

ENERGY STABLE AND CONSERVATIVE DYNAMICAL LOW-RANK APPROXIMATION FOR THE SU–OLSON PROBLEM*

LENA BAUMANN[†], LUKAS EINKEMMER[‡], CHRISTIAN KLINGENBERG[†], AND
JONAS KUSCH[§]

Abstract. Computational methods for thermal radiative transfer problems exhibit high computational costs and a prohibitive memory footprint when the spatial and directional domains are finely resolved. A strategy to reduce such computational costs is dynamical low-rank approximation (DLRA), which represents and evolves the solution on a low-rank manifold, thereby significantly decreasing computational and memory requirements. Efficient discretizations for the DLRA evolution equations need to be carefully constructed to guarantee stability while enabling mass conservation. In this work, we focus on the Su–Olson closure leading to a linearized internal energy model and derive a stable discretization through an implicit coupling of internal energy and particle density. Moreover, we propose a rank-adaptive strategy to preserve local mass conservation. Numerical results are presented which showcase the accuracy and efficiency of the proposed low-rank method compared to the solution of the full system.

Key words. thermal radiative transfer, Su–Olson closure, dynamical low-rank approximation, energy stability, mass conservation, rank adaptivity

MSC codes. 35L65, 35Q49, 65M12, 65M22

DOI. 10.1137/23M1586215

1. Introduction. Numerically solving the radiative transfer equations is a challenging task, especially due to the high dimensionality of the solution’s phase space. A common strategy to tackle this issue is to choose coarse numerical discretizations and mitigate numerical artifacts [23, 27, 32] which arise due to the insufficient resolution; see, e.g., [3, 15, 1, 24, 39]. Despite the success of these approaches in a large number of applications, the requirement of picking user-determined and problem dependent tuning parameters can render them impracticable. Another approach to deal with the problem’s high dimensionality is the use of model order reduction techniques. A reduced order method which is gaining a considerable amount of attention in the field of radiation transport is dynamical low-rank approximation (DLRA) [20] due to its ability to yield accurate solutions while not requiring an expensive offline training phase. DLRA’s core idea is to approximate the solution on a low-rank manifold and evolve it accordingly. Past work in the area of radiative transfer has focused on asymptotic-preserving schemes [10, 9], mass conservation [34], stable discretizations [21], imposing boundary conditions [22, 18], and implicit time discretizations [35].

*Submitted to the journal’s Software, High-Performance Computing, and Computational Science and Engineering section July 13, 2023; accepted for publication (in revised form) February 1, 2024; published electronically April 18, 2024.

<https://doi.org/10.1137/23M1586215>

Funding: The work of the first author was supported by the Würzburg Mathematics Center for Communication and Interaction (WMCCI) as well as the Stiftung der Deutschen Wirtschaft (Foundation of German Business). The work of the fourth author was supported by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) - 491976834.

[†]Department of Mathematics, University of Wuerzburg, Wuerzburg, DE-97074, Germany (lena.baumann@uni-wuerzburg.de, klingen@mathematik.uni-wuerzburg.de).

[‡]Numerical Analysis and Scientific Computing, University of Innsbruck, Innsbruck, A-6020, Austria (lukas.einkemmer@uibk.ac.at).

[§]Scientific Computing, Norwegian University of Life Sciences, Ås, NO-1432, Norway (jonas.kusch@nmbu.no).

A discontinuous Galerkin discretization of the DLRA evolution equations for thermal radiative transfer has been proposed in [5].

A key building block of efficient, accurate, and stable methods for DLRA is the construction of time integrators which are robust irrespective of small singular values in the solution [19]. Three integrators that move on the low-rank manifold while not being restricted by its curvature are the *projector-splitting* (PS) integrator [25], the *basis update & Galerkin* (BUG) integrator [8], and the *parallel* integrator [7]. Since the PS integrator evolves one of the required subflows backward in time, the BUG and parallel integrators are preferable for diffusive problems while facilitating the construction of stable numerical discretization for hyperbolic problems [21]. Moreover, the BUG integrator allows for a basis augmentation step [6], which can be used to construct conservative schemes for the Schrödinger equation [6] and the Vlasov–Poisson equations [14].

In this work we consider the thermal radiative transfer equations using the Su–Olson closure. This leads to a linearized internal energy model for which we propose an energy stable and mass conservative DLRA scheme. The main novelties of this paper are as follows:

- *A stable numerical scheme for thermal radiative transfer:* We show that a naive IMEX scheme fails to guarantee energy stability. To overcome this unphysical behavior we propose a scheme which advances radiation and internal energy implicitly in a coupled fashion. In addition, our novel analysis gives a classic hyperbolic CFL condition that enables us to operate up to a time step size of $\Delta t = \text{CFL} \cdot \Delta x$.
- *A mass conservative and rank-adaptive integrator:* We employ the basis augmentation step from [6] as well as an adaptation of the conservative truncation strategy from [14, 17] to guarantee local mass conservation and rank adaptivity. In contrast to [14, 17] we do not need to impose conservation through a modified L -step equation but solely use the basis augmentation strategy from [6].

Both of these properties are extremely important as they ensure key physical principles and allow us to choose an optimal time step size which reduces the computational effort. Moreover, we demonstrate numerical experiments which underline the derived stability and conservation properties of the proposed low-rank method while showing significantly reduced computational costs and memory requirements compared to the full-order system.

This paper is structured as follows: After the introduction in section 1, we review the background on thermal radiative transfer and dynamical low-rank approximation in section 2. In section 3 we present the evolution equations for the thermal radiative transfer equations when using the rank-adaptive BUG integrator. Section 4 discretizes the resulting equations in angle and space. The main method is presented in section 5, where a stable time discretization is proposed. We discuss local mass conservation of the scheme in section 6. Numerical experiments are demonstrated in section 7.

2. Background.

2.1. Thermal radiative transfer. In this work, we study radiation particles moving through and interacting with a background material. By absorbing particles, the material heats up and emits new particles which can in turn again interact with the background. This process is described by the thermal radiative transport equations

$$\begin{aligned} \frac{1}{c} \partial_t f(t, x, \mu) + \mu \partial_x f(t, x, \mu) &= \sigma(B(t, x) - f(t, x, \mu)), \\ \partial_t e(t, x) &= \sigma(\langle f(t, x, \cdot) \rangle_\mu - B(t, x)), \end{aligned}$$

where we omit boundary and initial conditions for now. This system can be solved for the particle density $f(t, x, \mu)$ and the internal energy $e(t, x)$ of the background medium. Here, $x \in D \subset \mathbb{R}$ is the spatial variable and $\mu \in [-1, 1]$ denotes the directional (or velocity) variable. The opacity σ encodes the rate at which particles are absorbed by the medium, and we use brackets $\langle \cdot \rangle_\mu, \langle \cdot \rangle_x$ to indicate an integration over the directional domain and the spatial domain, respectively. Moreover, the speed of light is denoted by c , and the black body radiation at the material temperature T is denoted by $B(T)$. It often is described by the Stefan–Boltzmann law

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

where $a = \frac{4\sigma_{\text{SB}}}{c}$ is the radiation density constant and σ_{SB} is the Stefan–Boltzmann constant. Different closures exist to determine a relation between the temperature T and the internal energy e . Following the ideas of Pomraning [37] and Su and Olson [38], we assume $e(T) = \alpha B(T)$. Without loss of generality we set $\alpha = 1$ and obtain

$$(2.1a) \quad \partial_t f(t, x, \mu) + \mu \partial_x f(t, x, \mu) = \sigma(B(t, x) - f(t, x, \mu)),$$

$$(2.1b) \quad \partial_t B(t, x) = \sigma(\langle f(t, x, \cdot) \rangle_\mu - B(t, x)).$$

We call this system the Su–Olson problem. It is a linear system for the particle density f and the internal energy B that is analytically solvable and serves as a common benchmark for numerical considerations [33, 30, 31, 28]. Note that we leave out the speed of light by doing a rescaling of time $\tau = t/c$ and in an abuse of notation use t to denote τ in the remainder. Constructing numerical schemes to solve the above equation is challenging. First, the potentially stiff opacity term has to be treated by an implicit time integration scheme. Second, for three-dimensional spatial domains the computational costs and memory requirements of finely resolved spatial and angular discretizations become prohibitive. To tackle the high dimensionality, we choose a dynamical low-rank approximation which we introduce in the following.

2.2. Dynamical low-rank approximation. The core idea of DLRA is to approximate the solution of a given equation $\partial_t f(t, x, \mu) = F(f(t, x, \mu))$ by a representation of the form

$$(2.2) \quad f(t, x, \mu) \approx \sum_{i,j=1}^r X_i(t, x) S_{ij}(t) V_j(t, \mu),$$

where the orthonormal functions $\{X_i : i = 1, \dots, r\}$ depend only on t and x and the orthonormal functions $\{V_j : j = 1, \dots, r\}$ depend only on t and μ . The number of basis functions is set to r and we call r the rank of this approximation. This terminology stems from the matrix setting for which the concept of DLRA has been introduced [20]. Then, (2.2) can be interpreted as a continuous analogue to the singular value decomposition for matrices. As representation (2.2) is not unique, we impose the Gauge conditions $\langle \dot{X}_i, X_j \rangle_x = 0$ and $\langle \dot{V}_i, V_j \rangle_\mu = 0$, from which we can conclude that $\{X_i\}$ and $\{V_j\}$ are uniquely determined for invertible $\mathbf{S} = (S_{ij}) \in \mathbb{R}^{r \times r}$ [20, 10, 13]. That is, we seek an approximation of f that for each time t lies in the manifold

$$\mathcal{M}_r = \left\{ f \in L^2(D \times [-1, 1]) : f(\cdot, x, \mu) = \sum_{i,j=1}^r X_i(\cdot, x) S_{ij}(\cdot) V_j(\cdot, \mu) \text{ with invertible} \right. \\ \left. \mathbf{S} = (S_{ij}) \in \mathbb{R}^{r \times r}, X_i \in L^2(D), V_j \in L^2([-1, 1]) \text{ and } \langle X_i, X_j \rangle_x = \delta_{ij}, \right. \\ \left. \langle V_i, V_j \rangle_\mu = \delta_{ij} \right\}.$$

Note that in the following we denote the full rank and the low-rank solutions as f . Let $f(t, \cdot, \cdot)$ be a path on \mathcal{M}_r . A formal differentiation of f with respect to t leads to

$$\dot{f}(t, \cdot, \cdot) = \sum_{i,j=1}^r \left(\dot{X}_i(t, \cdot) S_{ij}(t) V_j(t, \cdot) + X_i(t, \cdot) \dot{S}_{ij}(t) V_j(t, \cdot) + X_i(t, \cdot) S_{ij}(t) \dot{V}_j(t, \cdot) \right).$$

These functions restrict the solution dynamics onto the low-rank manifold \mathcal{M}_r and constitute the corresponding tangent space which under the Gauge conditions reads

$$\begin{aligned} \mathcal{T}_f \mathcal{M}_r = \left\{ \dot{f} \in L^2(D \times [-1, 1]) : \dot{f}(\cdot, x, \mu) = \sum_{i,j=1}^r \left(\dot{X}_i(\cdot, x) S_{ij}(\cdot) V_j(\cdot, \mu) \right. \right. \\ \left. \left. + X_i(\cdot, x) \dot{S}_{ij}(\cdot) V_j(\cdot, \mu) + X_i(\cdot, x) S_{ij}(\cdot) \dot{V}_j(\cdot, \mu) \right) \right. \\ \left. \text{with } \dot{S}_{ij} \in \mathbb{R}, \dot{X}_i \in L^2(D), \dot{V}_j \in L^2([-1, 1]) \text{ and } \langle \dot{X}_i, X_j \rangle_x = 0, \right. \\ \left. \langle \dot{V}_i, V_j \rangle_\mu = 0 \right\}. \end{aligned}$$

Having defined the low-rank manifold and its corresponding tangent space, we now wish to determine $f(t, \cdot, \cdot) \in \mathcal{M}_r$ such that $\partial_t f(t, \cdot, \cdot) \in \mathcal{T}_f \mathcal{M}_r$ and $\|\partial_t f(t, \cdot, \cdot) - F(f(t, \cdot, \cdot))\|_{L^2(D \times [-1, 1])}$ is minimized. That is, one wishes to determine f such that

$$(2.3) \quad \langle \partial_t f(t, \cdot, \cdot) - F(f(t, \cdot, \cdot)), \dot{f} \rangle_{x, \mu} = 0 \quad \text{for all } \dot{f} \in \mathcal{T}_f \mathcal{M}_r.$$

The orthogonal projector onto the tangent plane $\mathcal{T}_f \mathcal{M}_r$ can be explicitly given as

$$P(f)F(f) = \sum_{j=1}^r \langle V_j, F(f) \rangle_\mu V_j - \sum_{i,j=1}^r X_i \langle X_i V_j, F(f) \rangle_{x, \mu} V_j + \sum_{i=1}^r X_i \langle X_i, F(f) \rangle_x.$$

With this definition at hand, we can reformulate (2.3) as

$$\partial_t f(t, x, \mu) = P(f(t, x, \mu))F(f(t, x, \mu)).$$

To evolve the approximation of the solution in time according to the above equation is not trivial. Indeed standard time integration schemes suffer from the curvature of the low-rank manifold, which is proportional to the smallest singular value of the low-rank solution [20]. Three integrators which move along the manifold without suffering from its high curvature exist: the projector–splitting integrator [25], the BUG integrator [8], and the parallel integrator [7]. In this work, we will use the basis-augmented extension to the BUG integrator [6], which we explain in the following.

The rank-adaptive BUG integrator [6] updates and augments the bases $\{X_i\}, \{V_j\}$ in parallel in the first two steps. In the third step, a Galerkin step is performed for the augmented bases followed by a truncation step to a new rank r_1 . In detail, to evolve the approximation of the distribution function from $f(t_0, x, \mu) = \sum_{i,j=1}^r X_i^0(x) S_{ij}^0 V_j^0(\mu)$ at time t_0 to $f(t_1, x, \mu) = \sum_{i,j=1}^{r_1} X_i^1(x) S_{ij}^1 V_j^1(\mu)$ at time $t_1 = t_0 + \Delta t$ the integrator performs the following steps.

K-step. Write $K_j(t, x) = \sum_{i=1}^r X_i(t, x) S_{ij}(t)$. Then we obtain the representation $f(t, x, \mu) = \sum_{j=1}^r K_j(t, x) V_j^0(\mu)$ with $\{V_j^0\}$ kept fixed in this step. The basis functions $X_i^0(x)$ with $i = 1, \dots, r$ are updated by solving the partial differential equation

$$\partial_t K_j(t, x) = \left\langle V_j^0, F \left(\sum_{k=1}^r K_k(t, x) V_k^0 \right) \right\rangle_\mu, \quad K_j(t_0, x) = \sum_{i=1}^r X_i^0(x) S_{ij}^0,$$

and applying Gram–Schmidt to $[K_j(t_1, x), X_i^0] = \sum_{i=1}^{2r} \widehat{X}_i^1(x) R_{ij}^1$. Then, the updated and augmented basis in physical space consists of $\widehat{X}_i^1(x)$ with $i = 1, \dots, 2r$. Note that R_{ij}^1 is discarded after this step. Compute $\widehat{M}_{ki} = \langle \widehat{X}_k^1, X_i^0 \rangle_x$.

L-step. Write $L_i(t, \mu) = \sum_{j=1}^r S_{ij}(t) V_j(t, \mu)$. Then we obtain the representation $f(t, x, \mu) = \sum_{i=1}^r X_i^0 L_i(t, \mu)$ with $\{X_i^0\}$ kept fixed in this step. The basis functions $V_j^0(\mu)$ with $j = 1, \dots, r$ are updated by solving the partial differential equation

$$\partial_t L_i(t, \mu) = \left\langle X_i^0, F \left(\sum_{\ell=1}^r X_\ell^0 L_\ell(t, \mu) \right) \right\rangle_x, \quad L_i(t_0, \mu) = \sum_{j=1}^r S_{ij}^0 V_j^0(\mu),$$

and applying Gram–Schmidt to $[L_i(t_1, \mu), V_j^0(\mu)] = \sum_{j=1}^{2r} \widehat{V}_j^1(\mu) R_{ij}^2$. Then, the updated and augmented basis in velocity space consists of $\widehat{V}_j^1(\mu)$ with $j = 1, \dots, 2r$. Note that R_{ij}^2 is discarded after this step. Compute $\widehat{N}_{\ell j} = \langle \widehat{V}_\ell^1, V_j^0 \rangle_\mu$.

S-step. Update S_{ij}^0 with $i, j = 1, \dots, r$ to \widehat{S}_{ij}^1 with $i, j = 1, \dots, 2r$ by solving the ordinary differential equation

$$\dot{\widehat{S}}_{ij}(t) = \left\langle \widehat{X}_i^1 \widehat{V}_j^1, F \left(\sum_{\ell, k=1}^{2r} \widehat{X}_\ell^1 \widehat{S}_{\ell k}(t) \widehat{V}_k^1 \right) \right\rangle_{x, \mu}, \quad \widehat{S}_{ij}(t_0) = \sum_{k, \ell=1}^r \widehat{M}_{ik} S_{k\ell}^0 \widehat{N}_{j\ell}.$$

Truncation. Let \widehat{S}_{ij}^1 be the entries of the matrix $\widehat{\mathbf{S}}^1$. Compute the singular value decomposition of $\widehat{\mathbf{S}}^1 = \widehat{\mathbf{P}} \widehat{\mathbf{\Sigma}} \widehat{\mathbf{Q}}^\top$ with $\widehat{\mathbf{\Sigma}} = \text{diag}(\sigma_j)$. Given a tolerance ϑ , choose the new rank $r_1 \leq 2r$ as the minimal number such that

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

Let \mathbf{S}^1 with entries S_{ij}^1 be the $r_1 \times r_1$ diagonal matrix with the r_1 largest singular values and let \mathbf{P}^1 with entries P_{ij}^1 and \mathbf{Q}^1 with entries Q_{ji}^1 contain the first r_1 columns of $\widehat{\mathbf{P}}$ and $\widehat{\mathbf{Q}}$, respectively. Set $X_i^1(x) = \sum_{i=1}^{2r} \widehat{X}_i^1(x) P_{ij}^1$ for $i = 1, \dots, r_1$ and $V_j^1(\mu) = \sum_{j=1}^{2r} \widehat{V}_j^1(\mu) Q_{ji}^1$ for $j = 1, \dots, r_1$.

The updated approximation of the solution after one time step is then given by $f(t_1, x, \mu) = \sum_{i,j=1}^{r_1} X_i^1(x) S_{ij}^1 V_j^1(\mu)$. Note that we are not limited to augmenting with the old basis, which we will use to construct our scheme.

3. Dynamical low-rank approximation for Su–Olson. Let us now derive the evolution equations of the rank-adaptive BUG integrator for system (2.1), i.e., the partial differential equations appearing in the K - and L -steps and the ordinary differential equation for the S -step. To simplify notation, all derivations are performed for one spatial and one directional variable. However, the derivation trivially extends to higher dimensions. We start with considering the evolution equations for the low-rank approximation of the particle density (2.1a).

K-step. Write $K_j(t, x) = \sum_{i=1}^r X_i(t, x) S_{ij}(t)$. Then we have the representation $f(t, x, \mu) = \sum_{j=1}^r K_j(t, x) V_j^0(\mu)$ for the low-rank approximation of the solution. Again $\{V_j^0\}$ denotes the set of orthonormal basis functions for the velocity space that shall be kept fixed in this step. Inserting this representation of f into (2.1a) and projecting onto $V_k^0(\mu)$ gives the partial differential equation

$$(3.1) \quad \partial_t K_k(t, x) = - \sum_{j=1}^r \partial_x K_j(t, x) \langle V_k^0, \mu V_j^0 \rangle_\mu + \sigma (B(t, x) \langle V_k^0 \rangle_\mu - K_k(t, x)).$$

L-step. Write $L_i(t, \mu) = \sum_{j=1}^r S_{ij}(t) V_j(t, \mu)$. Then we have the representation $f(t, x, \mu) = \sum_{i=1}^r X_i^0(x) L_i(t, \mu)$ for the low-rank approximation of the solution. Again $\{X_i^0\}$ denotes the set of spatial orthonormal basis functions that shall be kept fixed in this step. Inserting this representation of f into (2.1a) and projecting onto $X_k^0(x)$ yields the partial differential equation

$$(3.2) \quad \partial_t L_k(t, \mu) = -\mu \sum_{i=1}^r \left\langle X_k^0, \frac{d}{dx} X_i^0 \right\rangle_x L_i(t, \mu) + \sigma \left(\langle X_k^0, B(t, \cdot) \rangle_x - L_k(t, \mu) \right).$$

Lastly, we derive the augmented Galerkin step of the rank-adaptive BUG integrator. We denote the time updated spatial basis augmented with X_i^0 as \widehat{X}_i^1 . The augmented directional basis \widehat{V}_i^1 is constructed in the corresponding way. Then, the augmented Galerkin step is constructed according to the following step.

S-step. We use the initial condition $\widehat{S}_{ij}(t_0) = \sum_{\ell, k=1}^r \langle \widehat{X}_i^1 X_\ell^0 \rangle_x S_{\ell k}(t_0) \langle \widehat{V}_j^1 V_k^0 \rangle_\mu$ and approximate the solution f as $f(t, x, \mu) = \sum_{i,j=1}^{2r} \widehat{X}_i^1(x) \widehat{S}_{ij}(t) \widehat{V}_j^1(\mu)$. Inserting this representation into (2.1a) and testing against \widehat{X}_k^1 and \widehat{V}_ℓ^1 gives the ordinary differential equation

$$(3.3) \quad \dot{\widehat{S}}_{k\ell}(t) = - \sum_{i,j=1}^{2r} \left\langle \widehat{X}_k^1, \frac{d}{dx} \widehat{X}_i^1 \right\rangle_x \widehat{S}_{ij}(t) \langle \widehat{V}_\ell^1, \mu \widehat{V}_j^1 \rangle_\mu + \sigma \left(\langle \widehat{X}_k^1, B(t, \cdot) \rangle_x \langle \widehat{V}_\ell^1 \rangle_\mu - \widehat{S}_{k\ell}(t) \right),$$

from which we get the augmented quantity $\widehat{S}_{ij}(t)$. Inserting all augmented low-rank factors into (2.1b) leads to the partial differential equation

$$(3.4) \quad \partial_t B(t, x) = \sigma \left(\sum_{i,j=1}^{2r} \widehat{X}_i^1(x) \widehat{S}_{ij}(t) \langle \widehat{V}_j^1 \rangle_\mu - B(t, x) \right).$$

Before repeating this process and evolving the subequations further in time, we truncate the augmented quantities to a new rank r_1 using a suitable truncation strategy.

4. Angular and spatial discretization. Having derived the K -, L -, and S -steps of the rank-adaptive BUG integrator, we can now proceed with discretizing in angle and space. For the angular discretization, we use the modal representations

$$V_j^0(\mu) \simeq \sum_{n=0}^{N-1} V_{nj}^0 P_n(\mu), \quad \widehat{V}_j^1(\mu) \simeq \sum_{n=0}^{N-1} \widehat{V}_{nj}^1 P_n(\mu), \quad L_i(t, \mu) \simeq \sum_{n=0}^{N-1} L_{ni}(t) P_n(\mu),$$

where P_n are the normalized Legendre polynomials. Note that in the following, we use Einstein's sum convention when not stated otherwise to ensure compactness of notation. Let us define the matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ with entries $A_{mn} := \langle P_m, \mu P_n \rangle_\mu$. Then we can rewrite $\langle V_k^0, \mu V_j^0 \rangle_\mu = V_{km}^0 A_{mn} V_{jn}^0$. The evolution equations with angular discretization then read

(4.1a)

$$\partial_t K_k(t, x) = -\partial_x K_j(t, x) V_{nj}^0 A_{mn} V_{mk}^0 + \sigma \left(B(t, x) V_{0k}^0 - K_k(t, x) \right),$$

(4.1b)

$$\dot{L}_{mk}(t) = - \left\langle X_k^0, \frac{d}{dx} X_i^0 \right\rangle_x L_{ni}(t) A_{mn} + \sigma \left(\langle X_k^0, B(t, \cdot) \rangle_x \delta_{m0} - L_{mk}(t) \right),$$

$$(4.1c) \quad \dot{\widehat{S}}_{k\ell}(t) = - \left\langle \widehat{X}_k^1, \frac{d}{dx} \widehat{X}_i^1 \right\rangle_x S_{ij}(t) \widehat{V}_{nj}^1 A_{mn} \widehat{V}_{m\ell}^1 + \sigma \left(\langle \widehat{X}_k^1, B(t, \cdot) \rangle_x \widehat{V}_{0\ell}^1 - \widehat{S}_{k\ell}(t) \right).$$

For the angular discretization of (3.4) we get

$$(4.1d) \quad \partial_t B(t, x) = \sigma \left(\widehat{X}_i^1(x) \widehat{S}_{ij}(t) \widehat{V}_{0j}^1 - B(t, x) \right).$$

To derive a spatial discretization we choose a spatial grid $x_1 < \dots < x_{n_x}$ with equidistant spacing Δx . The solution in a given cell p is then approximated by

$$\begin{aligned} X_{pk}(t) &\approx \frac{1}{\Delta x} \int_{x_p}^{x_{p+1}} X_k(t, x) dx, & K_{pk}(t) &\approx \frac{1}{\Delta x} \int_{x_p}^{x_{p+1}} K_k(t, x) dx, \\ B_p(t) &\approx \frac{1}{\Delta x} \int_{x_p}^{x_{p+1}} B(t, x) dx. \end{aligned}$$

Spatial derivatives are approximated and stabilized through the tridiagonal stencil matrices $\mathbf{D}^x \approx \partial_x$ and $\mathbf{D}^{xx} \approx \frac{1}{2} \Delta x \partial_{xx}$ with entries

$$D_{p,p\pm 1}^x = \frac{\pm 1}{2\Delta x}, \quad D_{p,p}^{xx} = -\frac{1}{\Delta x}, \quad D_{p,p\pm 1}^{xx} = \frac{1}{2\Delta x}.$$

Applying the matrix $\mathbf{D}^x \in \mathbb{R}^{n_x \times n_x}$ corresponds to a first order and the stabilization matrix $\mathbf{D}^{xx} \in \mathbb{R}^{n_x \times n_x}$ to a second order central differencing scheme. Moreover, from now on we assume periodic boundary conditions. Recall the symmetric matrix \mathbf{A} . It is diagonalizable in the form $\mathbf{A} = \mathbf{Q}\mathbf{M}\mathbf{Q}^\top$ with \mathbf{Q} orthogonal and $\mathbf{M} = \text{diag}(\sigma_1, \dots, \sigma_n)$. We define matrix $|\mathbf{A}|$ as $|\mathbf{A}| = \mathbf{Q}|\mathbf{M}|\mathbf{Q}^\top$. We then obtain the spatially and angular discretized matrix ODEs

$$(4.2a) \quad \begin{aligned} \dot{K}_{pk}(t) &= -D_{qp}^x K_{pj}(t) V_{nj}^0 A_{mn} V_{mk}^0 + D_{qp}^{xx} K_{pj}(t) V_{nj}^0 |A|_{mn} V_{mk}^0 \\ &\quad + \sigma \left(B_p(t) V_{0k}^0 - K_{pk}(t) \right), \end{aligned}$$

$$(4.2b) \quad \begin{aligned} \dot{L}_{mk}(t) &= -A_{mn} L_{ni}(t) X_{pi}^0 D_{qp}^x X_{qk}^0 + |A|_{mn} L_{ni}(t) X_{pi}^0 D_{qp}^{xx} X_{qk}^0 \\ &\quad + \sigma \left(\delta_{m0} B_p(t) X_{pk}^0 - L_{mk}(t) \right), \end{aligned}$$

$$(4.2c) \quad \begin{aligned} \dot{S}_{kl}(t) &= -\widehat{X}_{pk}^1 D_{pq}^x \widehat{X}_{qi}^1 \widehat{S}_{ij}(t) \widehat{V}_{nj}^1 A_{mn} \widehat{V}_{ml}^1 + \widehat{X}_{pk}^1 D_{pq}^{xx} \widehat{X}_{qi}^1 \widehat{S}_{ij}(t) \widehat{V}_{nj}^1 |A|_{mn} \widehat{V}_{ml}^1 \\ &\quad + \sigma \left(\widehat{X}_{pk}^1 B_p(t) \widehat{V}_{0l}^1 - \widehat{S}_{kl}(t) \right). \end{aligned}$$

Lastly, we obtain from (4.1d) for the internal energy B the spatially discretized equation

$$(4.2d) \quad \dot{B}_p(t) = \sigma \left(\widehat{X}_{pi}^1 \widehat{S}_{ij}(t) \widehat{V}_{0j}^1 - B_p(t) \right) = \sigma \left(u_{p0}^1(t) - B_p(t) \right),$$

where we use the notation $\widehat{X}_{pi}^1 \widehat{S}_{ij}(t) \widehat{V}_{mj}^1 =: u_{pm}^1(t)$. We can now show that the semi-discrete time-dependent system (4.2) is energy stable. For this, let us first give a definition of the total energy of the system.

DEFINITION 4.1 (total energy). *Let the matrix $\mathbf{u}^1(t) \in \mathbb{R}^{n_x \times N}$ with low-rank entries $u_{pm}^1(t) = \widehat{X}_{pi}^1 \widehat{S}_{ij}(t) \widehat{V}_{mj}^1$ denote the angularly and spatially discretized approximation of the solution of (2.1a) and $\mathbf{B}(t) \in \mathbb{R}^{n_x}$ be the spatially discretized approximation of the solution of (2.1b). Then we call*

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

with $\|\cdot\|_F$ denoting the Frobenius and $\|\cdot\|_E$ denoting the Euclidean norm, the total energy of the system (4.2).

Further, we note the following properties of the chosen spatial stencil matrices, which we write down denoting all sums explicitly.

LEMMA 4.2 (summation by parts). *Let $y, z \in \mathbb{R}^{n_x}$ with indices $p, q = 1, \dots, n_x$. In addition, we set $y_0 = y_{n_x}$ and $y_{n+1} = y_1$, for z respectively, due to the periodic boundary conditions. Then the stencil matrices fulfill the following properties:*

$$\sum_{p,q=1}^{n_x} y_p D_{pq}^x z_q = - \sum_{p,q=1}^{n_x} z_p D_{pq}^x y_q, \quad \sum_{p,q=1}^{n_x} z_p D_{pq}^x z_q = 0, \quad \sum_{p,q=1}^{n_x} y_p D_{pq}^{xx} z_q = \sum_{p,q=1}^{n_x} z_p D_{pq}^{xx} y_q.$$

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

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

Then, $\sum_{p,q=1}^{n_x} z_p D_{pq}^{xx} z_q = - \sum_{p=1}^{n_x} (\sum_{q=1}^{n_x} D_{pq}^+ z_q)^2$.

Proof. The assertions follow directly by plugging in the definitions of the stencil matrices and rearranging the sums of the products in an adequate way:

$$\begin{aligned} \sum_{p,q=1}^{n_x} y_p D_{pq}^x z_q &= \frac{1}{2\Delta x} \sum_{p=1}^{n_x} y_p (z_{p+1} - z_{p-1}) = -\frac{1}{2\Delta x} \sum_{p=1}^{n_x} z_p (y_{p+1} - y_{p-1}) \\ &= - \sum_{p,q=1}^{n_x} z_p D_{pq}^x y_q, \\ \sum_{p,q=1}^{n_x} z_p D_{pq}^x z_q &= - \sum_{p,q=1}^{n_x} z_p D_{pq}^x z_q = 0, \\ \sum_{p,q=1}^{n_x} y_p D_{pq}^{xx} z_q &= -\frac{1}{\Delta x} \sum_{p=1}^{n_x} y_p z_p + \frac{1}{2\Delta x} \sum_{p=1}^{n_x} y_p (z_{p+1} + z_{p-1}) \\ &= -\frac{1}{\Delta x} \sum_{p=1}^{n_x} z_p y_p + \frac{1}{2\Delta x} \sum_{p=1}^{n_x} z_p (y_{p+1} + y_{p-1}) = \sum_{p,q=1}^{n_x} z_p D_{pq}^{xx} y_q, \\ \sum_{p,q=1}^{n_x} z_p D_{pq}^{xx} z_q &= -\frac{1}{\Delta x} \sum_{p=1}^{n_x} z_p^2 + \frac{1}{2\Delta x} \sum_{p=1}^{n_x} z_p (z_{p+1} + z_{p-1}) \\ &= -\frac{1}{2\Delta x} \sum_{p=1}^{n_x} (z_p^2 - 2z_p z_{p+1} + z_{p+1}^2) = -\frac{1}{2\Delta x} \sum_{p=1}^{n_x} (z_p - z_{p+1})^2 \\ &= - \sum_{p=1}^{n_x} \left(\sum_{q=1}^{n_x} D_{pq}^+ z_q \right)^2. \quad \square \end{aligned}$$

With these properties at hand, we can now show dissipation of the total energy.

THEOREM 4.3. *The semidiscrete time-continuous system consisting of (4.2) is energy stable; that is, $\dot{E}(t) \leq 0$.*

Proof. Let us start from the S -step in (4.2c):

$$\begin{aligned} \hat{S}_{k\ell}(t) &= -\hat{X}_{pk}^1 D_{pq}^x \hat{X}_{qi}^1 \hat{S}_{ij}(t) \hat{V}_{nj}^1 A_{mn} \hat{V}_{m\ell}^1 + \hat{X}_{pk}^1 D_{pq}^{xx} \hat{X}_{qi}^1 \hat{S}_{ij}(t) \hat{V}_{nj}^1 |A|_{mn} \hat{V}_{m\ell}^1 \\ &\quad + \sigma \left(\hat{X}_{pk}^1(x) B_p(t) \hat{V}_{0\ell}^1 - \hat{S}_{k\ell}(t) \right). \end{aligned}$$

We multiply with $\hat{X}_{\alpha k}^1 \hat{V}_{\beta \ell}^1$, where $\alpha = 1, \dots, n_x$ and $\beta = 0, \dots, N-1$, sum over k and ℓ , and introduce the projections $P_{\alpha p}^{X,1} = \hat{X}_{\alpha k}^1 \hat{X}_{pk}^1$ and $P_{m\beta}^{V,1} = \hat{V}_{m\ell}^1 \hat{V}_{\beta \ell}^1$. With the notation $\hat{X}_{qi}^1 \hat{S}_{ij}(t) \hat{V}_{nj}^1 = u_{qn}^1(t)$ we get

$$\begin{aligned} \dot{u}_{\alpha\beta}^1(t) = & -P_{\alpha p}^{X,1} D_{pq}^x u_{qn}^1(t) A_{mn} P_{m\beta}^{V,1} + P_{\alpha p}^{X,1} D_{pq}^{xx} u_{qn}^1(t) |A|_{mn} P_{m\beta}^{V,1} \\ & + \sigma \left(P_{\alpha p}^{X,1} B_p(t) \delta_{0m} P_{m\beta}^{V,1} - u_{\alpha\beta}^1(t) \right). \end{aligned}$$

Next, we multiply with $u_{\alpha\beta}^1(t)$ and sum over α and β . Note that

$$P_{\alpha p}^{X,1} u_{\alpha\beta}^1(t) = u_{p\beta}^1(t) \quad \text{and} \quad P_{m\beta}^{V,1} u_{p\beta}^1(t) = u_{pm}^1(t).$$

This leads to

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \|\mathbf{u}^1(t)\|_F^2 = & -u_{pm}^1(t) D_{pq}^x u_{qn}^1(t) A_{mn} + u_{pm}^1(t) D_{pq}^{xx} u_{qn}^1(t) |A|_{mn} \\ & + \sigma \left(u_{pm}^1(t) B_p(t) \delta_{0m} - \|\mathbf{u}^1(t)\|^2 \right). \end{aligned}$$

Recall that we can write $\mathbf{A} = \mathbf{Q}\mathbf{M}\mathbf{Q}^\top$ with $\mathbf{M} = \text{diag}(\sigma_1, \dots, \sigma_N)$. Inserting this representation gives

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \|\mathbf{u}^1(t)\|_F^2 = & -u_{pm}^1(t) D_{pq}^x u_{qn}^1(t) Q_{nk} \sigma_k Q_{mk} + u_{pm}^1(t) D_{pq}^{xx} u_{qn}^1(t) Q_{nk} |\sigma_k| Q_{mk} \\ & + \left(u_{pm}^1(t) B_p(t) \delta_{0m} - \|\mathbf{u}^1(t)\|^2 \right) \\ = & -\sigma_k \tilde{u}_{pk}^1(t) D_{pq}^x \tilde{u}_{qk}^1(t) + |\sigma_k| \tilde{u}_{pk}^1(t) D_{pq}^{xx} \tilde{u}_{qk}^1(t) \\ & + \left(u_{pm}^1(t) B_p(t) \delta_{0m} - \|\mathbf{u}^1(t)\|^2 \right), \end{aligned}$$

where $\tilde{u}_{pk}^1(t) = u_{pm}^1(t) Q_{mk}$. With the properties of the stencil matrices we get

$$(4.3) \quad \frac{1}{2} \frac{d}{dt} \|\mathbf{u}^1(t)\|_F^2 = - \left(D_{pq}^+ \tilde{u}_{qm}^1(t) |A|_{mn}^{1/2} \right)^2 + \sigma \left(u_{p0}^1(t) B_p(t) - \|\mathbf{u}^1(t)\|_F^2 \right).$$

Next we consider (4.2d). Multiplication with $B_p(t)$ and summation over p gives

$$(4.4) \quad \frac{1}{2} \frac{d}{dt} \|\mathbf{B}(t)\|_E^2 = \sigma \left(u_{p0}^1(t) B_p(t) - \|\mathbf{B}(t)\|_E^2 \right).$$

For the total energy of the system it holds that $E(t) = \frac{1}{2} \|\mathbf{u}^1(t)\|_F^2 + \frac{1}{2} \|\mathbf{B}(t)\|_E^2$. Adding the evolution equations (4.3) and (4.4), we get

$$\begin{aligned} \frac{d}{dt} E(t) = & - \left(D_{pq}^+ \tilde{u}_{qm}^1(t) |A|_{mn}^{1/2} \right)^2 + \sigma \left(u_{p0}^1(t) B_p(t) - \|\mathbf{u}^1(t)\|_F^2 \right) \\ & + \sigma \left(u_{p0}^1(t) B_p(t) - \|\mathbf{B}(t)\|_E^2 \right) \\ = & - \left(D_{pq}^+ \tilde{u}_{qm}^1(t) |A|_{mn}^{1/2} \right)^2 - \sigma \left((u_{p0}^1(t) - B_p(t))^2 + (u_{pm}^1(t))^2 (1 - \delta_{m0}) \right), \end{aligned}$$

where we rewrote $\|\mathbf{B}(t)\|_E^2 = B_p(t)^2$ and $\|\mathbf{u}^1(t)\|_F^2 = (u_{pm}^1(t))^2$. This expression is strictly negative, which means that E is dissipated in time. Hence, the system is energy stable. \square

5. Time discretization. Our goal is to construct a conservative DLRA scheme which is energy stable under a sharp time step restriction. Constructing time discretization schemes which preserve the energy dissipation shown in Theorem 4.3 while not suffering from the potentially stiff opacity term is not trivial. In fact a naive IMEX time discretization potentially will increase the total energy, which we demonstrate in the following.

5.1. Naive time discretization. We start from system (4.2), which still depends continuously on the time t . For the time discretization we choose a naive IMEX Euler scheme where we perform a splitting of the internal energy and radiation transport equation. That is, we use an explicit Euler step for the transport part of the evolution equations, treat the internal energy B explicitly, and use an implicit Euler step for the radiation absorption term. Note that the scheme describes the evolution from time t_0 to time $t_1 = t_0 + \Delta t$ but holds for all further time steps equivalently. This yields the fully discrete scheme

$$(5.1a) \quad K_{pk}^1 = K_{pk}^0 - \Delta t D_{qp}^x K_{pj}^0 V_{nj}^0 A_{mn} V_{mk}^0 + \Delta t D_{qp}^{xx} K_{pj}^0 V_{nj}^0 |A|_{mn} V_{mk}^0 \\ + \sigma \left(\Delta t B_p^0 V_{0k}^0 - \Delta t K_{pk}^1 \right),$$

$$(5.1b) \quad L_{mk}^1 = L_{mk}^0 - \Delta t X_{qk}^0 D_{qp}^x X_{pi}^0 L_{ni}^0 A_{mn} + \Delta t X_{qk}^0 D_{qp}^{xx} X_{pi}^0 L_{ni}^0 |A|_{mn} \\ + \sigma \left(\Delta t X_{pk}^0 B_p^0 \delta_{m0} - \Delta t L_{mk}^1 \right).$$

We perform a QR-decomposition of the quantities $[K_{pk}^1, X_{pk}^0]$ and $[L_{pk}^1, V_{pk}^0]$ to obtain the augmented and time updated bases \widehat{X}_{pk}^1 and \widehat{V}_{pk}^1 according to the rank-adaptive BUG integrator [6]. Lastly, we perform a Galerkin step for the augmented bases according to

$$(5.1c) \quad \widetilde{S}_{kl}^1 = \widetilde{S}_{kl}^0 - \Delta t \widehat{X}_{pk}^1 D_{pq}^x \widehat{X}_{qi}^1 \widetilde{S}_{ij}^0 \widehat{V}_{nj}^1 A_{mn} \widehat{V}_{m\ell}^1 + \Delta t \widehat{X}_{pk}^1 D_{pq}^{xx} \widehat{X}_{qi}^1 \widetilde{S}_{ij}^0 \widehat{V}_{nj}^1 |A|_{mn} \widehat{V}_{m\ell}^1 \\ + \sigma \left(\Delta t \widehat{X}_{pk}^1 B_p^0 \widehat{V}_{0\ell}^1 - \Delta t \widetilde{S}_{kl}^1 \right),$$

where $\widetilde{S}_{kl}^0 := \widehat{X}_{pk}^1 X_{pi}^0 S_{ij}^0 V_{nj}^0 \widehat{V}_{n\ell}^1$. The internal energy is then updated via

$$(5.1d) \quad B_p^1 = B_p^0 + \sigma \Delta t \left(\widehat{X}_{pi}^1 \widehat{S}_{ij}^1 \widehat{V}_{0j}^1 - B_p^1 \right).$$

However, this numerical method has the undesirable property that it can increase the total energy during a time step. In Theorem 5.1 we show this analytically. This behavior is, obviously, completely unphysical.

THEOREM 5.1. *Let $\mathbf{u}^0 \in \mathbb{R}^{n_x \times N}$ with entries $u_{pm}^0 = X_{pk}^0 S_{kl}^0 V_{m\ell}^0$ denote the angularly and spatially discretized low-rank approximation of the function f at time $t = t_0$, and let $\mathbf{u}^1 \in \mathbb{R}^{n_x \times N}$ with entries $u_{\alpha\beta}^1 = \widehat{X}_{\alpha k}^1 \widehat{S}_{kl}^1 \widehat{V}_{\beta\ell}^1$ denote the basis augmented angularly and spatially discretized low-rank approximation at time $t = t_1$ using the rank-adaptive BUG integrator. Further, $\mathbf{B}^0 \in \mathbb{R}^{n_x}$ shall denote the spatially discretized low-rank approximation of B at time $t = t_0$, and $\mathbf{B}^1 \in \mathbb{R}^{n_x}$ at time $t = t_1$, respectively. The total energy at time $t = t_0$ is denoted by E^0 and E^1 at time $t = t_1$, respectively. Then, there exist initial value pairs $(\mathbf{u}^0, \mathbf{B}^0)$ and time step sizes Δt such that the naive scheme (5.1) results in $(\mathbf{u}^1, \mathbf{B}^1)$, for which the total energy increases, i.e., for which $E^1 > E^0$.*

Proof. Let us multiply the S -step (5.1c) with $\widehat{X}_{\alpha k}^1 \widehat{V}_{\beta\ell}^1$ and sum over k and ℓ . Again we make use of the projections $P_{\alpha p}^{X,1} = \widehat{X}_{\alpha k}^1 \widehat{X}_{pk}^1$ and $P_{m\beta}^{V,1} = \widehat{V}_{m\ell}^1 \widehat{V}_{\beta\ell}^1$. With the definition of $\widetilde{S}_{k\ell}^0$ we obtain

$$(5.2) \quad u_{\alpha\beta}^1 = u_{pm}^0 - P_{\alpha p}^{X,1} \Delta t D_{pq}^x u_{qn}^0 A_{mn} P_{m\beta}^{V,1} + P_{\alpha p}^{X,1} \Delta t D_{pq}^{xx} u_{qn}^0 |A|_{mn} P_{m\beta}^{V,1} \\ + \sigma \left(\Delta t P_{\alpha p}^{X,1} B_p^0 \delta_{m0} P_{m\beta}^{V,1} - \Delta t u_{\alpha\beta}^1 \right).$$

Let us choose a constant solution in space, i.e., $B_p^1 = B^1$ and $u_{\alpha\beta}^1 = u^1 \delta_{\beta 0}$ for all spatial indices $p, \alpha = 1, \dots, n_x$. The scalar values B^1 and u^1 are chosen such that $B^1 = u^1 + \alpha$, where

$$0 < \alpha < \frac{\sigma \Delta t}{1 + \sigma \Delta t + \sigma^2 \Delta t^2 + \frac{1}{2} \sigma^3 \Delta t^3} u^1.$$

We can now verify that we obtain our chosen values for B_p^1 and $u_{\alpha\beta}^1$ after a single step of (5.2) when using the initial condition

$$(5.3a) \quad B_p^0 = B^1 + \sigma \Delta t \alpha = u^1 + \alpha(1 + \sigma \Delta t),$$

$$(5.3b) \quad u_{pm}^0 = (u^1 + \sigma \Delta t(u^1 - B_p^0)) \delta_{m0} = (u^1 - \sigma \Delta t \alpha(1 + \sigma \Delta t)) \delta_{m0}.$$

To show this, note that since the solution is constant in space, all terms containing the stencil matrices \mathbf{D}^x and \mathbf{D}^{xx} drop out, and we are left with

$$(5.4) \quad u_{\alpha\beta}^1 = u_{pm}^0 + \sigma \left(\Delta t P_{\alpha p}^{X,1} B_p^0 \delta_{m0} P_{m\beta}^{V,1} - \Delta t u_{\alpha\beta}^1 \right).$$

Since B_p^0 is constant in space and δ_{m0} lies in the span of our basis, we know that all projections in the above equation are exact. Plugging the initial values (5.3) into (5.4), we then directly obtain $u_{\alpha\beta}^1 = u^1 \delta_{\beta 0}$. Similarly, by plugging (5.3) into (5.1d), we obtain $B_p^1 = B^1$.

Then, we square both of the initial terms (5.3) to get

$$(B_p^0)^2 = (B^1)^2 + 2\sigma \Delta t \alpha B^1 + \sigma^2 \Delta t^2 \alpha^2 = (B^1)^2 + 2\sigma \Delta t \alpha (u^1 + \alpha) + \sigma^2 \Delta t^2 \alpha^2,$$

$$(u_{pm}^0)^2 = ((u^1)^2 - 2\sigma \Delta t \alpha u^1 (1 + \sigma \Delta t) + \sigma^2 \Delta t^2 \alpha^2 (1 + \sigma \Delta t)^2) \delta_{m0}.$$

Summing over p and m , adding these two terms, and multiplying with $\frac{1}{2}$ yields

$$E^1 = E^0 + \sigma^2 \Delta t^2 \alpha u^1 - \sigma \Delta t \alpha^2 - \frac{1}{2} \sigma^2 \Delta t^2 \alpha^2 - \frac{1}{2} \sigma^2 \Delta t^2 \alpha^2 (1 + \sigma \Delta t)^2.$$

Note that $E^1 > E^0$ if

$$\sigma \Delta t u^1 - \alpha - \frac{1}{2} \sigma \Delta t \alpha - \frac{1}{2} \sigma \Delta t \alpha (1 + \sigma \Delta t)^2 > 0.$$

Rearranging gives

$$\alpha < \frac{\sigma \Delta t}{1 + \sigma \Delta t + \sigma^2 \Delta t^2 + \frac{1}{2} \sigma^3 \Delta t^3} u^1.$$

This is exactly the domain α is chosen from. Hence, we have $E^1 > E^0$, which is the desired result. \square

5.2. Energy stable space-time discretization. We have seen that the naive scheme presented in (5.1) can increase the total energy in one time step. The main goal of this section is to construct a novel energy stable time integration scheme for which the corresponding analysis leads to a classic hyperbolic CFL condition that enables us to operate up to a time step size of $\Delta t = \text{CFL} \cdot \Delta x$. For constructing this energy stable scheme, we write the original equations in two parts followed by a basis augmentation and correction step.

In detail, we first solve

$$(5.5a) \quad K_{pk}^* = K_{pk}^0 - \Delta t D_{qp}^x K_{pj}^0 V_{nj}^0 A_{mn} V_{mk}^0 + \Delta t D_{qp}^{xx} K_{pj}^0 V_{nj}^0 |A|_{mn} V_{mk}^0,$$

$$(5.5b) \quad L_{mk}^* = L_{mk}^0 - \Delta t X_{qk}^0 D_{qp}^x X_{pi}^0 L_{ni}^0 A_{mn} + \Delta t X_{qk}^0 D_{qp}^{xx} X_{pi}^0 L_{ni}^0 |A|_{mn}.$$

We perform a QR-decomposition of the augmented quantities $\mathbf{X}^*\mathbf{R} = [\mathbf{K}^*, \mathbf{X}^0]$ and $\mathbf{V}^*\mathbf{R} = [\mathbf{L}^*, \mathbf{V}^0]$ to obtain the augmented and time updated bases \mathbf{X}^* and \mathbf{V}^* . Note that \mathbf{R} and $\tilde{\mathbf{R}}$ are discarded. With $\tilde{S}_{\alpha\beta}^0 = X_{j\alpha}^* X_{j\ell}^0 S_{\ell m}^0 V_{km}^0 V_{k\beta}^*$ we then solve the S -step equation

$$(5.5c) \quad S_{\alpha\beta}^* = \tilde{S}_{\alpha\beta}^0 - \Delta t X_{p\alpha}^* D_{pq}^x X_{qi}^* \tilde{S}_{ij}^0 V_{nj}^* A_{mn} V_{m\beta}^* + \Delta t X_{p\alpha}^* D_{pq}^{xx} X_{qi}^* \tilde{S}_{ij}^0 V_{nj}^* |A|_{mn} V_{m\beta}^*.$$

Second, we solve the coupled equations for the internal energy $\mathbf{B} \in \mathbb{R}^{n_x}$ and the quantity $\hat{\mathbf{u}}_0^1 = (\hat{u}_{j0}^1)_j \in \mathbb{R}^{n_x}$ to which we refer as the zeroth order moment according to

$$(5.5d) \quad \hat{u}_{j0}^1 = X_{j\ell}^0 S_{\ell m}^0 V_{0m}^0 - \Delta t D_{ji}^x X_{in}^* \tilde{S}_{nm}^0 V_{\ell m}^* A_{0\ell} + \Delta t D_{ji}^{xx} X_{in}^* \tilde{S}_{nm}^0 V_{\ell m}^* |A|_{0\ell} + \sigma \Delta t (B_j^1 - \hat{u}_{j0}^1),$$

$$(5.5e) \quad B_j^1 = B_j^0 + \sigma \Delta t (\hat{u}_{j0}^1 - B_j^1).$$

Following [21, section 6] we perform the opacity update only on $\mathbf{L} = \mathbf{V}^*\mathbf{S}^*$ according to

$$(5.5f) \quad L_{mk}^{*,\text{scat}} = \frac{1}{1 + \Delta t \sigma} L_{mk} \quad \text{for } k \neq 0$$

and perform a QR-decomposition $\mathbf{V}^{*,\text{scat}} \mathbf{S}^{*,\text{scat},\top} = \mathbf{L}^{*,\text{scat}}$ to retrieve the factorized basis $\mathbf{V}^{*,\text{scat}}$ and the coefficients from the matrix $\mathbf{S}^{*,\text{scat}}$. We then augment the basis matrices according to

$$(5.5g) \quad \tilde{\mathbf{X}}^1 = \text{qr}([\hat{\mathbf{u}}_0^1, \mathbf{X}^*]), \quad \tilde{\mathbf{V}}^1 = \text{qr}([\mathbf{e}_1, \mathbf{V}^{*,\text{scat}}]).$$

Third, the coefficient matrix is updated via

$$(5.5h) \quad \tilde{\mathbf{S}}^1 = \tilde{\mathbf{X}}^{1,\top} \mathbf{X}^* \mathbf{S}^{*,\text{scat}} \mathbf{V}^{*,\text{scat},\top} (\mathbf{I} - \mathbf{e}_1 \mathbf{e}_1^\top) \tilde{\mathbf{V}}^1 + \tilde{\mathbf{X}}^{1,\top} \hat{\mathbf{u}}_0^1 \mathbf{e}_{1,\top} \tilde{\mathbf{V}}^1 \in \mathbb{R}^{(2r+1) \times (2r+1)}.$$

Then, we obtain the updated solution $\tilde{\mathbf{X}}^1 \tilde{\mathbf{S}}^1 \tilde{\mathbf{V}}^{1,\top} \in \mathbb{R}^{n_x \times N}$. Lastly, we truncate this rank $2r+1$ solution to a new rank r_1 using a suited truncation strategy such as that proposed in [6] or the conservative truncation strategy of [14]. This finally gives the low-rank factors $\mathbf{X}^1, \mathbf{S}^1$, and \mathbf{V}^1 . We show that the given scheme is energy stable and start with the following lemma.

LEMMA 5.2. *Let us denote $u_{jk}^1 := \tilde{X}_{j\alpha}^1 \tilde{S}_{\alpha\beta}^1 \tilde{V}_{k\beta}^1$. Under the time step restriction $\Delta t \leq \Delta x$ it holds that*

$$(5.6) \quad \frac{\Delta t}{2} (D_{ji}^x u_{jk}^1 A_{k\ell} - D_{ji}^{xx} u_{jk}^1 |A|_{k\ell})^2 - \left(D_{ji}^+ u_{ik}^1 |A|_{k\ell}^{1/2} \right)^2 \leq 0.$$

Proof. Following [21], we employ a Fourier analysis, which allows us to write the stencil matrices $\mathbf{D}^{x,xx,+}$ in diagonal form. Let us define $\mathbf{E} \in \mathbb{C}^{n_x \times n_x}$ with entries

$$E_{k\alpha} = \sqrt{\Delta x} \exp(i\alpha\pi x_k), \quad k, \alpha = 1, \dots, n_x,$$

with $i \in \mathbb{C}$ being the imaginary unit. Then, the matrix \mathbf{E} is orthonormal, i.e., $\mathbf{E}\mathbf{E}^H = \mathbf{E}^H\mathbf{E} = \mathbf{I}$ (the uppercase H denotes the complex transpose), and it diagonalizes the stencil matrices:

$$(5.7) \quad \mathbf{D}^{x,xx,+} \mathbf{E} = \mathbf{E} \mathbf{\Lambda}^{x,xx,+}.$$

The matrices $\mathbf{A}^{x,xx,+}$ are diagonal with entries

$$\begin{aligned}\lambda_{\alpha,\alpha}^x &= \frac{1}{2\Delta x} (e^{i\alpha\pi\Delta x} - e^{-i\alpha\pi\Delta x}) = \frac{i}{\Delta x} \sin(\omega_\alpha), \\ \lambda_{\alpha,\alpha}^{xx} &= \frac{1}{2\Delta x} (e^{i\alpha\pi\Delta x} - 2 + e^{-i\alpha\pi\Delta x}) = \frac{1}{\Delta x} (\cos(\omega_\alpha) - 1), \\ \lambda_{\alpha,\alpha}^+ &= \frac{1}{\sqrt{2\Delta x}} (e^{i\alpha\pi\Delta x} - 1) = \frac{1}{\sqrt{2\Delta x}} (\cos(\omega_\alpha) + i \sin(\omega_\alpha) - 1),\end{aligned}$$

where we use $\omega_\alpha := \alpha\pi\Delta x$. Moreover, recall that we can write $\mathbf{A} = \mathbf{Q}\mathbf{M}\mathbf{Q}^\top$, where $\mathbf{M} = \text{diag}(\sigma_1, \dots, \sigma_N)$. We then have with $\hat{u}_{jk} = E_{j\ell} u_{\ell m} Q_{mk}$

$$\begin{aligned}\frac{\Delta t}{2} (D_{ji}^x u_{jk}^1 A_{k\ell} - D_{ji}^{xx} u_{jk}^1 |A|_{k\ell})^2 - \left(D_{ji}^+ u_{jk}^1 |A|_{k\ell}^{1/2} \right)^2 \\ = \frac{\Delta t}{2} \left| \lambda_{jj}^x \hat{u}_{jk}^1 \sigma_k - \lambda_{jj}^{xx} \hat{u}_{jk}^1 |\sigma_k| \right|^2 - \left| \lambda_{jj}^+ \hat{u}_{jk}^1 |\sigma_k|^{1/2} \right|^2 \\ \leq \left[\Delta t \left(\frac{|\sigma_k|^2}{\Delta x^2} \cdot |1 - \cos(\omega_j)| \right) - \frac{|\sigma_k|}{\Delta x} \cdot |1 - \cos(\omega_j)| \right] (\hat{u}_{jk}^1)^2.\end{aligned}$$

To ensure negativity, we must have

$$\Delta t \left(\frac{|\sigma_k|^2}{\Delta x^2} \cdot |1 - \cos(\omega_j)| \right) \leq \frac{|\sigma_k|}{\Delta x} \cdot |1 - \cos(\omega_j)|.$$

Hence, for $\Delta t \leq \frac{\Delta x}{|\sigma_k|}$, (5.6) holds. Since $|\sigma_k| \leq 1$, we have proven the lemma. \square

We can now show energy stability of the proposed scheme.

THEOREM 5.3. *Under the time step restriction $\Delta t \leq \Delta x$, the scheme (5.5) is energy stable; i.e.,*

$$(5.8) \quad \frac{1}{2} \|\mathbf{B}^1\|_E^2 + \frac{1}{2} \|\mathbf{X}^1 \mathbf{S}^1 \mathbf{V}^{1,\top}\|_F^2 \leq \frac{1}{2} \|\mathbf{B}^0\|_E^2 + \frac{1}{2} \|\mathbf{X}^0 \mathbf{S}^0 \mathbf{V}^{0,\top}\|_F^2.$$

Proof. First, we multiply (5.5e) with B_j^1 and sum over j . Then,

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

Let us note that

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

Hence,

$$(5.9) \quad \frac{1}{2} (B_j^1)^2 = \frac{1}{2} (B_j^0)^2 - \frac{1}{2} (B_j^1 - B_j^0)^2 + \sigma \Delta t \left(u_{j0}^1 B_j^1 - (B_j^1)^2 \right).$$

To obtain a similar expression for $(u_{jk}^1)^2$, we multiply (5.5c) with $X_{j\alpha}^* V_{k\beta}^*$ and sum over α and β . For simplicity of notation, let us define $u_{jk}^* := X_{j\alpha}^* S_{\alpha\beta}^* V_{k\beta}^*$ and $u_{jk}^0 := X_{j\alpha}^* \tilde{S}_{\alpha\beta}^0 V_{k\beta}^*$ as well as the projections $P_{jp}^X := X_{j\alpha}^* X_{p\alpha}^*$ and $P_{km}^V := V_{k\beta}^* V_{m\beta}^*$. Then, we obtain the system

$$(5.10) \quad u_{jk}^* = u_{jk}^0 - \Delta t P_{jp}^X D_{pq}^x u_{qn}^0 A_{mn} P_{km}^V + \Delta t P_{jp}^X D_{pq}^{xx} u_{qn}^0 |A|_{mn} P_{km}^V.$$

Next, we define $u_{jk}^1 := \tilde{X}_{j\alpha}^1 \tilde{S}_{\alpha\beta}^1 \tilde{V}_{k\beta}^1$ and note that by construction we have that

$$u_{jk}^1 = \frac{u_{jk}^* (1 - \delta_{k0})}{1 + \sigma \Delta t} + \hat{u}_{j0}^1 \delta_{k0}.$$

Hence, plugging in the schemes for u_{jk}^* and \widehat{u}_{j0}^1 , that is, (5.10) and (5.5d), we get

$$(1 + \sigma \Delta t) u_{jk}^1 = (u_{jk}^0 - \Delta t P_{jp}^X D_{pq}^x u_{qn}^0 A_{mn} P_{km}^V + \Delta t P_{jp}^X D_{pq}^{xx} u_{qn}^0 |A|_{mn} P_{km}^V) (1 - \delta_{k0}) \\ + \left(X_{j\ell}^0 S_{\ell m}^0 V_{0m}^0 - \Delta t D_{ji}^x X_{in}^* \widetilde{S}_{nm}^0 V_{\ell m}^* A_{0\ell} + \Delta t D_{ji}^{xx} X_{in}^* \widetilde{S}_{nm}^0 V_{\ell m}^* |A|_{0\ell} \right. \\ \left. + \sigma \Delta t B_j^1 \right) \delta_{k0}.$$

Let us note that $P_{km}^V P_{jp}^X u_{jk}^1 = u_{jk}^1$ for $k \neq 0$. Hence, multiplying the above equation with u_{jk}^1 and summing over j and k gives

$$\frac{1}{2} (u_{jk}^1)^2 = \frac{1}{2} (u_{jk}^0)^2 - \frac{1}{2} (u_{jk}^1 - u_{jk}^0)^2 - \Delta t u_{jk}^1 D_{ji}^x u_{i\ell}^0 A_{k\ell} + \Delta t u_{jk}^1 D_{ji}^{xx} u_{i\ell}^0 |A|_{k\ell} \\ + \sigma \Delta t u_{jk}^1 (B_j^1 \delta_{k0} - u_{jk}^1).$$

Let us now add the zero term $\Delta t u_{jk}^1 D_{ji}^x u_{i\ell}^1 A_{k\ell}$ and add and subtract the term $\Delta t u_{jk}^1 D_{ji}^{xx} u_{i\ell}^1 |A|_{k\ell}$. Then,

$$\frac{1}{2} (u_{jk}^1)^2 = \frac{1}{2} (u_{jk}^0)^2 - \frac{1}{2} (u_{jk}^1 - u_{jk}^0)^2 - \Delta t u_{jk}^1 D_{ji}^x (u_{i\ell}^0 - u_{i\ell}^1) A_{k\ell} \\ + \Delta t u_{jk}^1 D_{ji}^{xx} (u_{i\ell}^0 - u_{i\ell}^1) |A|_{k\ell} + \Delta t u_{jk}^1 D_{ji}^{xx} u_{i\ell}^1 |A|_{k\ell} \\ + \sigma \Delta t u_{jk}^1 (B_j^1 \delta_{k0} - u_{jk}^1).$$

In the following, we use Young's inequality, which states that for $a, b \in \mathbb{R}$ we have $a \cdot b \leq \frac{a^2}{2} + \frac{b^2}{2}$. We now apply this to the term

$$- \Delta t u_{jk}^1 D_{ji}^x (u_{i\ell}^0 - u_{i\ell}^1) A_{k\ell} + \Delta t u_{jk}^1 D_{ji}^{xx} (u_{i\ell}^0 - u_{i\ell}^1) |A|_{k\ell} \\ \leq \frac{1}{2} (u_{i\ell}^0 - u_{i\ell}^1)^2 + \frac{\Delta t^2}{2} (D_{ji}^x u_{jk}^1 A_{k\ell} - D_{ji}^{xx} u_{jk}^1 |A|_{k\ell})^2.$$

Hence, using $u_{jk}^1 D_{ji}^{xx} u_{i\ell}^1 |A|_{k\ell} = -(D_{ji}^+ u_{ik}^1 |A|_{k\ell}^{1/2})^2$, we get

$$\frac{1}{2} (u_{jk}^1)^2 \leq \frac{1}{2} (u_{jk}^0)^2 + \frac{\Delta t^2}{2} (D_{ji}^x u_{jk}^1 A_{k\ell} - D_{ji}^{xx} u_{jk}^1 |A|_{k\ell})^2 - \Delta t (D_{ji}^+ u_{ik}^1 |A|_{k\ell}^{1/2})^2 \\ (5.11) \quad + \sigma \Delta t u_{jk}^1 (B_j^1 \delta_{k0} - u_{jk}^1).$$

As for the continuous case, we add (5.11) and (5.9) to obtain a time update equation for $E^0 := \frac{1}{2} (u_{jk}^0)^2 + \frac{1}{2} (B_j^0)^2$:

$$E^1 \leq E^0 + \frac{\Delta t^2}{2} (D_{ji}^x u_{jk}^1 A_{k\ell} - D_{ji}^{xx} u_{jk}^1 |A|_{k\ell})^2 - \Delta t (D_{ji}^+ u_{ik}^1 |A|_{k\ell}^{1/2})^2 \\ + \sigma \Delta t (u_{j0}^1 B_j^1 - (u_{jk}^1)^2) - \frac{1}{2} (B_j^1 - B_j^0)^2 + \sigma \Delta t (u_{j0}^1 B_j^1 - (B_j^1)^2) \\ \leq E^0 + \frac{\Delta t^2}{2} (D_{ji}^x u_{jk}^1 A_{k\ell} - D_{ji}^{xx} u_{jk}^1 |A|_{k\ell})^2 - \Delta t (D_{ji}^+ u_{ik}^1 |A|_{k\ell}^{1/2})^2 \\ (5.12) \quad - \sigma \Delta t (B_j^1 - u_{jk}^1)^2 - \frac{1}{2} (B_j^1 - B_j^0)^2.$$

With Lemma 5.2 we have that

$$\frac{\Delta t}{2} (D_{ji}^x u_{jk}^1 A_{k\ell} - D_{ji}^{xx} u_{jk}^1 |A|_{k\ell})^2 - (D_{ji}^+ u_{ik}^1 |A|_{k\ell}^{1/2})^2 \leq 0$$

for $\Delta t \leq \Delta x$. Since the truncation step is designed to not alter the zero order moments, we conclude that $E^1 \leq E^0$ and the full scheme is energy stable under the time step restriction $\Delta t \leq \Delta x$. \square

6. Mass conservation. A drawback of dynamical low-rank approximation using the classical integrators introduced in section 1 is that the method does not preserve physical invariants. It has been shown in [12] that this problem can be overcome when using a modified L -step equation. On this basis, [14, 17] have presented conservative DLRA algorithms where they additionally introduced a conservative truncation step. In contrast to [14, 17], we do not need to consider a modified L -step equation due to the applied basis augmentation strategy from [6] but use the conservative truncation step. Then we can show that besides being energy stable, our scheme ensures local conservation of mass. The conservative truncation strategy works as follows:

1. Compute $\tilde{\mathbf{K}} = \tilde{\mathbf{X}}^1 \tilde{\mathbf{S}}^1$ and split it into two parts $\tilde{\mathbf{K}} = [\tilde{\mathbf{K}}^{\text{cons}}, \tilde{\mathbf{K}}^{\text{rem}}]$, where $\tilde{\mathbf{K}}^{\text{cons}}$ corresponds to the first and $\tilde{\mathbf{K}}^{\text{rem}}$ consists of the remaining columns of $\tilde{\mathbf{K}}$. Analogously, distribute $\tilde{\mathbf{V}}^1 = [\tilde{\mathbf{V}}^{\text{cons}}, \tilde{\mathbf{V}}^{\text{rem}}]$, where $\tilde{\mathbf{V}}^{\text{cons}}$ corresponds to the first and $\tilde{\mathbf{V}}^{\text{rem}}$ consists of the remaining columns of $\tilde{\mathbf{V}}$.
2. Derive $\mathbf{X}^{\text{cons}} = \tilde{\mathbf{K}}^{\text{cons}} / \|\tilde{\mathbf{K}}^{\text{cons}}\|$ and $\mathbf{S}^{\text{cons}} = \|\tilde{\mathbf{K}}^{\text{cons}}\|$.
3. Perform a QR-decomposition of $\tilde{\mathbf{K}}^{\text{rem}}$ to obtain $\tilde{\mathbf{K}}^{\text{rem}} = \tilde{\mathbf{X}}^{\text{rem}} \tilde{\mathbf{S}}^{\text{rem}}$.
4. Compute the singular value decomposition of $\tilde{\mathbf{S}}^{\text{rem}} = \mathbf{U} \mathbf{\Sigma} \mathbf{W}^T$ with $\mathbf{\Sigma} = \text{diag}(\sigma_j)$. Given a tolerance ϑ , choose the new rank $r_1 \leq 2r$ as the minimal number such that

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

Let \mathbf{S}^{rem} be the $r_1 \times r_1$ diagonal matrix with the r_1 largest singular values, and let \mathbf{U}^{rem} and \mathbf{W}^{rem} contain the first r_1 columns of \mathbf{U} and \mathbf{W} , respectively. Set $\mathbf{X}^{\text{rem}} = \tilde{\mathbf{X}}^{\text{rem}} \mathbf{U}^{\text{rem}}$ and $\mathbf{V}^{\text{rem}} = \tilde{\mathbf{V}}^{\text{rem}} \mathbf{W}^{\text{rem}}$.

5. Set $\hat{\mathbf{X}} = [\mathbf{X}^{\text{cons}}, \mathbf{X}^{\text{rem}}]$ and $\hat{\mathbf{V}} = [\mathbf{e}_1, \mathbf{V}^{\text{rem}}]$. Perform a QR-decomposition of $\hat{\mathbf{X}} = \mathbf{X}^1 \mathbf{R}^1$ and $\hat{\mathbf{V}} = \mathbf{V}^1 \mathbf{R}^2$.
6. Set

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

The updated solution at time $t_1 = t_0 + \Delta t$ is then given by $\mathbf{u}^1 = \mathbf{X}^1 \mathbf{S}^1 \mathbf{V}^{1,\top}$. Then, the scheme is conservative.

THEOREM 6.1. *The scheme (5.5) is locally conservative. That is, for the scalar flux at time t_n denoted by $\Phi_j^n = X_{j\ell}^n S_{\ell m}^n V_{0m}^n$, where $n \in \{0, 1\}$ and $u_{jk}^0 = X_{j\ell}^0 S_{\ell m}^0 V_{km}^0$, it fulfills the conservation law*

$$(6.1a) \quad \Phi_j^1 = \Phi_j^0 - \Delta t D_{ji}^x u_{i\ell}^0 A_{0\ell} + \Delta t D_{ji}^{xx} u_{i\ell}^0 |A|_{0\ell} + \sigma \Delta t (B_j^1 - \Phi_j^1),$$

$$(6.1b) \quad B_j^1 = B_j^0 + \sigma \Delta t (\Phi_j^1 - B_j^1).$$

Proof. The conservative truncation step is designed such that it does not alter the first column of $\tilde{\mathbf{X}}^1 \tilde{\mathbf{S}}^1 \tilde{\mathbf{V}}^{1,\top}$. Together with the basis augmentation (5.5g) and correction step (5.5f) we then know that

$$\Phi_j^1 = X_{j\ell}^1 S_{\ell m}^1 V_{0m}^1 = \tilde{X}_{j\ell}^1 \tilde{S}_{\ell m}^1 \tilde{V}_{0m}^1 = \hat{u}_{j0}^1.$$

Hence, with (5.5d) and (5.5e) we get that

$$\Phi_j^1 = X_{j\ell}^0 S_{\ell m}^0 V_{0m}^0 - \Delta t D_{ji}^x X_{in}^* \tilde{S}_{nm}^0 V_{\ell m}^* A_{0\ell} + \Delta t D_{ji}^{xx} X_{in}^* \tilde{S}_{nm}^0 V_{\ell m}^* |A|_{0\ell} + \sigma \Delta t (B_j^1 - \Phi_j^1),$$

$$B_j^1 = B_j^0 + \sigma \Delta t (\Phi_j^1 - B_j^1).$$

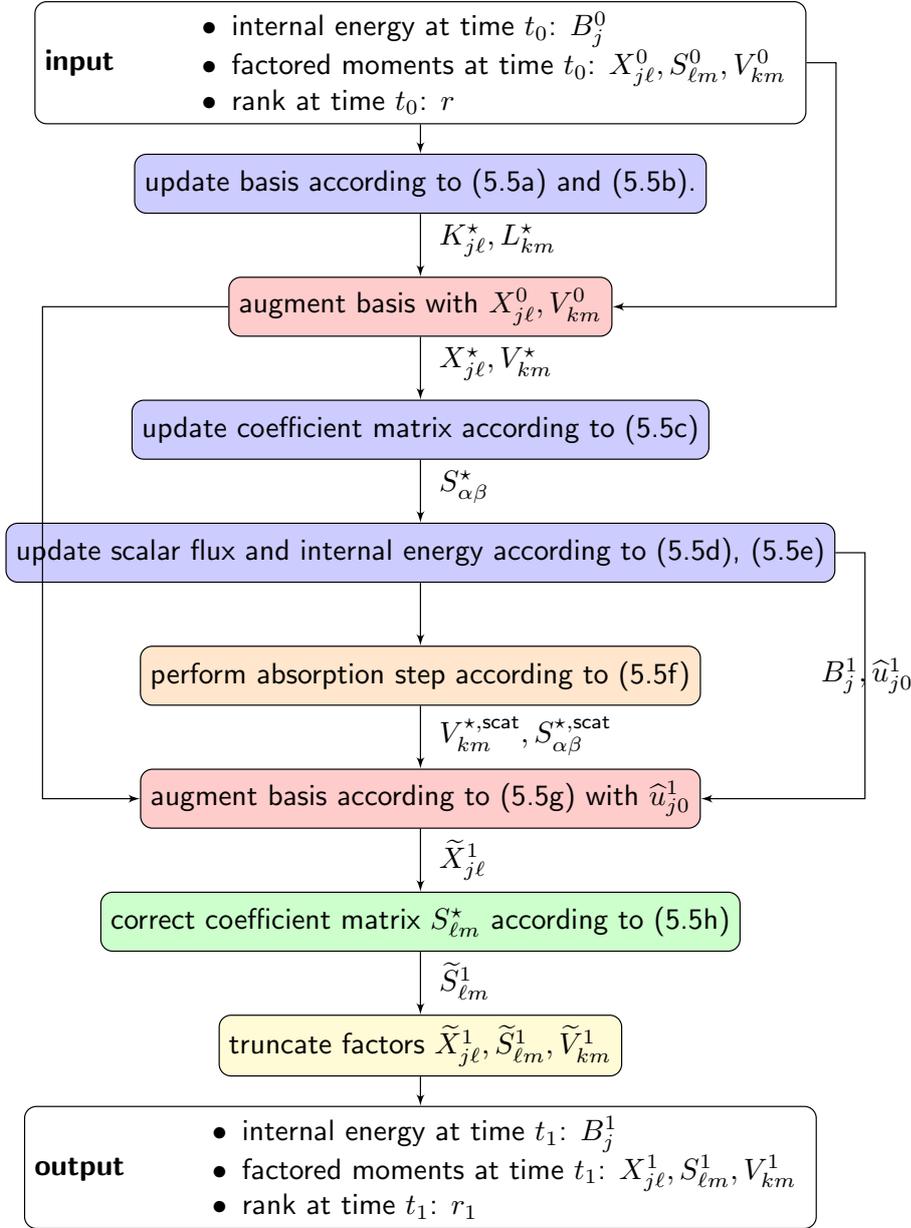


FIG. 1. Flowchart of the stable and conservative method (5.5).

Since the basis augmentation with \mathbf{X}^0 and \mathbf{V}^0 ensures $X_{j\ell}^0 S_{\ell m}^0 V_{km}^0 = X_{in}^* \tilde{S}_{nm}^0 V_{\ell m}^* = u_{i\ell}^0$, the local conservation law (6.1) holds. \square

Hence, equipped with a conservative truncation step, the energy stable algorithm presented in (5.5) conserves mass locally. To give an overview of the algorithm, we visualize the main steps in Figure 1.

7. Numerical results. In this section we give numerical results to validate the proposed DLRA algorithm. The source code to reproduce the presented numerical results is openly available; see [2].

7.1. 1D plane source. We consider the thermal radiative transfer equations as described in (2.1a) on the spatial domain $D = [-10, 10]$. As initial distribution we choose a cutoff Gaussian

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

with constant deviation $\sigma_{\text{IC}} = 0.03$. Particles are initially centered around $x = 1$ and move into all directions $\mu \in [-1, 1]$. The initial value for the internal energy is set to $B^0 = 1$ and we start computations with a rank of $r = 20$. The opacity σ is set to the constant value of 1. Note that this setting is an extension of the so-called *plane source* problem, which is a common test case for the radiative transfer equation [16]. In the context of dynamical low-rank approximation it has been studied in [6, 21, 34, 36]. We compare the solution of the full coupled-implicit system without DLRA, which reads

$$(7.1a) \quad u_{jk}^1 = u_{jk}^0 - \Delta t D_{ji}^x u_{i\ell}^0 A_{k\ell} + \Delta t D_{ji}^{xx} u_{i\ell}^0 |A|_{k\ell} + \sigma \Delta t (B_j^1 \delta_{k0} - u_{jk}^1),$$

$$(7.1b) \quad B_j^1 = B_j^0 + \sigma \Delta t (u_{j0}^1 - B_j^1),$$

to the presented energy stable mass conservative DLRA solution from (5.5). We refer to (7.1) as the full system. The total mass at any time t_n shall be defined as $m^n = \Delta x \sum_j (u_{j0}^n + B_j^n)$. As computational parameters we use $n_x = 1000$ cells in the spatial domain and $N = 500$ moments to represent the directional variable. The time step size is chosen as $\Delta t = \text{CFL} \cdot \Delta x$ with a CFL number of $\text{CFL} = 0.99$. In Figure 2 we present computational results for the solution $f(x, \mu)$, the scalar flux $\Phi = \langle f \rangle_\mu$, and the temperature T at the end time $t_{\text{end}} = 8$. Further, the evolution of the rank r in time, and the relative mass error $\frac{|m^0 - m^n|}{\|m^0\|}$ are shown. One can observe that the DLRA scheme captures well the behavior of the full system. For a chosen tolerance of $\vartheta = 10^{-1} \|\Sigma\|_2$ the rank increases up to $r = 24$ before it reduces again. The relative mass error is of order $\mathcal{O}(10^{-14})$. Hence, our proposed scheme is mass conservative up to machine precision.

7.2. 1D Su–Olson problem. For the next test problem we add a source term $Q(x)$ to the previously investigated equations leading to

$$\begin{aligned} \partial_t f(t, x, \mu) + \mu \partial_x f(t, x, \mu) &= \sigma (B(t, x) - f(t, x, \mu)) + Q(x), \\ \partial_t B(t, x) &= \sigma (\langle f(t, x, \cdot) \rangle_\mu - B(t, x)). \end{aligned}$$

In our example we use the source function $Q(x) = \chi_{[-0.5, 0.5]}(x)/a$ with $a = \frac{4\sigma_{\text{SB}}}{c}$ being the radiation constant. Again we consider the spatial domain $D = [-10, 10]$ and choose the initial condition

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

with constant deviation $\sigma_{\text{IC}} = 0.03$ and particles moving into all directions $\mu \in [-1, 1]$. The initial value for the internal energy is set to $B_0 = 50$, and the initial value for the rank is set to $r = 20$. The opacity σ is again chosen to have the constant value of 1. As computational parameters we use $n_x = 1000$ cells in the spatial domain and $N = 500$ moments to represent the directional variable. The time step size is chosen as $\Delta t = \text{CFL} \cdot \Delta x$ with a CFL number of $\text{CFL} = 0.99$. The isotropic source term generates radiation particles flying through and interacting with a background material. The interaction is driven by the opacity σ . In turn, particles heat up the material leading to a traveling temperature front, also called a *Marshak wave* [26]. Again this traveling heat wave can lead to the emission of new particles from the background material generating a particle wave. At a given time point $t_{\text{end}} = 3.16$

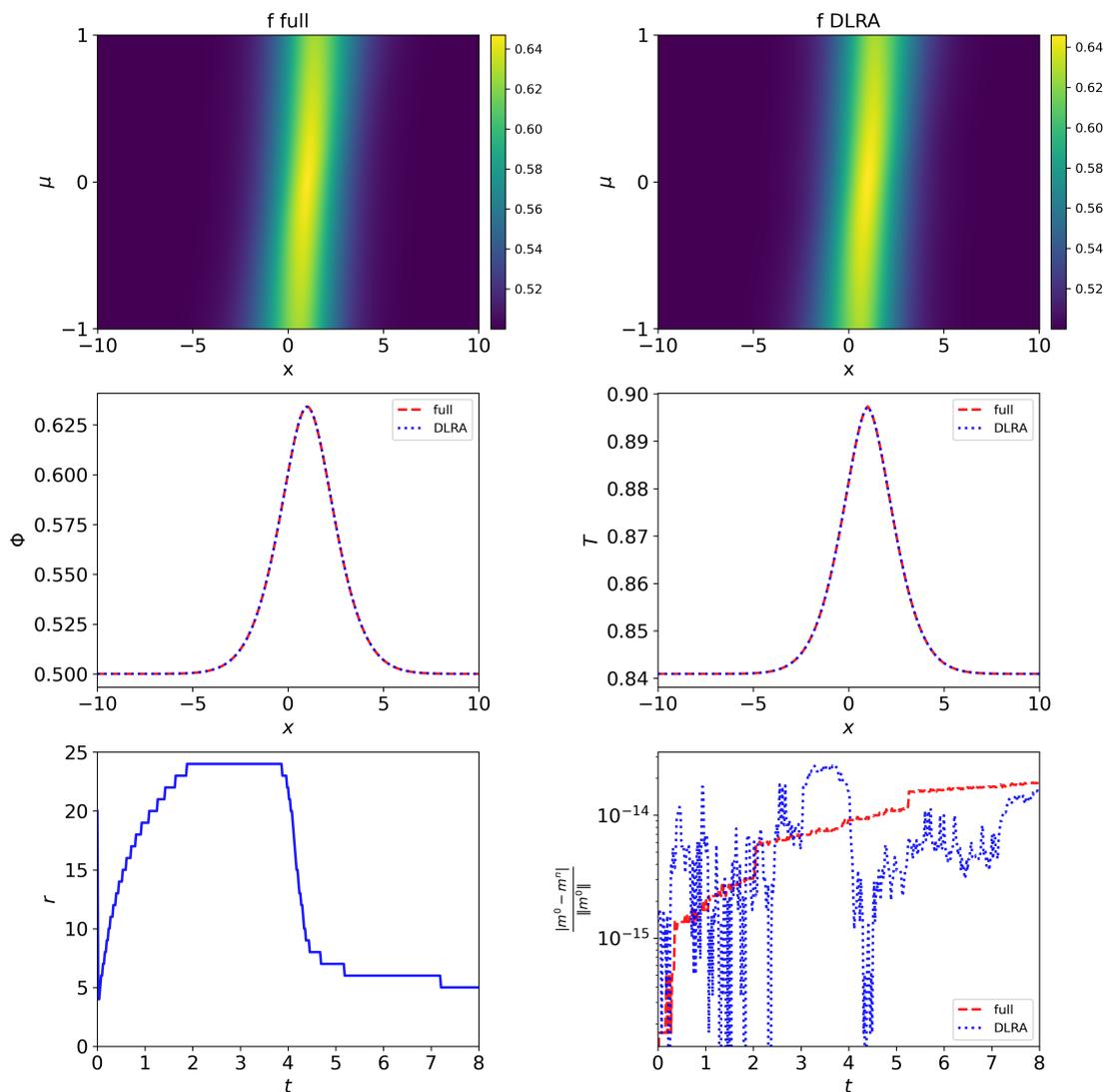


FIG. 2. Top row: Numerical results for the solution $f(x, \mu)$ of the plane source problem at time $t_{\text{end}} = 8$ computed with the full coupled-implicit system (left) and the DLRA system (right). Middle row: Traveling particle (left) and heat wave (right) for both the full system and the DLRA system. Bottom row: Evolution of the rank in time for the DLRA method (left) and relative mass error compared for both methods (right).

this wave can be seen in Figure 3, where we display numerical results for the solution $f(x, \mu)$, the scalar flux $\Phi = \langle f \rangle_\mu$, and the temperature T . We compare the solution of the full coupled-implicit system differing from (7.1) by an additional source term to the presented energy stable mass conservative DLRA solution from (5.5), where we have also added this source term. Further, the evolution of the rank in time is presented for a tolerance parameter of $\vartheta = 10^{-2} \|\Sigma\|_2$. Again we observe that the proposed DLRA scheme approximates well the behavior of the full system. In addition, a very low rank is sufficient to obtain accurate results. Note that due to the source term there is no mass conservation in this example.

7.3. 2D beam. To approve the computational benefits of the presented method we extend it to a two-dimensional setting. The set of equations becomes

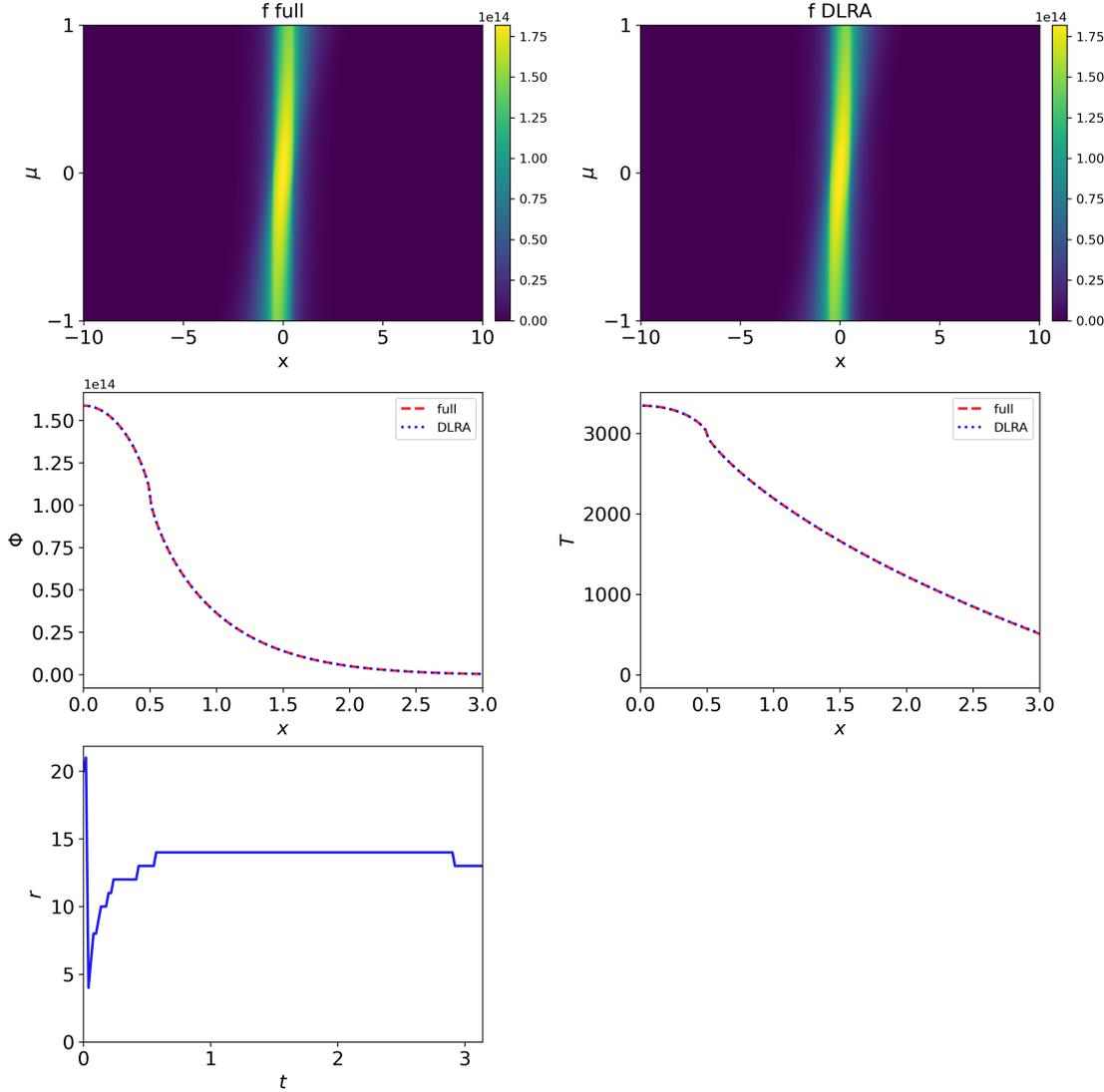


FIG. 3. Top row: Numerical results for the solution $f(x, \mu)$ of the Su-Olson problem at time $t_{\text{end}} = 3.16$ computed with the full coupled-implicit system (left) and the DLRA system (right). Middle row: Traveling particle (left) and heat wave (right) for both the full system and the DLRA system. Bottom row: Evolution of the rank in time for the DLRA method.

$$\begin{aligned} \partial_t f(t, \mathbf{x}, \boldsymbol{\Omega}) + \boldsymbol{\Omega} \cdot \nabla_{\mathbf{x}} f(t, \mathbf{x}, \boldsymbol{\Omega}) &= \sigma(B(t, \mathbf{x}) - f(t, \mathbf{x}, \boldsymbol{\Omega})), \\ \partial_t B(t, \mathbf{x}) &= \sigma(\langle f(t, \mathbf{x}, \cdot) \rangle_{\boldsymbol{\Omega}} - B(t, \mathbf{x})). \end{aligned}$$

For the numerical experiments let $\mathbf{x} = (x_1, x_2) \in [-1, 1] \times [-1, 1]$, $\boldsymbol{\Omega} = (\Omega_1, \Omega_2, \Omega_3) \in \mathcal{S}^2$, and $\sigma = 0.5$. The initial condition of the two-dimensional beam is given by

$$f(t=0, \mathbf{x}, \boldsymbol{\Omega}) = 10^6 \cdot \frac{1}{2\pi\sigma_x^2} \exp\left(-\frac{\|\mathbf{x}\|^2}{2\sigma_x^2}\right) \cdot \frac{1}{2\pi\sigma_\Omega^2} \exp\left(-\frac{(\Omega_1 - \Omega^*)^2 + (\Omega_3 - \Omega^*)^2}{2\sigma_\Omega^2}\right),$$

with $\Omega^* = \frac{1}{\sqrt{2}}$, $\sigma_x = \sigma_\Omega = 0.1$. The initial value for the internal energy is set to $B^0 = 1$, and the initial value for the rank is set to $r = 100$. The total mass at any time t_n shall be defined as $m^n = \Delta x_1 \Delta x_2 \sum_j (u_{j0}^n + B_j^n)$. We perform our computations on a spatial grid with $N_{\text{CellsX}} = 500$ points in x_1 and $N_{\text{CellsY}} = 500$ points in x_2 . For the angular basis we use again a modal approach, namely the spherical harmonics

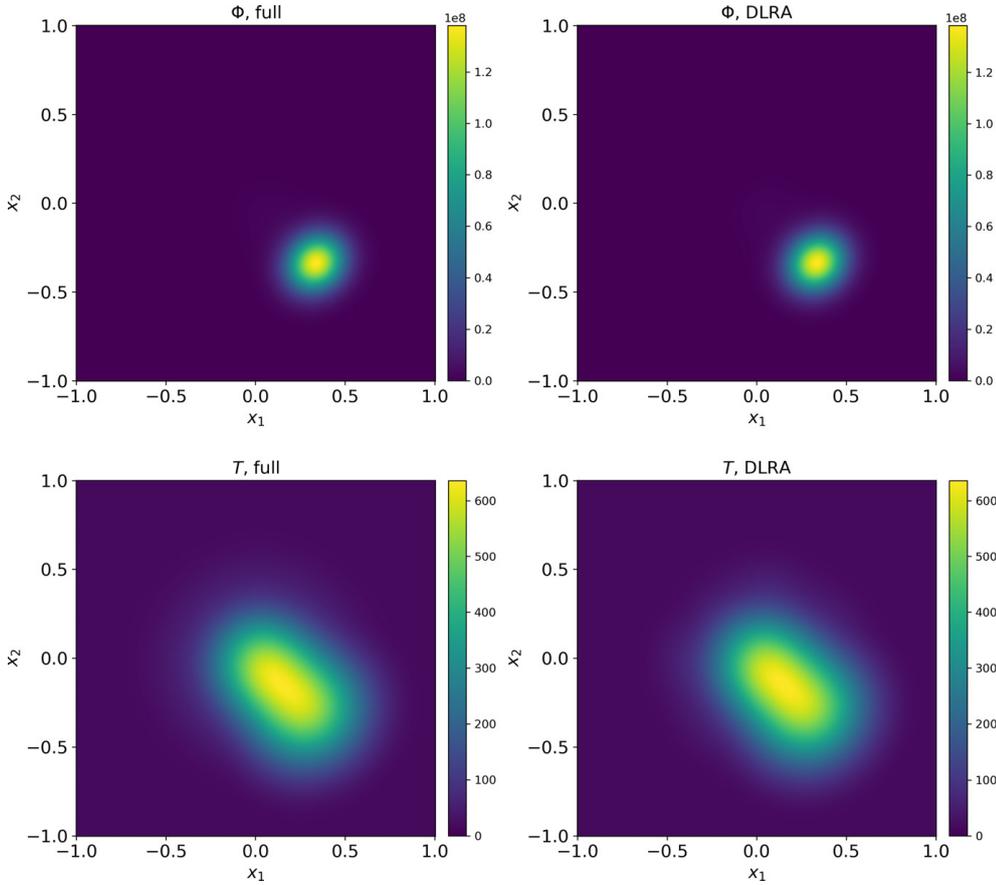


FIG. 4. Numerical results of the scalar flux and the temperature for the 2D beam example for the full coupled-implicit system (left) and the DLRA system (right) at the time $t = 0.5$.

(P_N) method. Technical details can be found in [4, 31, 29], whereas [36, 22] relate the method to dynamical low-rank approximation. The polynomial degree shall be chosen large enough such that the behavior is captured correctly but small enough to stay in a reasonable computational regime. An increasing order of unknowns usually leads to an increasing complexity and therefore to the need of a higher polynomial degree. For our example we use a polynomial degree of $n_{PN} = 29$ corresponding to 900 expansion coefficients in angle. The time step size is chosen as $\Delta t = CFL \cdot \Delta x$ with a CFL number of $CFL = 0.7$. We compare the solution of the two-dimensional full system corresponding to (7.1) to the two-dimensional DLRA solution corresponding to (5.5). The extension to two dimensions is straightforward. In Figure 4 we show numerical results for the scalar flux $\Phi = \int_{S^2} f(t, \mathbf{x}, \cdot) d\Omega$ and the temperature T at the time $t = 0.5$. We again observe the accuracy of the proposed DLRA scheme. For this setup the computational benefit of the DLRA method is significant as the run time compared to the solution of the full problem is reduced by a factor of approximately 8 from 20023 seconds to 2509 seconds. For the evolution of the rank r in time and the relative mass error $\frac{|m^0 - m^n|}{\|m^0\|}$ we consider a time interval up to $t = 1.5$. In Figure 5 one can observe that for a chosen tolerance parameter of $\vartheta = 5 \cdot 10^{-4} \|\Sigma\|_2$ the rank increases but does not approach its allowed maximal value of 100. Further, the relative mass error stagnates and the DLRA method shows its mass conservation property.

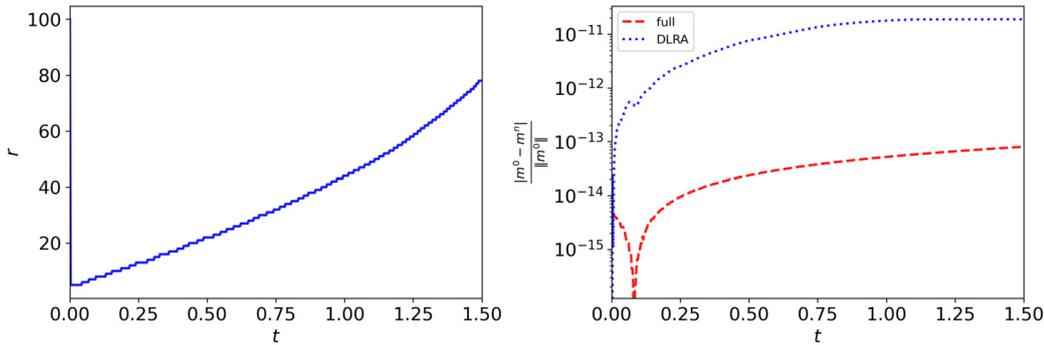


FIG. 5. Evolution of the rank in time for the 2D beam example for the DLRA method (left) and relative mass error compared for both methods (right) until a time of $t = 1.5$.

8. Conclusion and outlook. We have introduced an energy stable and mass conservative dynamical low-rank algorithm for the Su–Olson problem. The key points leading to these properties consist in treating both equations in a coupled-implicit way and using a mass conservative truncation strategy. Numerical examples both in 1D and 2D validate the accuracy of the DLRA method. Its efficiency compared to the solution of the full system can especially be seen in the two-dimensional setting. For future work, we propose to implement the parallel integrator of [7] to further enhance the efficiency of the DLRA method. Moreover, we expect to draw conclusions from this Su–Olson system to the Boltzmann–BGK system and the DLRA algorithm presented in [11] regarding stability and an appropriate choice of the size of the time step.

REFERENCES

- [1] I. ABU-SHUMAYS, *Angular quadratures for improved transport computations*, *Transp. Theory Stat. Phys.*, 30 (2001), pp. 169–204.
- [2] L. BAUMANN, L. EINKEMMER, C. KLINGENBERG, AND J. KUSCH, *Numerical test-cases for “Energy stable and conservative dynamical low-rank approximation for the Su–Olson problem,”* 2023, <https://github.com/JonasKu/publication-Energy-stable-and-conservative-dynamical-low-rank-approximation-for-the-Su-Olson-problem.git>.
- [3] T. CAMMINADY, M. FRANK, K. KÜPPER, AND J. KUSCH, *Ray effect mitigation for the discrete ordinates method through quadrature rotation*, *J. Comput. Phys.*, 382 (2019), pp. 105–123.
- [4] K. M. CASE AND P. F. ZWEIFEL, *Linear Transport Theory*, Addison-Wesley, Reading, MA, 1967.
- [5] G. CERUTI, M. FRANK, AND J. KUSCH, *Dynamical Low-Rank Approximation for Marshak Waves*, CRC 1173 Preprint 2022/76, Karlsruhe Institute of Technology, 2022.
- [6] G. CERUTI, J. KUSCH, AND C. LUBICH, *A rank-adaptive robust integrator for dynamical low-rank approximation*, *BIT*, 62 (2022), pp. 1149–1174.
- [7] G. CERUTI, J. KUSCH, AND C. LUBICH, *A Parallel Rank-Adaptive Integrator for Dynamical Low-Rank Approximation*, preprint, arXiv:2304.05660, 2023.
- [8] G. CERUTI AND C. LUBICH, *An unconventional robust integrator for dynamical low-rank approximation*, *BIT*, 62 (2022), pp. 23–44.
- [9] L. EINKEMMER, J. HU, AND J. KUSCH, *Asymptotic–Preserving and Energy Stable Dynamical Low-Rank Approximation*, preprint, arXiv:2212.12012, 2022.
- [10] L. EINKEMMER, J. HU, AND Y. WANG, *An asymptotic-preserving dynamical low-rank method for the multi-scale multi-dimensional linear transport equation*, *J. Comput. Phys.*, 439 (2021), 110353.
- [11] 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 (2021), pp. B1057–B1080, <https://doi.org/10.1137/21M1392772>.
- [12] L. EINKEMMER AND J. LON, *A mass, momentum, and energy conservative dynamical low-rank scheme for the Vlasov equation*, *J. Comput. Phys.*, 443 (2021), 110493.

- [13] L. EINKEMMER AND C. LUBICH, *A low-rank projector-splitting integrator for the Vlasov–Poisson equation*, SIAM J. Sci. Comput., 40 (2018), pp. B1330–B1360, <https://doi.org/10.1137/18M116383X>.
- [14] L. EINKEMMER, A. OSTERMANN, AND C. SCALONE, *A Robust and Conservative Dynamical Low-Rank Algorithm*, preprint, arXiv:2206.09374, 2022.
- [15] M. FRANK, J. KUSCH, T. CAMMINADY, AND C. D. HAUCK, *Ray effect mitigation for the discrete ordinates method using artificial scattering*, Nucl. Sci. Eng., 194 (2020), pp. 971–988.
- [16] B. D. GANAPOL, *Analytical Benchmarks for Nuclear Engineering Applications, Case Studies in Neutron Transport Theory*, NEA 6292, OECD, 2008.
- [17] W. GUO AND J.-M. QIU, *A Conservative Low Rank Tensor Method for the Vlasov Dynamics*, preprint, arXiv:2201.10397, 2022.
- [18] J. HU AND Y. WANG, *An adaptive dynamical low rank method for the nonlinear Boltzmann equation*, J. Sci. Comput., 92 (2022), 75.
- [19] E. KIERI, C. LUBICH, AND H. WALACH, *Discretized dynamical low-rank approximation in the presence of small singular values*, SIAM J. Numer. Anal., 54 (2016), pp. 1020–1038, <https://doi.org/10.1137/15M1026791>.
- [20] O. KOCH AND C. LUBICH, *Dynamical low-rank approximation*, SIAM J. Matrix Anal. Appl., 29 (2007), pp. 434–454, <https://doi.org/10.1137/050639703>.
- [21] J. KUSCH, L. EINKEMMER, AND G. CERUTI, *On the stability of robust dynamical low-rank approximations for hyperbolic problems*, SIAM J. Sci. Comput., 45 (2023), pp. A1–A24, <https://doi.org/10.1137/21M1446289>.
- [22] J. KUSCH AND P. STAMMER, *A robust collision source method for rank adaptive dynamical low-rank approximation in radiation therapy*, ESAIM Math. Model. Numer. Anal., 57 (2023), pp. 865–891.
- [23] K. D. LATHROP, *Ray effects in discrete ordinates equations*, Nucl. Sci. Eng., 32 (1968), pp. 357–369.
- [24] K. D. LATHROP, *Remedies for ray effects*, Nucl. Sci. Eng., 45 (1971), pp. 255–268.
- [25] C. LUBICH AND I. V. OSELEDETS, *A projector-splitting integrator for dynamical low-rank approximation*, BIT, 54 (2014), pp. 171–188.
- [26] R. E. MARSHAK, *Effect of radiation on shock wave behavior*, Phys. Fluids, 1 (1958), pp. 24–29.
- [27] K. A. MATHEWS, *On the propagation of rays in discrete ordinates*, Nucl. Sci. Eng., 132 (1999), pp. 155–180.
- [28] 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 (2008), pp. 7561–7586.
- [29] R. G. MCCLARREN AND C. D. HAUCK, *Robust and accurate filtered spherical harmonics expansions for radiative transfer*, J. Comput. Phys., 229 (2010), pp. 5597–5614.
- [30] 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 (2008), pp. 389–403.
- [31] 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 (2008), pp. 2864–2885.
- [32] J. MOREL, T. WAREING, R. LOWRIE, AND D. PARSONS, *Analysis of ray-effect mitigation techniques*, Nucl. Sci. Eng., 144 (2003), pp. 1–22.
- [33] 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 (2000), pp. 619–634.
- [34] 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 (2021), 110672.
- [35] Z. PENG AND R. G. MCCLARREN, *A sweep-based low-rank method for the discrete ordinate transport equation*, J. Comput. Phys., 473 (2023), 111748.
- [36] Z. PENG, R. G. MCCLARREN, AND M. FRANK, *A low-rank method for two-dimensional time-dependent radiation transport calculations*, J. Comput. Phys., 421 (2020), 109735.
- [37] G. C. POMRANING, *The non-equilibrium Marshak wave problem*, J. Quant. Spectrosc. Radiat. Transf., 21 (1979), pp. 249–261.
- [38] B. SU AND G. L. OLSON, *An analytical benchmark for non-equilibrium radiative transfer in an isotropically scattering medium*, Ann. Nucl. Energy, 24 (1997), pp. 1035–1055.
- [39] J. TENCER, *Ray effect mitigation through reference frame rotation*, J. Heat Transf., 138 (2016), 112701.