewritten and bounded as follows:

$$
\begin{aligned}
& \sigma \Delta t \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left( B_j^1 \right)^2 v_{jk}^1 \left( \sqrt{2} \delta_{k0} - v_{jk}^1 \right) + \sigma \Delta t \sum_{j=1}^{N_x} \left( B_j^1 \right)^2 \left( \sqrt{2} v_{j0}^1 - 2 \right) \\
& \leq \sigma \Delta t \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left( B_j^1 \right)^2 \left( - \left( v_{jk}^1 \right)^2 + 2\sqrt{2} v_{jk}^1 \delta_{k0} - 2\delta_{k0} \right) \\
& = - \sigma \Delta t \sum_{j=1}^{N_x} \sum_{k=0}^{N_\mu-1} \left( B_j^1 \right)^2 \left( v_{jk}^1 - \sqrt{2} \delta_{k0} \right)^2 \leq 0 .
\end{aligned}
$$

Hence, we have shown that under the time step restriction $\Delta t \leq \Delta x$ it holds $E^1 \leq E^0$, and the system is energy stable.

## 5 Dynamical low-rank approximation for the energy stable system

Having attained an energy stable discretization of the multiplicative Su-Olson problem, its practical implementation can still pose numerical challenges such as large memory demands and computational costs, especially in a higher-dimensional setting. To overcome these problems, we introduce the concept of dynamical low-rank approximation.

### 5.1 Background on DLRA

The method of dynamical low-rank approximation has originally been introduced in a semi-discrete time-continuous matrix setting [26], in which it will also be explained in this subsection. Let us consider the matrix differential equation

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

where $\mathbf{f}(t) \in \mathbb{R}^{N_x \times N_\mu}$ is the solution of the equation and $\mathbf{F}\left(t, \mathbf{f}(t)\right) : \mathbb{R}^{N_x \times N_\mu} \rightarrow \mathbb{R}^{N_x \times N_\mu}$ denotes its right-hand side. The dynamical low-rank approximation of $\mathbf{f}(t)$ is then given by

$$\mathbf{f}_r(t) = \mathbf{X}(t)\mathbf{S}(t)\mathbf{V}(t)^\top, \tag{15}$$

where $\mathbf{X}(t) \in \mathbf{R}^{N_x \times r}$ is the orthonormal basis in space, $\mathbf{V}(t) \in \mathbf{R}^{N_\mu \times r}$ the orthonormal basis in direction, and $\mathbf{S}(t) \in \mathbb{R}^{r \times r}$ the coupling or coefficient matrix of the rank $r$ approximation. All matrices of the form (15) constitute the manifold of low-rank matrices $\mathcal{M}_r$. The basis matrices $\mathbf{X}(t)$ and $\mathbf{V}(t)$, and the coefficient matrix $\mathbf{S}(t)$ shall then be evolved in time such that the minimization problem

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

is fulfilled at all times $t$, where $\mathcal{T}_{\mathbf{f}_r(t)}\mathcal{M}_r$ denotes the tangent space of the low-rank manifold $\mathcal{M}_r$ at $\mathbf{f}_r(t)$. In [26], it has been shown that solving this minimization problem is equivalent to projecting the right-hand side $\mathbf{F}\left(t, \mathbf{f}_r(t)\right)$ by means of an orthogonal projection $\mathbf{P}$ onto $\mathcal{T}_{\mathbf{f}_r(t)}\mathcal{M}_r$, and solving

$$\dot{\mathbf{f}}_r(t) = \mathbf{P}\left(\mathbf{f}_r(t)\right)\mathbf{F}\left(t, \mathbf{f}_r(t)\right). \tag{16}$$

For $\mathbf{f}_r = \mathbf{X}\mathbf{S}\mathbf{V}^\top$, the orthogonal projection $\mathbf{P}$ onto $\mathcal{T}_{\mathbf{f}_r(t)}\mathcal{M}_r$ is explicitly given as

$$\mathbf{P}(\mathbf{f}_r)\mathbf{F}\left(\mathbf{f}_r\right) = \mathbf{X}\mathbf{X}^\top \mathbf{F}\left(\mathbf{f}_r\right) - \mathbf{X}\mathbf{X}^\top \mathbf{F}\left(\mathbf{f}_r\right)\mathbf{V}\mathbf{V}^\top + \mathbf{F}\left(\mathbf{f}_r\right)\mathbf{V}\mathbf{V}^\top.$$

Different time integrators that make use of the special form of this orthogonal projection such as the projector-splitting [32], the (augmented) basis update & Galerkin (BUG) [7,5], and the parallel integrator [6] exist. They are able to evolve the solution on the low-rank manifold while not suffering from the stiffness of (16). In this paper,

we focus on the augmented basis update & Galerkin (BUG) integrator, that shall be explained in the following.

The BUG integrator first updates and augments the spatial basis $\mathbf{X}$ and the directional basis $\mathbf{V}$ in parallel. This leads to an increase from rank $r$ to rank $2r$. Note that we denote augmented quantities of rank $2r$ with hats. Next, a Galerkin step is conducted for the coefficient matrix $\mathbf{S}$ in the augmented setting, before in a last step all augmented quantities are truncated to a new rank $r_1 \leq 2r$. To be more specific, the BUG integrator evolves the low-rank solution $\mathbf{f}_r^0 = \mathbf{X}^0 \mathbf{S}^0 \mathbf{V}^{0,\top}$ at time $t_0$ to the time-updated low-rank solution $\mathbf{f}_r^1 = \mathbf{X}^1 \mathbf{S}^1 \mathbf{V}^{1,\top}$ at time $t_1 = t_0 + \Delta t$ as follows:

***K*-Step**: We fix the directional basis $\mathbf{V}^0$ at time $t_0$ and introduce the notation $\mathbf{K}(t) = \mathbf{X}(t)\mathbf{S}(t)$. Then, we update the spatial basis from $\mathbf{X}^0$ to $\widehat{\mathbf{X}}^1 \in \mathbb{R}^{N_x \times 2r}$ by solving the PDE

$$\dot{\mathbf{K}}(t) = \mathbf{F}\left(t, \mathbf{K}(t)\mathbf{V}^{0,\top}\right)\mathbf{V}^0, \quad \mathbf{K}(t_0) = \mathbf{X}^0\mathbf{S}^0,$$

and determining $\widehat{\mathbf{X}}^1$ as an orthonormal basis of the augmented matrix $[\mathbf{K}(t_1), \mathbf{X}^0] \in \mathbb{R}^{N_x \times 2r}$, e.g. by QR-decomposition. We store $\widehat{\mathbf{M}} = \widehat{\mathbf{X}}^{1,\top}\mathbf{X}^0 \in \mathbb{R}^{2r \times r}$.

***L*-Step**: We fix the spatial basis $\mathbf{X}^0$ at time $t_0$ and introduce the notation $\mathbf{L}(t) = \mathbf{V}(t)\mathbf{S}(t)^\top$. Then, we update the directional basis from $\mathbf{V}^0$ to $\widehat{\mathbf{V}}^1 \in \mathbb{R}^{N_\mu \times 2r}$ by solving the PDE

$$\dot{\mathbf{L}}(t) = \mathbf{F}\left(t, \mathbf{X}^0 \mathbf{L}(t)^\top\right)^\top \mathbf{X}^0, \quad \mathbf{L}(t_0) = \mathbf{V}^0\mathbf{S}^{0,\top},$$

and determining $\widehat{\mathbf{V}}^1$ as an orthonormal basis of the augmented matrix $[\mathbf{L}(t_1), \mathbf{V}^0] \in \mathbb{R}^{N_\mu \times 2r}$, e.g. by QR-decomposition. We store $\widehat{\mathbf{N}} = \widehat{\mathbf{V}}^{1,\top}\mathbf{V}^0 \in \mathbb{R}^{2r \times r}$.

***S*-step**: We update the coupling matrix from $\mathbf{S}^0 \in \mathbb{R}^{r \times r}$ to $\widehat{\mathbf{S}}^1 \in \mathbb{R}^{2r \times 2r}$ by solving the ODE

$$\dot{\widehat{\mathbf{S}}}(t) = \widehat{\mathbf{X}}^{1,\top}\mathbf{F}\left(t, \widehat{\mathbf{X}}^1\widehat{\mathbf{S}}(t)\widehat{\mathbf{V}}^{1,\top}\right)\widehat{\mathbf{V}}^1, \quad \widehat{\mathbf{S}}(t_0) = \widehat{\mathbf{M}}\mathbf{S}^0\widehat{\mathbf{N}}^\top.$$

**Truncation**: We compute the singular value decomposition of $\widehat{\mathbf{S}}^1 = \widehat{\mathbf{P}}\mathbf{\Sigma}\widehat{\mathbf{Q}}^\top$ with $\mathbf{\Sigma} = \mathrm{diag}(\sigma_j)$. The new rank $r_1 \leq 2r$ is chosen such that for a prescribed tolerance parameter $\vartheta$ it holds

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

We set $\mathbf{S}^1 \in \mathbb{R}^{r_1 \times r_1}$ to be the matrix containing the $r_1$ largest singular values. For the update of the spatial and the directional basis we introduce the matrices $\mathbf{P}^1 \in \mathbb{R}^{2r \times r_1}$ and $\mathbf{Q}^1 \in \mathbb{R}^{2r \times r_1}$ containing the first $r_1$ columns of $\widehat{\mathbf{P}}$ and $\widehat{\mathbf{Q}}$, respectively, and set $\mathbf{X}^1 = \widehat{\mathbf{X}}^1\mathbf{P}^1 \in \mathbb{R}^{N_x \times r_1}$ and $\mathbf{V}^1 = \widehat{\mathbf{V}}^1\mathbf{Q}^1 \in \mathbb{R}^{N_\mu \times r_1}$.

Altogether, this gives the time-updated low-rank approximation $\mathbf{f}_r^1 = \mathbf{X}^1 \mathbf{S}^1 \mathbf{V}^{1,\top}$ at time $t_1 = t_0 + \Delta t$. Note that in the following, in an abuse of notation, we will write $\mathbf{f}$ instead of $\mathbf{f}_r$ and call this the low-rank solution of the considered problem.

5.2 DLRA scheme for multiplicative Su-Olson

In this subsection, the DLRA method is applied to the energy stable conservative form (11) of the Su-Olson problem to evolve $\mathbf{v}^0 = \left(v_{jk}^0\right)$ to $\mathbf{v}^1 = \left(v_{jk}^1\right)$. First note that for the derivation of the low-rank scheme we rewrite equations (11). In (11a), we bring all terms containing $v_{jk}^1$ to the left-hand side and divide by $1 + \sigma\Delta t$. Further, we multiply (11b) with $\frac{1}{B_j^0}$. This gives the system

$$\frac{B_j^1}{B_j^0}v_{jk}^1 = \frac{1}{1+\sigma\Delta t}v_{jk}^0 - \frac{\Delta t}{1+\sigma\Delta t}\sum_{i=1}^{N_x}\sum_{\ell=0}^{N_\mu-1}\frac{1}{B_j^0}D_{ji}^x B_i^0 v_{i\ell}^0 A_{k\ell} \tag{17a}$$

$$+ \frac{\Delta t}{1+\sigma\Delta t}\sum_{i=1}^{N_x}\sum_{\ell=0}^{N_\mu-1}\frac{1}{B_j^0}D_{ji}^{xx}B_i^0 v_{i\ell}^0 |A|_{k\ell} + \frac{\sqrt{2}\sigma\Delta t}{1+\sigma\Delta t}\frac{B_j^1}{B_j^0}\delta_{k0},$$

$$\frac{B_j^1}{B_j^0} = 1 + \sigma\Delta t\frac{B_j^1}{B_j^0}\left(\sqrt{2}v_{j0}^1 - 2\right). \tag{17b}$$

We then apply DLRA to this set of equations as follows. In a first step, we want to update $v_{jk}^0 = \sum_{m,n=1}^r X_{jm}^0 S_{mn}^0 V_{kn}^0$ to $\frac{B_j^1}{B_j^0}v_{jk}^* = \sum_{m,n=1}^{3r}\widehat{\widehat{X}}_{jm}^*\widehat{\widehat{S}}_{mn}^*\widehat{\widehat{V}}_{kn}^*$ for $k \neq 0$. We introduce the notation $K_{jn}^0 = \sum_{m=1}^r X_{jm}^0 S_{mn}^0$ and solve the $K$-step equation

$$K_{jp}^* = \frac{1}{1+\sigma\Delta t}K_{jp}^0 - \frac{\Delta t}{1+\sigma\Delta t}\frac{1}{B_j^0}\sum_{i=1}^{N_x}D_{ji}^x B_i^0 \sum_{n=1}^r K_{in}^0 \sum_{k,\ell=0}^{N_\mu-1}V_{\ell n}^0 A_{k\ell}V_{kp}^0 \tag{18a}$$

$$+ \frac{\Delta t}{1+\sigma\Delta t}\frac{1}{B_j^0}\sum_{i=1}^{N_x}D_{ji}^{xx}B_i^0\sum_{n=1}^r K_{im}^0\sum_{k,\ell=0}^{N_\mu-1}V_{\ell n}^0 |A|_{k\ell}V_{kp}^0.$$

We derive the updated basis $\widehat{\mathbf{X}}^*$ of rank $2r$ from $\widehat{\mathbf{X}}^* = \mathrm{qr}\left([\mathbf{K}^*, \mathbf{X}^0]\right)$. Moreover, we perform an additional basis augmentation step according to

$$\widehat{\widehat{\mathbf{X}}}^* = \mathrm{qr}\left([\widehat{\mathbf{X}}^*, \frac{1}{\mathbf{B}^0}\odot\mathbf{D}^x\left(\mathbf{B}^0\odot\mathbf{X}^0\right)]\right), \tag{18b}$$

which ensures the exactness of the corresponding projection operator in the proof of energy stability of the DLRA scheme. Here, the symbol $\odot$ stands for a pointwise multiplication and the vector $\frac{1}{\mathbf{B}^0}\in\mathbb{R}^{N_x}$ is defined to contain the element $\frac{1}{B_j}$ for each $j = 1, ..., N_x$. In addition, we compute and store $\widehat{\widehat{\mathbf{M}}} = \widehat{\widehat{\mathbf{X}}}^{*,\top}\mathbf{X}^0$. Note that quantities of rank $2r$ are denoted with one single hat and quantities of rank $3r$ with double hats.

The $L$-step can be computed in parallel with the $K$-step. We introduce the notation $L_{mk}^0 = \sum_{n=1}^r S_{nm}^0 V_{nk}^0$ and solve

$$L_{kp}^* = \frac{1}{1+\sigma\Delta t}L_{kp}^0 - \frac{\Delta t}{1+\sigma\Delta t}\sum_{\ell=0}^{N_\mu-1}\sum_{m=1}^r A_{\ell k}L_{\ell m}^0\sum_{i=1}^{N_x}X_{im}^0 B_i^0\sum_{j=1}^{N_x}D_{ij}^x\frac{1}{B_j^0}X_{jp}^0 \tag{18c}$$

$$+ \frac{\Delta t}{1+\sigma\Delta t}\sum_{\ell=0}^{N_\mu-1}\sum_{m=1}^r |A|_{\ell k}L_{\ell m}^0\sum_{i=1}^{N_x}X_{im}^0 B_i^0\sum_{j=1}^{N_x}D_{ij}^{xx}\frac{1}{B_j^0}X_{jp}^0.$$

We derive the updated basis $\widehat{\mathbf{V}}^*$ of rank $2r$ from $\widehat{\mathbf{V}}^* = \mathrm{qr}\left([\mathbf{L}^*, \mathbf{V}^0]\right)$. Moreover, we perform an additional basis augmentation step according to

$$\widehat{\widehat{\mathbf{V}}}^* = \mathrm{qr}\left([\widehat{\mathbf{V}}^*, \mathbf{A}^\top \mathbf{V}^0]\right), \tag{18d}$$

that again ensures the exactness of the corresponding projection operator and will be made clear in the proof of energy stability of the DLRA scheme in the next subsection. In addition, we compute and store $\widehat{\widehat{\mathbf{N}}} = \widehat{\widehat{\mathbf{V}}}^{*,\top} \mathbf{V}^0$.

For the $S$-step we use the information from the $K$- and $L$-step, set $\widetilde{S}^0_{mn} = \sum_{j,k=1}^r \widehat{\widehat{M}}_{mj} S^0_{jk} \widehat{\widehat{N}}_{nk}$, and solve

$$\widehat{\widehat{S}}^*_{qp} = \frac{1}{1+\sigma\Delta t} \widetilde{S}^0_{qp} - \frac{\Delta t}{1+\sigma\Delta t} \sum_{j=1}^{N_x} \widehat{\widehat{X}}^*_{jq} \frac{1}{B^0_j} \sum_{i=1}^{N_x} D^x_{ji} B^0_i \sum_{m,n=1}^{3r} \widehat{\widehat{X}}^*_{im} \widetilde{S}^0_{mn} \sum_{k,\ell=0}^{N_\mu - 1} \widehat{\widehat{V}}^*_{\ell n} A_{k\ell} \widehat{\widehat{V}}^*_{kp} \tag{18e}$$

$$+ \frac{\Delta t}{1+\sigma\Delta t} \sum_{j=1}^{N_x} \widehat{\widehat{X}}^*_{jq} \frac{1}{B^0_j} \sum_{i=1}^{N_x} D^{xx}_{ji} B^0_i \sum_{m,n=1}^{3r} \widehat{\widehat{X}}^*_{im} \widetilde{S}^0_{mn} \sum_{k,\ell=0}^{N_\mu - 1} \widehat{\widehat{V}}^*_{\ell n} |A|_{k\ell} \widehat{\widehat{V}}^*_{kp}.$$

In the next step, we consider the equations for $k = 0$. In this case, the equations for $\frac{B^1_j}{B^0_j} v^1_{j0}$ and $\frac{B^1_j}{B^0_j}$ couple and we solve the system

$$\frac{B^1_j}{B^0_j} v^1_{j0} = \frac{1}{1+\sigma\Delta t} \sum_{m,n=1}^r X^0_{jm} S^0_{mn} V^0_{kn}$$

$$- \frac{\Delta t}{1+\sigma\Delta t} \sum_{i=1}^{N_x} \sum_{\ell=0}^{N-1} \frac{1}{B^0_j} D^x_{ji} B^0_i \sum_{m,n=1}^{3r} \widehat{\widehat{X}}^*_{im} \widetilde{S}^0_{mn} \widehat{\widehat{V}}^*_{\ell n} A_{0\ell} \tag{18f}$$

$$+ \frac{\Delta t}{1+\sigma\Delta t} \sum_{i=1}^{N_x} \sum_{\ell=0}^{N-1} \frac{1}{B^0_j} D^{xx}_{ji} B^0_i \sum_{m,n=1}^{3r} \widehat{\widehat{X}}^*_{im} \widetilde{S}^0_{mn} \widehat{\widehat{V}}^*_{\ell n} |A|_{0\ell} + \frac{\sqrt{2}\sigma\Delta t}{1+\sigma\Delta t} \frac{B^1_j}{B^0_j},$$

$$\frac{B^1_j}{B^0_j} = 1 + \sigma\Delta t \frac{B^1_j}{B^0_j} \left(\sqrt{2} v^1_{j0} - 2\right). \tag{18g}$$

From (18f) and (18g) we can then retrieve $\mathbf{v}^1_0 = \left(v^1_{j0}\right)$ and $\mathbf{B}^1 = \left(B^1_j\right)$. We use the latter to divide out the factor $\frac{B^1_j}{B^0_j}$ in the low-rank representation of $\frac{B^1_j}{B^0_j} v^*_{jk} = \sum_{m,n=1}^{3r} \widehat{\widehat{X}}^*_{jm} \widehat{\widehat{S}}^*_{mn} \widehat{\widehat{V}}^*_{kn}$. We perform the transformation step

$$K^{*,\mathrm{trans}}_{jp} = \frac{B^0_j}{B^1_j} K^*_{jp}. \tag{18h}$$

From a QR-decomposition we then obtain $\widehat{\widehat{\mathbf{X}}}^{*,\mathrm{trans}} \widehat{\widehat{\mathbf{S}}}^{*,\mathrm{trans}} = \mathrm{qr}\left(\mathbf{K}^{*,\mathrm{trans}}\right)$. Then, we perform an additional basis augmentation step according to

$$\widehat{\widehat{\mathbf{X}}}^1 = \mathrm{qr}\left([\mathbf{v}^1_0, \widehat{\widehat{\mathbf{X}}}^{*,\mathrm{trans}}]\right), \quad \widehat{\widehat{\mathbf{V}}}^1 = \mathrm{qr}\left([\mathbf{e}_1, \widehat{\widehat{\mathbf{V}}}^*]\right), \tag{18i}$$

where we add $\mathbf{v}_0^1$ to the updated spatial low-rank basis since the mass of the system is given by the zeroth order moment $\mathbf{v}_0^1$. In the directional basis, we add $\mathbf{e}_1 \in \mathbb{R}^{N_\mu}$, denoting the first unit vector in $\mathbb{R}^{N_\mu}$. Again, this ensures mass conservation of the proposed low-rank scheme. Then, we have to adjust the coefficient matrix $\widehat{\widehat{\mathbf{S}}}^{*,\mathrm{trans}}$ correspondingly as

$$\widehat{\widehat{\mathbf{S}}}^1 = \widehat{\widehat{\mathbf{X}}}^{1,\top} \widehat{\widehat{\mathbf{X}}}^{*,\mathrm{trans}} \widehat{\widehat{\mathbf{S}}}^{*,\mathrm{trans}} \widehat{\widehat{\mathbf{V}}}^{*} \left( \mathbf{I} - \mathbf{e}_1 \mathbf{e}_1^\top \right) \widehat{\widehat{\mathbf{V}}}^1 + \widehat{\widehat{\mathbf{X}}}^{1,\top} \mathbf{v}_0^1 \mathbf{e}_1^\top \widehat{\widehat{\mathbf{V}}}^1 \in \mathbb{R}^{(3r+1) \times (3r+1)}. \tag{18j}$$

In a last step, we truncate the augmented quantities $\widehat{\widehat{\mathbf{X}}}^1, \widehat{\widehat{\mathbf{S}}}^1$ and $\widehat{\widehat{\mathbf{V}}}^1$ back to a new rank $r_1$, using the truncation strategy described in [5] or an adjusted truncation method inspired by [19] that ensures conservation of mass and will be given in a following subsection. Altogether, we obtain the updated low-rank factors $\mathbf{X}^1, \mathbf{S}^1$ and $\mathbf{V}^1$ such that $v_{jk}^1 = \sum_{m,n=1}^{r_1} X_{jm}^1 S_{mn}^1 V_{kn}^1$. The structure of the DLRA scheme is visualized in Figure 1.

5.3 Energy stability of the proposed low-rank scheme

We can then show that the proposed DLRA scheme preserves the energy stability of the full system.

**Theorem 3** *Under the time step restriction $\Delta t \leq \Delta x$, the fully discrete DLRA scheme* (18) *is energy stable, i.e. it holds $E^1 \leq E^0$.*

*Proof* We start with the internal energy $\mathbf{B}$ and multiply (18g) with $B_j^0 B_j^1$. This leads to

$$\left(B_j^1\right)^2 = B_j^0 B_j^1 + \sigma \Delta t \left(B_j^1\right)^2 \left(\sqrt{2} v_{j0}^1 - 2\right)$$

Analogously to the proof of Theorem 2, we rewrite the product $B_j^0 B_j^1$ and sum over $j$, giving the relation

$$\frac{1}{2} \sum_{j=1}^{N_x} \left(B_j^1\right)^2 = \frac{1}{2} \sum_{j=1}^{N_x} \left(B_j^0\right)^2 - \frac{1}{2} \sum_{j=1}^{N_x} \left(B_j^1 - B_j^0\right)^2 + \sigma \Delta t \sum_{j=1}^{N_x} \left(B_j^1\right)^2 \left(\sqrt{2} v_{j0}^1 - 2\right). \tag{19}$$

In the next step, we multiply (18e) with $\widehat{\widehat{X}}_{\alpha q}^* \widehat{\widehat{V}}_{\beta p}^*$ and sum over $q$ and $p$. We introduce the projection operators $P_{\alpha j}^{X^*} = \sum_{q=1}^{3r} \widehat{\widehat{X}}_{\alpha q}^* \widehat{\widehat{X}}_{jq}^*$ and $P_{k\beta}^{V^*} = \sum_{p=1}^{3r} \widehat{\widehat{V}}_{kp}^* \widehat{\widehat{V}}_{\beta p}^*$ and the notation $v_{\alpha\beta}^* := \sum_{p,q=1}^{3r} \widehat{\widehat{X}}_{\alpha q}^* \widehat{\widehat{S}}_{qp}^* \widehat{\widehat{V}}_{\beta p}^*$ and obtain

$$v_{\alpha\beta}^* = \frac{1}{1+\sigma\Delta t} v_{\alpha\beta}^0 - \frac{\Delta t}{1+\sigma\Delta t} \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} P_{\alpha j}^{X^*} \frac{1}{B_j^0} D_{ji}^x B_i^0 v_{i\ell}^0 A_{k\ell} P_{k\beta}^{V^*} \tag{20}$$

$$+ \frac{\Delta t}{1+\sigma\Delta t} \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} P_{\alpha j}^{X^*} \frac{1}{B_j^0} D_{ji}^{xx} B_i^0 v_{i\ell}^0 |A|_{k\ell} P_{k\beta}^{V^*}.$$
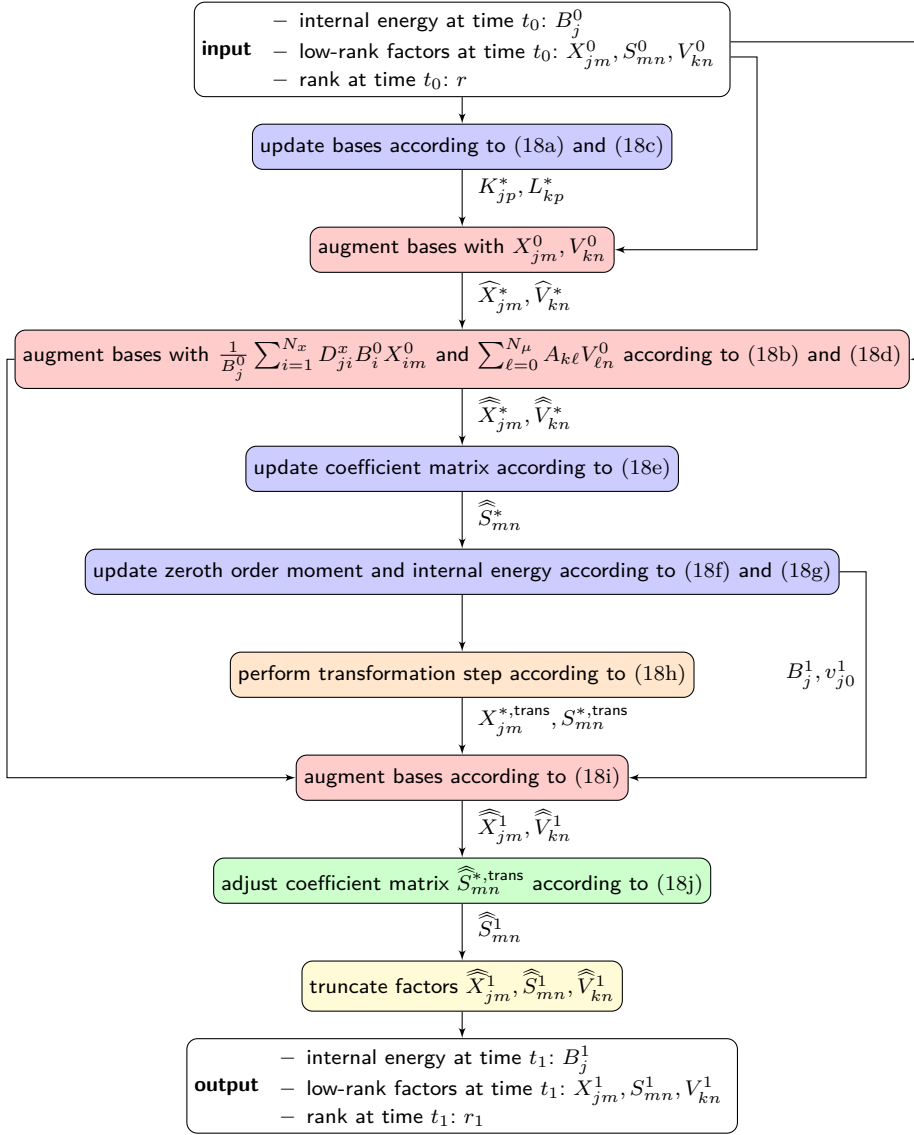
Fig. 1: Flowchart of the stable and conservative method (18).

Further, we denote $\widehat{\widehat{v}}^1_{\alpha\beta} := \sum_{p,q=1}^{3r+1} \widehat{\widehat{X}}^1_{\alpha q} \widehat{\widehat{S}}^1_{qp} \widehat{\widehat{V}}^1_{\beta p}$. From equation (18j), we have that

$$\frac{B^1_\alpha}{B^0_\alpha} \widehat{\widehat{v}}^1_{\alpha\beta} = v^*_{\alpha\beta} \left(1 - \delta_{\beta 0}\right) + \frac{B^1_\alpha}{B^0_\alpha} v_{\alpha 0} \delta_{\beta 0}.$$

We plug in (20) and (18f) and get

$$
\begin{aligned}
\frac{B_\alpha^1}{B_\alpha^0} \widehat{\widetilde{v}}^1_{\alpha\beta} \left(1 + \sigma \Delta t\right) = &\left( v^0_{\alpha\beta} - \Delta t \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} P^{X^*}_{\alpha j} \frac{1}{B_j^0} D^x_{ji} B_i^0 v^0_{i\ell} A_{k\ell} P^{V^*}_{k\beta} \right. \\
&+ \left. \Delta t \sum_{i,j=1}^{N_x} \sum_{k,\ell=0}^{N_\mu-1} P^{X^*}_{\alpha j} \frac{1}{B_j^0} D^{xx}_{ji} B_i^0 v^0_{i\ell} |A|_{k\ell} P^{V^*}_{k\beta} \right) \left(1 - \delta_{\beta0}\right) \\
&+ \left( v^0_{\alpha\beta} - \Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \frac{1}{B_\alpha^0} D^x_{\alpha i} B_i^0 v^0_{i\ell} A_{0\ell} \right. \\
&+ \left. \Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \frac{1}{B_\alpha^0} D^{xx}_{\alpha i} B_i^0 v^0_{i\ell} |A|_{0\ell} + \sqrt{2} \sigma \Delta t \frac{B_\alpha^1}{B_\alpha^0} \right) \delta_{\beta0}.
\end{aligned}
$$

We now use the fact that we have augmented the spatial basis according to (18b) and (18d). This allows us to write any function $h_i^0 \in \mathrm{span}\left(X_i^0\right)$ and $\tilde{h}_\ell^0 \in \mathrm{span}\left(V_\ell^0\right)$ as

$$
\sum_{i,j=1}^{N_x} P^{X^*}_{\alpha j} \frac{1}{B_j^0} D^x_{ji} B_i^0 h_i^0 = \frac{1}{B_\alpha^0} \sum_{i=1}^{N_x} D^x_{\alpha i} B_i^0 h_i^0 \quad \text{and} \quad \sum_{k,\ell=0}^{N_\mu-1} \tilde{h}_\ell^0 A_{k\ell} P^{V^*}_{k\beta} = \sum_{\ell=0}^{N_\mu-1} \tilde{h}_\ell^0 A_{\beta\ell}.
$$

The basis augmentations as well as the property of the projection operators hence enables us to obtain a representation of the form

$$
\frac{B_\alpha^1}{B_\alpha^0} \widehat{\widetilde{v}}^1_{\alpha\beta} \left(1 + \sigma \Delta t\right) = F \left(1 - \delta_{\beta0}\right) + F \delta_{\beta0} + \sqrt{2} \sigma \Delta t \frac{B_\alpha^1}{B_\alpha^0} \delta_{\beta0}
$$

with

$$
F = v^0_{\alpha\beta} - \Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \frac{1}{B_\alpha^0} D^x_{\alpha i} B_i^0 v^0_{i\ell} A_{0\ell} + \Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \frac{1}{B_\alpha^0} D^{xx}_{\alpha i} B_i^0 v^0_{i\ell} |A|_{0\ell}.
$$

We use that on the right-hand side $F\delta_{\beta0}$ cancels out such that we are left with

$$
\begin{aligned}
\frac{B_\alpha^1}{B_\alpha^0} \widehat{\widetilde{v}}^1_{\alpha\beta} \left(1 + \sigma \Delta t\right) = &v^0_{\alpha\beta} - \Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \frac{1}{B_\alpha^0} D^x_{\alpha i} B_i^0 v^0_{i\ell} A_{0\ell} \\
&+ \Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} \frac{1}{B_\alpha^0} D^{xx}_{\alpha i} B_i^0 v^0_{i\ell} |A|_{0\ell} \delta_{\beta0} + \sqrt{2} \sigma \Delta t \frac{B_\alpha^1}{B_\alpha^0} \delta_{\beta0}.
\end{aligned}
$$

In the next step, we multiply with $B_\alpha^1 B_\alpha^0 \widehat{\widehat{v}}_{\alpha\beta}^1$, sum over $\alpha$ and $\beta$, rearrange, and introduce the notation $\widehat{\widehat{u}}_{\alpha\beta}^1 = B_\alpha^1 \widehat{\widehat{v}}_{\alpha\beta}^1$. This leads to

$$
\sum_{\alpha=1}^{N_x} \sum_{\beta=0}^{N_\mu-1} \left( \widehat{\widehat{u}}_{\alpha\beta}^1 \right)^2 = \sum_{\alpha=1}^{N_x} \sum_{\beta=0}^{N_\mu-1} u_{\alpha\beta}^0 \widehat{\widehat{u}}_{\alpha\beta}^1 - \Delta t \sum_{i,\alpha=1}^{N_x} \sum_{\ell,\beta=0}^{N_\mu-1} \widehat{\widehat{u}}_{\alpha\beta}^1 D_{\alpha i}^x u_{i\ell}^0 A_{\beta\ell}
$$

$$
+ \Delta t \sum_{i,\alpha=1}^{N_x} \sum_{\ell,\beta=0}^{N_\mu-1} \widehat{\widehat{u}}_{\alpha\beta}^1 D_{\alpha i}^{xx} v_{i\ell}^0 |A|_{\beta\ell}
$$

$$
+ \sigma \Delta t \sum_{\alpha=1}^{N_x} \sum_{\beta=0}^{N_\mu-1} \left( B_\alpha^1 \right)^2 \widehat{\widehat{v}}_{\alpha\beta}^1 \left( \sqrt{2}\delta_{\beta 0} - \widehat{\widehat{v}}_{\alpha\beta}^1 \right).
$$

Inserting the relation

$$
\sum_{\alpha=1}^{N_x} \sum_{\beta=0}^{N_\mu-1} u_{\alpha\beta}^0 \widehat{\widehat{u}}_{\alpha\beta}^1 = \sum_{\alpha=1}^{N_x} \sum_{\beta=0}^{N_\mu-1} \left( \frac{1}{2} \left( \widehat{\widehat{u}}_{\alpha\beta}^1 \right)^2 + \frac{1}{2} \left( u_{\alpha\beta}^0 \right)^2 - \frac{1}{2} \left( \widehat{\widehat{u}}_{jk}^1 - u_{jk}^0 \right)^2 \right)
$$

then gives

$$
\frac{1}{2} \sum_{\alpha=1}^{N_x} \sum_{\beta=0}^{N_\mu-1} \left( \widehat{\widehat{u}}_{\alpha\beta}^1 \right)^2 = \frac{1}{2} \sum_{\alpha=1}^{N_x} \sum_{\beta=0}^{N_\mu-1} \left( u_{\alpha\beta}^0 \right)^2 - \frac{1}{2} \sum_{\alpha=1}^{N_x} \sum_{\beta=0}^{N_\mu-1} \left( \widehat{\widehat{u}}_{jk}^1 - u_{jk}^0 \right)^2
$$

$$
- \Delta t \sum_{i,\alpha=1}^{N_x} \sum_{\ell,\beta=0}^{N_\mu-1} \widehat{\widehat{u}}_{\alpha\beta}^1 D_{\alpha i}^x u_{i\ell}^0 A_{\beta\ell} + \Delta t \sum_{i,\alpha=1}^{N_x} \sum_{\ell,\beta=0}^{N_\mu-1} \widehat{\widehat{u}}_{\alpha\beta}^1 D_{\alpha i}^{xx} v_{i\ell}^0 |A|_{\beta\ell}
$$

$$
+ \sigma \Delta t \sum_{\alpha=1}^{N_x} \sum_{\beta=0}^{N_\mu-1} \left( B_\alpha^1 \right)^2 \widehat{\widehat{v}}_{\alpha\beta}^1 \left( \sqrt{2}\delta_{\beta 0} - \widehat{\widehat{v}}_{\alpha\beta}^1 \right).
$$

We estimate this expression as in the proof of Theorem 2 and add it with equation (19). Analogously to the proof of Theorem 2, and as the truncation step does not alter the zeroth order moment, we obtain that the DLRA scheme is energy stable under the time step restriction $\Delta t \leq \Delta x$.

5.4 Mass conservation

In addition, the DLRA scheme (18) can be shown to be mass conservative when using a suitable truncation strategy. We follow the ideas in [19,22,14] and adjust the truncation step such that it conserves the zeroth order moment. Different from [19] and as explained in [14], we do not need to adjust the $L$-step equation due to the usage of the augmented BUG integrator from [5]. Starting from the augmented quantities $\widehat{\widehat{\mathbf{X}}}^1$, $\widehat{\widehat{\mathbf{S}}}^1$ and $\widehat{\widehat{\mathbf{V}}}^1$, the conservative truncation strategy then works as follows:

1. We set $\widehat{\widehat{\mathbf{K}}}^1 = \widehat{\widehat{\mathbf{X}}}^1 \widehat{\widehat{\mathbf{S}}}^1$ and split it into two parts $\widehat{\widehat{\mathbf{K}}}^1 = [\widehat{\widehat{\mathbf{K}}}^{1,\mathrm{cons}}, \widehat{\widehat{\mathbf{K}}}^{1,\mathrm{rem}}]$, where $\widehat{\widehat{\mathbf{K}}}^{1,\mathrm{cons}}$ corresponds to the first and $\widehat{\widehat{\mathbf{K}}}^{1,\mathrm{rem}}$ to the remaining columns of $\widehat{\widehat{\mathbf{K}}}^1$. Analogously, we split $\widehat{\widehat{\mathbf{V}}}^1$ into $\widehat{\widehat{\mathbf{V}}}^1 = [\widehat{\widehat{\mathbf{V}}}^{1,\mathrm{cons}}, \widehat{\widehat{\mathbf{V}}}^{1,\mathrm{rem}}]$, where $\widehat{\widehat{\mathbf{V}}}^{1,\mathrm{cons}}$ corresponds to the first and $\widehat{\widehat{\mathbf{V}}}^{1,\mathrm{rem}}$ to the remaining columns of $\widehat{\widehat{\mathbf{V}}}^1$.

2. We compute $\widehat{\widehat{\mathbf{X}}}^{1,\text{cons}} = \widehat{\widehat{\mathbf{K}}}^{1,\text{cons}}/\|\widehat{\widehat{\mathbf{K}}}^{1,\text{cons}}\|$ and $\widehat{\widehat{\mathbf{S}}}^{1,\text{cons}} = \|\widehat{\widehat{\mathbf{K}}}^{1,\text{cons}}\|$.

3. We perform a QR-decomposition of $\widehat{\widehat{\mathbf{K}}}^{1,\text{rem}}$ to obtain $\widehat{\widehat{\mathbf{K}}}^{1,\text{rem}} = \widehat{\widehat{\mathbf{X}}}^{1,\text{rem}}\widehat{\widehat{\mathbf{S}}}^{1,\text{rem}}$.

4. We compute the singular value decomposition of $\widehat{\widehat{\mathbf{S}}}^{1,\text{rem}} = \widehat{\widehat{\mathbf{P}}}\boldsymbol{\Sigma}\widehat{\widehat{\mathbf{Q}}}^\top$ with $\boldsymbol{\Sigma} = \text{diag}(\sigma_j)$. The new rank $r_1 \leq 3r+1$ is chosen such that for a prescribed tolerance parameter $\vartheta$ it holds

$$\left( \sum_{j=r_1+1}^{3r+1} \sigma_j^2 \right)^{1/2} \leq \vartheta.$$

We set $\mathbf{S}^{1,\text{rem}} \in \mathbb{R}^{r_1 \times r_1}$ to be the matrix containing the $r_1$ largest singular values. For the update of the spatial and the directional basis we introduce the matrices $\widehat{\widehat{\mathbf{P}}}^{\text{rem}} \in \mathbb{R}^{(3r+1)\times r_1}$ and $\widehat{\widehat{\mathbf{Q}}}^{\text{rem}} \in \mathbb{R}^{(3r+1)\times r_1}$ containing the first $r_1$ columns of $\widehat{\widehat{\mathbf{P}}}$ and $\widehat{\widehat{\mathbf{Q}}}$, respectively, and set $\mathbf{X}^{1,\text{rem}} = \widehat{\widehat{\mathbf{X}}}^{1,\text{rem}}\widehat{\widehat{\mathbf{P}}}^{\text{rem}} \in \mathbb{R}^{N_x \times r_1}$ and $\mathbf{V}^{1,\text{rem}} = \widehat{\widehat{\mathbf{V}}}^{1,\text{rem}}\widehat{\widehat{\mathbf{Q}}}^{\text{rem}} \in \mathbb{R}^{N_\mu \times r_1}$.

5. We set $\widetilde{\mathbf{X}}^1 = [\widehat{\widehat{\mathbf{X}}}^{1,\text{cons}}, \mathbf{X}^{1,\text{rem}}]$ and $\widetilde{\mathbf{V}}^1 = [\mathbf{e}_1, \mathbf{V}^{1,\text{rem}}]$ and perform a QR-decomposition to obtain $\widetilde{\mathbf{X}}^1 = \mathbf{X}^1\mathbf{R}^1$ and $\widetilde{\mathbf{V}}^1 = \mathbf{V}^1\mathbf{R}^2$, respectively.

6. We compute

$$\mathbf{S}^1 = \mathbf{R}^1 \begin{bmatrix} \widehat{\widehat{\mathbf{S}}}^{1,\text{cons}} & 0 \\ 0 & \mathbf{S}^{1,\text{rem}} \end{bmatrix} \mathbf{R}^{2,\top}.$$

This leads to the updated solution $\mathbf{v}^1 = \mathbf{X}^1\mathbf{S}^1\mathbf{V}^{1,\top}$ after one time step at time $t_1 = t_0 + \Delta t$.

In order to show local mass conservation for the proposed DLRA scheme, we translate the macroscopic quantities given in Definition 1 to the fully discretized setting.

**Definition 4 (Fully discretized macroscopic quantities)** The *mass* of the fully discretized multiplicative Su-Olson problem at time $t_0$ is defined as

$$\rho_j^0 = \sqrt{2}B_j^0 v_{j0}^0 + B_j^0.$$

The *momentum* is given as

$$u_j^0 = \sqrt{2}\sum_{j=1}^{N_x} B_j v_{j\ell}^0 A_{0\ell}.$$

For $t_1 = t_0 + \Delta t$ the definitions shall hold analogously.

We can then show that the DLRA algorithm together with the conservative truncation strategy fulfills the following local conservation law.

**Theorem 4** *The DLRA scheme* (18) *together with the conservative truncation strategy is locally mass conservative, i.e. it fulfills the local conservation law*

$$\frac{1}{\Delta t}\left( \sqrt{2}B_j^1\Phi_j^1 + B_j^1 - \left( \sqrt{2}B_j^0\Phi_j^0 + B_j^0 \right) \right)$$
$$= -\sqrt{2}\sum_{i=1}^{N_x}\sum_{\ell=0}^{N_\mu-1} D_{ji}^x B_i^0 v_{i\ell}^0 A_{0\ell} + \sqrt{2}\sum_{i=1}^{N_x}\sum_{\ell=0}^{N_\mu-1} D_{ji}^{xx} B_i^0 v_{i\ell}^0 |A|_{0\ell},$$

where $\Phi_j^0 = \sum_{m,n=1}^{r} X_{jm}^0 S_{mn}^0 V_{0n}^0$, $\Phi_j^1 = \sum_{m,n=1}^{r_1} X_{jm}^1 S_{mn}^1 V_{0n}^1$ and $v_{jk}^0 = X_{jm}^0 S_{mn}^0 V_{kn}^0$. *This is a discretization of the continuous local conservation law given in* (5).

*Proof* We know that the conservative truncation strategy is designed such that it does not alter the zeroth order moment, i.e. it holds $\sum_{m,n=1}^{3r} \widehat{\widehat{X}}_{jm}^1 \widehat{\widehat{S}}_{mn}^1 \widehat{\widehat{V}}_{0m}^1 = v_{j0}^1$. In addition, we know from the basis augmentation (18i) and the correction step (18j) that $\sum_{m,n=1}^{3r} \widehat{\widehat{X}}_{jm}^1 \widehat{\widehat{S}}_{mn}^1 \widehat{\widehat{V}}_{0m}^1 = \sum_{m,n=1}^{r} X_{jm}^1 S_{mn}^1 V_{0n}^1$. Combining both, we have

$$\Phi_j^1 = \sum_{m,n=1}^{r_1} X_{jm}^1 S_{mn}^1 V_{0n}^1 = \sum_{m,n=1}^{3r} \widehat{\widehat{X}}_{jm}^1 \widehat{\widehat{S}}_{mn}^1 \widehat{\widehat{V}}_{0m}^1 = v_{j0}^1.$$

We insert this relation into the coupled equations (18f) and (18g). We multiply (18f) with $\sqrt{2}$, rearrange it and multiply both equations with $B_j^0$. This gives

$$\sqrt{2}B_j^1 \Phi_j^1 = \sqrt{2}B_j^0 \Phi_j^0 - \sqrt{2}\Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} D_{ji}^x B_i^0 \sum_{m,n=1}^{3r} \widehat{\widehat{X}}_{im}^* \widetilde{S}_{mn}^0 \widehat{\widehat{V}}_{\ell n}^* A_{0\ell} \qquad (21a)$$

$$+ \sqrt{2}\Delta t \sum_{i=1}^{N_x} \sum_{\ell=0}^{N_\mu-1} D_{ji}^{xx} B_i^0 \sum_{m,n=1}^{3r} \widehat{\widehat{X}}_{im}^* \widetilde{S}_{mn}^0 \widehat{\widehat{V}}_{\ell n}^* |A|_{0\ell} + \sigma \Delta t B_j^1 \left( 2 - \sqrt{2}\Phi_j^1 \right),$$

$$B_j^1 = B_j^0 + \sigma \Delta t B_j^1 \left( \sqrt{2}\Phi_j^1 - 2 \right). \qquad (21b)$$

Due to the basis augmentations with $\mathbf{X}^0$ and $\mathbf{V}^0$ from the augmented BUG integrator, we also know that

$$\sum_{m,n=1}^{3r} \widehat{\widehat{X}}_{im}^* \widetilde{S}_{mn}^0 \widehat{\widehat{V}}_{\ell n}^* = \sum_{m,n=1}^{r} X_{im}^0 S_{mn}^0 V_{\ell n}^0 = v_{i\ell}^0.$$

We insert this into (21a), add equation (21a) and (21b), and rearrange. This leads to the local conservation law (4), ensuring the local conservation of mass.

## 6 Numerical results

In this section, we compare the solution of the DLRA scheme (18) to the solution of the full equations (17) to underline the efficiency and accuracy of the proposed method. We give different test examples in one space and one directional dimension that validate our theoretical results.

### 6.1 1D plane source

We start with the one-dimensional plane source test case which is a common test example for thermal radiative transfer [20, 21, 1, 42]. We consider the spatial domain $D = [-10, 10]$ and the directional domain $[-1, 1]$. The initial condition shall be chosen as the cutoff Gaussian

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

where the constant deviation is given as $\sigma_{\mathrm{IC}} = 0.03$. The traveling particles are initially centered around $x = 1$ and move into all directions $\mu \in [-1, 1]$. The initial internal energy is set to $B^0 = 1$ and the opacity to the constant value $\sigma = 1$. For the low-rank computations we use an initial rank of $r = 20$. The total mass $m^n$ at time $t_n$ can be derived as $m^n = \Delta x \sum_j \left( \sqrt{2} B_j^n v_{j0}^n + B_j^n \right)$. As computational parameters we use $N_x = 1000$ grid points in the spatial domain and $N_\mu = 500$ moments for the approximation in the directional domain. The time step size is determined by $\Delta t = \mathrm{CFL} \cdot \Delta x$ with $\mathrm{CFL} = 0.99$. In Figure 2 we compare the solution of the DLRA scheme (DLRA) with the solution of the full system (full). We observe that the solution $f(x, \mu)$ as well as the scalar flux $\Phi = \frac{1}{\sqrt{2}} \langle f \rangle_\mu$ and the dimensionless temperature $T = \sqrt[4]{B}$ are captured well by the DLRA scheme. For a chosen tolerance parameter of $\vartheta = 10^{-1} \|\mathbf{\Sigma}\|_2$ the rank $r$ increases up to $r = 23$ before it significantly reduces again. The relative mass error $\frac{|m^0 - m^n|}{\|m^0\|}$ is of order $\mathcal{O}(10^{-13})$, i.e. the DLRA scheme is mass conservative up to machine precision. These results confirm our theoretical results.
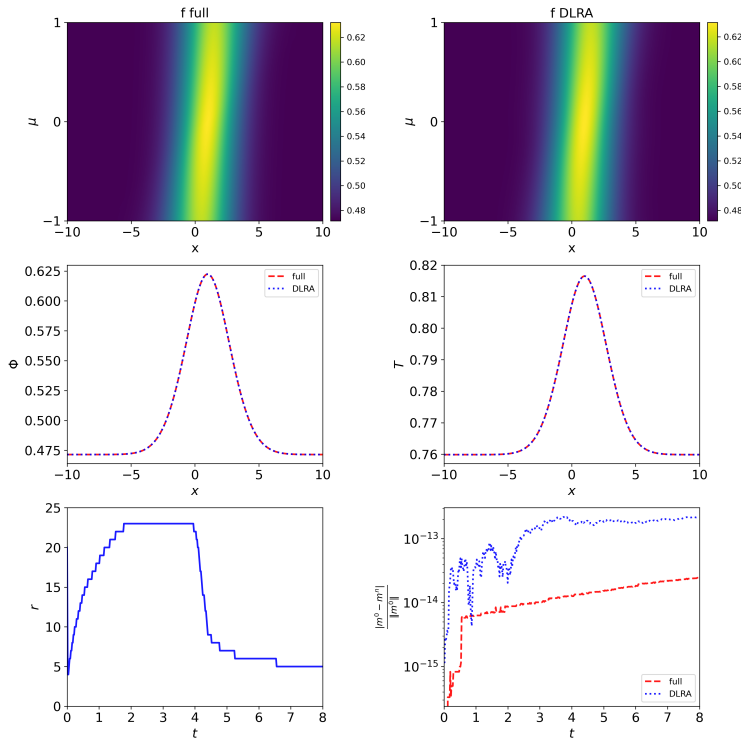


Fig. 2: Top row: Numerical results for the solution $B(x)g(x, \mu)$ of the plane source problem at time $t_{\mathrm{end}} = 8$ computed with the full solver (left) and the multiplicative DLRA scheme (right). Middle row: Scalar flux $\Phi$ (left) and temperature $T$ (right) for both the full system and the multiplicative DLRA scheme. Bottom row: Evolution of the rank in time for the multiplicative DLRA method (left) and relative mass error for both methods (right).

6.2 1D external source

In a second test example we take an external source term $Q(x)$ into account. It shall be added to the conservative form of the Su-Olson system as follows

$$\partial_t g(t,x,\mu) = -\frac{\mu}{B(t,x)} \partial_x \left(B(t,x)g(t,x,\mu)\right) + \sigma \left(1 - g(t,x,\mu)\right) - \frac{g(t,x,\mu)}{B(t,x)} \partial_t B(t,x)$$
$$+ \frac{Q(x)}{B(t,x)},$$
$$\partial_t B(t,x) = \sigma B(t,x) \left(\langle g(t,x,\mu)\rangle_\mu - 2\right).$$

This source term again generates radiation that moves through and interacts with the background material which in turn heats up and itself emits particles. The resulting travelling temperature wave is called a *Marshak wave* [33]. Again, we compare the solution of the full equations (17) (full) with the presented low-rank scheme (18) (DLRA). Both schemes are now adjusted such that the additional source term is taken into account. For our numerical example we choose the function $Q(x) = \chi_{[-0.5,0.5]}(x)/a$ with $a = \frac{4\sigma_{\mathrm{SB}}}{c}$ being the radiation constant. We again consider the spatial domain $D = [-10, 10]$ and the directional domain $[-1, 1]$. The initial condition shall be chosen as

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

where the constant deviation is given as $\sigma_{\mathrm{IC}} = 0.03$ and the particles move into all directions $\mu \in [-1, 1]$. The initial internal energy is set to $B^0 = 50$ and the opacity to the constant value $\sigma = 1$. For the low-rank computations we use an initial rank of $r = 20$. As computational parameters we use $N_x = 1000$ grid points in the spatial domain and $N_\mu = 500$ moments for the approximation in the directional domain. The time step size is determined by $\Delta t = \mathrm{CFL} \cdot \Delta x$ with $\mathrm{CFL} = 0.99$. Figure 3 then shows the numerical results for the solution $f(x,\mu)$, for the scalar flux $\Phi = \frac{1}{\sqrt{2}} \langle f \rangle_\mu$, and for the dimensionless temperature $T = \sqrt[4]{B}$, computed with both solvers. We again observe that the DLRA scheme captures the solution of the full system. The rank $r$ increases up to a value of $r = 23$ for a chosen tolerance parameter of $\vartheta = 10^{-3}\|\mathbf{\Sigma}\|_2$. Due to the additional source term, there is no conservation of mass in this test example.

## 7 Conclusion and outlook

We have presented a DLRA discretization for the multiplicative Su-Olson problem that is energy stable and mass conservative. To achieve both of these features, an additional basis augmentation in the augmented BUG integrator combined with an adjusted truncation step are performed. This enables us to give a mathematically rigorous stability analysis. Numerical test examples confirm the theoretical properties and validate the accuracy and computational advantages of the DLRA scheme. However, the extension of the considered stability analysis from a linear to a non-linear problem, for example the isothermal Boltzmann-BGK equation treated in [12], poses additional challenges as the general theoretical setting is significantly more difficult.
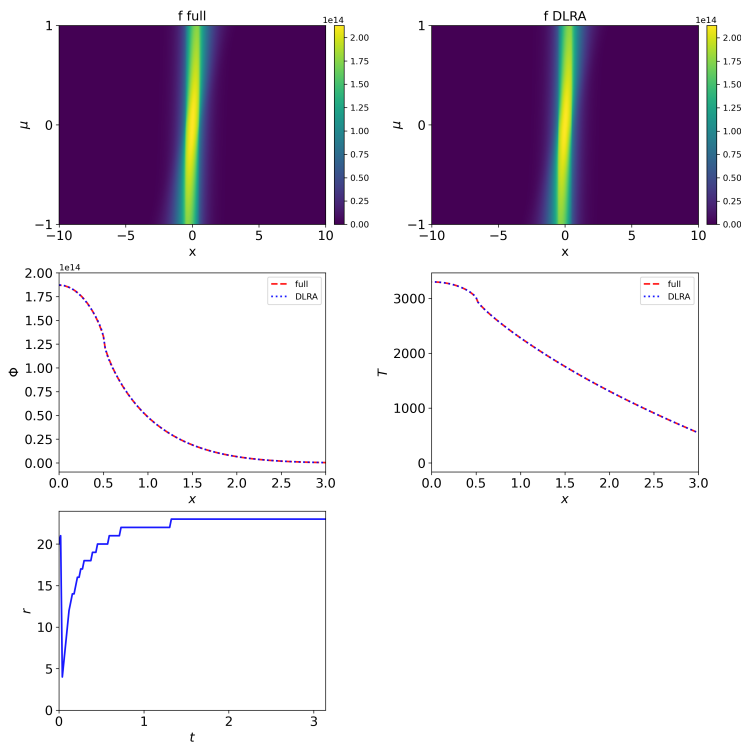
Fig. 3: Top row: Numerical results for the solution $B(x)g(x,\mu)$ of the Su-Olson problem at time $t_{\text{end}} = 3.16$ computed with the full solver (left) and the multiplicative DLRA scheme (right). Middle row: Scalar flux $\Phi$ (left) and temperature $T$ (right) for both the full system and the multiplicative DLRA scheme. Bottom row: Evolution of the rank in time for the multiplicative DLRA method.

For example, to our knowledge, the theoretical framework we have used here is not available in the non-linear case. Nevertheless, the analysis performed in this paper provides valuable insight into the choice of a suitable space discretization and stabilization when considering a multiplicative splitting of the distribution function. This splitting approach can be extremely useful for the construction of DLRA schemes for more complicated problems, which we consider future work.

**Acknowledgements**

## Declarations

Not applicable. The authors have no relevant financial or non-financial interests to disclose.

## References

1. L. Baumann, L. Einkemmer, C. Klingenberg, and J. Kusch. Energy stable and conservative dynamical low-rank approximation for the Su-Olson problem. *SIAM J. Sci. Comput.*, 46(2):B137–B158, 2024.
2. L. Baumann, L. Einkemmer, C. Klingenberg, and J. Kusch. A stable multiplicative dynamical low-rank discretization for the linear Boltzmann-BGK equation. *arXiv preprint arXiv:2411.06844*, 2024.
3. G. I. Bell and S. Glasstone. *Nuclear Reactor Theory*. Van Nostrand Reinhold Company, New York, 1970.
4. G. Ceruti, L. Einkemmer, J. Kusch, and C. Lubich. A robust second-order low-rank BUG integrator based on the midpoint rule. *BIT Numer. Math.*, 64(30), 2024.
5. G. Ceruti, J. Kusch, and C. Lubich. A rank-adaptive robust integrator for dynamical low-rank approximation. *BIT Numer. Math.*, 62:1149–1174, 2022.
6. G. Ceruti, J. Kusch, and C. Lubich. A parallel rank-adaptive integrator for dynamical low-rank approximation. *SIAM J. Sci. Comput.*, 46(3):B205–B228, 2024.
7. G. Ceruti and C. Lubich. An unconventional robust integrator for dynamical low-rank approximation. *BIT Numer. Math.*, 62(1):23–44, 2022.
8. S. Dargaville, A. Buchan, R. Smedley-Stevenson, P. Smith, and C. Pain. Angular adaptivity with spherical harmonics for Boltzmann transport. *J. Comput. Phys.*, 397:108846, 2019.
9. Z. Ding, L. Einkemmer, and Q. Li. Dynamical low-rank integrator for the linear Boltzmann equation: error analysis in the diffusion limit. *SIAM J. Numer. Anal.*, 59(4):2254–2285, 2021.
10. L. Einkemmer. A low-rank algorithm for weakly compressible flow. *SIAM J. Sci. Comput.*, 41(5):A2795–A2814, 2019.
11. L. Einkemmer, J. Hu, and J. Kusch. Asymptotic-preserving and energy stable dynamical low-rank approximation. *SIAM J. Numer. Anal.*, 62(1):73–92, 2024.
12. 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 J. Sci. Comput.*, 43(5):B1057–B1080, 2021.
13. L. Einkemmer and I. Joseph. A mass, momentum, and energy conservative dynamical low-rank scheme for the Vlasov equation. *J. Comput. Phys.*, 443:110493, 2021.
14. L. Einkemmer, J. Kusch, and S. Schotthöfer. Conservation properties of the augmented basis update & Galerkin integrator for kinetic problems. *arXiv preprint arXiv:2311.06399*, 2023.
15. L. Einkemmer and C. Lubich. A low-rank projector-splitting integrator for the Vlasov-Poisson equation. *SIAM J. Sci. Comput.*, 40(5):B1330–B1360, 2018.
16. L. Einkemmer and C. Lubich. A quasi-conservative dynamical low-rank algorithm for the Vlasov equation. *SIAM J. Sci. Comput.*, 41(5):B1061–B1081, 2019.
17. L. Einkemmer, J. Mangott, and M. Prugger. A low-rank complexity reduction algorithm for the high-dimensional kinetic chemical master equation. *J. Comput. Phys.*, 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. *J. Comput. Phys.*, 403:109063, 2020.
19. L. Einkemmer, A. Ostermann, and C. Scalone. A robust and conservative dynamical low-rank algorithm. *J. Comput. Phys.*, 484:112060, 2023.
20. B. Ganapol, R. Baker, and J. Dahl. Homogeneous infinite media time-dependent analytical benchmarks. *Los Alamos Natl. Lab.*, 2001.
21. B. D. Ganapol. Analytical benchmarks for nuclear engineering applications. Case studies in neutron transport theory. *Nucl. energy agency, Organ. for econ. coop. and dev.*, 2008.
22. W. Guo and J.-M. Qiu. A conservative low rank tensor method for the Vlasov dynamics. *SIAM J. Sci. Comput.*, 46(1):A232–A263, 2024.

23. E. Hairer and G. Wanner. *Solving Ordinary Differential Equations II: Stiff and Differential-Algebraic Problems*. Springer Berlin, Heidelberg, 1996.
24. J. Hu and Y. Wang. An adaptive dynamical low rank method for the nonlinear Boltzmann equation. *J. Sci. Comput.*, 92:75, 2022.
25. E. Kieri, C. Lubich, and H. Walach. Discretized dynamical low-rank approximation in the presence of small singular values. *SIAM J. Numer. Anal.*, 54(2):1020–1038, 2016.
26. O. Koch and C. Lubich. Dynamical low-rank approximation. *SIAM J. Matrix Anal. Appl.*, 29(2):434–454, 2007.
27. J. Kusch. Second-order robust parallel integrators for dynamical low-rank approximation. *arXiv preprint arXiv:2403.02834*, 2024.
28. 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.
29. J. Kusch, B. Whewell, R. McClarren, and M. Frank. A low-rank power iteration scheme for neutron transport criticality problems. *J. Comput. Phys.*, 470:111587, 2022.
30. V. M. Laboure, R. G. McClarren, and C. D. Hauck. Implicit filtered $P_N$ for high-energy density thermal radiation transport using discontinuous Galerkin finite elements. *J. Comput. Phys.*, 321:624—-643, 2016.
31. R. J. LeVeque. *Finite Volume Methods for Hyperbolic Problems*. Cambridge University Press, Cambridge, 2002.
32. C. Lubich and I. V. Oseledets. A projector-splitting integrator for dynamical low-rank approximation. *BIT Numer. Math.*, 54(1):171–188, 2014.
33. R. E. Marshak. Effect of radiation on shock wave behavior. *Phys. Fluids*, 1:24–29, 1958.
34. R. G. McClarren, T. M. Evans, R. B. Lowrie, and J. D. Densmore. Semi-implicit time integration for $P_N$ thermal radiative transfer. *J. Comput. Phys.*, 227:7561–7586, 2008.
35. R. G. McClarren and C. D. Hauck. Robust and accurate filtered spherical harmonics expansions for radiative transfer. *J. Comput. Phys.*, 229:5597–5614, 2010.
36. R. G. McClarren, J. P. Holloway, and T. A. Brunner. Analytic $P_1$, solutions for time-dependant, thermal radiative transfer in several geometries. *J. Quant. Spectrosc. Radiat. Transf.*, 109:389–403, 2008.
37. R. G. McClarren, J. P. Holloway, and T. A. Brunner. On solutions to the $P_n$ equations for thermal radiative transfer. *J. Comput. Phys.*, 227:2864–2885, 2008.
38. G. L. Olson, L. H. Auer, and M. L. Hall. Diffusion, $P_1$, and other approximate forms of radiation transport. *J. Quant. Spectrosc. Radiat. Transf.*, 62:619–634, 2000.
39. C. Patwardhan, M. Frank, and J. Kusch. Asymptotic-preserving and energy stable dynamical low-rank approximation for thermal radiative transfer equations. *Multiscale Model. Simul.*, 23(1):278–312, 2025.
40. Z. Peng and R. G. McClarren. A high-order/low-order (HOLO) algorithm for preserving conservation in time-dependent low-rank transport calculations. *J. Comput. Phys.*, 447:110672, 2021.
41. Z. Peng and R. G. McClarren. A sweep-based low-rank method for the discrete ordinate transport equation. *J. Comput. Phys.*, 473:111748, 2023.
42. Z. Peng, R. G. McClarren, and M. Frank. A low-rank method for two-dimensional time-dependent radiation transport calculations. *J. Comput. Phys.*, 421:109735, 2020.
43. G. C. Pomraning. The non-equilibrium marshak wave problem. *J. Quant. Spectrosc. Radiat. Transf.*, 21:249–261, 1979.
44. M. Prugger, L. Einkemmer, and C. F. Lopez. A dynamical low-rank approach to solve the chemical master equation for biological reaction networks. *J. Comput. Phys.*, 489:112250, 2023.
45. B. Su and G. L. Olson. An analytical benchmark for non-equilibrium radiative transfer in an isotropically scattering medium. *Ann. Nucl. Energy*, 24(13):1035–1055, 1997.
46. P. Yin, E. Endeve, C. D. Hauck, and S. R. Schnake. Towards dynamical low-rank approximation for neutrino kinetic equations. Part I: Analysis of an idealized relaxation model. *Math. Comput.*, 2024.
47. W. Zheng and R. G. McClarren. Moment closures based on minimizing the residual of the $P_N$ angular expansion in radiation transport. *J. Comput. Phys.*, 314:682—-699, 2016.