

Randomized low-rank Runge-Kutta methods*

Hei Yin Lam[†]

Gianluca Ceruti[‡]

Daniel Kressner[†]

September 11, 2024

Abstract

This work proposes and analyzes a new class of numerical integrators for computing low-rank approximations to solutions of matrix differential equation. We combine an explicit Runge-Kutta method with repeated randomized low-rank approximation to keep the rank of the stages limited. The so-called generalized Nyström method is particularly well suited for this purpose; it builds low-rank approximations from random sketches of the discretized dynamics. In contrast, all existing dynamical low-rank approximation methods are deterministic and usually perform tangent space projections to limit rank growth. Using such tangential projections can result in larger error compared to approximating the dynamics directly. Moreover, sketching allows for increased flexibility and efficiency by choosing structured random matrices adapted to the structure of the matrix differential equation. Under suitable assumptions, we establish moment and tail bounds on the error of our randomized low-rank Runge-Kutta methods. When combining the classical Runge-Kutta method with generalized Nyström, we obtain a method called Rand RK4, which exhibits fourth-order convergence numerically – up to the low-rank approximation error. For a modified variant of Rand RK4, we also establish fourth-order convergence theoretically. Numerical experiments for a range of examples from the literature demonstrate that randomized low-rank Runge-Kutta methods compare favorably with two popular dynamical low-rank approximation methods, in terms of robustness and speed of convergence.

1 Introduction

In this work, we aim at approximating the solution $A(t)$ to large-scale matrix differential equations of the form

$$\dot{A}(t) = F(A(t)), \quad A(0) = A_0 \in \mathbb{R}^{m \times n}. \quad (1)$$

In many situations of practical interest, an autonomous ordinary differential equation can be naturally viewed as such a matrix differential equation. Examples include applications in physics [13, 22], uncertainty quantification [2], and machine learning [21]. For large m , n , the solution of (1) becomes expensive; in fact, it may not even be possible to store the entire matrix $A(t)$ explicitly. To circumvent this limitation, model order reduction techniques can be employed. An increasingly popular approach is based on exploiting (approximate) low-rank structure of $A(t)$, which arises, for example, from smoothness properties of the underlying physical system. In particular, dynamical low-rank approximation [15] approximates $A(t)$ by evolving matrices $Y(t)$ on the manifold \mathcal{M}_r of rank- r matrices. As only the rank- r factors of $Y(t)$ need to be stored, this already reduces memory requirements significantly when $r \ll m, n$.

*This work was supported by the SNSF research project *Fast algorithms from low-rank updates*, grant number: 200020_178806.

[†]École Polytechnique Fédérale de Lausanne (EPFL), Institute of Mathematics, Switzerland. hysan.lam@epfl.ch, daniel.kressner@epfl.ch.

[‡]Department of Mathematics, University of Innsbruck, Innsbruck, Austria. gianluca.ceruti@uibk.ac.at.

By the Dirac-Frenkel variational principle, the matrix $Y(t)$ is obtained by solving the differential equation

$$\dot{Y}(t) = P_r(Y(t))F(Y(t)), \quad Y(0) = Y_0 \in \mathcal{M}_r, \quad (2)$$

where $P_r(Y(t))$ denotes the orthogonal projection onto $T_{Y(t)}\mathcal{M}_r$, the tangent space of \mathcal{M}_r at $Y(t)$. To also achieve a reduction of computational cost, one needs to exploit the low-rank structure of $Y(t)$ when integrating (2). As shown in [15], this can be achieved by rewriting (2) as a system of differential equations for the rank- r factors of Y . However, directly integrating this system with standard explicit time integration methods often leads to poor approximation results, unless (very) small time step sizes are used. This is caused by the additional stiffness introduced by small singular values of $Y(t)$. To address this issue, special integrators have been proposed that are robust to the presence of small singular values and allow for much larger step sizes. These methods include the projected splitting integrator [19], projection methods [14, 9] as well as Basis Update & Galerkin (BUG) integrators [6, 7, 8]. Under the assumption

$$\|F(Y) - P_r(Y)F(Y)\|_F \leq \tilde{\epsilon}, \text{ for all } Y \in \mathcal{M}_r \cap \{\text{suitable neighbourhood of } A(t)\} \quad (3)$$

all these methods exhibit at least first-order convergence up to $\mathcal{O}(\tilde{\epsilon})$, both theoretically and numerically. The mid-point BUG [6], a variant of the parallel integrator [18], and the projected Runge–Kutta methods [14] are the only provable second-order integrators up to $\mathcal{O}(\tilde{\epsilon})$. Projected Runge–Kutta methods can also achieve higher order.

Assumption (3), which says that $F(Y)$ is nearly contained in the tangent space, is arguably a strong assumption. For small $\tilde{\epsilon} > 0$ it implies, at least for short times, that $A(t)$ can be well approximated by a rank- r matrix, but the reverse is not true. In particular, it is possible that $A(t)$ can be well approximated by a rank- r matrix even if (3) is not satisfied with small $\tilde{\epsilon}$. According to [14], this can occur, for instance, when the manifold \mathcal{M}_r close to $A(t)$ has tiny, high-frequency wiggles, or when the range and co-range of Y are contained in the orthogonal complement of the range and co-range of $F(Y)$, respectively [1]. Concrete examples are given in Section 4. When Assumption (3) is not satisfied with small $\tilde{\epsilon}$, the use of tangent space projections in numerical methods bears the danger of introducing unacceptably high errors.

In this work, we develop low-rank time integration methods for (1) that do not rely on (3) but only require $A(t)$ to admit accurate low-rank approximations. Our approach is based on the notion of projected integrators [11, Ch. IV.4], which first perform a standard time integration step and then project back to the manifold. For the manifold \mathcal{M}_r , the efficiency of projected integrators is impaired by the occurrence of high-rank matrices, e.g., during the intermediate stages of a Runge–Kutta method. In [14], this issue is addressed by repeatedly applying tangent space projection, which limits the rank to $2r$ at the expense of having to impose (3).

In this work, we take a novel approach to avoid the high ranks encountered by projected integrators. Our approach is based on performing randomized low-rank approximation, which uses random sketches instead of tangent space projections. For a *constant* matrix B , such randomized approaches have been studied intensively during the last 1-2 decades, including the popular randomized SVD [12] and the (generalized) Nyström approximation [20, 24]. The generalized Nyström approximation utilizes two sketches $B\Omega$ and $\Psi^T B$ to approximately capture the range and co-range of B , where Ω and Ψ are random matrices with the number of columns chosen to be slightly larger than r .

To the best of our knowledge, this is the first work that proposes and analyzes randomized low-rank approximation methods for time integration. The randomized low-rank Runge–Kutta (RK) methods proposed in this work combine explicit RK methods with randomized low-rank approximation. Our analysis applies to any randomized low-rank approximation satisfying the moment assumptions defined in Section 2.3, but for simplicity we will restrict our algorithmic considerations to the generalized Nyström approximation. Assuming that the dynamics generated by F preserve rank- r matrices approximately, we derive a probabilistic result that

establishes a convergence order (up to the level of rank- r approximation error) based on the so-called stage order of the underlying RK method. This matches the order established in [14] for projected RK methods. However, unlike the results in [14], our numerical experiments indicate that randomized low-rank RK methods actually achieve the usual convergence order of the RK method, which can be significantly higher. For the randomized low-rank RK method based on RK 4, we also establish order 4 theoretically when allowing for modest intermediate rank increases in the stages. This compares favorably to order 2 implied by the techniques from in [14].

The remainder of the paper is organized as follows. In Section 2, after providing some preliminaries, we propose and analyze an idealized projection method, which assumes that the exact flow of (1) is given. We show that applying randomized low-rank approximation causes, with high probability, little to no harm to time integration. In Section 3, we propose a practical method that uses an RK method to approximate the exact flow and provide error analysis based on the stage order. Furthermore, we prove that if we allow rank increase in the intermediate stages then the classical RK method combined with randomized low-rank approximation can still achieve convergence order 4, up to the level of the low-rank approximation error. Finally, in Section 4, we provide a range of numerical experiments that confirm the theoretical results and demonstrate the robust convergence of randomized RK methods.

2 Preliminaries and an idealized randomized projection method

In this section, we provide preliminaries on (randomized) low-rank approximation and introduce an idealized randomized projection method that allows one to study the impact of randomization in an isolated fashion. In the following, $\|\cdot\|_F$ denotes the Frobenius norm of a (constant) matrix and $\|Y\|_{L_q} = (\mathbb{E}[\|Y\|_F^q])^{\frac{1}{q}}$ denotes the L_q norm, for some $q \geq 1$, of a random matrix Y .

2.1 Assumptions

Our analysis will be based on the following three assumptions on F . The first two assumptions are the same as in [6, 14, 18], while the third one is a modification of the usual low-rank approximability assumption in dynamical low-rank approximation.

Assumption 1. We assume that F is Lipschitz continuous, that is, there is a Lipschitz constant $L > 0$ such that

$$\|F(X) - F(Y)\|_F \leq L\|X - Y\|_F \quad \text{for all } X, Y \in \mathbb{R}^{m \times n}. \quad (4)$$

By the Picard-Lindelöf theorem, this implies that the solution of (1) exists and it is unique for some finite time interval.

Assumption 2. Let Φ_F^t denote the exact flow of F , that is, given the solution $A(t)$ of (1), we have that $A(t) = \Phi_F^t(A_0)$. For a method of order τ , we assume that the first τ derivatives

$$\frac{d^{j+1}}{dt^{j+1}} \Phi_F^t(Y), \quad j = 0, 1, \dots, \tau, \quad (5)$$

have uniformly bounded norm for all $Y \in \mathbb{R}^{m \times n}$. This assumption is needed when, e.g., performing local error analysis of higher-order methods [10, Chapter II.1].

Assumption 3. To ensure low-rank approximability, we assume for every $h \leq h_0$ that

$$\|\Phi_F^h(Y) - [\Phi_F^h(Y)]_r\|_F \leq C_M h \epsilon \quad \text{for all } Y \in \mathcal{M}_r, \quad (6)$$

where C_M is a constant that depends on L and h_0 only. Here and in the following, we use $[\cdot]_r$ to denote a best rank- r approximation of a matrix. Assumption (6) replaces the assumption

(3), usually made in the analysis of dynamical low-rank approximation [13, 14]; see Section 2.2 below for a more detailed discussion.

In this paper, we assume that (4), (5) and (6) hold *globally*, because this greatly simplifies the analysis and the results. As in (3), it is common to impose such assumptions only in a neighbourhood of the exact solution $A(t)$. When using randomized techniques, there is always a tiny but non-zero probability that the approximation leaves *any* neighborhood. In Remark 6, we discuss how our analysis can be modified to account for this effect, resulting in (slightly) weaker results.

2.2 Low-rank approximability

In this section, we discuss the relation between the assumption (3) on the tangent space projection and the low-rank assumption (6).

First of all, (3) implies (6). To see this, suppose that F satisfies (3). Then, by [14, Lemma 1], there exists $h_0 > 0$ such that

$$\|\Phi_F^h(Y) - \Phi_{P_r F}^h(Y)\|_F \leq \tilde{\epsilon} \int_0^h e^{Ls} ds \leq e^{Lh} h \tilde{\epsilon}, \quad \forall 0 \leq h \leq h_0.$$

Because of $\Phi_{P_r F}^h(Y) \in \mathcal{M}_r$, this implies

$$\|\Phi_F^h(Y) - \llbracket \Phi_F^h(Y) \rrbracket_r\|_F \leq \|\Phi_F^h(Y) - \Phi_{P_r F}^h(Y)\|_F \leq e^{Lh} h \tilde{\epsilon}.$$

Hence, (6) is satisfied with $C_M = e^{Lh_0}$ and $\epsilon = \tilde{\epsilon}$.

Assumption (6) does not necessarily imply (3), that is, the existence of a good low-rank approximation does not require the tangent space projection error to be small. In other words, Assumption (6) is weaker. This is demonstrated by the following example.

Example 1. Consider the rank-2 approximation of the differential equation

$$\dot{Y}(t) = F(Y(t)) = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & -10 \end{pmatrix} Y(t) + \begin{pmatrix} 0 & 0 & 0 \\ 0 & 10^{-5}e^{-1} & 0 \\ 0 & 0 & 0 \end{pmatrix}, \quad Y(0) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 10^{-6}e \end{pmatrix}.$$

We have that

$$\Phi_F^h(Y) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 10^{-6}(e^{10h-1} - e^{-1}) & 0 \\ 0 & 0 & 10^{-6}e^{1-10h} \end{pmatrix}$$

admits an excellent rank-2 approximation at $h = 1$: $\|\Phi_F^1(Y) - \llbracket \Phi_F^1(Y) \rrbracket_2\|_F \leq 1.24 \times 10^{-10}$. On the other hand, the tangent space projection

$$\Phi_{P_r F}^h(Y) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 10^{-6}e^{1-10h} \end{pmatrix}.$$

results in a rank-2 matrix with a much larger error: $\|\Phi_F^1(Y) - \Phi_{P_r F}^1(Y)\|_F \geq 0.008$.

2.3 Randomized low-rank approximation

Conceptually, a rank- r approximation is a map $\mathcal{R} : \mathbb{R}^{m \times n} \rightarrow \mathcal{M}_r$. When randomization is used, the map \mathcal{R} is random, usually due to the use of random matrices for sketching. We measure the quality of \mathcal{R} through moments, which will later be used to derive concentration inequalities. We say that \mathcal{R} satisfies the moment assumption for $q \geq 1$ if

$$\|\mathcal{R}(Z) - Z\|_{L_q} \leq C_{\mathcal{R}} \|Z - \llbracket Z \rrbracket_r\|_F, \quad (7)$$

holds for fixed but arbitrary $Z \in \mathbb{R}^{m \times n}$, with a constant $C_{\mathcal{R}}$ being independent of Z .

Randomized low-rank approximations that utilize Gaussian random matrices for sketching usually satisfy the moment assumption (7). In the following, we will establish this fact for the generalized Nyström method [20, 24], which proceeds as follows. Given oversampling parameters $p, \ell \in \mathbb{N}$ and random matrices $\Omega \in \mathbb{R}^{n \times (r+p)}$, $\Psi \in \mathbb{R}^{m \times (r+p+\ell)}$, generalized Nyström constructs a rank- r approximation of Z by first performing an oblique projection onto $\text{span}(Z\Omega)$ and then truncating to rank r :

$$Z \approx \llbracket Z\Omega(\Psi^T Z\Omega)^\dagger \Psi^T Z \rrbracket_r := \mathcal{N}(Z), \quad (8)$$

where $(\cdot)^\dagger$ denotes the pseudoinverse. If $\Psi^T Z\Omega \in \mathbb{R}^{(r+p+\ell) \times (r+p)}$ has full column rank, we have the equivalent expression

$$\mathcal{N}(Z) = Q \llbracket (\Psi^T Q)^\dagger \Psi^T Z \rrbracket_r, \quad (9)$$

where Q is an orthonormal basis of $\text{span}(Z\Omega)$, computed by, e.g., a QR factorization, $Z\Omega = QR$. For dense and unstructured matrices Ω and Ψ , computing the sketches $Z\Omega$ and $\Psi^T Z$ usually dominates the overall computational cost. In particular, this is true when Ω and Ψ are Gaussian random matrices, i.e., their entries are independent standard normal Gaussian random variables. The following theorem shows that generalized Nyström satisfies the moment assumption (7) in this case.

Theorem 2. *Suppose that $\Omega \in \mathbb{R}^{n \times (r+p)}$ and $\Psi \in \mathbb{R}^{m \times (r+p+\ell)}$ are independent standard Gaussian matrices with $p, \ell \geq 4$. Setting $q = \min\{p, \ell\}$, it holds for $Z \in \mathbb{R}^{m \times n}$ that*

$$\|\mathcal{N}(Z) - Z\|_{L_q} = (\mathbb{E}[\|\mathcal{N}(Z) - Z\|_F^q])^{\frac{1}{q}} \leq C_{\mathcal{N}} \|Z - \llbracket Z \rrbracket_r\|_F$$

with $C_{\mathcal{N}} = 1 + 2\sqrt{(1+r+p)(1+r)}$.

Proof. By the triangle inequality

$$\begin{aligned} \|\mathcal{N}(Z) - Z\|_{L_q} &\leq \|\llbracket Z\Omega(\Psi^T Z\Omega)^\dagger \Psi^T Z \rrbracket_r - Z\Omega(\Psi^T Z\Omega)^\dagger \Psi^T Z\|_{L_q} + \|Z\Omega(\Psi^T Z\Omega)^\dagger \Psi^T Z - Z\|_{L_q} \\ &\leq \|\llbracket Z \rrbracket_r - Z\Omega(\Psi^T Z\Omega)^\dagger \Psi^T Z\|_{L_q} + \|Z\Omega(\Psi^T Z\Omega)^\dagger \Psi^T Z - Z\|_{L_q} \\ &\leq \|\llbracket Z \rrbracket_r - Z\|_{L_q} + 2\|Z\Omega(\Psi^T Z\Omega)^\dagger \Psi^T Z - Z\|_{L_q}. \end{aligned}$$

To bound the second term, we follow the proof of [16, Theorem 11]. Considering an orthonormal basis Q of $Z\Omega$, one obtains that

$$\begin{aligned} \mathbb{E}[\|Z - Z\Omega(\Psi^T Z\Omega)^\dagger \Psi^T Z\|_F^q] &\leq (1+r+p)^{q/2} \mathbb{E}[\|(I - QQ^T)Z\|_F^q] \\ &\leq (1+r+p)^{q/2} (\mathbb{E}^p[\|(I - QQ^T)Z\|_F]^q) \\ &= (1+r+p)^{q/2} \left(\mathbb{E}^{p/2}[\|(I - QQ^T)Z\|_F^2] \right)^{q/2} \\ &\leq (1+r+p)^{q/2} ((1+r)\|Z - \llbracket Z \rrbracket_r\|_F^2)^{q/2}. \end{aligned}$$

Therefore, $\|\mathcal{N}(Z) - Z\|_{L_q} \leq \|Z - \llbracket Z \rrbracket_r\|_F + 2\sqrt{(1+r+p)(1+r)}\|Z - \llbracket Z \rrbracket_r\|_F = C_{\mathcal{N}}\|Z - \llbracket Z \rrbracket_r\|_F$. \square

To keep our developments concrete, we will always use generalized Nyström instead of an abstract randomized method \mathcal{R} in the rest of the paper. However, the theoretical results remain valid for any \mathcal{R} satisfying (7).

2.4 Idealized randomized projection method

Following [11, Ch. IV.4], a rank- r approximation to the solution $A((i+1)h)$ is obtained by combining exact integration with rank- r truncation

$$Y_{i+1} = \llbracket \Phi_F^h(Y_i) \rrbracket_r, \quad Y_0 = \llbracket A_0 \rrbracket_r.$$

This is an idealized integrator because the exact flow Φ_F^h still needs to be approximated in order to obtain a practical method. Replacing rank- r truncation by the generalized Nyström method one gets the *idealized randomized projection method*

$$Y_{i+1} = \mathcal{N}_i(\Phi_F^h(Y_i)) = \llbracket \Phi_F^h(Y_i) \Omega_i (\Psi_i \Phi_F^h(Y_i) \Omega_i)^\dagger \Psi_i^T \Phi_F^h(Y_i) \rrbracket_r. \quad Y_0 = \mathcal{N}_0(A_0). \quad (10)$$

The subscript i of \mathcal{N} is used to emphasize that the generalized Nyström method is used with different (independent) $\Omega_i \in \mathbb{R}^{n \times (r+p)}$, $\Psi_i \in \mathbb{R}^{m \times (r+p+\ell)}$ in every time step.

2.4.1 Error analysis

In this section, we provide an error analysis of the idealized method (10) when Ω_i, Ψ_i are independent standard Gaussian matrices. The following theorem provides a bound on the L_q norm of the error. The proof basically follows from the proof of [14, Theorem 2], with the Frobenius norm replaced by the L_q norm. It is included for convenience, because similar arguments will be used again below.

Theorem 3. *With the assumptions stated in Section 2.2 and assuming $\|\llbracket A_0 \rrbracket_r - A_0\|_F \leq \delta$ holds for the initial data, the method (10) with independent standard Gaussian matrices $\Omega_i \in \mathbb{R}^{n \times (r+p)}$, $\Psi_i \in \mathbb{R}^{m \times (r+p+\ell)}$ and oversampling parameters $p, \ell \geq 4$ satisfies the error estimate*

$$\|Y_N - A(Nh)\|_{L_q} \leq C(\delta + \epsilon)$$

for $q = \min\{p, \ell\}$ on a finite time-interval $0 \leq Nh \leq T$ for every $0 < h \leq h_0$. The constant C only depends on L, T, h_0, C_M , and C_N .

Proof. We first note that \mathcal{N}_i is stochastically independent of $\Phi_F^h(Y_i)$. By the the law of total expectation, Theorem 2 and Assumption (6), we get

$$\begin{aligned} \|\mathcal{N}_i(\Phi_F^h(Y_i)) - \Phi_F^h(Y_i)\|_{L_q} &= \left(\mathbb{E} \left[\mathbb{E} [\|\mathcal{N}_i(\Phi_F^h(Y_i)) - \Phi_F^h(Y_i)\|_F^q | Y_i] \right] \right)^{1/q} \\ &\leq \left(\mathbb{E} [C_N^q \|\Phi_F^h(Y_i) - \llbracket \Phi_F^h(Y_i) \rrbracket_r\|_F^q] \right)^{1/q} \leq C_N C_M \epsilon h. \end{aligned} \quad (11)$$

To bound the L_q norm of the global error, we follow the proof of [14, Theorem 2] and use a telescoping sum:

$$\|Y_N - A(Nh)\|_{L_q} = \left\| \sum_{i=1}^N (\Phi_F^{(N-i)h}(Y_i) - \Phi_F^{(N-i+1)h}(Y_{i-1})) + \Phi_F^{Nh}(Y_0) - \Phi_F^{Nh}(A_0) \right\|_{L_q} \leq \sum_{i=0}^N E_i,$$

where we define

$$\begin{aligned} E_i &= \|\Phi_F^{(N-i)h}(Y_i) - \Phi_F^{(N-i)h}(\Phi_F^h(Y_{i-1}))\|_{L_q}, \quad i = 1, \dots, N, \\ E_0 &= \|\Phi_F^{Nh}(Y_0) - \Phi_F^{Nh}(A_0)\|_{L_q}. \end{aligned}$$

The Lipschitz continuity of F implies $\|\Phi_F^t(X) - \Phi_F^t(Y)\|_F \leq e^{Lt}\|X - Y\|_F$. Therefore, by the law of total expectation and (11),

$$\begin{aligned}
E_i &= \left(\mathbb{E} \left[\mathbb{E} \left[\|\Phi_F^{(N-i)h}(Y_i) - \Phi_F^{(N-i)h}(\Phi_F^h(Y_{i-1}))\|_F^q | Y_{i-1} \right] \right] \right)^{\frac{1}{q}} \\
&\leq e^{Lh(N-i)} \left(\mathbb{E} \left[\mathbb{E} \left[\|Y_i - \Phi_F^h(Y_{i-1})\|_F^q | Y_{i-1} \right] \right] \right)^{\frac{1}{q}} \\
&= e^{Lh(N-i)} \left(\mathbb{E} \left[\mathbb{E} \left[\|\mathcal{N}_{i-1}(\Phi_F^h(Y_{i-1})) - \Phi_F^h(Y_{i-1})\|_F^q | Y_{i-1} \right] \right] \right)^{\frac{1}{q}} \\
&\leq C_N \cdot C_M e^{Lh(N-i)} \epsilon h.
\end{aligned} \tag{12}$$

In summary,

$$\|Y_N - A(Nh)\|_{L_q} \leq C_N e^{LNh} \delta + C_N \cdot C_M \epsilon \sum h e^{Lh(N-i)},$$

which yields the desired result by bounding the sum by an integral, as in [14, P.80]. \square

The Markov inequality turns the moment bound of Theorem 3 into tail bounds for the approximation error.

Corollary 4. *With the assumptions and notation stated in Theorem 3, the error estimate*

$$\|Y_N - A(Nh)\|_F \leq C\eta(\delta + \epsilon),$$

holds for any $\eta \geq 1$ with probability at least $1 - \eta^{-q}$.

Proof. By Markov's inequality and Theorem 3,

$$\Pr\{\|Y_N - A(Nh)\|_F \geq C\eta(\delta + \epsilon)\} \leq \left(\frac{[\mathbb{E}\|Y_N - A(Nh)\|_F^q]^{1/q}}{C\eta(\delta + \epsilon)} \right)^q \leq \frac{1}{\eta^q}.$$

\square

Corollary 4 states that the generalized Nyström method produces an error on the level of $\epsilon + \delta$ with high probability. Under the assumptions stated in Section 2.2, the same type of error bound is obtained when using exact rank- r truncations.

3 Randomized low-rank Runge-Kutta methods

To turn the idealized projection method (10) into a practical method, we need to combine it with a time-integration method, e.g., a RK method. However, directly replacing the exact flow Φ_F by a RK method will result in high ranks in the intermediate stages. To mitigate this issue, we also apply the generalized Nyström method to these intermediate stages. To make this idea specific, let us consider a general explicit Runge-Kutta method with s stages applied to the matrix differential equation (1):

$$\begin{aligned}
\tilde{Z}_j &= A_i + h \sum_{l=1}^{j-1} a_{jl} F(\tilde{Z}_l), \quad j = 1, \dots, s, \\
A_{i+1} &= A_i + h \sum_{j=1}^s b_j F(\tilde{Z}_j).
\end{aligned} \tag{13}$$

Performing generalized Nyström in the intermediate stages yields our *Randomized low-rank RK method*:

$$\begin{aligned} Z_j &= Y_i + h \sum_{l=1}^{j-1} a_{jl} F(\mathcal{N}_l(Z_l)), \quad j = 1, \dots, s, \\ Y_{i+1} &= \mathcal{N}_{s+1} \left(Y_i + h \sum_{j=1}^s b_j F(\mathcal{N}_j(Z_j)) \right). \end{aligned} \tag{14}$$

Note that the index of \mathcal{N} is now used to emphasize the use of different (independent) random matrices Ω_j, Ψ_j for different stages. Across different time steps, the random matrices are, of course, also independently drawn. In (14), the rank of Z_j increases linearly with respect to j . However, there is no need to construct and store Z_j explicitly. For all subsequent purposes, we only need access to the Nyström approximation of Z_j , and thus it suffices to compute and store the sketches $Z_j \Omega_j$ and $\Psi_j^T Z_j$. In fact, the method (14) is equivalent to

$$\begin{cases} Z_j \Omega_j = Y_i \Omega_j + h \sum_{l=1}^{j-1} a_{jl} F(\mathcal{N}_l(Z_l)) \Omega_j \\ \Psi_j^T Z_j = \Psi_j^T Y_i + h \sum_{l=1}^{j-1} a_{jl} \Psi_j^T F(\mathcal{N}_l(Z_l)), \end{cases} \quad j = 1, \dots, s, \tag{15}$$

$$Y_{i+1} = \mathcal{N}_{s+1} \left(Y_i + h \sum_{j=1}^s b_j F(\mathcal{N}_j(Z_j)) \right).$$

We will use the expression (9) to evaluate $\mathcal{N}_j(Z_j)$ and, for this purposes, only needs $\Psi_j^T Z_j$, Ψ_j , and $Z_j \Omega_j$ (or, rather, the orthogonal factor of its QR decomposition).

Algorithm 1 contains the pseudo-code of the Randomized low-rank RK method. It closely follows (15), except that we also precompute the sketches of $F([\hat{Z}_j(\Psi_j^T \hat{Z}_j)^\dagger \tilde{Z}_j]_r)$ as they are needed in subsequent stages.

3.1 Implementation aspects and cost

We now consider the efficient implementation and cost of a time step performed by Algorithm 1. To simplify the discussion, we assume $m = n$ and $p = \ell = \mathcal{O}(r) \ll n$. We let c_n denote the cost of multiplying a vector of length n with Ω or Ψ . In the worst case, when Ω, Ψ are unstructured dense random matrices, $c_n = \mathcal{O}(nr)$. The use of structured random matrices can lead to lower c_n . For example, using the Subsampled Randomized Fourier Transform (SRFT) [12] for sketching reduces c_n to $\mathcal{O}(n \log(r))$.

Every (large) $n \times n$ matrix occurring in the algorithm is represented in factored form $U \Sigma V^T \in \mathbb{R}^{n \times n}$, where U, V are tall matrices (not necessarily orthonormal) and Σ is a small square matrix (not necessarily diagonal) of size equal to the rank of the matrix.

The evaluation of $[\hat{Z}_j(\Psi_j^T \hat{Z}_j)^\dagger \tilde{Z}_j]_r$ in Algorithm 1 requires $\mathcal{O}(nr^2 + jnr)$ operations for computing \hat{Z}_j, \tilde{Z}_j and another $\mathcal{O}(nr^2)$ operations for computing the factored form of the matrix.

The cost of applying F to a rank- r matrix and obtaining a factored form of the result strongly depends on the nature of F . We will denote this cost by c_F and the resulting rank by r_F . In some cases (see Section 4.1 for an example), $c_F = \mathcal{O}(nr)$ and $r_F = \mathcal{O}(r)$. In cases that lead to large r_F , the use of random sketches gives flexibility to exploit structure. For example, consider the case that $F(A)$ contains Hadamard products of the matrix A with itself, originating from, e.g., quadratic nonlinearities in the underlying partial differential equation. Then although the rank of the matrix F_j in Algorithm 1 is much larger than r , its factored representation has rich Kronecker product structure, which can be exploited when sketching F_j with Khatri-Rao

Algorithm 1 Randomized low-rank Runge-Kutta method with s stages

Input: Differential equation (1) defined by F and initial condition $A_0 \in \mathbb{R}^{m \times n}$. Target rank r , oversampling parameters p, ℓ , step size $h > 0$, number of steps $N \geq 0$.

Output: Approximation $Y_N \in \mathcal{M}_r$ of $A(Nh)$.

Draw ind. random matrices $\Omega \in \mathbb{R}^{n \times (r+p)}$, $\Psi \in \mathbb{R}^{m \times (r+p+\ell)}$.

$\hat{Y}_0 = A_0 \Omega$, $\tilde{Y}_0 = \Psi^T A_0$

for $i = 0, \dots, N - 1$ **do**

Draw ind. random matrices $\Omega_j \in \mathbb{R}^{n \times (r+p)}$, $\Psi_j \in \mathbb{R}^{m \times (r+p+\ell)}$ for $j = 1, 2, \dots, (s + 1)$.

for $j = 1, \dots, s$ **do**

$\hat{Z}_j = \llbracket \hat{Y}_i (\Psi^T \hat{Y}_i)^\dagger \tilde{Y}_i \rrbracket_r \Omega_j + h \sum_{l=1}^{j-1} a_{jl} \hat{K}_{lj}$ ▷ Evaluate $\llbracket \hat{Y}_i (\Psi^T \hat{Y}_i)^\dagger \tilde{Y}_i \rrbracket_r$ using (9)

$\tilde{Z}_j = \Psi_j^T \llbracket \hat{Y}_i (\Psi^T \hat{Y}_i)^\dagger \tilde{Y}_i \rrbracket_r + h \sum_{l=1}^{j-1} a_{jl} \tilde{K}_{lj}$

$F_j = F(\llbracket \hat{Z}_j (\Psi_j^T \hat{Z}_j)^\dagger \tilde{Z}_j \rrbracket_r)$ ▷ Evaluate $\llbracket \hat{Z}_j (\Psi_j^T \hat{Z}_j)^\dagger \tilde{Z}_j \rrbracket_r$ using (9)

for $q = j + 1, \dots, s + 1$ **do** ▷ Pre-compute sketches of F for other stages

$\hat{K}_{jq} = F_j \Omega_q$

$\tilde{K}_{jq} = \Psi_q^T F_j$

end for

end for

$\hat{Y}_{i+1} = Y_i \Omega_{s+1} + h \sum_{j=1}^s b_j \hat{K}_{j,s+1}$

$\tilde{Y}_{i+1} = \Psi_{s+1}^T Y_i + h \sum_{j=1}^s b_j \tilde{K}_{j,s+1}$

Set $\Psi = \Psi_{s+1}$

end for

Return $Y_N = \llbracket \hat{Y}_N (\Psi \hat{Y}_N)^\dagger \tilde{Y}_N \rrbracket_r$. ▷ Compute $\llbracket \hat{Y}_N (\Psi \hat{Y}_N)^\dagger \tilde{Y}_N \rrbracket_r$ using (9)

products of random matrices [4, 17]. As another example, when F has block structure, one can use the Block SRFT [3] to benefit from parallel computing. In Section 4.1.1, we test numerically the possibility to speed up the computation of \hat{K}_{jq} and \tilde{K}_{jq} by using the same random matrices. Although there is little justification, this variant appears to lead to an accuracy comparable to Algorithm 1. When not using any of the tricks mentioned above, the computation of each \hat{K}_{jq} and \tilde{K}_{jq} requires $\mathcal{O}(r_F c_n + r_F^2 r + nr_F r) = \mathcal{O}(r_F c_n + nr_F r)$ operations.

In summary, the complexity of the j th stage is $\mathcal{O}(jnr + nr^2 + c_F + (s + 1 - j)(r_F c_n + nr_F r))$, and thus one step of the Randomized RK method has a total complexity of

$$\mathcal{O}(s^2(r_F c_n + nr_F r) + snr^2 + sc_F). \quad (16)$$

It can be seen that the computation of \hat{K}_{jq} , \tilde{K}_{jq} is a dominating part of the cost.

The linearity of the sketches with respect to the data (sometimes also called a streaming property), makes the generalized Nyström method a preferred choice for the randomized low-rank approximation in Algorithm 1. It allows one to only store and work with sketches of F_j , $j = 1, \dots, s$ when computing subsequent stages. In contrast, the randomized SVD uses an orthogonal projection that is not linear in the data and, in turn, the full rank- r_F factorizations of the potentially high-rank matrices F_j need to be stored. Another disadvantages of the randomized SVD is that its direct application to the sums appearing in the Runge-Kutta stages is quite expensive, requiring $\mathcal{O}(ns^2 r_F^2)$ operations.

3.2 Comparison to Projected Runge-Kutta method

The Projected Runge-Kutta method (PRK) from [14] is closely related to our proposed method (14). It proceeds by performing the time stepping

$$\begin{aligned} Z_j &= Y_i + h \sum_{l=1}^{j-1} a_{jl} P_r(\mathcal{R}(Z_l)) F(\mathcal{R}(Z_l)), \quad j = 1, \dots, s, \\ Y_{i+1} &= \mathcal{R}\left(Y_i + h \sum_{j=1}^s b_j P_r(\mathcal{R}(Z_j)) F(\mathcal{R}(Z_j))\right). \end{aligned} \tag{17}$$

Here, \mathcal{R} denotes a retraction to the manifold \mathcal{M}_r and a common choice is the truncated SVD. The tangent space projection P_r in (17) reduces the rank of $F(\mathcal{R}(Z_l))$ to $2r$, which can make the subsequent application of \mathcal{R} significantly cheaper. Our method uses sketching instead of tangent space projection, which has two potential advantages: (1) As discussed in Section 2.2, the tangent space projection could introduce significant error, even when the solution admits a good rank- r approximation. In this situation, we expect our method (14) to be more accurate. This expectation is confirmed by the numerical experiments in Section 4. (2) As discussed in Section 3.1 above, sketching gives additional flexibility in the choice of random matrices to exploit structure in F applied to a rank- r matrix. Tangent space projection does not offer this flexibility.

One iteration of (17) needs to apply retraction to the matrices Z_j , which have rank at most $2jr$ for $j = 1 \dots s$. When using the truncated SVD to perform retraction, the cost is dominated by the QR factorizations of the $n \times 2jr$ factors of Z_j . In summary, the total complexity for one time step of PRK (17) is

$$\mathcal{O}\left(\underbrace{s^3 nr^2}_{\text{total retraction cost}} + \underbrace{snr_F r}_{s \times \text{application of } P_r} + \underbrace{sc_F}_{s \times \text{evaluation of } F} \right)$$

Compared to (16), we see that the term $s^2 nr_F r$ is reduced to $snr_F r$, which only becomes relevant when $r_F > sr$. This is because PRK can reuse computations related to tangent space projections across different stages, while the randomized low-rank RK method uses different sketches for every stage and can thus not reuse computations. As mentioned above, this issue can be mitigated by reusing random matrices for sketching, at the expense of theoretical justification. However, as s is typically very small (e.g., 1, 2 or 4), this issue may not be too relevant in practice.

3.3 Error analysis

With high probability, our method achieves at least the same qualitative error behavior that has been established for PRK. To see this, we follow [14] and consider the stage orders $\gamma_1, \dots, \gamma_s$, which are defined as the local errors of the stages \tilde{Z}_j in the standard RK method (13): For every $h \leq h_0$,

$$\|\tilde{Z}_j - \phi_F^{c_j h}(A_j)\|_F \leq C_L h^{\gamma_j+1}, \quad j = 1, \dots, s, \tag{18}$$

with $c_j = a_{j1} + a_{j2} + \dots + a_{j,j-1}$. We then obtain the following result, which corresponds to Theorem 6 in [14].

Theorem 5. *Consider the randomized low-rank RK method (14) utilizing independent standard Gaussian matrices, oversampling parameters $p, \ell \geq 4$, and an explicit s -stage RK method of order τ with stage orders $\gamma_1 \leq \gamma_2 \leq \dots \leq \gamma_s$. Denote*

$$\gamma = \begin{cases} \min(\tau, \gamma_2 + 1) & \text{if } b_2 \neq 0, \\ \min(\tau, \gamma_3 + 1, \gamma_2 + 2) & \text{if } b_2 = 0. \end{cases}$$

Then with the assumptions stated in Theorem 3, the global error is bounded for $q = \min\{p, \ell\}$ by

$$\|Y_N - A(Nh)\|_{L_q} \leq C(\delta + \epsilon + h^\gamma)$$

on the finite time interval $0 \leq Nh \leq T$, for all $h \leq h_0$. The constant C depends only on $L, T, h_0, s, C_L, C_N, \max_{ij} |a_{ij}|$ and $\max_i |b_i|$. In particular, for any $\eta \geq 1$, it holds for fixed h and N that

$$\Pr \{ \|Y_N - A(Nh)\|_F \geq C\eta(\epsilon + h^\gamma + \delta) \} \leq \frac{1}{\eta^q}.$$

Proof. The result of this theorem essentially follows from replacing the Frobenius norm in the proof of [14, Theorem 6] by the L_q norm and performing some additional minor modifications. For completeness, we have included the proof in the appendix. \square

Theorem 5 above shows that the randomized RK methods with the coefficients given by the following RK methods of order 1, 2 and 3, enjoy the usual convergence order up to $\mathcal{O}(\epsilon)$:

- RK 1 (Euler): $b_1 = 1$,
- RK 2 (Heun's method): $a_{21} = 1, b_1 = b_2 = \frac{1}{2}$,
- RK 3 (Heun's third-order method): $a_{21} = a_{32} = \frac{1}{2}, a_{43} = 1, b_1 = b_4 = \frac{1}{6}, b_2 = b_3 = \frac{1}{3}$.

Unfortunately, for RK 4 we only obtain order 2 from Theorem 5. Section 3.4 investigates this combination further.

Remark 6. As already noted, the global nature of the three assumptions in Section 2.1 can be quite limiting. If we modify these assumptions (4), (5) and (6) such that they only hold in a neighborhood of the exact solution $A(t)$ for $0 \leq t \leq T$, we need to additionally ensure that Y_i and the intermediate stages remain in the neighborhood. Due to the presence of randomness, this complicates the analysis and yields slightly worse results. In the following, we sketch which modifications need to be performed in order to localize the assumptions.

Suppose we want to ensure $E_i := \|Y_i - A(ih)\|_F \leq M$ for every $i = 1, \dots, N-1$ to use the properties (4), (5) and (6) locally. (We simplify the discussion by ignoring the probability that the intermediate stages leave the neighborhood; it is easy to adapt the approach outlined here by additionally requiring $\|Z_j - \tilde{Z}_j\|_F \leq M'$ for $j = 1, \dots, s$, which will yield the same convergence rate with a different constant and a somewhat higher failure probability, increased by a factor s .) To proceed, we can utilize the failure probability estimate of Theorem 5 to conclude

$$\Pr \{ E_i > M \mid \cap_{j=0}^{i-1} \{ E_j \leq M \} \} \leq \left(\frac{C(\epsilon + h^\gamma + \delta)}{M} \right)^q.$$

Using conditional probability and Bernoulli's inequality,

$$\begin{aligned} \Pr \{ \cap_{j=0}^{N-1} \{ E_j \leq M \} \} &= \Pr \{ E_{N-1} \leq M \mid \cap_{j=0}^{N-2} \{ E_j \leq M \} \} \Pr \{ \cap_{j=0}^{N-2} \{ E_j \leq M \} \} \\ &\geq \left(1 - \left(\frac{C(\epsilon + h^\gamma + \delta)}{M} \right)^q \right) \Pr \{ \cap_{j=0}^{N-2} \{ E_j \leq M \} \} \\ &\geq \left(1 - \left(\frac{C(\epsilon + h^\gamma + \delta)}{M} \right)^q \right)^{(N-1)} \\ &\geq 1 - (N-1) \left(\frac{C(\epsilon + h^\gamma + \delta)}{M} \right)^q. \end{aligned}$$

Similarly, we have

$$\Pr \{ \|Y_N - A(Nh)\|_F \geq C\eta(\epsilon + h^\gamma + \delta) \} \leq \frac{1}{\eta^q} + (N-1) \left(\frac{C(\epsilon + h^\gamma + \delta)}{M} \right)^q.$$

The additional second term is not large when M is large and/or $\epsilon + h^\gamma + \delta$ is small. To further quantify how this term affects the order of convergence, we substitute $\eta = \frac{M}{h^{1/q}}$ and assume $C(\epsilon + h^\gamma + \delta) \leq 1$, $h \leq 1$, leading to

$$\Pr \left\{ \|Y_N - A(Nh)\|_F \geq CM(\epsilon + h^\gamma + \delta)h^{-\frac{1}{q}} \right\} \leq \frac{h}{M^q} + (N-1) \left(\frac{C(\epsilon + h^\gamma + \delta)}{M} \right)^q \leq \frac{N}{M^q}. \quad (19)$$

The presence of the additional factor $h^{-\frac{1}{q}}$ reduces the convergence order by $\frac{1}{q}$. Because of $q = \min\{p, \ell\}$, even modest choices of the oversampling parameters p, ℓ imply that this potential order loss is negligible. \checkmark

3.4 Error analysis of randomized low-rank Runge-Kutta 4 method

Although our numerical experiments indicate convergence order 4 for the randomized RK method based on RK 4, it appears to be difficult to establish this order theoretically. In the following, we establish order 4 when the intermediate stages are oversampled. Let us emphasize that this oversampling is only performed for theoretical purposes; in practice, it does not seem to be needed.

Concretely, we plug the coefficients of the classical RK 4 method into (14) and oversample the intermediate stage as follows:

$$\begin{aligned} \hat{Z}_1 &= \mathcal{N}_1^{15r}(Y_i) \\ \hat{Z}_2 &= \mathcal{N}_2^{15r}\left(Y_i + \frac{h}{2}F(\hat{Z}_1)\right) \\ \hat{Z}_3 &= \mathcal{N}_3^{15r}\left(Y_i + \frac{h}{2}F(\hat{Z}_2)\right) \\ \hat{Z}_4 &= \mathcal{N}_4^{15r}\left(Y_i + hF(\hat{Z}_3)\right) \\ Y_{i+1} &= \mathcal{N}_5\left(Y_i + \frac{h}{6}\left(F(\hat{Z}_1) + 2F(\hat{Z}_2) + 2F(\hat{Z}_3) + F(\hat{Z}_4)\right)\right). \end{aligned} \quad (20)$$

Here, \mathcal{N}_5 refers to the usual generalized Nyström method (8) with target rank r , while \mathcal{N}_i^{15r} , for $i = 1, \dots, 4$, refers to the generalized Nyström method with the increased target rank $15r$. As the method leaves \mathcal{M}_r , we need to impose a stronger assumption on low-rank approximability. For $h \leq h_0$, we assume that

$$\|\Phi_F^h(Y) - \llbracket \Phi_F^h(Y) \rrbracket_k\|_F \leq C_M h \epsilon, \quad \forall Y \in \mathcal{M}_k, \quad \forall r \leq k \leq 15r. \quad (21)$$

We start our analysis of (20) with a result on the low-rank approximability of F implied by (21).

Lemma 7. *Assuming that F satisfies (21), let k be any integer such that $r \leq k \leq 15r$. Then*

$$\|F(Y) - \llbracket F(Y) \rrbracket_{2k}\|_F \leq C_M \epsilon, \quad \forall Y \in \mathcal{M}_k.$$

Proof. By the definition of the flow, $F(Y)$ is the time derivative of $\Phi_F^t(Y)$. Thus, for any $\gamma > 0$, there exists $h > 0$ such that

$$\left\| \frac{\Phi_F^h(Y) - \Phi_F^0(Y)}{h} - F(Y) \right\|_F \leq \gamma.$$

Because $\llbracket \Phi_F^{h_j}(Y) \rrbracket_k - \Phi_F^0(Y)$ has rank at most $2k$, it follows that

$$\begin{aligned} \|F(Y) - \llbracket F(Y) \rrbracket_{2k}\|_F &\leq \left\| F(Y) - \frac{\llbracket \Phi_F^{h_j}(Y) \rrbracket_k - \Phi_F^0(Y)}{h_j} \right\|_F \\ &\leq \left\| F(Y) - \frac{\Phi_F^{h_j}(Y) - \Phi_F^0(Y)}{h_j} \right\|_F + \left\| \frac{\Phi_F^{h_j}(Y) - \llbracket \Phi_F^{h_j}(Y) \rrbracket_k}{h_j} \right\|_F \\ &\leq \gamma + C_M \epsilon. \end{aligned}$$

The result of the lemma is obtained by taking $\gamma \rightarrow 0$. \square

The following auxiliary result helps to bound the local error of (20).

Lemma 8. *Let $Z \in \mathcal{M}_r, X \in \mathbb{R}^{m \times n}, \alpha \geq 0$, and $k \leq 14r$. Then*

$$\|Z + \alpha X - \llbracket Z + \alpha X \rrbracket_{15r}\|_F \leq \alpha \|X - \llbracket B \rrbracket_k\|_F$$

holds for any $B \in \mathbb{R}^{m \times n}$.

Proof. Using that $Z + \alpha \llbracket B \rrbracket_k$ has rank at most $15r$, the result follows from

$$\|Z + \alpha X - \llbracket Z + \alpha X \rrbracket_{15r}\|_F \leq \|Z + \alpha X - \llbracket Z + \alpha \llbracket B \rrbracket_k \rrbracket_{15r}\|_F = \alpha \|X - \llbracket B \rrbracket_k\|_F.$$

\square

We are now in the position to establish a local error estimate for (20).

Lemma 9. *Suppose that the assumptions stated in Section 2.1 and (21) hold. Given $Y_i \in \mathcal{M}_r$, one step of the method (20) with independent standard Gaussian matrices and oversampling parameters $p, \ell \geq 4$ satisfies the local error estimate*

$$\|\Phi_F^h(Y_i) - Y_{i+1}\|_{L_q} \leq C(h\epsilon + h^5),$$

for $q = \min\{p, \ell\}$ and all $0 < h \leq h_0$. The constant C depends only on L, T, h_0, C_M , and C_N .

Proof. The proof proceeds by bounding the moments of the differences between the stages \hat{Z}_j of (20) and the stages \tilde{Z}_j of the classic RK4 applied to Y_i , as defined in (13). For $j = 1$, $Y_i \in \mathcal{M}_r$ implies

$$\|\hat{Z}_1 - \tilde{Z}_1\|_{L_q} = \|\mathcal{N}_1^{15r}(Y_i) - Y_i\|_{L_q} = 0.$$

For $j = 2$, we set $Y_{i,1} := Y_i + \frac{h}{2}F(\hat{Z}_1)$

$$\begin{aligned} \|\hat{Z}_2 - \tilde{Z}_2\|_{L_q} &= \left\| \mathcal{N}_2^{15r}(Y_{i,1}) - Y_i - \frac{h}{2}F(\tilde{Z}_1) \right\|_{L_q} \\ &\leq \left\| \mathcal{N}_2^{15r}(Y_{i,1}) - \mathcal{N}_2^{15r}(Y_i + \frac{h}{2}\llbracket F(\tilde{Z}_1) \rrbracket_{2r}) \right\|_{L_q} \\ &\quad + \left\| \mathcal{N}_2^{15r}(Y_i + \frac{h}{2}\llbracket F(\tilde{Z}_1) \rrbracket_{2r}) - Y_i - \frac{h}{2}F(\tilde{Z}_1) \right\|_{L_q} \\ &= \left\| \mathcal{N}_2^{15r}(Y_{i,1}) - \mathcal{N}_2^{15r}(Y_i + \frac{h}{2}\llbracket F(\tilde{Z}_1) \rrbracket_{2r}) \right\|_{L_q} + \frac{h}{2} \|\llbracket F(\tilde{Z}_1) \rrbracket_{2r} - F(\tilde{Z}_1)\|_{L_q} \\ &\leq \left\| \mathcal{N}_2^{15r}(Y_{i,1}) - \mathcal{N}_2^{15r}(Y_i + \frac{h}{2}\llbracket F(\tilde{Z}_1) \rrbracket_{2r}) \right\|_{L_q} + \frac{h}{2}C_M \epsilon \end{aligned} \tag{22}$$

The second equality above holds because $Y_i + \frac{h}{2}\llbracket F(\tilde{Z}_1)\rrbracket_{2r}$ has rank at most $3r$, and therefore $Y_i + \frac{h}{2}\llbracket F(\tilde{Z}_1)\rrbracket_{2r} = \mathcal{N}_2^{15r}(Y_i + \frac{h}{2}\llbracket F(\tilde{Z}_1)\rrbracket_{2r})$ holds almost surely. The last inequality follows from Lemma 7. With similar reasoning, we obtain the following bound for the first term in (22):

$$\begin{aligned} \left\| \mathcal{N}_2^{15r}(Y_{i,1}) - \mathcal{N}_2^{15r}\left(Y_i + \frac{h}{2}\llbracket F(\tilde{Z}_1)\rrbracket_{2r}\right) \right\|_{L_q} &= \left\| \mathcal{N}_2^{15r}(Y_{i,1}) - \left(Y_i + \frac{h}{2}\llbracket F(\tilde{Z}_1)\rrbracket_{2r}\right) \right\|_{L_q} \\ &\leq \left\| \mathcal{N}_2^{15r}(Y_{i,1}) - \left(Y_i + \frac{h}{2}F(\tilde{Z}_1)\right) \right\|_{L_q} + \frac{h}{2}C_M\epsilon. \end{aligned}$$

Using that $\hat{Z}_1 = \tilde{Z}_1 = Y_i$ holds almost surely, the law of total expectation, and Theorem 2 give

$$\begin{aligned} \left\| \mathcal{N}_2^{15r}(Y_{i,1}) - \left(Y_i + \frac{h}{2}F(\tilde{Z}_1)\right) \right\|_{L_q} &= \left(\mathbb{E} \left[\mathbb{E} \left(\left\| \mathcal{N}_2^{15r}(Y_{i,1}) - \left(Y_i + \frac{h}{2}F(\tilde{Z}_1)\right) \right\|_{L_q}^q \middle| \hat{Z}_1 \right) \right] \right)^{\frac{1}{q}} \\ &\leq C_{\mathcal{N}} \left\| Y_{i,1} - \llbracket Y_{i,1} \rrbracket_{15r} \right\|_{L_q} \\ &\leq C_{\mathcal{N}} \left\| Y_i + \frac{h}{2}F(Y_i) - \llbracket Y_i + \frac{h}{2}F(Y_i) \rrbracket_{3r} \right\|_{L_q} \\ &\leq C_{\mathcal{N}} \left\| Y_i + \frac{h}{2}F(Y_i) - \left(Y_i + \frac{h}{2}\llbracket F(Y_i)\rrbracket_{2r}\right) \right\|_{L_q} \leq \frac{C_{\mathcal{N}}h}{2}C_M\epsilon. \end{aligned}$$

In summary, we have proven

$$\|\hat{Z}_2 - \tilde{Z}_2\|_{L_q} \leq \frac{(C_{\mathcal{N}} + 2)C_M}{2}h\epsilon. \quad (23)$$

Note that we also have the following inequality by Lemma 7:

$$\|\tilde{Z}_2 - \llbracket \tilde{Z}_2 \rrbracket_{3r}\|_{L_q} \leq \left\| Y_i + \frac{h}{2}F(\tilde{Z}_1) - Y_i - \frac{h}{2}\llbracket F(\tilde{Z}_1)\rrbracket_{2r} \right\|_{L_q} \leq \frac{h}{2}C_M\epsilon. \quad (24)$$

For $j = 3$, we set $Y_{i,2} := Y_i + \frac{h}{2}F(\hat{Z}_2)$ and apply an analogous reasoning:

$$\begin{aligned} \|\hat{Z}_3 - \tilde{Z}_3\|_{L_q} &= \left\| \mathcal{N}_3^{15r}(Y_{i,2}) - Y_i - \frac{h}{2}F(\tilde{Z}_2) \right\|_{L_q} \\ &\leq \left\| \mathcal{N}_3^{15r}(Y_{i,2}) - \mathcal{N}_3^{15r}\left(Y_i + \frac{h}{2}\llbracket F(\llbracket \tilde{Z}_2 \rrbracket_{3r}) \rrbracket_{6r}\right) \right\|_{L_q} \\ &\quad + \left\| \mathcal{N}_3^{15r}\left(Y_i + \frac{h}{2}\llbracket F(\llbracket \tilde{Z}_2 \rrbracket_{3r}) \rrbracket_{6r}\right) - Y_i - \frac{h}{2}F(\tilde{Z}_2) \right\|_{L_q} \\ &= \left\| \mathcal{N}_3^{15r}(Y_{i,2}) - \mathcal{N}_3^{15r}\left(Y_i + \frac{h}{2}\llbracket F(\llbracket \tilde{Z}_2 \rrbracket_{3r}) \rrbracket_{6r}\right) \right\|_{L_q} + \frac{h}{2}\|F(\tilde{Z}_2) - \llbracket F(\llbracket \tilde{Z}_2 \rrbracket_{3r}) \rrbracket_{6r}\|_{L_q} \\ &= \left\| \mathcal{N}_3^{15r}(Y_{i,2}) - \left(Y_i + \frac{h}{2}\llbracket F(\llbracket \tilde{Z}_2 \rrbracket_{3r}) \rrbracket_{6r}\right) \right\|_{L_q} + \frac{h}{2}\|F(\tilde{Z}_2) - \llbracket F(\llbracket \tilde{Z}_2 \rrbracket_{3r}) \rrbracket_{6r}\|_{L_q} \\ &\leq \left\| \mathcal{N}_3^{15r}(Y_{i,2}) - (Y_{i,2}) \right\|_{L_q} + \frac{h}{2}\|F(\hat{Z}_2) - \llbracket F(\llbracket \tilde{Z}_2 \rrbracket_{3r}) \rrbracket_{6r}\|_{L_q} \\ &\quad + \frac{h}{2}\|F(\tilde{Z}_2) - \llbracket F(\llbracket \tilde{Z}_2 \rrbracket_{3r}) \rrbracket_{6r}\|_{L_q}. \end{aligned}$$

The first term of the last inequality can be bounded using Lemma 8:

$$\left\| \mathcal{N}_3^{15r}(Y_{i,2}) - (Y_{i,2}) \right\|_{L_q} \leq C_{\mathcal{N}} \left\| Y_{i,2} - \llbracket Y_{i,2} \rrbracket_{15r} \right\|_{L_q} \leq \frac{h}{2}C_{\mathcal{N}}\|F(\hat{Z}_2) - \llbracket F(\llbracket \tilde{Z}_2 \rrbracket_{3r}) \rrbracket_{6r}\|_{L_q}.$$

Therefore, by (23) and (24),

$$\begin{aligned}
\|\hat{Z}_3 - \tilde{Z}_3\|_{L_q} &\leq \frac{(C_N + 1)h}{2} \|F(\hat{Z}_2) - \llbracket F(\llbracket \tilde{Z}_2 \rrbracket_{3r}) \rrbracket_{6r}\|_{L_q} + \frac{h}{2} \|F(\tilde{Z}_2) - \llbracket F(\llbracket \tilde{Z}_2 \rrbracket_{3r}) \rrbracket_{6r}\|_{L_q} \\
&\leq \frac{(C_N + 2)h}{2} (\|F(\hat{Z}_2) - F(\tilde{Z}_2)\|_{L_q} + \|F(\tilde{Z}_2) - \llbracket F(\llbracket \tilde{Z}_2 \rrbracket_{3r}) \rrbracket_{6r}\|_{L_q} \\
&\quad + \|F(\llbracket \tilde{Z}_2 \rrbracket_{3r}) - \llbracket F(\llbracket \tilde{Z}_2 \rrbracket_{3r}) \rrbracket_{6r}\|_{L_q}) \\
&\leq \frac{h}{2} (C_N + 2) (L \frac{(C_N + 2)C_M}{2} h\epsilon + L \frac{h}{2} C_M \epsilon + C_M \epsilon) \leq C_3 h \epsilon.
\end{aligned}$$

Also, as for $j = 2$,

$$\|\tilde{Z}_3 - \llbracket \tilde{Z}_3 \rrbracket_{7r}\|_{L_q} \leq \left\| Y_i + \frac{h}{2} F(\tilde{Z}_2) - Y_i - \frac{h}{2} \llbracket F(\llbracket \tilde{Z}_2 \rrbracket_{3r}) \rrbracket_{6r} \right\|_{L_q} \leq \frac{h}{2} [C_M \epsilon + L \frac{h}{2} C_M \epsilon].$$

Finally, for $j = 4$, we set $Y_{i,3} := Y_i + hF(\hat{Z}_3)$ and obtain

$$\begin{aligned}
\|\hat{Z}_4 - \tilde{Z}_4\|_{L_q} &= \left\| \mathcal{N}_4^{15r}(Y_{i,3}) - \hat{Y}_i - hF(\tilde{Z}_3) \right\|_{L_q} \\
&\leq \left\| \mathcal{N}_4^{15r}(Y_{i,3}) - \mathcal{N}_4^{15r}(Y_i + h\llbracket F(\llbracket \tilde{Z}_3 \rrbracket_{7r}) \rrbracket_{14r}) \right\|_{L_q} \\
&\quad + \left\| \mathcal{N}_4^{15r}(Y_i + h\llbracket F(\llbracket \tilde{Z}_3 \rrbracket_{7r}) \rrbracket_{14r}) - Y_i - hF(\tilde{Z}_3) \right\|_{L_q} \\
&\leq \left\| \mathcal{N}_4^{15r}(Y_{i,3}) - (Y_i + h\llbracket F(\llbracket \tilde{Z}_3 \rrbracket_{7r}) \rrbracket_{14r}) \right\|_{L_q} + h \left\| F(\tilde{Z}_3) - \llbracket F(\llbracket \tilde{Z}_3 \rrbracket_{7r}) \rrbracket_{14r} \right\|_{L_q} \\
&\leq \left\| \mathcal{N}_4^{15r}(Y_{i,3}) - (Y_{i,3}) \right\|_{L_q} + h \left\| F(\hat{Z}_3) - \llbracket F(\llbracket \tilde{Z}_3 \rrbracket_{7r}) \rrbracket_{14r} \right\|_{L_q} \\
&\quad + h \left\| F(\tilde{Z}_3) - \llbracket F(\llbracket \tilde{Z}_3 \rrbracket_{7r}) \rrbracket_{14r} \right\|_{L_q}.
\end{aligned}$$

The first term of the last inequality can once again be bounded using Lemma 8:

$$\left\| \mathcal{N}_4^{15r}(Y_{i,3}) - (Y_{i,3}) \right\|_{L_q} \leq C_N \|Y_{i,3} - \llbracket Y_{i,3} \rrbracket_{15r}\|_{L_q} \leq hC_N \|F(\hat{Z}_3) - \llbracket F(\llbracket \tilde{Z}_3 \rrbracket_{7r}) \rrbracket_{14r}\|_{L_q}.$$

Hence,

$$\begin{aligned}
\|\hat{Z}_4 - \tilde{Z}_4\|_{L_q} &\leq h(C_N + 1) \|F(\hat{Z}_3) - \llbracket F(\llbracket \tilde{Z}_3 \rrbracket_{7r}) \rrbracket_{14r}\|_{L_q} + h \|F(\tilde{Z}_3) - \llbracket F(\llbracket \tilde{Z}_3 \rrbracket_{7r}) \rrbracket_{14r}\|_F \\
&\leq h(C_N + 2) (\|F(\hat{Z}_3) - F(\tilde{Z}_3)\|_{L_q} + \|F(\tilde{Z}_3) - \llbracket F(\llbracket \tilde{Z}_3 \rrbracket_{7r}) \rrbracket_{14r}\|_F \\
&\quad + \|F(\llbracket \tilde{Z}_3 \rrbracket_{7r}) - \llbracket F(\llbracket \tilde{Z}_3 \rrbracket_{7r}) \rrbracket_{14r}\|_F) \\
&\leq h(C_N + 2) (LC_3 h \epsilon + L \frac{h}{2} [C_M \epsilon + L \frac{h}{2} C_M \epsilon] + C_M \epsilon) \leq C_4 h \epsilon.
\end{aligned}$$

Collecting the obtained bounds for the stages and using Lipschitz continuity, we have

$$\begin{aligned}
&\left\| Y_i + \frac{h}{6} (F(\hat{Z}_1) + 2F(\hat{Z}_2) + 2F(\hat{Z}_3) + F(\hat{Z}_4)) \right. \\
&\quad \left. - Y_i - \frac{h}{6} (F(\tilde{Z}_1) + 2F(\tilde{Z}_2) + 2F(\tilde{Z}_3) + F(\tilde{Z}_4)) \right\|_{L_q} \leq C_5 h^2 \epsilon. \tag{25}
\end{aligned}$$

Recall that one step of the classic RK 4 method has error $\mathcal{O}(h^5)$. With the true solution $\Phi_F^h(Y_i)$ satisfying the low-rank approximability assumption (21), we find that one step of RK 4 satisfies

$$\begin{aligned}
&\left\| Y_i + \frac{h}{6} (F(\tilde{Z}_1) + 2F(\tilde{Z}_2) + 2F(\tilde{Z}_3) + F(\tilde{Z}_4)) \right. \\
&\quad \left. - \llbracket Y_i + \frac{h}{6} (F(\tilde{Z}_1) + 2F(\tilde{Z}_2) + 2F(\tilde{Z}_3) + F(\tilde{Z}_4)) \rrbracket_r \right\|_F \leq C_6 (h^5 + h\epsilon). \tag{26}
\end{aligned}$$

Using Theorem 2 and the inequalities (25), (26), we obtain

$$\begin{aligned}
& \left\| Y_{i+1} - Y_i - \frac{h}{6}(F(\hat{Z}_1) + 2F(\hat{Z}_2) + 2F(\hat{Z}_3) + F(\hat{Z}_4)) \right\|_{L_q} \\
& \leq C_{\mathcal{N}} \left\| Y_i + \frac{h}{6}(F(\hat{Z}_1) + 2F(\hat{Z}_2) + 2F(\hat{Z}_3) + F(\hat{Z}_4)) \right. \\
& \quad \left. - \mathbb{E} \left[Y_i + \frac{h}{6}(F(\hat{Z}_1) + 2F(\hat{Z}_2) + 2F(\hat{Z}_3) + F(\hat{Z}_4)) \right]_r \right\|_{L_q} \\
& \leq C_{\mathcal{N}} \left\| Y_i + \frac{h}{6}(F(\hat{Z}_1) + 2F(\hat{Z}_2) + 2F(\hat{Z}_3) + F(\hat{Z}_4)) \right. \\
& \quad \left. - \mathbb{E} \left[Y_i + \frac{h}{6}(F(\tilde{Z}_1) + 2F(\tilde{Z}_2) + 2F(\tilde{Z}_3) + F(\tilde{Z}_4)) \right]_r \right\|_{L_q} \\
& \leq C_{\mathcal{N}}(C_5 h^2 \epsilon + C_6(h^5 + h\epsilon)).
\end{aligned}$$

The result of the lemma is concluded from the following inequality:

$$\begin{aligned}
\|\Phi_F^h(Y_i) - Y_{i+1}\|_{L_q} & \leq \|\Phi_F^h(Y_i) - Y_i - \frac{h}{6}(F(\tilde{Z}_1) + 2F(\tilde{Z}_2) + 2F(\tilde{Z}_3) + F(\tilde{Z}_4))\|_{L_q} \\
& \quad + C_5 h^2 \epsilon + C_{\mathcal{N}}(C_5 h^2 \epsilon + C_6(h^5 + h\epsilon)) \leq \mathcal{O}(h^5 + h\epsilon).
\end{aligned}$$

□

Finally, the following theorem establishes order 4 with respect to h of the modified randomized low-rank RK 4 method (20). This provides some theoretical explanation for the convergence order 4 we observe for the randomized low-rank RK 4 method (without intermediate rank increases).

Theorem 10. *Suppose that the assumptions stated in Section 2.1 and (21) hold. Under the assumption (21) and the assumptions stated in section 2.1. The the global error of the scheme (20) with independent standard Gaussian matrices and oversampling parameters $p, \ell \geq 4$ satisfies for $q = \min\{p, \ell\}$ the bound*

$$\|Y_N - A(Nh)\|_{L_q} \leq C(\epsilon + h^4 + \delta),$$

on the finite time-interval $0 \leq nh \leq T$ and for every $0 < h \leq h_0$. The constant C depends only on L, T, h_0, C_M and $C_{\mathcal{N}}$. In particular, for any $\eta \geq 1$, it holds for fixed h and N that

$$\Pr\{\|Y_N - A(Nh)\|_F \geq C\eta(\epsilon + h^4 + \delta)\} \leq \frac{1}{\eta^q}.$$

Proof. By Lemma 9 we have

$$\left(\mathbb{E}(\mathbb{E}[\|\Phi_F^h(Y_i) - Y_{i+1}\|_F^q | Y_i]) \right)^{1/q} \leq \left(\mathbb{E}[C^q(h^5 + h\epsilon)^q] \right)^{1/q} = C(h^5 + h\epsilon).$$

Substituting this bound into (12) proves the first statement of the theorem. The second statement is proved by Markov's inequality, as in the proof of Corollary 4. □

4 Numerical Experiments

In the following numerical experiments, we verify the accuracy of randomized low-rank RK methods, Algorithm 1, which will be denoted as Rand RK. In particular, we consider Rand RK1 (Euler), Rand RK2 and Rand RK4, based upon the Euler method, Heun's method and the classical RK 4 method, respectively. For convenience, we recall the corresponding coefficients:

- Rand RK1 (Euler): $b_1 = 1$

- Rand RK2: $a_{21} = 1, b_1 = b_2 = \frac{1}{2}$,
- Rand RK4: $a_{21} = a_{32} = \frac{1}{2}, a_{43} = 1, b_1 = b_4 = \frac{1}{6}, b_2 = b_3 = \frac{1}{3}$.

We have implemented the generalized Nystöm method following [23]; we always use Gaussian random matrices and the oversampling parameters $p = \ell = \max\{2, 0.1r\}$ for sketching. Although the generalized Nyström method is usually numerically stable [20], it may exhibit instabilities, especially when $\Psi^T Z \Omega$ is numerically rank deficient. We have never observed such an instability in our numerical experiments, but let us point out that a numerically safer variant with regularization is described in [20]. Note that our implementation of Rand RK4 does *not* use oversampling in the intermediate stage, that is, we have implemented Algorithm 1 with the coefficients of RK 4 instead of (20).

The error is measured by the Frobenius norm between the reference solution and the approximation. The reference solution is obtained by solving the full matrix differential equation with MATLAB ODE45 using the tolerances `{'RelTol', 1e-10, 'AbsTol', 1e-10}`. Because our methods involve randomization, we report the mean approximation error as well as the spread between the largest and smallest errors for 10 independent random trials (indicated by lower/upper horizontal lines in the graph).

In the first two experiments, we also report the error of our implementation of the projected RK s method [14] for $s = 1, 2, 4$ as well as the projector splitting integrator [19] with the sub-steps computed by MATLAB's ODE45 using the same tolerance parameters as above. The initial value for these methods is obtained by applying the truncated SVD to the initial matrix A_0 .

All experiments have been performed in Matlab (version 2023a) on a Macbook Pro with an Apple M1 Pro processor. The code used to perform the experiments and produce the figures can be found at https://github.com/hysanlam/rand_RK.

4.1 Lyapunov matrix differential equation

As a first simple experiment, we approximate the solution of a Lyapunov matrix differential equation [25, Section 6.1], which takes the form

$$\dot{A}(t) = LA(t) + A(t)L + \alpha \frac{C}{\|C\|_F}, \quad A(0) = A_0,$$

where $A(t) \in \mathbb{R}^{n \times n}$, $L \in \mathbb{R}^{n \times n}$, $C \in \mathbb{R}^{n \times n}$, $\alpha \geq 0$ and $t \in [0, T]$. We set $n = 128$ and use the symmetric matrix $L = \text{diag}(1, -2, 1) \in \mathbb{R}^{128 \times 128}$. For setting the entries of the source term C and initial matrix A_0 , we follow a construction similar to the one used in [6, section 5.1]:

$$C_{ij} = \sum_{k=1}^{11} 10^{-(k-1)} \cdot e^{-k(x_i^2 + y_j^2)},$$

$$(A_0)_{ij} = \sum_{k=1}^{20} b_k \cdot \sin(kx_i) \sin(ky_j), \quad \text{with } b_k = \begin{cases} 1, & \text{if } k = 1 \\ 5e^{-(7+0.5(k-2))}, & \text{if } k > 1, \end{cases}$$

where (x_i, y_j) , for $i, j = 1, \dots, 128$, are uniform discretization points of the square $[-\pi, \pi] \times [-\pi, \pi]$. The final time is set to $T = 1$ and first consider $\alpha = 1$. Figure 1 displays the singular values of the reference solution at $t = T$, and the error of the approximation obtained by Rand Euler and Rand RK4 with different ranks. We can see that Rand Euler achieves first-order convergence in time, while Rand RK4 achieves fourth-order convergence in time until it reaches the level of low rank approximation error. Moreover, both methods demonstrate robust behavior despite randomness; among these 10 trials, the maximum error is only at most three times the empirical mean.

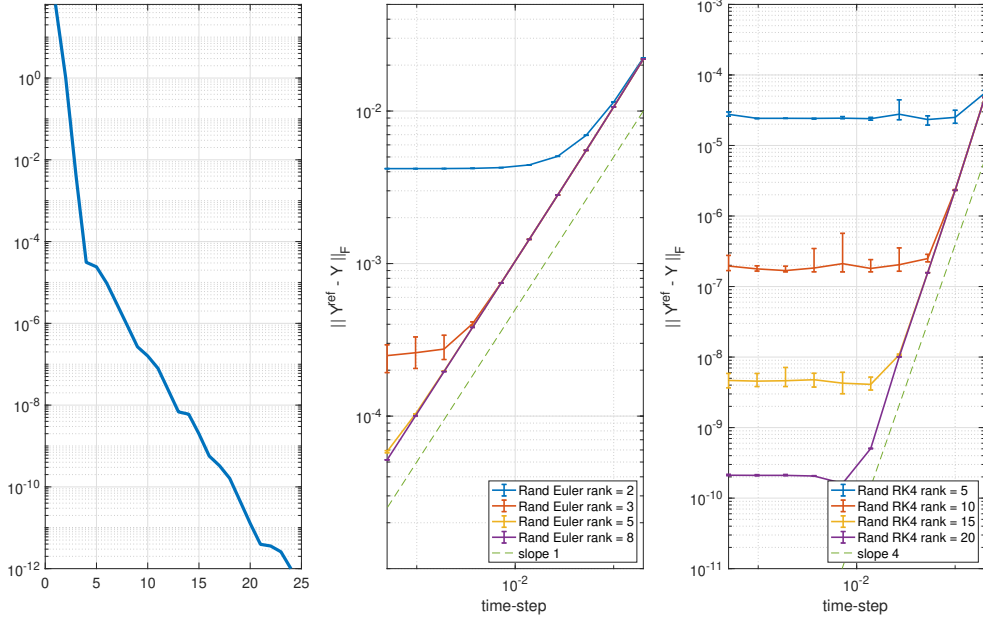


Figure 1: Lyapunov matrix differential equation with $\alpha = 1$. The singular values of the reference solution at time $T = 1$ together with the approximation errors of the numerical approximation obtained via the Rand Euler and Rand RK4 for different ranks and time-step sizes.

We now fix the rank to $r = 10$ and compare the accuracy of Rand Euler, Rand RK2, and Rand RK4 with other methods, including PRK1, PRK2, PRK4, and projector splitting for $\alpha = 10^{-5}$ as well as $\alpha = 1$. In Table 1, we report the average and maximum values of the tangential projection error $\|F(Y_i) - P_r(Y_i)F(Y_i)\|_F$, where Y_i is the approximation computed at the i th time step by PRK 2 with $h = 5 \times 10^{-3}$. For $\alpha = 1$ this error is much larger than $\alpha = 10^{-5}$. However, when the tangential projection error is large, the accuracy and convergence behavior of PRK and projector splitting are not guaranteed. This is indeed observed in Figure 2. For $\alpha = 10^{-5}$, PRK 1, PRK 2, and projector splitting exhibit the expected order of convergence. PRK 4 seems to even exhibit fourth-order convergence initially but this quickly deteriorates for smaller h . When $\alpha = 1$, all these methods only show first-order convergence. On the other hand, the randomized methods remain robust: Rand Euler, Rand RK2, and Rand RK4 exhibit first, second, and fourth-order convergence, respectively, for both choices of α .

α	10^{-5}	1
average $\ F(Y_i) - P_r(Y_i)F(Y_i)\ _F$	1.3553×10^{-7}	4.0144×10^{-4}
max $\ F(Y_i) - P_r(Y_i)F(Y_i)\ _F$	1×10^{-5}	0.7932

Table 1: Lyapunov matrix differential equation with $\alpha = 10^{-5}$ and $\alpha = 1$. Average and maximum $\|F(Y_i) - P_r(Y_i)F(Y_i)\|_F$ for the approximation Y_i at the i th time step of PRK 2 with $h = 5 \times 10^{-3}$.

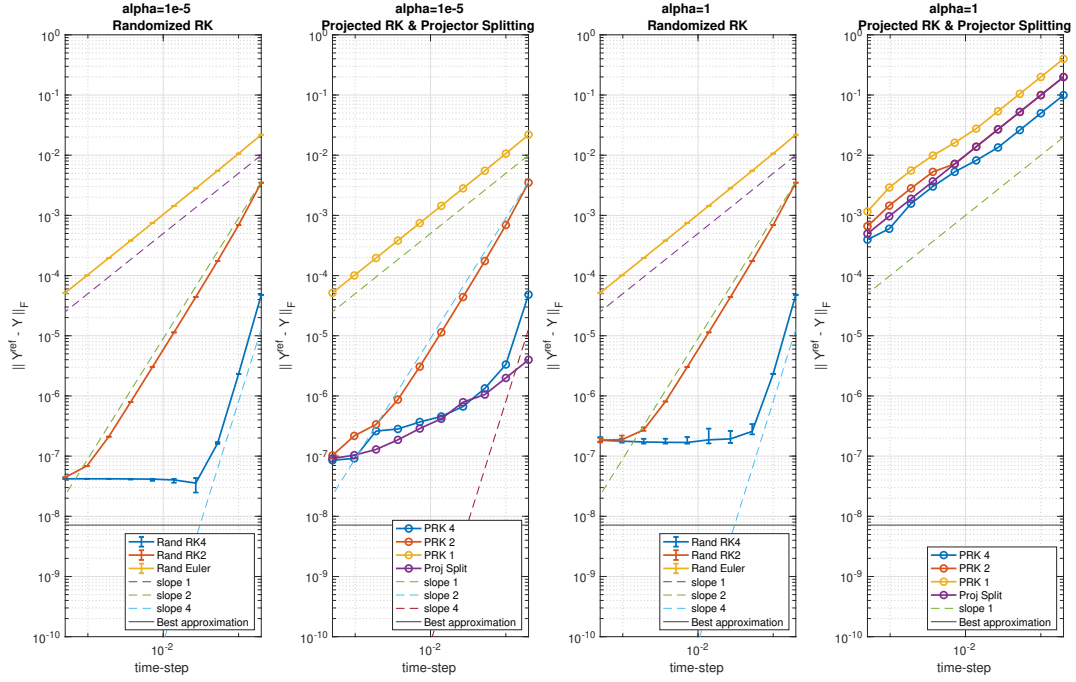


Figure 2: Lyapunov matrix differential equation with $\alpha = 10^{-5}$ and $\alpha = 1$. Comparison of absolute approximation errors for different low-rank integrators using rank $r = 10$.

4.1.1 Speeding up by constant random matrices

As indicated in Section 3.2, the cost of computing the sketches $\hat{K}_{jq}, \tilde{K}_{jq}$ in Algorithm 1 can be divided by nearly a factor s when choosing the same random matrices across the stages in one time step. That is, we draw two independent Gaussian matrices Ω_1 and Ψ_1 , and set $\Omega_1 = \Omega_2 = \dots = \Omega_{s+1}$, $\Psi_1 = \Psi_2 = \dots = \Psi_{s+1}$. In this experiment, we ran 10 random trials to solve the differential Lyapunov equation for $\alpha = 1$ and compare the described constant choice with the standard (independent) choice of random matrices in Algorithm 1. We tested Rand RK2 and Rand RK4 with different ranks and plotted the obtained errors in Figure 3. For this example, we observe comparable performance when using either different or the same random matrices across the stages in a time step. Although this modification is computationally attractive and likely the preferred way to run Algorithm 1, we are unable to establish an error bound due to the lack of stochastic independence of the stages.

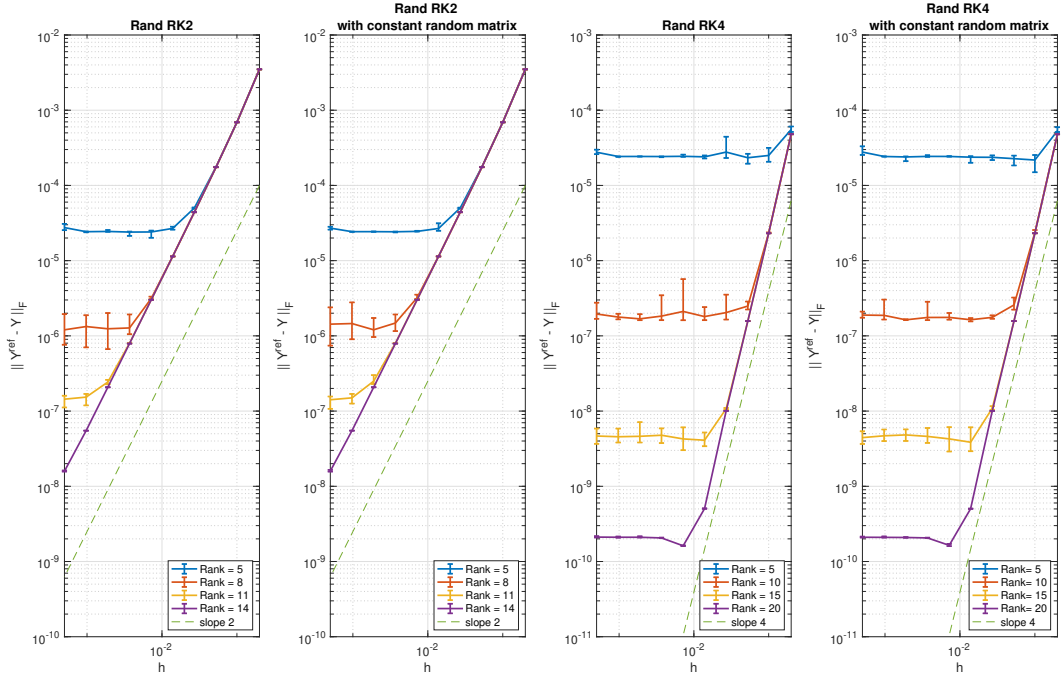


Figure 3: Lyapunov matrix differential equation with $\alpha = 1$. Comparison of absolute approximation errors for Rand RK2, Rand RK2 with different (independent) random matrices across stages vs. Rand RK2 and Rand RK4 with the same random matrices across stages.

4.2 Non-linear Schrödinger equation

We now consider the non-linear Schrödinger equation from [14, Section 5.3], where $A : [0, T] \rightarrow \mathbb{C}^{n \times n}$ evolves according to

$$\dot{A}(t) = i\left[\frac{1}{2}(BA + AB) + \alpha|A|^2A\right]. \quad (27)$$

The cubic nonlinearity $|A|^2A$ is taken element-wise and $B = \text{diag}(1, 0, 1)$. We choose $n = 100$, $T = 5$ and the initial data

$$(A_0)_{ij} = \exp\left(-\frac{(i-60)^2}{100} - \frac{(j-50)^2}{100}\right) + \exp\left(-\frac{(i-50)^2}{100} - \frac{(j-40)^2}{100}\right).$$

In this example, we aim at computing approximations of rank up to 30. To ensure that A_0 has rank at least 30 (making sure it satisfies the assumptions of PRK), we perturb A_0 by taking the full SVD of A_0 and setting the singular values 3, 4, \dots , 32 to 10^{-9} .

First, we set $\alpha = 0.3$. The singular values of the reference solution at $t = T$, as well as the approximation errors by Rand Euler and Rand RK4 with different ranks are shown in Figure 4. We again observe that Rand Euler achieves first-order convergence in time, while Rand RK4 achieves fourth-order convergence in time until it reaches the level of low-rank approximation error. In this example, the maximum error of both methods is also very close to the empirical mean, deviating by less than twice the mean.

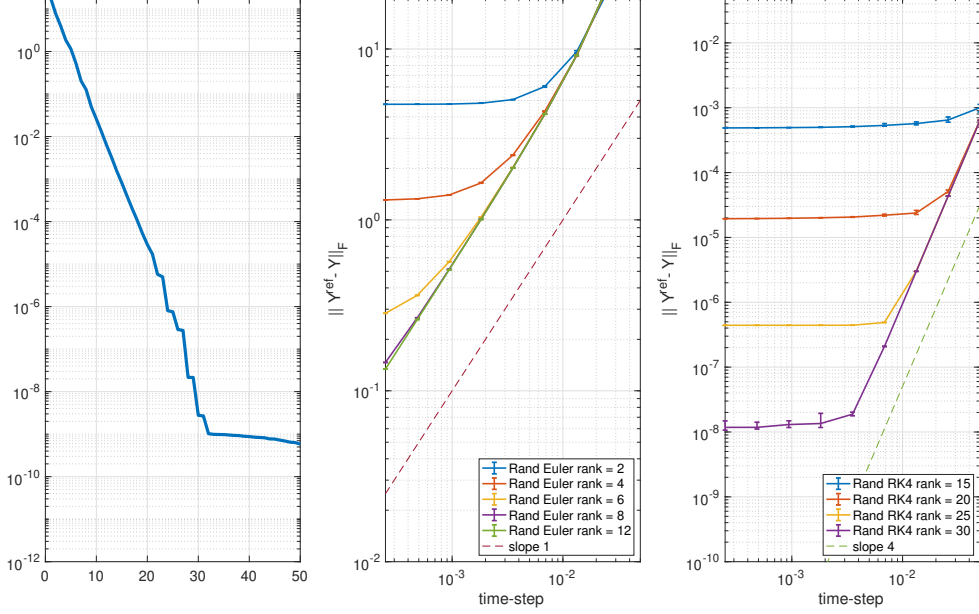


Figure 4: Non-linear Schrödinger equation with $\alpha = 0.3$. Singular values of the reference solution at time $T = 5$ together with the approximation errors of the numerical approximation obtained by Rand Euler and Rand RK4 for different ranks and time-step sizes.

Now we fix the rank to $r = 30$ and compare Rand RK with PRK and projector splitting for $\alpha = 3 \times 10^{-4}$ and $\alpha = 0.3$ in Figure 5. For small $\alpha = 3 \times 10^{-4}$, we observe that PRK 1 and PRK 2 exhibit the correct order of convergence. Again, one observes that the initially visible fourth-order convergence of PRK 4 quickly deteriorates as h decreases. For $\alpha = 0.3$, we see that the order of PRK 2 decreases to 1 when h is small, and PRK 4 only shows first-order convergence as well. Again, this is likely due to the large tangential projection error $\|P_r(Y_i)F(Y_i) - F(Y_i)\|_F$; see Table 2. On the other hand, all the Randomized RK methods exhibit robust convergence of the expected order. Surprisingly and for reasons unclear to us, the projector splitting method provides very accurate results for this example.

α	3e-4	3e-1
average $\ F(Y_i) - P_r(Y_i)F(Y_i)\ _F$	3.0145×10^{-8}	2.9315×10^{-5}
max $\ F(Y_i) - P_r(Y_i)F(Y_i)\ _F$	5.8059×10^{-4}	0.5806

Table 2: Average and maximum $\|F(Y_i) - P_r(Y_i)F(Y_i)\|_F$ for the approximation Y_i at the i th time step of PRK 2 with $h = 2.5 \times 10^{-4}$.

4.3 Discrete Schrödinger equation in imaginary time

In this example, we aim at approximating the solution of the discrete Schrödinger equation in imaginary time from [8]:

$$\dot{A}(t) = -H[A(t)], \quad A(0) = A_0, \quad t \in [0, 0.5] \quad (28)$$

where

$$H[A(t)] = -\frac{1}{2}(DA(t) + A(t)D) + V_{\cos}A(t)V_{\cos} \in \mathbb{R}^{n \times n},$$

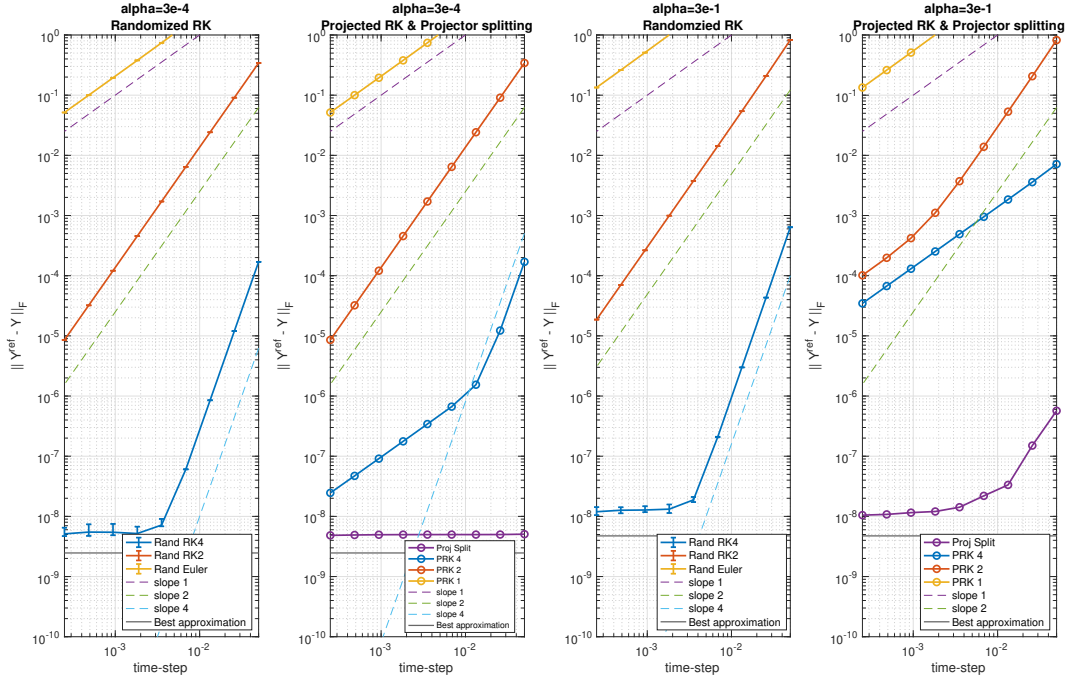


Figure 5: Non-linear Schrödinger equation with $\alpha = 3 \times 10^{-4}$ and $\alpha = 3 \times 10^{-1}$. Comparison of absolute approximation errors measured in Frobenius norm for different integrators for rank-30.

$D = \text{diag}(-1, 2, -1)$ is the discrete 1D Laplace, and V_{\cos} is the diagonal matrix with diagonal entries $1 - \cos(2j\pi/n)$ for $j = -n/2, \dots, n/2 - 1$. We choose $n = 512$ and an initial value A_0 that is randomly generated with prescribed singular values 10^{-i} , $i = 1, \dots, 512$. In Figure 6, we first plot the singular values of the reference solution at $T = 0.5$. Next, we plot the error of Rand Euler, Rand RK2, and Rand RK4 with a rank of 40. The final plot shows the error of Rand RK4 with various ranks. Once again, we observe the expected order of convergence until it reaches the level of low-rank approximation error.

4.4 Allen-Cahn equation

Following [5, Section 5.3], we consider the matrix differential equation arising from discretizing the Allen-Cahn equation via finite differences:

$$\dot{A} = \epsilon(LA + AL) + A - A^3,$$

with initial data

$$(A_0)_{ij} = \frac{[e^{-\tan^2(x_i)} + e^{-\tan^2(y_j)}] \sin(x_i) \sin(y_j)}{1 + e^{|\csc(-x_i/2)|} + e^{|\csc(-y_i/2)|}},$$

where A^3 is to be understood element-wise and $(x_i, y_j) \in [0, 2\pi]^2$, with $i, j = 1, \dots, 256$, are uniform discretization points. The matrix L is the one-dimensional finite-difference stencil, $\epsilon = 0.01$ and the time interval is $[0, 10]$. For this example, we apply Rand RK4 with rank 2 and time step size $h = 10^{-3}$. We plot the contours of the results at $t = 1, 3, 5, 7, 10$ in Figure 7. Additionally, we calculate the difference between the reference solution and normalize the difference by the Frobenius norm of the reference solution. We observe that Rand RK 4 accurately captures the contours despite using a very low rank.

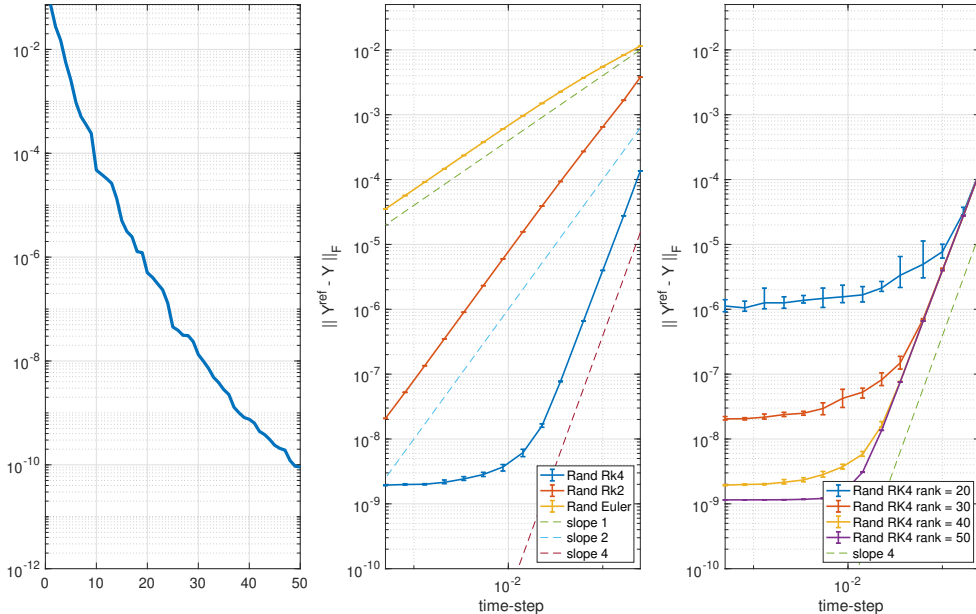


Figure 6: Discrete Schrödinger equation in imaginary time. The singular values of the reference solution at time $T = 0.5$, and rank-40 approximation errors for Rand Euler, Rand RK2 and Rand RK4. Also, the approximation errors for different ranks using Rand RK4.

5 Conclusions

In this work, we have proposed randomized low-rank Runge-Kutta methods. The analysis and numerical experiments clearly demonstrate the great potential of these methods to constitute an attractive alternative to existing dynamical low-rank methods. To fully realize this potential, further work is needed. On the practical side, rank adaptivity, parallelization, the preservation of quantities, and the combination with splitting/exponential integrators belong to the aspects that need to be studied in order to match the progress on dynamical low-rank methods achieved during the last decade. Also, the use of structured random matrices for efficiently addressing matrix differential equations with nonlinearities merits exploration. On the theoretical side, our analysis has raised a number of open questions. In particular, it would be important to develop a more direct approach for establishing the full order of randomized low-rank Runge-Kutta methods.

A Appendix: Proof of Theorem 5

The proof of Theorem 5 follows from modifying the proof of [14, Theorem 6] appropriately. We first provide an bound that relates the stages of the standard RK method (13) with the ones of the randomized RK method (14).

Theorem 11. *With the assumptions stated in Theorem 3, suppose that $\gamma_1 \leq \gamma_2 \leq \dots \leq \gamma_s$ are the stage orders of the explicit RK method (13) and denote $\tilde{\gamma}_j = \min(\gamma_j, \gamma_2 + 1)$. Then the differences between the stages \tilde{Z}_j of the standard RK method (13) with $A_i = Y_i$ and the stages Z_j of the randomized RK method (14) are bounded by*

$$\|Z_j - \tilde{Z}_j\|_{L_q} \leq \begin{cases} 0 & \text{if } j = 1, 2, \\ C(h^2\epsilon + h^{\gamma_2+2}) & \text{if } j = 3, 4, \dots, s, \end{cases} \quad (29)$$

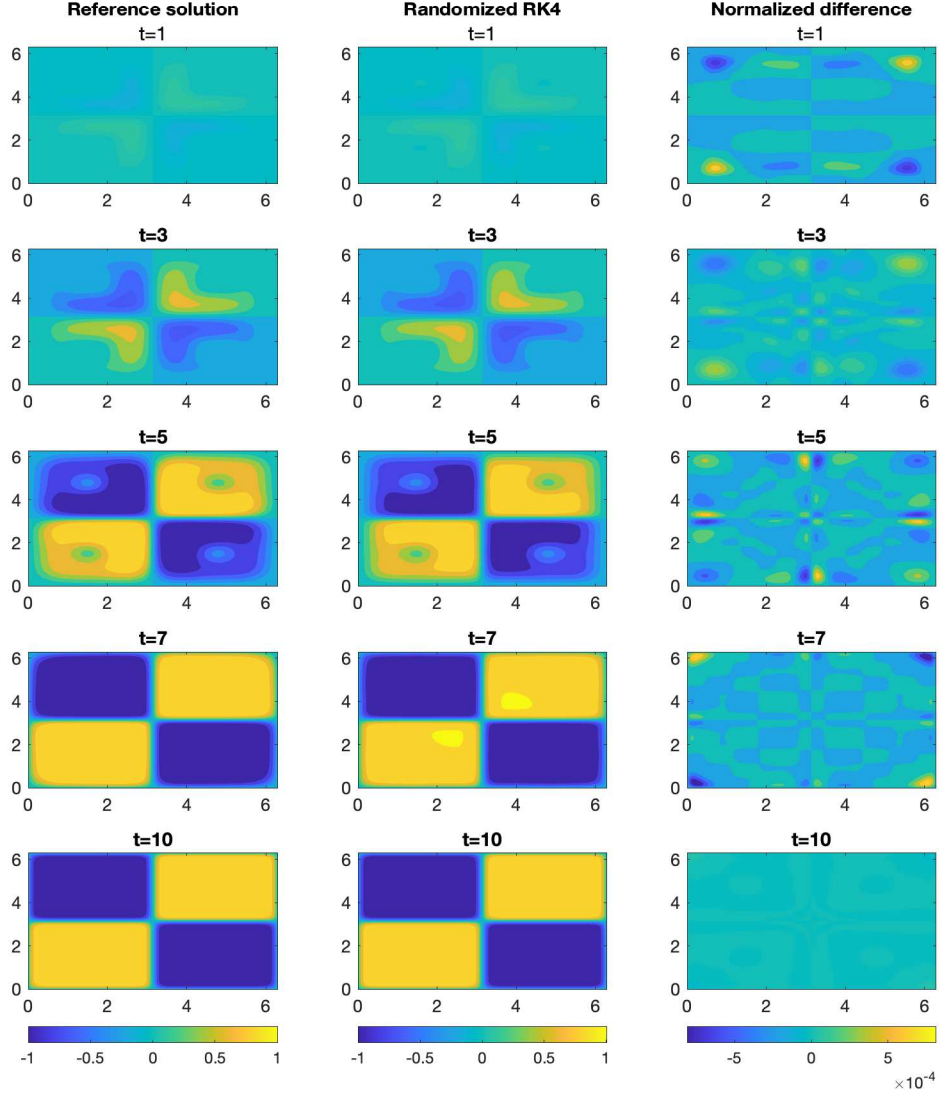


Figure 7: Allen-Cahn equation. The contour of the reference solution and the rank-2 approximation obtained by Rand RK4 with $h = 10^{-3}$.

$$\|F(\mathcal{N}_j(Z_j)) - F(\tilde{Z}_j)\|_{L^q} \leq \begin{cases} 0 & \text{if } j = 1, \\ C(h\epsilon + h^{\tilde{\gamma}_j+1}) & \text{if } j = 2, 3, \dots, s, \end{cases} \quad (30)$$

for any $0 \leq h \leq h_0$, where $q = \min\{p, \ell\}$ and the constant C depends only on $L, T, C_{\mathcal{N}}, h_0, s$ and $\max_{i,j} |a_{ij}|$.

Proof. For $j = 1$, we have $Y_i = Z_1 = \tilde{Z}_1$ and, therefore,

$$F(\mathcal{N}_1(Z_1)) = F(Z_1) = F(\tilde{Z}_1),$$

holds almost surely, because Y_i has rank at most r . For $j \geq 2$, we proceed by induction and assume that the statement of the theorem holds up to $j - 1$. Using the induction hypothesis

for (30), we obtain that

$$\begin{aligned} \|Z_j - \tilde{Z}_j\|_{L_q} &\leq h \sum_{l=1}^{j-1} |a_{jl}| \cdot \|F(\mathcal{N}_l(Z_l)) - F(\tilde{Z}_l)\|_{L_q} \\ &\leq \begin{cases} 0 & \text{if } j = 2, \\ C_A C h (sh\epsilon + h^{\tilde{\gamma}_2+1} + \dots + h^{\tilde{\gamma}_{j-1}+1}) & \text{if } j = 3, \dots, s. \end{cases} \end{aligned}$$

Using the assumption on the ordering of the stage orders and $\tilde{\gamma}_2 = \gamma_2$, it follows that

$$\|Z_j - \tilde{Z}_j\|_{L_q} \leq \begin{cases} 0 & \text{if } j = 1, 2, \\ C_Z h (h\epsilon + h^{\gamma_2+1}) & \text{if } j = 3, \dots, s, \end{cases} \quad (31)$$

which shows (29). To establish (30), we note that Z_j is independent from Ω_j and Ψ_j , which allows us to use the law of total expectation and Theorem 2 to conclude that

$$\begin{aligned} \|F(\mathcal{N}_j(Z_j)) - F(\tilde{Z}_j)\|_{L_q} &\leq L \|\mathcal{N}_j(Z_j) - \tilde{Z}_j\|_{L_q} \\ &\leq L \left((\mathbb{E}\{\mathbb{E}[\|\mathcal{N}_j(Z_j) - Z_j\|_F^q | Z_j]\})^{1/q} + \|Z_j - \tilde{Z}_j\|_{L_q} \right) \\ &\leq L \left(C_{\mathcal{N}} \|\llbracket Z_j \rrbracket_r - Z_j\|_{L_q} + \|Z_j - \tilde{Z}_j\|_{L_q} \right) \\ &\leq LC_{\mathcal{N}} \|\llbracket \tilde{Z}_j \rrbracket_r - Z_j\|_{L_q} + L \|Z_j - \tilde{Z}_j\|_{L_q} \\ &\leq LC_{\mathcal{N}} \|\llbracket \tilde{Z}_j \rrbracket_r - \tilde{Z}_j\|_F + L(C_{\mathcal{N}} + 1) \|Z_j - \tilde{Z}_j\|_{L_q}. \end{aligned}$$

While the second term of the last inequality is bounded by (31), we bound the first term by

$$\begin{aligned} \|\llbracket \tilde{Z}_j \rrbracket_r - \tilde{Z}_j\|_F &\leq \|\llbracket \phi_F^{c_j h}(Y_i) \rrbracket_r - \phi_F^{c_j h}(Y_i)\|_F + \|\phi_F^{c_j h}(Y_i) - \tilde{Z}_j\|_F \\ &\leq \|\llbracket \phi_F^{c_j h}(Y_i) \rrbracket_r - \phi_F^{c_j h}(Y_i)\|_F + C_L h^{\gamma_j+1} \\ &\leq (C_M \epsilon h + C_L h^{\gamma_j+1}), \end{aligned}$$

where we recall that the coefficient c_j was used in the definition (18) of stage order. In summary, we have

$$\begin{aligned} \|F(\mathcal{N}_j(Z_j)) - F(\tilde{Z}_j)\|_{L_q} &\leq LC_{\mathcal{N}}(C_M \epsilon h + C_L h^{\gamma_j+1}) + L(C_{\mathcal{N}} + 1) C_Z h (h\epsilon + h^{\gamma_2+1}) \\ &\leq C_F h (\epsilon + h^{\gamma_j} + h^{\gamma_2+1}). \end{aligned}$$

This concludes the proof of (30) using the definition of $\tilde{\gamma}_j$. \square

Proof of Theorem 3. By the triangular inequality, the local error satisfies

$$\|Y_{i+1} - \phi_F^h(Y_i)\|_{L_q} \leq \|Y_{i+1} - \tilde{Y}_{i+1}\|_{L_q} + \|\tilde{Y}_{i+1} - \phi_F^h(Y_i)\|_{L_q}, \quad (32)$$

where $\tilde{Y}_{i+1} = Y_i + h \sum_{j=1}^s b_j F(\tilde{Z}_j)$. Using that Ω_{s+1} and Ψ_{s+1} are independent of $\Omega_1, \dots, \Omega_s$ and

Ψ_1, \dots, Ψ_s , Theorem 2 yields

$$\begin{aligned}
& \|Y_{i+1} - \tilde{Y}_{i+1}\|_{L_q} = \left\| \mathcal{N}_{s+1} \left(Y_i + h \sum_{j=1}^s b_j F(\mathcal{N}_j(Z_j)) \right) - \tilde{Y}_{i+1} \right\|_{L_q} \\
& \leq C_{\mathcal{N}} \left\| \left[Y_i + h \sum_{j=1}^s b_j F(\mathcal{N}_j(Z_j)) \right]_r - Y_i - h \sum_{j=1}^s b_j F(\mathcal{N}_j(Z_j)) \right\|_{L_q} + \left\| h \sum_{j=1}^s b_j (F(\mathcal{N}_j(Z_j)) - F(\tilde{Z}_j)) \right\|_{L_q} \\
& \leq C_{\mathcal{N}} \left\| [\tilde{Y}_{i+1}]_r - Y_i - h \sum_{j=1}^s b_j F(\mathcal{N}_j(Z_j)) \right\|_{L_q} + \left\| h \sum_{j=1}^s b_j (F(\mathcal{N}_j(Z_j)) - F(\tilde{Z}_j)) \right\|_{L_q} \\
& \leq C_{\mathcal{N}} \left\| [\tilde{Y}_{i+1}]_r - \tilde{Y}_{i+1} \right\|_{L_q} + (1 + C_{\mathcal{N}}) \left\| h \sum_{j=1}^s b_j (F(\mathcal{N}_j(Z_j)) - F(\tilde{Z}_j)) \right\|_{L_q} \\
& \leq C \left(\epsilon h + h^{\tau+1} + h \sum_{j=1}^s |b_j| \| (F(\mathcal{N}_j(Z_j)) - F(\tilde{Z}_j)) \|_{L_q} \right).
\end{aligned}$$

Hence, the local error is bounded by

$$\|Y_{i+1} - \phi_F^h(Y_i)\|_{L_q} \leq \|Y_{i+1} - \tilde{Y}_{i+1}\|_{L_q} + C_1 h^{\tau+1} \leq Ch(\epsilon + h^\gamma + h^\tau).$$

This is turned into a bound on the global error using as in the proof of Theorem 3, which yields the bounds on the L_q norm claimed in the statement of Theorem 5. The tail bound is obtained using Markov's inequality; see Corollary 4. \square

References

- [1] Daniel Appelö and Yingda Cheng. Robust implicit Adaptive Low Rank Time-Stepping Methods for Matrix Differential Equations. *arXiv preprint arXiv:2402.05347*, 2024.
- [2] Hessam Babae, Minseok Choi, Themistoklis P. Sapsis, and George Em Karniadakis. A robust bi-orthogonal/dynamically-orthogonal method using the covariance pseudo-inverse with application to stochastic flow problems. *J. Comput. Phys.*, 344:303–319, 2017.
- [3] Oleg Balabanov, Matthias Beaupère, Laura Grigori, and Victor Lederer. Block subsampled randomized Hadamard transform for Nyström approximation on distributed architectures. In *Proceedings of the 40th International Conference on Machine Learning*, 2023.
- [4] Zvonimir Bujanović, Luka Grubišić, Daniel Kressner, and Hei Yin Lam. Subspace embedding with random Khatri-Rao products and its application to eigensolvers. *arXiv preprint arXiv:2405.11962*, 2024.
- [5] Benjamin Carrel and Bart Vandereycken. Projected exponential methods for stiff dynamical low-rank approximation problems. *arXiv preprint arXiv:2312.00172*, 2023.
- [6] Gianluca Ceruti, Lukas Einkemmer, Jonas Kusch, and Christian Lubich. A robust second-order low-rank BUG integrator based on the midpoint rule. *BIT*, 64(3):Paper No. 30, 2024.
- [7] Gianluca Ceruti, Jonas Kusch, and Christian Lubich. A parallel rank-adaptive integrator for dynamical low-rank approximation. *SIAM J. Sci. Comput.*, 46(3):B205–B228, 2024.
- [8] Gianluca Ceruti and Christian Lubich. An unconventional robust integrator for dynamical low-rank approximation. *BIT*, 62(1):23–44, 2022.

- [9] Aaron Charous and Pierre F. J. Lermusiaux. Dynamically orthogonal Runge-Kutta schemes with perturbative retractions for the dynamical low-rank approximation. *SIAM J. Sci. Comput.*, 45(2):A872–A897, 2023.
- [10] E. Hairer, S. P. Nørsett, and G. Wanner. *Solving ordinary differential equations. I*, volume 8 of *Springer Series in Computational Mathematics*. Springer-Verlag, Berlin, 1987.
- [11] Ernst Hairer, Christian Lubich, and Gerhard Wanner. *Geometric numerical integration*, volume 31 of *Springer Series in Computational Mathematics*. Springer, Heidelberg, 2010.
- [12] N. Halko, P. G. Martinsson, and J. A. Tropp. Finding structure with randomness: probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Rev.*, 53(2):217–288, 2011.
- [13] Emil Kieri, Christian Lubich, and Hanna Walach. Discretized dynamical low-rank approximation in the presence of small singular values. *SIAM J. Numer. Anal.*, 54(2):1020–1038, 2016.
- [14] Emil Kieri and Bart Vandereycken. Projection methods for dynamical low-rank approximation of high-dimensional problems. *Comput. Methods Appl. Math.*, 19(1):73–92, 2019.
- [15] Othmar Koch and Christian Lubich. Dynamical low-rank approximation. *SIAM J. Matrix Anal. Appl.*, 29(2):434–454, 2007.
- [16] Daniel Kressner and Hei Yin Lam. Randomized low-rank approximation of parameter-dependent matrices. *Numer. Linear Algebra Appl.*, page e2576, 2024.
- [17] Daniel Kressner and Lana Periša. Recompression of Hadamard products of tensors in Tucker format. *SIAM J. Sci. Comput.*, 39(5):A1879–A1902, 2017.
- [18] Jonas Kusch. Second-order robust parallel integrators for dynamical low-rank approximation. *arXiv preprint arXiv:2403.02834*, 2024.
- [19] Christian Lubich and Ivan V. Oseledets. A projector-splitting integrator for dynamical low-rank approximation. *BIT*, 54(1):171–188, 2014.
- [20] Yuji Nakatsukasa. Fast and stable randomized low-rank matrix approximation. *arXiv preprint arXiv:2009.11392*, 2020.
- [21] Steffen Schotthöfer, Emanuele Zangrando, Jonas Kusch, Gianluca Ceruti, and Francesco Tudisco. Low-rank lottery tickets: finding efficient low-rank neural networks via matrix differential equations. *Advances in Neural Information Processing Systems*, 35:20051–20063, 2022.
- [22] Andrea Trombettoni and Augusto Smerzi. Discrete solitons and breathers with dilute bose-einstein condensates. *Physical Review Letters*, 86(11):2353, 2001.
- [23] Joel A. Tropp, Alp Yurtsever, Madeleine Udell, and Volkan Cevher. Fixed-rank approximation of a positive-semidefinite matrix from streaming data. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 1225–1234, 2017.
- [24] Joel A. Tropp, Alp Yurtsever, Madeleine Udell, and Volkan Cevher. Practical sketching algorithms for low-rank matrix approximation. *SIAM J. Matrix Anal. Appl.*, 38(4):1454–1485, 2017.
- [25] André Uschmajew and Bart Vandereycken. Geometric methods on low-rank matrix and tensor manifolds. In *Handbook of variational methods for nonlinear geometric data*, pages 261–313. Springer, Cham, 2020.