SIAM J. SCI. COMPUT. Vol. 45, No. 2, pp. A872–A897

# DYNAMICALLY ORTHOGONAL RUNGE–KUTTA SCHEMES WITH PERTURBATIVE RETRACTIONS FOR THE DYNAMICAL LOW-RANK APPROXIMATION<sup>\*</sup>

AARON CHAROUS<sup>†</sup> AND PIERRE F. J. LERMUSIAUX<sup>†</sup>

Abstract. Whether due to the sheer size of a computational domain, the fine resolution required. or the multiples scales and stochasticity of the dynamics, the dimensionality of a system must often be reduced so that problems of interest become computationally tractable. In this paper, we develop retractions for time-integration schemes that efficiently and accurately evolve the dynamics of a system's low-rank approximation. Through differential geometry, we analyze the error incurred at each time-step due to the high-order curvature of the manifold of fixed-rank matrices. We first obtain a novel, explicit, computationally inexpensive set of algorithms that we refer to as perturbative retractions and show that the set converges to an ideal retraction that projects optimally and exactly to the manifold of fixed-rank matrices by reducing what we define as the projection-retraction error. Furthermore, each perturbative retraction itself exhibits high-order convergence to the best lowrank approximation of the full-rank solution. Using perturbative retractions, we then develop a new class of integration techniques that we refer to as dynamically orthogonal Runge–Kutta (DORK) schemes. DORK schemes integrate along the nonlinear manifold, updating the subspace upon which we project the system's dynamics as it is integrated. Through numerical test cases, we demonstrate our schemes for matrix addition, real-time data compression, and deterministic and stochastic partial differential equations. We find that DORK schemes are highly accurate by incorporating knowledge of the dynamic, nonlinear manifold's high-order curvature, and they are computationally efficient by limiting the growing rank needed to represent the evolving dynamics.

**Key words.** retraction, stochastic dynamical systems, reduced-order modeling, fixed-rank matrix manifold, dynamical low-rank approximation, curvature, dynamically orthogonal equations, Riemannian matrix optimization

 $\mathbf{MSC}$  codes. 54C15, 65F55, 53B21, 15A23, 57Z05, 57Z20, 57Z25, 60G60, 65C30, 65M12, 65M22, 65C20, 81S22, 94A08, 53A07, 35R60

DOI. 10.1137/21M1431229

1. Introduction. Simulation needs will always outstrip current computational resources. As computing power grows, so too does our desire to solve larger and larger problems. The curse of dimensionality limits the possibility of computing exact solutions to high-dimensional problems, so obtaining sufficiently accurate approximate solutions via optimal reduced-order modeling is essential.

In this paper, we develop a perturbative methodology to evolve a high-order lowrank approximation, X(t), in time that approximates a full-rank system state,  $\mathfrak{X}(t)$ . First studied in the context of matrix initial value problems, this approach is called the *dynamical low-rank approximation* [34]. More precisely, we seek X(t) such that at all fixed times t, X(t) is the best approximation to  $\mathfrak{X}(t)$ . That is,

(1.1) 
$$X(t) = \underset{\tilde{X} \in \mathscr{M}_r}{\arg\min} \|\tilde{X}(t) - \mathfrak{X}(t)\|,$$

Submitted to the journal's Methods and Algorithms for Scientific Computing section July 6, 2021; accepted for publication (in revised form) November 9, 2022; published electronically April 27, 2023.

https://doi.org/10.1137/21M1431229

**Funding:** This work was partially supported by the Office of Naval Research under grants N00014-19-1-2693 (IN-BDA) and N00014-19-1-2664 (Task Force Ocean, DEEP-AI).

<sup>&</sup>lt;sup>†</sup>Mechanical Engineering, Center for Computational Science and Engineering, Massachusetts Institute of Technology, Cambridge, MA 02139 USA (acharous@mit.edu, pierrel@mit.edu).

where  $\mathscr{M}_r$  denotes the manifold of matrices of r. More realistically, we are only concerned with discrete values of t,  $\{t_i\}_i$ , with  $\Delta t \equiv t_{i+1} - t_i$ . In this paper, we assume  $X(t_i), \mathfrak{X}(t_i) \in \mathbb{R}^{m \times n}$ , and we use the Frobenius norm. We also assume  $\operatorname{rank}(\mathfrak{X}(t_i)) \geq r$ for all  $t_i$  to ensure that the solution to (1.1) exists. Because we are interested in highdimensional problems, m and n are assumed to be very large; maybe  $\mathfrak{X}(t_i)$  cannot be stored (either on a local hard drive or in RAM), so we seek a compressed form of  $\mathfrak{X}(t_i)$ for both easy storage and computationally inexpensive matrix operations. As such, we restrict  $\operatorname{rank}(X) = r$  with  $r \ll m, n$ , so X admits a low-rank representation  $UZ^T$ , where  $U \in \mathbb{R}^{m \times r}, Z \in \mathbb{R}^{n \times r}$ . With this, we only evolve and store (m+n)r values for each time  $t_i$  rather than mn values.

Our present goal is to obtain new retractions—mappings that ensure the lowrank approximation remains on  $\mathcal{M}_r$ —for numerical integrators so that, as  $\Delta t \to 0$ ,  $X(t_i)$  converges with high order to the best possible approximation to  $\mathfrak{X}(t_i)$ . Previous works (e.g., [34, 21, 20, 51, 32, 55, 56, 10]) already provide robust algorithms for approximate solutions to this problem. For a thorough review on retractions, see [2]. But, most existing algorithms only promise first-order convergence in time to the best low-rank approximation. One notable exception is the projector-splitting integrator, which may yield arbitrarily high-order convergence to the best approximation by taking symmetric compositions using the adjoint method [51]. Furthermore, it has been shown to preserve fixed-point iteration convergence rates under some regularity conditions [35]. However, the projector-splitting integrator, along with the randomized SVD [24] and truncated SVD algorithms, do not preserve mode continuity. As pointed out in [22], unlike the retractions we develop, the error bounds for the randomized SVD provided in [24] do not suggest convergence to the best low-rank approximation, i.e., the truncated SVD itself; instead, the error between the full-rank solution  $\mathfrak{X}(t_i)$  and its rank-r approximation  $X(t_i)$  is bounded as a multiple of  $\sigma_{r+1}$ , the largest singular of  $\mathfrak{X}(t_i)$  not captured by a rank-r approximation, which could be large. Depicted in Figure 1, we show that after randomly updating a point on the low-rank manifold, our new perturbative retractions only slightly change the vectors in U, whereas other retractions in the literature change them completely since they only preserve the subspace's span. Mode continuity [20, 47, 67] is most useful for the reduced-order evolution of dynamical systems in time, especially in the context of uncertainty quantification, which further motivates the need for new retractions.

Downloaded 05/06/23 to 24.61.23.148 by Pierre Lermusiaux (pierrel@mit.edu). Redistribution subject to SIAM license or copyright; see https://epubs.siam.org/terms-privacy

Another set of integration schemes, the projected Runge–Kutta methods [33], appears to retain the order of accuracy of classical Runge–Kutta schemes up to the Dirac–Frenkel model closure error  $\varepsilon_{DF}$  (to be defined in section 3), which the projector-splitting integrator is also susceptible to. However, the projected Runge– Kutta methods do not incorporate knowledge of the high-order curvature of the low-



FIG. 1. Starting from a randomly initialized point  $U_0 Z_0^T$  on a rank-10 manifold, we show a subset of the updated subspaces  $U_1$  after randomly updating the point via addition. That is, we apply different retractions to  $U_0 Z_0^T + \Delta t \overline{\mathscr{S}}$  (where  $\overline{\mathscr{S}}$  is a random matrix and  $\Delta t = 0.25$ ) and investigate how close  $U_0$  and  $U_1$  appear.

rank manifold. Furthermore, errors are incurred by departing and reprojecting onto the low-rank manifold. We show that integrating along the low-rank manifold provides more accurate integration schemes.

In this paper, we derive a set of retractions (to be defined in section 3) and integration schemes that theoretically guarantee high-order convergence of X(t) to the best low-rank approximation of  $\mathfrak{X}(t)$ . Moreover, these retractions are explicit and computationally inexpensive and have a fixed number of iterations in order to obtain a designated order of convergence. Within time-steps, our new Runge-Kutta schemes dynamically update the subspace upon which we project the system's dynamics, resulting in improved accuracy and reduced computational cost. The framework we propose also allows for not only the derivation of new integration schemes but also the adaptation of preexisting full-rank integration schemes in a low-rank setting. This may allow for new low-rank integration schemes that preserve physical quantities of interest, e.g., energy, or other application-specific properties such as symplecticity.

We motivate our work with two examples. First, suppose we are given  $\mathfrak{X}(t_i)$  at all  $t_i$ , but we would like to only store  $X(t_i)$ . We could compute  $X(t_i)$  via the truncated SVD at each fixed time as it is the best low-rank approximation in the Frobenius norm [18, 54, 63]. However, this is often prohibitively expensive: from [66, p. 237], the SVD may be computed via Golub-Kahan bidiagonalization in about  $4mn^2 - \frac{4}{3}n^3$  flops. If, however, we can compute one SVD of  $\mathfrak{X}(t_i)$  at only the first time instant, can we somehow predict  $X(t_i)$  at future times with some cheaper approximate the SVD, we accomplish this feat, avoiding the computational burden of computing an SVD at each time-step, and the efficacy of our algorithm is demonstrated in the real-time data compression example given in section 7. More broadly, we show that we may apply our perturbative retractions to problems where dynamics are data-driven rather than given by a PDE as studied previously (see, e.g., [34, 67, 20]).

Second, suppose we have a dynamical system, stochastic or deterministic, and we are given initial conditions  $\mathfrak{X}(0)$ . We cannot afford to evolve  $\mathfrak{X}(t)$ , so we must somehow evolve a low-rank approximation X(t). However, the differential equations that describe the system's dynamics are given for  $\mathfrak{X}(t)$ , not X(t). How can we obtain the differential equations for X(t) that achieve high-order accuracy with respect to the full space? Furthermore, are there numerical schemes that reduce the error between the numerically integrated X(t) and the theoretical best approximation to  $\mathfrak{X}(t)$ ? By using high-order classical time-integration schemes coupled with our perturbative retractions with high-order corrections, we show how to advance X efficiently and accurately in section 7.

In section 2, we review the dynamic Karhunen–Loéve expansion as the mathematical foundation for what follows. Section 3 gives definitions of several matrix spaces and terminology referenced throughout. In section 4, we introduce the key errors and outline the problem setup. In section 5, we derive *perturbative retractions*, the first main contribution of this paper. We build off of these retractions in section 6 to derive the second main contribution of this paper, *dynamically orthogonal Runge–Kutta* (DORK) schemes, which offer a new family of schemes that integrate along a dynamic, nonlinear manifold. We demonstrate the efficiency of these retractions and schemes in section 7 with examples of low-rank matrix addition, matrix differential equations, real-time data compression, partial differential equations, and stochastic differential equations, including comparisons with existing schemes. Finally, we summarize and discuss our results and provide future research directions in section 8. For further details on some of the results presented next, we refer to [11].

A874

Downloaded 05/06/23 to 24.61.23.148 by Pierre Lermusiaux (pierrel@mit.edu). Redistribution subject to SIAM license or copyright; see https://epubs.siam.org/terms-privacy

2. Reduced-order modeling preliminaries. To further motivate the (spatially) discrete problem of interest, we consider reduced-order methods for (spatially) continuous problems. For stochastic dynamics, one main approach is the polynomial chaos expansion, where a fixed polynomial basis parameterizes the stochastic space [71, 72, 8, 73, 74]. Another approach is the Karhunen–Loéve (KL) expansion or proper orthogonal decomposition, which is a data-driven decomposition yielding the best low-rank approximation in total mean square error [52, 6, 36, 29, 49]. We start our analysis from the *dynamic* KL expansion [34, 39, 61]. Given a square-integrable stochastic process  $\Phi(x,t;\omega)$  (analogous to our discrete  $\mathfrak{X}$  previously) defined over a probability space  $(\Omega, \mathcal{F}, \mathbf{P})$ , where  $x \in \mathcal{D} \subseteq \mathbb{R}^m$  denotes our spatial variable and  $t \in [0,T]$  time.  $\Omega$  is the sample space with  $\omega$  denoting a simple event,  $\mathcal{F}$  is the  $\sigma$ -algebra, and  $\mathbf{P}$  is the probability measure.  $\Phi$  may be decomposed as follows,

$$\Phi(x,t;\omega) = \mathbb{E}[\Phi(x,t;\omega)] + \sum_{i=1}^{\infty} \varphi_i(x,t)\zeta_i(t;\omega) \,.$$

Here,  $\mathbb{E}$  denotes the expectation operator,  $\varphi_i$  the spatial modes, and  $\zeta_i$  the stochastic processes or coefficients. The mean, modes, and coefficients are allowed to evolve in time. At each fixed time t, the modes form an orthonormal basis for  $\mathcal{D}$ , and the  $\zeta_i$  are zero-mean, uncorrelated random variables. Truncating the series to a finite number of modes and coefficients with the greatest variance would then yield the best possible approximation at each time. This is, of course, assuming that we know  $\Phi$  for all t.

As an aside that will become relevant later, we note that this decomposition is the stochastic analog to the separation of variables technique to solve partial differential equations [48] and for function approximation [78]. That is, instead of expressing a deterministic function as  $f(x_1, x_2) = \sum_{i=1}^{\infty} g_i(x_1)h_i(x_2)$ , we allow  $x_2$  to denote a simple event in the event space (denoted  $\omega \in \Omega$ ).

Now we consider an initial value problem where  $\Phi(x,0;\omega)$  is known, and the dynamics of each stochastic realization of  $\Phi$  are given by partial differential equations,

$$\frac{\partial \Phi}{\partial t} = \mathscr{L}_c(\Phi, x, t; \omega)$$

where  $\mathscr{L}_c$  may be any differential operator, nonlinear or linear, stochastic or deterministic. Following [61], we seek to evolve the mean, modes, and coefficients so that we evolve a truncated approximation to  $\Phi$  (analogous to our discrete X previously) without reconstructing the series at each time. We impose a gauge condition known as the dynamically orthogonal (DO) condition:  $\langle \varphi_i, \frac{\partial \varphi_j}{\partial t} \rangle = 0$  for all i, j, where  $\langle \cdot, \cdot \rangle$ denotes a spatial inner product over  $\mathcal{D}$ . The DO condition eliminates redundant degrees of freedom, ensures the modes  $\phi_i$  remain orthonormal in time, and decouples our system so that we can write the DO equations [61, 21] (see section SM1).

After intrusively deriving the DO equations for a particular differential operator  $\mathscr{L}_c$ , one may spatially discretize these equations directly (see, e.g., [67, 20]) and solve the reduced-order model numerically. Another approach is to discretize the dynamically orthogonal equations themselves and then insert a spatially discretized differential operator  $\mathscr{L}$ , yielding a nonintrusive approach [21]. In a slightly different setting, this was first analyzed in [34], where the *dynamical low-rank approximation* was proposed to solve time-dependent matrix initial value problems. The connection between the dynamical low-rank approximation and the DO equations was made in [21]; the DO equations can be thought of as instantaneously projecting the full-rank dynamics onto a manifold of fixed-rank matrices (also called the low-rank manifold). That is, for a discrete system state  $\mathfrak{X}$  with dynamics given by

AARON CHAROUS AND PIERRE F. J. LERMUSIAUX

(2.1) 
$$\frac{d\mathfrak{X}}{dt} = \mathscr{L}(\mathfrak{X}, t; \omega)$$

A876

with initial conditions  $\mathfrak{X}(0) = \mathfrak{X}_0$ , the dynamics of the low-rank approximation X are

(2.2) 
$$\frac{dX}{dt} = \mathscr{P}_{\mathcal{T}_X \mathscr{M}_r}[\mathscr{L}(\mathfrak{X}, t; \omega)] \approx \mathscr{P}_{\mathcal{T}_X \mathscr{M}_r}[\mathscr{L}(X, t; \omega)]$$

with initial conditions  $X(0) = \mathscr{M}_r \mathfrak{X}_0$ , where  $\mathscr{P}_{\mathcal{T}_X \mathscr{M}_r}$  denotes the operator that projects onto the tangent space of the low-rank manifold at X, and  $\mathscr{P}_{\mathscr{M}_r}$  denotes the projection onto the low-rank manifold. Again, since  $\mathfrak{X}$  is typically unknown, (2.2) is an approximation, but it is instantaneously optimal by the Dirac–Frenkel principle. Intuitively, we project the dynamics onto the tangent space because this yields the best possible approximation of  $\mathscr{L}$  accessible to a low-rank solution. This is derived in [21, 11]. What we really seek, however, is not the evolution equation for X, as this would require reconstructing X at every time-step; we seek evolution equations for Uand Z so that we may evolve X implicitly. Adopting Newton's notation for differentiation, we note that  $X = UZ^T \Rightarrow \dot{X} = U\dot{Z}^T + \dot{U}Z^T$ . From this, one can develop a bijective map between  $\dot{X}$  and  $(\dot{U}, \dot{Z})$  given U, Z [21]. As such, we may equivalently write the following matrix differential equations,

(2.3)  
$$\dot{U} = \mathscr{P}_{U}^{\perp} \mathscr{L}(UZ^{T}, t; \omega) Z(Z^{T}Z)^{-1}, \\ \dot{Z} = \mathscr{L}(UZ^{T}, t; \omega)^{T}U,$$

where  $\mathscr{P}_U^{\perp} = I - UU^T$ ; the columns of U,  $u_i$ , correspond to  $\varphi_i$ ; and the columns of Z,  $z_i$ , correspond to realizations of  $\zeta_i$ . The discrete DO condition is now  $\dot{U}^T U = 0$ . These equations (2.3) will be revisited in section 5. In contrast to the original equations (SM1.1, SM1.2, SM1.3), (2.3) do not explicitly extract the mean of  $\Phi$ ; the mean may be included by adding a column in U and a column of ones in Z [20].

Though our motivation has been stochastic dynamics, we remark that this discrete decomposition can be used in the deterministic setting. Instead of thinking of  $\Omega$  as a stochastic event space, we may replace it with another physical space  $\tilde{\mathcal{D}}$ . The expectation operator can be thought of as an inner product weighted by a probability measure; so in the deterministic case, the expectation operator just becomes an inner product over  $\tilde{\mathcal{D}}$ . We have already mentioned that the proper orthogonal decomposition is the stochastic analog of separation-of-variables, so this is not a large intellectual leap. As we will see in this paper, the interpretation of  $\mathcal{D}$  and  $\Omega$  is unimportant to our analysis and may be abstracted away.

Note that this  $UZ^T$  parameterization is simply one choice; other common matrix parameterizations include  $USV^T$  as in [34], and though different parameterizations will yield different evolution equations, the overarching mathematics is the same. In fact, setting  $Z \leftarrow VS^T$  or  $[V, S^T] \leftarrow qr(Z)$  (where qr denotes algorithmic implementation of the QR decomposition) defines a mapping between the two parameterizations.

**3. Definitions.** Before proceeding, we formally define some mathematical structures and terminology, following the notation of [21].

DEFINITION 3.1 (low-rank manifold).  $\mathcal{M}_r = \{A \in \mathbb{R}^{m \times n} : \operatorname{rank}(A) = r\}$ 

DEFINITION 3.2 (Stiefel manifold).  $\mathcal{V}_{m,r} = \{A \in \mathbb{R}^{m \times r} : A^T A = I\}$ 

DEFINITION 3.3.  $\mathbb{R}^{n \times r}_* = \{A \in \mathbb{R}^{n \times r} : \operatorname{rank}(A) = r\}$ 

Note that any matrix  $X \in \mathscr{M}_r$  may be written as  $X = UZ^T$  with  $U \in \mathcal{V}_{m,r}$  and  $Z \in \mathbb{R}^{n \times r}_*$ —simply consider the SVD of  $X, X = U\Sigma V^T$ , setting  $Z = V\Sigma^T$ . This

The discrete DO equations, as well as alternative equations that paramaterize the dynamical low-rank approximation, instantaneously project the full-rank dynamics of the system onto the tangent space of the low-rank manifold. So, it is useful to be able to characterize the tangent space of  $\mathcal{M}_r$ . Any matrix in the tangent space of  $\mathcal{M}_r$  at point  $X = UZ^T$  may be written as  $\delta_X = U\delta_Z^T + \delta_U Z^T$  [68].

DEFINITION 3.4. The tangent space of  $\mathcal{M}_r$  at  $X = UZ^T$  is defined as follows:

$$\mathcal{T}_X \mathscr{M}_r = \left\{ U \delta_Z^T + \delta_U Z^T : \delta_U \in \mathbb{R}^{m \times r}, \delta_Z \in \mathbb{R}^{n \times r} \right\}$$

DEFINITION 3.5. The tangent space to the Stiefel manifold at U is as follows:

$$\mathcal{T}_U \mathcal{V}_{m,r} = \left\{ \delta_U \in \mathbb{R}^{m \times r} : \delta_U^T U + U^T \delta_U = 0 \right\} = \left\{ \delta_U \in \mathbb{R}^{m \times r} : U^T \delta_U \in \mathrm{so}(r) \right\} \,.$$

Above, so(r) denotes the set of skew-symmetric, real  $r \times r$  matrices.

Hence, the discrete DO condition  $\dot{U}^T U = 0$  ensures that  $\dot{U} \in \mathcal{T}_U \mathcal{V}_{m,r}$ . In fact, we can define the DO space  $\mathcal{U}_{m,r}$  as a more restrictive set of matrices such that  $U \in \mathcal{U}_{m,r}$ .

DEFINITION 3.6 (DO space).  $\mathcal{U}_{m,r} = \{\delta_U \in \mathbb{R}^{m \times r} : U^T \delta_U = 0\} \subset \mathcal{T}_U \mathcal{V}_{m,r}.$ 

With this, a more restrictive parameterization of the tangent space of  $\mathcal{M}_r$  can be written as follows (see [21, p. 517]), replacing Definition 3.4.

DEFINITION 3.7. The tangent space  $\mathcal{T}_X \mathscr{M}_r$  of  $\mathscr{M}_r$  at  $X = UZ^T$  admits the following representation, leading to a unique parameterization for each matrix in  $\mathcal{T}_X \mathscr{M}_r$ :

$$\mathcal{T}_X \mathscr{M}_r = \{ U \delta_Z^T + \delta_U Z^T : \delta_U \in \mathcal{U}_{m,r}, \delta_Z \in \mathbb{R}^{n \times r} \}$$

A so-called retraction, in simple terms, maps a matrix in the affine tangent space back to the low-rank manifold. An extended retraction maps a matrix in the embedding Euclidean space back to the low-rank manifold. From [3], we write out a formal definition for extended retractions.

DEFINITION 3.8. An extended retraction  $\mathcal{R}: \mathcal{E} \to \mathcal{M}_r$  is a mapping such that

- 1.  $\mathcal{R}$  is defined and smooth on a neighborhood of the zero section in  $\mathcal{TM}_r$ ,
- 2.  $\mathcal{R}_X(0) = X \quad \forall X \in \mathscr{M}_r,$
- 3.  $\frac{d}{dt}\mathcal{R}_X(t\xi)\Big|_{t=0} = \mathscr{P}_{\mathcal{T}_X\mathscr{M}_r}\xi$  for all  $X \in \mathscr{M}_r$  and  $\xi \in \mathcal{E}$ .

For a nonextended retraction  $\mathcal{R}: \mathcal{T}\mathscr{M}_r \to \mathscr{M}_r$ , condition 3 is modified such that  $\xi \in \mathcal{T}_X \mathscr{M}_r$ . We use the same notation for retractions and extended retractions since they accomplish the same feat, and whether the retraction itself is extended or not may be determined implicitly by the argument.

Finally, we define three errors due to closure.

DEFINITION 3.9 (normal closure error). The normal closure error,  $\varepsilon_{\mathcal{N}}$ , is defined as the difference between a full-rank state and its best low-rank approximation,

$$\varepsilon_{\mathcal{N}} \equiv \mathfrak{X} - \mathscr{P}_{\mathscr{M}_r}(\mathfrak{X}).$$

This error is normal to the tangent space of the low-rank manifold at  $\mathscr{P}_{\mathcal{M}_r}(\mathfrak{X})$ [11, 20].

DEFINITION 3.10 (dynamical model closure error). The dynamical model closure error,  $\varepsilon_D$ , is the difference in dynamics by applying the dynamical model to the full-rank state  $\mathfrak{X}$  and its approximation X:

$$\varepsilon_D \equiv \mathscr{L}(\mathfrak{X}) - \mathscr{L}(X).$$

Approximating  $\mathscr{L}(\mathfrak{X})$  as  $\mathscr{L}(X)$  is often useful when only X is available or affordable for computation. Assuming  $\mathscr{L}$  is Lipschitz continuous, one can then bound how error accumulates over time (see [21, 34]). In contrast to definition 3.9, the dynamical model closure error measures the local error of the system's time derivative.

DEFINITION 3.11 (Dirac–Frenkel time-dependent variational principle and Dirac– Frenkel model closure error). From [15, 50], the Dirac–Frenkel time-dependent variational principle is the use of the low-rank approximation when calculating the system dynamics and then projecting those dynamics onto the tangent space of an approximation manifold. This induces what we call the "Dirac–Frenkel model closure error":

$$\varepsilon_{DF} \equiv \mathscr{L}(\mathfrak{X}) - \mathscr{P}_{\mathcal{T}_{\mathcal{X}}\mathcal{M}_{r}}\mathscr{L}(X).$$

The Dirac–Frenkel model closure error is the difference that arises (at each time-step) from applying the dynamics to the full-rank state and applying the dynamics to the reduced-rank state and projecting the result onto the tangent space of the low-rank manifold, as prescribed by the Dirac–Frenkel time-dependent variational principle.

4. Projection-retraction error. With the above definitions, we now formally describe the problem we seek to solve. In the continuous-time limit where  $\Delta t \rightarrow 0$ , the equations of (2.3) give the instantaneously optimal approximation of the full-rank dynamics up to dynamical model closure error (Definition 3.10) [21]. That is, if we could integrate (2.3), we would obtain the exact dynamical rank-r approximation to (2.1). However, this is challenging: they are coupled, nonlinear matrix differential equations,  $\mathscr{L}$  depends implicitly on U, Z, and time, and importantly, over  $\Delta t$ , the true solution can leave the tangent space. Existing numerical schemes include projector-splitting integrators [51], but here we take an alternative approach as in [33, 14]. Instead of integrating (2.3), we return to the original (2.1). We split the problem in two: integrating the system in time and retracting back to the manifold.

For the time integration, we assume that the operation

(4.1) 
$$\overline{\mathscr{L}}(\bullet;\omega) \equiv \frac{1}{\Delta t} \int_{t_n}^{t_{n+1}} \mathscr{L}(\bullet,t;\omega) dt + \mathcal{O}(\Delta t^k)$$

is given.  $\overline{\mathscr{L}}$  may be obtained via a *k*th-order time-integration scheme (e.g., Euler, Runge–Kutta, leapfrog, Crank–Nicolson), which includes some time-integration error. We may also assume that any error due to spatial discretization is built into  $\overline{\mathscr{L}}$ . On the other hand,  $\overline{\mathscr{L}}$  may be given a priori as in the first motivating example where  $\mathfrak{X}$  is given. The first argument of  $\overline{\mathscr{L}}$  denoted with • refers to either the low-rank X(which would induce the dynamical model closure error) or full-rank  $\mathfrak{X}$ , whichever is available; the retractions we develop are agnostic to the errors induced by time integration, spatial discretization, or closure errors.  $\overline{\mathscr{L}}$  may be thought of as the fullrank direction in which we would like to travel, and it may leave the tangent space of the low-rank manifold. Let us denote  $\mathfrak{X}(t_n)$  as  $\mathfrak{X}_n$  and  $X(t_n)$  as  $X_n = U_n Z_n^T$ . By definition, we have that

$$\mathfrak{X}_{n+1} = \mathfrak{X}_n + \Delta t \overline{\mathscr{L}}$$

is a consistent kth-order time-integration scheme for (2.1). Because we seek the best low-rank approximation X, we have that (4.2) is a consistent integrator for X(t),

(4.2) 
$$X_{n+1}^* = \mathscr{P}_{\mathscr{M}_r}(\mathfrak{X}_n + \Delta t \overline{\mathscr{P}})$$

(4.3) 
$$\approx \mathscr{P}_{\mathscr{M}_r}(X_n + \Delta t \overline{\mathscr{L}}),$$

where  $X_{n+1}^*$  denotes the best possible approximation of  $\mathfrak{X}_{n+1}$  given our knowledge at time  $t_n$ . Of course we may not know  $\mathfrak{X}_n$ , so we approximate it with  $X_n$  in (4.3), which induces our normal closure error (see Definition 3.9).

Finally, in the limit of  $\Delta t \to 0$ ,  $\mathscr{P}_{\mathscr{M}_r} \to \mathscr{P}_{\mathcal{T}_X \mathscr{M}_r}$ , and so we may rewrite (4.3) as

(4.4) 
$$X_{n+1}^* = \mathscr{P}_{\mathcal{T}_{X_n}\mathscr{M}_r} \left( X_n + \Delta t \overline{\mathscr{L}} \right) + \mathcal{O}(\Delta t^2) = X_n + \Delta t \mathscr{P}_{\mathcal{T}_{X_n}\mathscr{M}_r} \overline{\mathscr{L}} + \mathcal{O}(\Delta t^2)$$

We remark that although  $X_{n+1}^*$ , the best low-rank approximation at time  $t_{n+1}$  given present information is completely characterized by a *particular* choice of the  $\mathcal{O}(\Delta t^2)$ term (defined by  $(\mathscr{M}_r - \mathscr{P}_{\mathcal{T}_{X_n}\mathcal{M}_r})(X_n + \Delta t \overline{\mathscr{L}}))$ , and any choice of the  $\mathcal{O}(\Delta t^2)$  will result in a consistent integrator for the system; it just may not be the best possible time-integration scheme. With the time integration settled, we now proceed to the retractions which approximate  $\mathscr{M}_r$  by explicitly obtaining the  $\mathcal{O}(\Delta t^2)$  term in (4.4).

In the literature [3], a second-order retraction is defined as a retraction whose second-order error belongs to the normal space of  $\mathcal{M}_r$  at X, which is a property shared with geodesics. From [21], we know the DO equations (and any other paramterization of the dynamical low-rank approximation) instantaneously apply the truncated SVD to the dynamics as in (4.3). Hence, we will develop retractions that approximate  $\mathcal{P}_{\mathcal{M}_r}$ rather than the geodesics of  $\mathcal{M}_r$  in order to form consistent and accurate integrators.

Given  $\overline{\mathscr{L}}, U_n, Z_n$ , we seek  $U_{n+1}, Z_{n+1}$ . Because we are restricted to the low-rank manifold, we must somehow correct  $\overline{\mathscr{L}}$  to stay on the low-rank manifold: this is the idea of a retraction. If we integrate (2.3), we can define

(4.5) 
$$\overline{U} \equiv \frac{U_{n+1} - U_n}{\Delta t} \approx \frac{1}{\Delta t} \int_{t_n}^{t_{n+1}} \mathscr{P}_U^{\perp} \mathscr{L}Z(Z^T Z)^{-1} dt$$

(4.6) 
$$\overline{Z} \equiv \frac{Z_{n+1} - Z_n}{\Delta t} \approx \frac{1}{\Delta t} \int_{t_n}^{t_{n+1}} \mathscr{L}^T U dt.$$

That is, we define  $\overline{U}\overline{Z}$  as the scaled difference between states at  $t_{n+1}$  and  $t_n$ . If U and Z were the exact best dynamic approximation over  $\Delta t$ , the approximate equalities before the integrals would become exact equalities.

As an example, suppose we define  $\overline{U} = U|_{t \equiv t_n} \equiv U_n$  and  $\overline{Z} = Z|_{t=t_n} \equiv Z_n$ , where U and Z are given as  $\mathscr{P}_U^{\perp} \overline{\mathscr{L}} Z(Z^T Z)^{-1}$  and  $\overline{\mathscr{L}}^T U$ , respectively. This will lead to the forward Euler retraction. Writing out some algebra, we obtain the following:

$$\begin{aligned} X_{n+1} &= U_{n+1} Z_{n+1}^T = (U_n + \Delta t \vec{U}) (Z_n + \Delta t \vec{Z})^T \\ &= U_n Z_n^T + \Delta t \left[ U_n \dot{Z}_n^T + \dot{U}_n Z_n^T \right] + \Delta t^2 \dot{U}_n \dot{Z}_n^T \\ &= \underbrace{X_n + \Delta t \mathscr{P}_{\mathcal{T}_{X_n} \mathscr{M}_r} \widetilde{\mathscr{L}}}_{\text{consistent integrator}} + \underbrace{\Delta t^2 \dot{U}_n \dot{Z}_n^T}_{\text{retraction}}, \end{aligned}$$

where the last equation arises from the definitions of the tangent space (Definition 3.7) and consistent integrator (4.4). Even if  $\overline{\mathscr{L}}$  were accurate up to high order, we would still have a new term that is  $\mathcal{O}(\Delta t^2)$ , which must be taken into account.

Finally, we present the goal of this paper. Because  $\Delta t > 0$ , it is insufficient to simply project the dynamics onto the affine tangent space of the low-rank manifold: we must consider the high-order curvature by approximately projecting onto the lowrank manifold itself. That is, we seek an efficient retraction or extended retraction (Definition 3.8), such that

(4.7) 
$$X_{n+1} \equiv \mathcal{R}_{X_n}(\Delta t \overline{\mathscr{L}}) = \mathscr{P}_{\mathscr{M}_r}(X_n + \Delta t \overline{\mathscr{L}}) + \mathcal{O}(\Delta t^k).$$

This allows us to write the error which we wish to correct, the *projection-retraction* error,  $\varepsilon_{pr}$ , defined below as the difference between our chosen retraction and the projection of the full-rank integral of the dynamics over  $\Delta t$ ,

(4.8) 
$$\varepsilon_{pr} \equiv \mathcal{R}_{X_n}(\Delta t \overline{\mathscr{L}}) - \mathscr{P}_{\mathscr{M}_r}(X_n + \Delta t \overline{\mathscr{L}}).$$

Henceforth, a "retraction with an *n*th-order correction" or a "retraction of *n*th order" will denote a retraction with  $\varepsilon_{pr} = \mathcal{O}(\Delta t^{n+1})$ .

Figure 2(a) graphically depicts where this projection-retraction error arises. After projecting the dynamics,  $\overline{\mathscr{G}}$ , onto the affine tangent space, we must somehow get back to the low-rank manifold. It is important to note that in the continuous-time limit, the magnitude of the retraction is exactly zero. The retraction scales as  $\mathcal{O}(\Delta t^2)$ , whereas the projection of the dynamics scales as  $\mathcal{O}(\Delta t)$ , so as  $\Delta t \to 0$ , the retraction may be ignored completely. But because we take noninfinitesimally small time-steps, we have a new type of error that originates from the high-order curvature of  $\mathcal{M}_r$ . From another viewpoint, from [68], we have that  $\mathcal{M}_r$  is a  $C^{\infty}$  smooth embedded submanifold of  $\mathbb{R}^{m \times n}$ . So as  $\Delta t \to 0$ ,  $\mathcal{M}_r \to \mathcal{P}_{T_X \mathcal{M}_r}$ , and hence  $\varepsilon_{pr} \to 0$ . In words,



FIG. 2. In a low-dimensional space, Figure 2(a) illustrates how a retraction maps a matrix, projected onto the affine tangent space at  $X_n$ , to the low-rank manifold. We show two possible retractions to points  $X_{n+1}$  and  $X_{n+1}^*$ , both of order  $\mathcal{O}(\Delta t^2)$ , which arise due to the high-order curvature of the low-rank manifold. These retractions are particular realizations of the  $\mathcal{O}(\Delta t^2)$ term in (4.4), and their difference defines the projection-retraction error. We do not show the error that may have already accumulated up to time  $t_n$ , which may give an inaccurate  $X_n$  acting as the initial condition in (4.3). The dynamical closure error (Definition 3.10), the spatial discretization error, and numerical integration error in (4.1) are built into our illustration of  $\overline{\mathscr{L}}$ . Figure 2(b) depicts the projection onto the low-rank manifold of the full-rank integral of the dynamics over  $\Delta t$ , *i.e.*,  $\mathscr{M}_r(X_n + \Delta t \overline{\mathscr{L}})$ . Note how the residual—equivalent to the normal closure error  $\varepsilon_N$  assuming no prior error accumulation—is orthogonal to the affine tangent space of the low-rank manifold at the retracted point.

Downloaded 05/06/23 to 24.61.23.148 by Pierre Lermusiaux (pierrel@mit.edu). Redistribution subject to SIAM license or copyright; see https://epubs.siam.org/terms-privacy

the low-rank manifold is arbitrarily well approximated by its affine tangent space as  $\Delta t \rightarrow 0$ . However, its high-order curvature induces error for  $\Delta t > 0$ .

We close this section with a summary of possible errors throughout the time integration. First is the dynamical model closure error,  $\varepsilon_D$ . Second is numerical error from both spatially discretizing the differential operator  $\mathscr{L}_c$  and approximating  $\overline{\mathscr{L}} \approx \frac{1}{\Delta t} \int \mathscr{L} dt$  with a classical numerical integration scheme. Third is the projection-retraction error  $\varepsilon_{pr}$ . The dynamical model closure is almost always inevitable but may be alleviated by increasing the rank of the solution. The second error may be reduced by using more accurate discretization and time-integration schemes. The third error is what we focus on correcting in a computationally efficient manner. Note that the normal closure error  $\varepsilon_N$  and Dirac–Frenkel model closure error  $\varepsilon_{DF}$  (see Definitions 3.9–3.11) are less relevant to our analysis since we are concerned with convergence to the best low-rank approximation rather than convergence to the full-rank solution.

5. Perturbative retractions. In this section, we develop the first main result, retractions that converge to  $\mathscr{P}_{\mathcal{M}_r}$  with high order. To begin, we analyze how the classic DO equations (2.3) were derived. The DO equations orthogonally project the system's dynamics onto the tangent space of  $\mathscr{M}_r$  defined at the current point. That is, we seek  $\overline{U}$  and  $\overline{Z}$  such that the residual between the full-rank direction,  $\overline{\mathscr{Q}}$ , and its tangent space projection,  $\overline{U}Z^T + U\overline{Z}^T$ , is minimized. This condition may be expressed two ways mathematically. First, as in [21], we seek  $\overline{U} \in \mathcal{U}_{m,r}, \overline{Z} \in \mathbb{R}^{n \times r}$  with  $U \in \mathcal{V}_{m,r}, Z \in \mathbb{R}^{n \times r}_*$  such that

$$\overline{U}, \overline{Z} = \underset{\substack{\tilde{U} \in \mathcal{U}_{m,r}, \\ \tilde{Z} \in \mathbb{R}^{n \times r}}{\arg \min} \| \tilde{U}Z^T + U\tilde{Z}^T - \overline{\mathscr{L}} \|.$$

Equivalently, as in [34], we seek  $\overline{U} \in \mathcal{U}_{m,r}$  and  $\overline{Z} \in \mathbb{R}^{n \times r}$  such that

$$\langle \vec{U}Z^T + U\vec{Z}^T - \overline{\mathscr{L}}, \delta_U Z^T + U\delta_Z^T \rangle = 0$$

for all  $\delta_U \in \mathcal{U}_{m,r}$  and  $\delta_Z \in \mathbb{R}^{n \times r}$ , where we use the Frobenius inner product above. The latter formulation is a Galerkin condition insisting that the residual is orthogonal to every possible matrix in the affine tangent space at  $X = UZ^T$ .

We will continue with the latter formulation. There are two issues if we want highorder convergence to  $\mathscr{P}_{\mathscr{M}_r}$ . First, the residual does not include the  $\mathcal{O}(\Delta t^2)$  retraction term; it is completely ignored. Second, the residual should be orthogonal to the affine tangent space at the new, retracted point rather than the affine tangent space at the original point. For a graphical comparison, contrast Figures 2(a) and (b). These two corrections induce the following novel formulation for the exact projection operator after an  $\mathcal{O}(\Delta t)$  perturbation:

(5.1) 
$$\langle \overline{U}Z^T + U\overline{Z}^T + \Delta t\overline{U}\overline{Z}^T - \overline{\mathscr{L}}, (U + \Delta t\overline{U})\delta_Z^T + \delta_U (Z + \Delta t\overline{Z})^T \rangle = 0.$$

This new condition (5.1) yields the following coupled, nonlinear matrix equations (see supplementary material section SM2 for a derivation):

(5.2) 
$$(U + \Delta t \overline{U})^T \Delta t \overline{\mathscr{L}} = \Delta t^2 \overline{U}^T \overline{U} (Z + \Delta t \overline{Z})^T + \Delta t \overline{Z}^T,$$
$$\mathscr{P}_U^{\perp} \Delta t \overline{\mathscr{L}} (Z + \Delta t \overline{Z}) = \Delta t \overline{U} (Z + \Delta t \overline{Z})^T (Z + \Delta t \overline{Z}).$$

## A882 AARON CHAROUS AND PIERRE F. J. LERMUSIAUX

In contrast to the classic DO equations (2.3), the new projective DO equations (5.2) are nonlinear and coupled. Of course, (5.2) may be solved via the truncated SVD, but this would be extremely costly. In principle, one could solve (5.2) using a generic iterative method such as Newton–Raphson, but this would require a large matrix inversion at each step. Alternatively, an alternating least squares approach, which is proposed for projected implicit methods in [33], may be considered. However, it is unclear how many iterations would be needed to preserve high-order convergence, and we will see that applying perturbation theory yields linear equations that guarantee high-order convergence in a fixed number of steps.

THEOREM 5.1 (perturbative retractions). Given  $U \in \mathcal{V}_{m,r}$ ,  $Z \in \mathbb{R}^{n \times r}_*$ ,  $\overline{\mathscr{Q}} \in \mathbb{R}^{m \times n}$ , and  $\Delta t \in \mathbb{R}$ , the nth-order solutions to (5.2) for  $n = 1, \ldots, 4$  are given by

(5.3) 
$$\Delta t \overline{U}^{(n)} = \sum_{i=1}^{n} \Delta t^{i} \dot{u}_{i}, \quad \Delta t \overline{Z}^{(n)} = \sum_{i=1}^{n} \Delta t^{i} \dot{\mathcal{Z}}_{i}$$

where  $\{\dot{\mathcal{U}}_i\}_{i=1}^n, \{\dot{\mathcal{Z}}_i\}_{i=1}^n$ , are given by the following linear equations. Consequently,  $(U + \Delta t \overline{U}^{(n)})(Z + \Delta t \overline{Z}^{(n)})^T$  converges as  $\mathcal{O}(\Delta t^{n+1})$  to  $\mathscr{P}_{\mathscr{M}_r}(UZ^T + \Delta t \overline{\mathscr{L}})$ .

$$\begin{cases} \dot{u}_{1} = \mathscr{P}_{U}^{\perp} \overline{\mathscr{L}} Z \left( Z^{T} Z \right)^{-1}, \\ \dot{z}_{1} = \overline{\mathscr{L}}^{T} U. \\ \end{cases}$$

$$\begin{cases} \dot{u}_{2} = \left[ \mathscr{P}_{U}^{\perp} \overline{\mathscr{L}} \dot{z}_{1} - \dot{u}_{1} (Z^{T} \dot{z}_{1} + \dot{z}_{1}^{T} Z) \right] (Z^{T} Z)^{-1}, \\ \dot{z}_{2} = \left( \overline{\mathscr{L}}^{T} - Z \dot{u}_{1}^{T} \right) \dot{u}_{1}. \\ \end{cases}$$

$$\begin{cases} \dot{u}_{3} = \left[ \mathscr{P}_{U}^{\perp} \overline{\mathscr{L}} \dot{z}_{2} - \dot{u}_{2} (Z^{T} \dot{z}_{1} + \dot{z}_{1}^{T} Z) - \dot{u}_{1} (Z^{T} \dot{z}_{2} + \dot{z}_{2}^{T} Z + \dot{z}_{1}^{T} \dot{z}_{1}) \right] (Z^{T} Z)^{-1}, \\ \dot{z}_{3} = \overline{\mathscr{L}}^{T} \dot{u}_{2} - Z (\dot{u}_{1}^{T} \dot{u}_{2} + \dot{u}_{2}^{T} \dot{u}_{1}) - \dot{z}_{1} \dot{u}_{1}^{T} \dot{u}_{1}. \\ \end{cases}$$

$$\begin{cases} \dot{u}_{4} = \left[ \mathscr{P}_{U}^{\perp} \overline{\mathscr{L}} \dot{z}_{3} - \dot{u}_{3} \left( Z^{T} \dot{z}_{1} + \dot{z}_{1}^{T} Z \right) - \dot{u}_{2} \left( Z^{T} \dot{z}_{2} + \dot{z}_{2}^{T} Z + \dot{z}_{1}^{T} \dot{z}_{1} \right) - \dot{u}_{1} \left( Z^{T} \dot{z}_{2} + \dot{z}_{1}^{T} \dot{z}_{1} \right) \right) \\ - \dot{u}_{1} \left( Z^{T} \dot{z}_{3} + \dot{z}_{3}^{T} Z + \dot{z}_{2}^{T} \dot{z}_{1} + \dot{z}_{1}^{T} \dot{z}_{2} \right) \right] (Z^{T} Z)^{-1}, \\ \dot{z}_{4} = \overline{\mathscr{L}}^{T} \dot{u}_{3} - Z \left( \dot{u}_{1}^{T} \dot{u}_{3} + \dot{u}_{2}^{T} \dot{u}_{2} + \dot{u}_{3}^{T} \dot{u}_{1} \right) \\ - \dot{z}_{1} \left( \dot{u}_{2}^{T} \dot{u}_{1} + \dot{u}_{1} \dot{u}_{2} \right) - \dot{z}_{2} \dot{u}_{1}^{T} \dot{u}_{1}. \end{cases}$$

*Proof.* The proof is completed by substituting (5.3) into (5.2) and grouping terms by  $\Delta t$ . We relegate the details to supplementary material section SM3.

Here, we note a few interesting properties of the retractions. First, we appreciate how auspicious it is that  $\dot{u}_i$  and  $\dot{z}_i$  can be solved for explicitly, due to the decoupling of (5.2) via the DO condition. Second, the only nonlinear terms appear as lower degrees of  $\dot{u}_i$  and  $\dot{z}_i$ , giving linear equations at each step. Third, the classic DO solution is naturally recovered in the first-order perturbative retraction. Fourth,  $\vec{U} \in \mathcal{U}_{m,r}$ because each  $\dot{u}_i \in \mathcal{U}_{m,r}$  due to the  $\mathscr{P}_U^{\perp}$  operator acting on any new terms. Note that related work on perturbations to the truncated SVD has been done concurrently [70], but we note significant differences in the supplementary material section SM6.

The retractions also admit a nice geometric interpretation. The second-order perturbation builds on the first, inducing a quadratic  $\overline{\mathscr{L}}$  term. Similarly, the third-order perturbation includes a cubic  $\overline{\mathscr{L}}$  term, and so on. So, one could interpret the *n*thorder perturbative retraction as implicitly projecting  $\overline{\mathscr{L}}$  onto an *n*th-order polynomial

approximation of  $\mathcal{M}_r$ . That is, the first-order perturbation projects onto the (linear) tangent space, the second-order perturbation onto a quadratic approximation of  $\mathcal{M}_r$ , the third-order perturbation onto a cubic approximation of  $\mathcal{M}_r$ , etc. As such, we are encoding high-order curvature information into our time-integration scheme.

The power of these retractions is that we can obtain an arbitrarily high order of convergence at a relatively low cost. Though we provide up to fourth-order corrections, we may calculate as high-order terms of the perturbation series as we desire. For each additional term of the perturbation series we compute, we obtain an additional order of accuracy, and once we fix the number of terms, p, in the perturbation series to compute, the given perturbative retraction will preserve the convergence of a pth-order time integrator (assuming  $\varepsilon_N$ ,  $\varepsilon_D$ , and  $\varepsilon_{DF}$  are negligible) when considering the global error.

It is direct to show that this method converges linearly in the iterative sense. The rate of convergence,  $q^*$ , for a sequence  $\{a_k\}$  that converges to a is defined as [26, 59]

$$q^* = \sup_{q} \left\{ q : \lim_{k \to \infty} \frac{|a_{k+1} - a|}{|a_k - a|^q} = 0 \right\}.$$

Letting *a* be the exact projection (4.3) and  $a_k$  the *k*th-order perturbative retraction, we have that,  $\frac{|a_{k+1}-a|}{|a_k-a|^q}\Big|_{q=1} = \mathcal{O}(\Delta x)$ , indicating linear convergence.

Each perturbative retraction is quite cheap (only requiring the inversion of an  $r \times r$  matrix), but the cost grows as the correction order increases. We consider the case where  $\overline{\mathscr{S}}$  is of rank  $r_L$  and, for simplicity, m = n. The computational complexity of our new retractions as well as common techniques such as the full SVD, the truncated SVD (TSVD) as in [2, 69], the randomized SVD (RSVD) [24], and the projector-splitting integrator [51] are provided in Table 1. These retractions are most efficient when  $r_L \ll m$ , but even in the worst case where  $\overline{\mathscr{S}}$  is full-rank, i.e.,  $r_L = m$ , the perturbative retractions are still only quadratic in m, meaning they are cheap when compared to a full-rank simulation.

We also provide the results of a timing study comparing these algorithms with reorthonormalization (aside from the full SVD since it offers no advantage over the truncated SVD) in Table 2. We apply a retraction  $\mathcal{R}_X(\Delta t \overline{\mathscr{L}})$ , where  $m = n = 10^4$  and  $\Delta t = 0.25$ . In the first case, we employ r = 10,  $r_L = 100$ , and in the second, r = 25,  $r_L = 500$ . The matrix  $X = UZ^T$  is formed by randomly sampling U and Z from a normal distribution, and then U is orthonormalized.  $\overline{\mathscr{L}}$  is formed the same way with-

TABLE 1 Asymptotic complexity of preexisting and our new retractions.

SVD	TSVD	RSVD	ProjSplit.	Perturbative
$\mathcal{O}(m^3)$	$\mathcal{O}(m(r+r_L)^2)$	$\mathcal{O}(mr(r_L+r))$	$\mathcal{O}(mr(r+r_L))$	$\mathcal{O}(mr(r+r_L)+r^3)$

T/

Wall-clock time (in milliseconds) of different retractions using the timeit function of MATLAB.

	TSVD	RSVD 1 Iter.	RSVD 2 Iters.	ProjSplit
$r = 10, r_L = 100$	23.14	12.83	17.19	5.22
$r = 25, r_L = 500$	325.27	56.38	74.12	21.50
	1 <sup>st</sup> Pert.	2 <sup>nd</sup> Pert.	3 <sup>rd</sup> Pert.	4 <sup>th</sup> Pert.
$r = 10, r_L = 100$	6.77	7.56	10.66	13.06
$r = 25, r_L = 500$	26.97	32.30	44.47	60.91

out orthonormalization, and both X and  $\overline{\mathscr{L}}$  are normalized to have Frobenius norm one. We implement the stabilized randomized subspace iteration of [24, Algorithm 4.4] with an oversampling parameter of 5. By changing the number of iterations in the randomized SVD, we can directly compare to the high-order perturbative retractions, where one iteration is analogous to our second-order perturbative retraction, and two iterations are analogous to our fourth-order perturbative retraction. The projector-splitting integrator is most similar to our first-order perturbative retraction. We compare errors in section 7. We see that for large  $r_L$ , the truncated SVD becomes prohibitively expensive due to its quadratic dependence on  $r_L$ . Compared to the randomized SVD, the perturbative retractions seem to be slightly faster, and compared to the projector-splitting integrator, the first-order perturbative retraction is slightly slower but still competitive.

There are, however, two issues with the perturbative retractions. First, we would like  $U^{n+1} \in \mathcal{V}_{m,r}$  given  $U^n \in \mathcal{V}_{m,r}$ . In the continuous-time case, the DO condition  $\dot{U}^T U = 0$  ensures that  $U(t) \in \mathcal{V}_{m,r}$  for all time. However, in the discrete-time case, the DO condition  $\overline{U}^T U = 0$  is only first-order accurate, as shown below.

$$U_{n+1}^T U_{n+1} = (U_n + \Delta t \vec{U})^T (U_n + \Delta t \vec{U}) = I + \Delta t^2 \vec{U}^T \vec{U}.$$

This issue is addressed in [47], which involves a very cheap reorthonormalization procedure. An algorithm is given in supplementary material section SM4.

Second, the perturbation series may not always converge, causing overshoots. In analogy to Taylor series, if large time-steps are taken, lower-order truncations can be more accurate than those that include higher-order corrections. This is demonstrated in subsection 7.1. To ensure convergence of the series (5.3), we define a nondimensional hyperparameter  $\varepsilon \ll 1$  and insist on a necessary condition for convergence below.

$$\max\left(\Delta t^{i}\frac{\|\dot{u}_{i}\|}{\|U\|},\Delta t^{i}\frac{\|\dot{\mathcal{Z}}_{i}\|}{\|Z\|}\right)<\varepsilon.$$

If this condition is violated for any *i*, the *i*th and higher-order corrections are not used. This gives the adaptive perturbative algorithm 5.1, letting  $\beta$  denote the maximum order correction allowed for a given  $\Delta t$ . We note that the truncated and randomized SVD are not susceptible to overshoot.

6. Dynamically orthogonal Runge–Kutta schemes. We now develop a new class of integration schemes for the dynamical low-rank approximation, building off of our perturbative retractions. In the previous section, we treated the dynamics of the system  $\overline{\mathscr{L}}$  as a fixed matrix while retracting, meaning we would first integrate the dynamics in the full embedding Euclidean space before retracting back to the low-rank manifold. Obtaining a high-order integration scheme in the full space, however, is often costly as the rank may grow quickly, increasing the cost of the retraction and function evaluation. In this section, we show how to integrate along the nonlinear manifold of fixed rank, entirely avoiding integrating into the full embedding space. In doing so, we capture the dynamic effects of manifold's high-order curvature and drastically reduce the computational cost of integration. We refer to these novel integration methods as dynamically orthogonal Runge–Kutta (DORK) schemes.

First, let us be precise about the system we seek to integrate. As in section 5, we consider schemes in the form of (4.7). But now, we specify that our dynamics are

 $\begin{array}{l} \textbf{Algorithm 5.1. Adaptive perturbative retraction.} \\ \hline \textbf{Input: } U_0 \in \mathcal{V}_{m,r}, Z_0 \in \mathbb{R}_*^{n \times r}, \overline{\mathscr{L}} \in \mathbb{R}^{m \times n}, \Delta t \in \mathbb{R}, \varepsilon \in \mathbb{R}, \beta \in \mathbb{N} \\ \textbf{Output: } U_1 \in \mathcal{V}_{m,r}, Z_1 \in \mathbb{R}^{n \times r} \\ 1: U_1 = U_0, \quad Z_1 = Z_0, \quad \dot{u}_0 = 0, \quad \dot{Z}_0 = 0, \quad \alpha_0 = \|Z_0\|, \quad \alpha = 0, \quad i = 1 \\ 2: \textbf{ while } i \leq \beta \textbf{ do} \\ 3: \quad \text{Compute } \dot{\mathcal{U}}_i \text{ and } \dot{\mathcal{Z}}_i \text{ from (5.3).} \\ 4: \quad \alpha = \frac{\Delta t^i}{\alpha_0} \max\left(\|\dot{\mathcal{U}}_i\|, \|\dot{\mathcal{Z}}_i\|\right) \\ 5: \quad \textbf{ if } \alpha > \varepsilon \textbf{ then} \end{array}$ 

6: Break. Go to line 11.

7: end if

8:

9:  $i \leftarrow i+1$ 10: **end while** 

integrated along the low-rank manifold, i.e.,

11:  $U_1, Z_1 \leftarrow \texttt{reorthonormalize}(U_1, Z_1)$ 

 $U_1 \leftarrow U_1 + \Delta t^i \dot{\mathcal{U}}_i, \qquad Z_1 \leftarrow Z_1 + \Delta t^i \dot{\mathcal{Z}}_i$ 

(6.1) 
$$\overline{\mathscr{L}} = \frac{1}{\Delta t} \int_{t_n}^{t_{n+1}} \mathscr{L}(\mathscr{P}_{\mathscr{M}_r}(X(t)), t; \omega) dt + \mathcal{O}(\Delta t^k).$$

At discrete times,  $X(t_n) \in \mathcal{M}_r$  by definition; it is our low-rank approximation of the system. However, when integrating between those times, X(t) may depart the low-rank manifold. In section 5,  $\overline{\mathscr{L}}$  was implicitly defined as in (6.1), except without the  $\mathscr{M}_r$  operator, hence letting the argument of  $\mathscr{L}$  depart the manifold; both are valid integration schemes and are just alternative approximations related to the dynamical model closure error (see Definition 3.10).

The key step in constructing the DORK schemes is to define  $\overline{\mathscr{L}}$  as a perturbation series itself, i.e.,

(6.2) 
$$\Delta t \overline{\mathscr{B}} \equiv \sum_{i=1}^{k} \Delta t^{i} \overline{\mathcal{L}}_{i} \,.$$

Observe that the partial sum  $\overline{\mathscr{Q}}^{(j)} \equiv \sum_{i=1}^{j} \Delta t^{i-1} \overline{\mathcal{L}}_i$  is a *j*th-order integration scheme. Obtaining this perturbation series for an arbitrary integration scheme is always feasible. As a concrete example, consider a second-order integration scheme. We may write  $\overline{\mathscr{Q}}^{(1)} = \overline{\mathcal{L}}_1$  as any first-order integration scheme for the same system, and  $\overline{\mathcal{L}}_2 = (\overline{\mathscr{Q}}^{(2)} - \overline{\mathscr{Q}}^{(1)})/\Delta t$ . One can proceed similarly for higher-order schemes, and we may generally write  $\overline{\mathcal{L}}_i = (\overline{\mathscr{Q}}^{(i)} - \overline{\mathscr{Q}}^{(i-1)})/\Delta t$ , and  $\overline{\mathscr{Q}}^{(0)} = 0$ . Similar to classic Runge–Kutta schemes, the choice of  $\{\overline{\mathcal{L}}_i\}_i$  is here not unique due to the plethora of choices of intermediate integration schemes for the partial sums  $\overline{\mathscr{Q}}^{(j)}$ . Furthermore, constructing a computationally efficient DORK scheme requires thought to minimize the number of function evaluations and retractions in the scheme.

To obtain DORK integration schemes, we thus solve the projective DO equations (5.2) assuming a perturbation series in  $\overline{U}, \overline{Z}$ , and now also in  $\overline{\mathscr{L}}$ . This yields the following schemes.

$$\begin{aligned} & \left\{ \begin{array}{l} \dot{u}_{1} = \mathscr{P}_{U}^{\perp} \overline{\mathcal{L}}_{1} Z \left( Z^{T} Z \right)^{-1} \\ \dot{z}_{1} = \overline{\mathcal{L}}_{1}^{T} U \\ \\ & \left\{ \begin{array}{l} \dot{u}_{2} = \left[ \mathscr{P}_{U}^{\perp} \left( \overline{\mathcal{L}}_{2} Z + \overline{\mathcal{L}}_{1} \dot{z}_{1} \right) - \dot{u}_{1} (Z^{T} \dot{z}_{1} + \dot{z}_{1}^{T} Z) \right] (Z^{T} Z)^{-1} \\ \dot{z}_{2} = \overline{\mathcal{L}}_{2}^{T} U + \left( \overline{\mathcal{L}}_{1}^{T} - Z \dot{u}_{1}^{T} \right) \dot{u}_{1} \\ \\ & \left\{ \begin{array}{l} \dot{u}_{3} = \left[ \mathscr{P}_{U}^{\perp} \left( \overline{\mathcal{L}}_{3} Z + \overline{\mathcal{L}}_{2} \dot{z}_{1} + \overline{\mathcal{L}}_{1} \dot{z}_{2} \right) - \dot{u}_{2} (Z^{T} \dot{z}_{1} + \dot{z}_{1}^{T} Z) \\ & - \dot{u}_{1} (Z^{T} \dot{z}_{2} + \dot{z}_{2}^{T} Z + \dot{z}_{1}^{T} \dot{z}_{1}) \right] (Z^{T} Z)^{-1} \\ \dot{z}_{3} = \overline{\mathcal{L}}_{3}^{T} U + \overline{\mathcal{L}}_{2}^{T} \dot{u}_{1} + \overline{\mathcal{L}}_{1}^{T} \dot{u}_{2} - Z (\dot{u}_{1}^{T} \dot{u}_{2} + \dot{u}_{2}^{T} \dot{u}_{1}) - \dot{z}_{1} \dot{u}_{1}^{T} \dot{u}_{1} \\ \\ & \left\{ \begin{array}{l} \dot{u}_{4} = \left[ \mathscr{P}_{U}^{\perp} \left( \overline{\mathcal{L}}_{4} Z + \overline{\mathcal{L}}_{3} \dot{z}_{1} + \overline{\mathcal{L}}_{2} \dot{z}_{2} + \overline{\mathcal{L}}_{1} \dot{z}_{3} \right) \\ & - \dot{u}_{3} \left( Z^{T} \dot{z}_{1} + \dot{z}_{1}^{T} Z \right) - \dot{u}_{2} \left( Z^{T} \dot{z}_{2} + \dot{z}_{2}^{T} Z + \dot{z}_{1}^{T} \dot{z}_{1} \right) \\ & - \dot{u}_{3} \left( Z^{T} \dot{z}_{1} + \dot{z}_{1}^{T} Z \right) - \dot{u}_{2} \left( Z^{T} \dot{z}_{2} + \dot{z}_{2}^{T} Z + \dot{z}_{1}^{T} \dot{z}_{1} \right) \\ & - \dot{u}_{4} \left( Z^{T} \dot{z}_{3} + \dot{z}_{3}^{T} \dot{z}_{1} + \dot{z}_{2}^{T} \dot{z}_{1} + \dot{z}_{1}^{T} \dot{z}_{2} \right) \right] \left( Z^{T} Z \right)^{-1} \\ \dot{z}_{4} = \overline{\mathcal{L}}_{4}^{T} U + \overline{\mathcal{L}}_{3}^{T} \dot{u}_{1} + \overline{\mathcal{L}}_{2}^{T} \dot{u}_{2} + \dot{z}_{1}^{T} \dot{u}_{3} \\ & - Z \left( \dot{u}_{1}^{T} \dot{u}_{3} + \dot{u}_{2}^{T} \dot{u}_{2} + \dot{u}_{3}^{T} \dot{u}_{1} \right) - \dot{z}_{1} \left( \dot{u}_{2}^{T} \dot{u}_{1} + \dot{u}_{1} \dot{u}_{2} \right) - \dot{z}_{2} \dot{u}_{1}^{T} \dot{u}_{1} \end{aligned} \right\} \right\}$$

As compared to (5.4)–(5.7), we see additional terms involving the high-order corrections to  $\overline{\mathscr{L}}$ . The  $\overline{\mathcal{L}}_i$  terms are projected onto the subspaces  $\dot{\mathcal{U}}_i$ , which we interpret as updating the subspace onto which we are projecting the system's dynamics as we integrate. We also note that the adaptive retraction in Algorithm 5.1 may be used with the schemes (6.3)–(6.6) to ensure the series converges.

To further explain the DORK schemes, we go over a second-order scheme in depth. Consider Heun's method for the full-rank system starting from point  $\mathfrak{X}_0$  at time  $t_0$ .

$$\begin{split} k_1 &= \mathscr{L}(\mathfrak{X}_0, t_0; \omega), \\ \hat{\mathfrak{X}}_1 &= \mathfrak{X}_0 + \Delta t k_1, \\ k_2 &= \mathscr{L}(\hat{\mathfrak{X}}_1, t_0 + \Delta t; \omega), \\ \mathfrak{X}_1 &= \mathfrak{X}_0 + \frac{\Delta t}{2} \left( k_1 + k_2 \right). \end{split}$$

For our low-rank system with  $\overline{\mathscr{L}}$  defined as in (6.1), Heun's method is as follows:

$$\begin{split} k_1 &= \mathscr{L}(X_0, t_0; \omega), \\ \hat{X}_1 &= \mathscr{P}_{\mathscr{M}_r} \left( X_0 + \Delta t k_1 \right), \\ k_2 &= \mathscr{L}(\hat{X}_1, t_0 + \Delta t; \omega), \\ X_1 &= \mathscr{P}_{\mathscr{M}_r} \left( X_0 + \frac{\Delta t}{2} \left( k_1 + k_2 \right) \right). \end{split}$$

Because Heun's method is only second order, we can approximate  $\mathscr{P}_{\mathscr{M}_r}$  with a retraction of first order when computing  $\hat{X}_1$  (since the resulting error will be  $\mathcal{O}(\Delta t^3)$  once  $k_2$  is multiplied by  $\Delta t$ ) and a retraction of second order when computing  $X_1$ .

If we stop our analysis here and project  $k_1$  and  $k_2$  onto their respective tangent spaces, this algorithm appears as the projected Runge–Kutta scheme from [33]. However, that scheme does not update the subspace U as it integrates, and we show

Algorithm 6.1. Second-order DORK scheme.

 $\begin{array}{l} \overline{\mathbf{Input: } U_0 \in \mathcal{V}_{m,r}, \, Z_0 \in \mathbb{R}^{n \times r}_*, \, \mathscr{L}(X,t;\omega) : \mathbb{R}^{m \times n} \times \mathbb{R} \times \Omega \to \mathbb{R}^{m \times n}, t_0 \in \mathbb{R}, \, \omega \in \Omega, \\ \Delta t \in \mathbb{R} \\ \mathbf{Output: } U_1 \in \mathcal{V}_{m,r}, \, Z_1 \in \mathbb{R}^{n \times r} \\ 1: \, k_1 = \overline{\mathcal{L}}_1 = \mathscr{L}(U_0 Z_0^T, t_0; \omega) \\ 2: \, \hat{U}_1 \hat{Z}_1^T = \mathcal{R}_{U_0 Z_0^T}(\Delta t k_1) \text{ with retraction defined by (5.4)} \\ 3: \, k_2 = \mathscr{L}(\hat{U}_1 \hat{Z}_1^T, t_0 + \Delta t; \omega) \\ 4: \, \overline{\mathcal{L}}_2 = \frac{1}{2\Delta t} \, (k_2 - k_1) \\ 5: \, U_1, Z_1 = \mathcal{R}_{U_0 Z_0^T}(\Delta t \overline{\mathcal{L}}_1 + \Delta t^2 \overline{\mathcal{L}}_2) \text{ with retraction defined by (6.3)-(6.4)} \\ 6: \, U_1, Z_1 \leftarrow \text{reorthonormalize}(U_1, Z_1) \end{array}$ 

in section 7 that using the DORK schemes greatly reduces its error. Furthermore, the projected Runge–Kutta schemes project the dynamics onto the tangent space of the low-rank manifold in each substep to reduce the cost of integration; otherwise, the rank of the DLRA would grow quickly before applying a retraction as the final step, making such as scheme computationally expensive. In the DORK schemes, an additional projection of the dynamics onto the tangent space is not required as the projection is inherently built into the method. The DORK and projected Runge–Kutta schemes are compared in the experiments of subsection 7.2.

To proceed, we let  $\overline{\mathcal{L}}_1 = k_1$ . Then, we have that  $\overline{\mathscr{L}}^{(2)} = \frac{1}{2}(k_1 + k_2)$ , and so  $\overline{\mathcal{L}}_2 = \frac{1}{2\Delta t}(k_2 - k_1)$ . We may now apply (6.3)–(6.4) and then reorthonormalize. This scheme is succinctly written in Algorithm 6.1, and we provide third- and fourth-order algorithms SM8.1 and SM8.2 derived from embedded (2,3) and (3,4) Runge–Kutta–Fehlberg schemes [19, 65] in the supplementary material section SM8.

## 7. Numerical experiments.

7.1. Matrix addition. We begin with the simple case of adding two matrices, X and  $\Delta tL$ , and then retracting back onto the low-rank manifold. Though simple, this test case demonstrates high-order local error convergence, and low-rank matrix addition is a building block of time-integration schemes. X and L are normalized such that they have Frobenius norm one, so  $\Delta t$  controls their relative scaling. In particular,  $X, L \in \mathbb{R}^{10,000 \times 10,000}$ ,  $\operatorname{rank}(X) = 10$ ,  $X = UZ^T$ , and the entries of U, Z are independently chosen from a standard normal distribution, after which U is orthonormalized via the orth function of MATLAB (which uses the left singular vectors of a given matrix).  $L = L_U L_Z^T$  has rank 100 and the entries of  $L_U$  and  $L_Z$  are sampled independently from a standard normal distribution.

We analyze the Frobenius norm of the projection-retraction error,  $\varepsilon_{pr}$  (4.8), the difference between  $\mathcal{R}_X(\Delta tL)$  and  $\mathscr{P}_{\mathcal{M}_r}(X + \Delta tL)$ , which examines the convergence to the best low-rank approximation. Figure 3(a) depicts this error for different order, and we see  $\mathcal{O}(\Delta t^{n+1})$  convergence to the best approximation for the *n*th-order perturbative retraction. For large  $\Delta t$ , we see the perturbative retractions with higher-order corrections overshoot, but our adaptive method with hyperparameter  $\varepsilon = 0.1$  rectifies this. In Figure 3(b), we compare perturbative retractions with the randomized SVD (see section 5, where we compare algorithm runtimes for details on implementation) and the projector-splitting integrator. The randomized SVD with one iteration converges at third order, as does the second-order perturbative retraction, but we see our second-order perturbative retraction is more accurate by about an order of magnitude. We see similar behavior for the two-iteration randomized SVD and our fourth-order perturbative retraction. The projector-splitting integrator, as expected,



FIG. 3. The perturbative retractions exhibit high-order convergence to the best approximation. The adaptive algorithm avoids overshoot. Our perturbative retractions are competitive with the projector-splitting integrator and randomized SVD.

converges at second-order and has nearly the same projection-retraction error as our first-order perturbative retraction.

7.2. Matrix differential equations. In this example, we consider global error convergence, so we expect one order lower convergence than in subsection 7.1. That is, we expect to see  $\mathcal{O}(\Delta t^n)$  convergence for an *n*th-order perturbative retraction. We investigate systems of coupled linear oscillators  $\mathfrak{X} \in \mathbb{R}^{26 \times 26}$  with differential equation  $\ddot{\mathfrak{X}} = -\Omega^2 \mathfrak{X}$  and initial conditions  $\mathfrak{X}(0) = \mathfrak{X}_0$ ,  $\dot{\mathfrak{X}}(0) = \dot{\mathfrak{X}}_0$ . We choose  $\Omega = \text{diag}(\omega_1, \omega_1, \omega_2, \omega_2, \dots, \omega_{13}, \omega_{13})$  with each  $\omega_i$  independently chosen from a standard normal distribution. Letting  $R(t) \in \mathbb{R}^{26 \times 26}$  be a tridiagonal filled with  $2 \times 2$  block rotation matrices  $R_1, \dots, R_{13}$  such that

$$R_i(t) = \begin{bmatrix} \cos(\omega_i t) & -\sin(\omega_i t) \\ \sin(\omega_i t) & \cos(\omega_i t) \end{bmatrix},$$

then  $R(t) = \text{diag}(R_1, \dots, R_{13})$ .  $Q \in \mathbb{R}^{26 \times 26}$  is formed by orthonormalizing a matrix with independent samples from a uniform distribution. Last,  $S \in \mathbb{R}^{26 \times 26}$  is a diagonal matrix with nonincreasing entries  $S_{ii} = 100 + 10z_i$  for  $i \le 16$  and  $S_{ii} = 10^{-3 - (i-17)/9}$  for i > 16, where  $z_i$  are realizations of independent, standard normal random variables. It is simple to verify that the exact solution is given by  $\mathfrak{X}(t) = R(t)QS$  with  $\mathfrak{X}_0 = R(0)QS$ and  $\mathfrak{X}_0 = R(0)QS$ . In constructing such a solution, we ensure that the singular values of the full-rank solution are fixed by S and is well approximated by a rank-16 matrix, which we will call X(t), our low-rank approximation. We initialize X(0) and X(0)by taking the truncated SVD of  $\mathfrak{X}_0$  and  $\mathfrak{X}(0)$ . To integrate the system, we convert the system of second-order ODEs to a system of first-order ODEs by expanding our state space to  $\begin{bmatrix} X^T & \dot{X}^T \end{bmatrix}^T$ . Then, we implement two integrations schemes. First, we use a fourth-order Runge–Kutta–Fehlberg integrator [19, 65] to calculate  $\overline{\mathscr{L}}$ , and we retract back to the manifold after each time-step using our perturbative retractions (5.4)-(5.7). In fixing the order of integration, we are able to investigate the role of projection-retraction error alone. This scheme equates to integrating in the full Euclidean embedding space. Second, we use our DORK schemes (6.3)-(6.6) with

Copyright (c) by SIAM. Unauthorized reproduction of this article is prohibited.

### DORK SCHEMES WITH PERTURBATIVE RETRACTIONS



FIG. 4. Here we show the error from numerically integrating systems of linear oscillators. The DORK schemes are more accurate than the projected Runge–Kutta except at first order since we always use a fourth-order integrator for the projected Runge–Kutta schemes. In contrast, the DORK integration order efficiently adapts with the order of retraction.

TABLE 3
Normalized errors of the second-order projected Runge–Kutta schemes with various retractions
(TSVD, RSVD 2 Iters., and ProjSplit.) and of the second-order DORK scheme. All integration
schemes except for DORK have the same error because the dominating error arises from departing
the manifold and reprojecting, which is what DORK corrects.

Tunen 9

$N_T$	TSVD	RSVD 2 Iters.	ProjSplit.	2 <sup>nd</sup> -Order DORK
50	$2.96 \cdot 10^{-2}$	$2.96 \cdot 10^{-2}$	$2.96 \cdot 10^{-2}$	$2.64 \cdot 10^{-2}$
134	$4.00 \cdot 10^{-3}$	$4.00 \cdot 10^{-3}$	$4.00 \cdot 10^{-3}$	$3.59 \cdot 10^{-3}$
968	$7.58 \cdot 10^{-5}$	$7.58 \cdot 10^{-5}$	$7.58 \cdot 10^{-5}$	$6.79 \cdot 10^{-5}$

Algorithms 6.1, SM8.1, and SM8.2. In both cases, we compare X(t) with  $\mathfrak{X}(t)$  in Figure 4 using the Frobenius norm normalized by  $||X_0||$ .

Figure 4(a) shows that the perturbative retractions with high-order corrections perform orders of magnitude better in reducing error accumulation. Furthermore, we see that the DORK schemes are markedly better than the projected Runge–Kutta schemes even though the DORK schemes use a Runge–Kutta scheme of the same order as the retraction, e.g., a second-order DORK scheme uses a second-order Runge–Kutta scheme. However, the projected Runge–Kutta schemes are all using fourth-order Runge–Kutta schemes. This explains why the first-order projected Runge–Kutta scheme performs better than the first-order DORK scheme, but it is also remarkable that the second- and third-order DORK schemes outperform their projected Runge– Kutta counterparts that are integrating at a higher order. This indicates that the dominating error in this problem is  $\varepsilon_{pr}$ .

In Figure 4(b), we see  $\mathcal{O}(\Delta t^k)$  convergence for both the projected Runge–Kutta schemes with kth-order perturbative expansion and the kth-order DORK schemes. We also use the second-order symmetrized projector-splitting integrator as well as the 2-iteration randomized SVD and truncated SVD as retractions for a second-order projected Runge–Kutta scheme, Heun's method. In Table 3, we compare the normalized errors of the second-order projected Runge–Kutta scheme to our second-order DORK scheme as a function of the number of points in time  $N_T$  we integrate over, keeping the time interval  $t \in [0, 10]$  fixed. Because the truncated SVD, randomized SVD, and projector-splitting integrators all perform to the same degree of accuracy,

we infer that the projection-retraction error does not dominate for these schemes. The second-order DORK scheme, however, outperforms these techniques, which indicates that updating the subspace as we integrate is crucial for highly accurate integration.

7.3. Real-time data compression. Imagine live-streaming high-definition (for example, 4k) video at 60 frames per second over a low-bandwidth data stream. Realistically, data would be sent in packets allowing for greater compression, but we consider the case where we send each frame as we record it. To compress the data on the fly, ideally, we could take the truncated SVD of each frame, but this is costly at such a rapid rate. Instead, we propose taking the truncated SVD for the first frame (inducing a fixed but acceptable time lag) and then retracting on the difference between the frames to obtain the next low-rank frame. That is, we consider  $\overline{\mathscr{R}} = (\text{Frame}_i - \text{Frame}_{i-1})/\Delta t$ , and then we retract back onto the low-rank manifold. Our choice of retraction drastically affects how much our solution drifts from the best approximation due to varying magnitudes of  $\varepsilon_{pr}$ .

We note that this choice of  $\overline{\mathscr{L}}$  is not optimal and only used to show the benefits of the high-order perturbative retractions, particularly when full-rank derivative information is unavailable. In this case, when our dynamics are driven only by data, we may instead choose  $\overline{\mathscr{L}} = (\operatorname{Frame}_i - X_{i-1})/\Delta t$ , which corrects for the drift. If there were no numerical error, in the continuous-time limit, these two choices of  $\overline{\mathscr{L}}$  would yield equivalent low-rank approximations since  $\operatorname{Frame}_{i-1} - X_{i-1}$  would exist in the normal space of the  $\mathscr{M}_r$  at  $X_{i-1}$  (i.e.,  $\mathscr{P}_{\mathcal{T}_{X_{i-1}}\mathscr{M}_r}(\operatorname{Frame}_{i-1} - X_{i-1}) = 0$ ). Since such idealities are unrealistic, both cases are presented in Figure 5. We find that the retraction drifts far from the truncated SVD, whereas the adaptive retraction preserves much more detail.



(a) Truncated SVD of frame



(b) Adaptive retraction applied with  $\overline{\mathscr{L}} = (\text{Frame}_i - X_{i-1})/\Delta t$ 



(c) Retraction of first order applied with  $\overline{\mathscr{D}} = (\operatorname{Frame}_i - \operatorname{Frame}_{i-1})/\Delta t$ 



(d) Adaptive retraction applied with  $\overline{\mathscr{L}} = (\operatorname{Frame}_i - \operatorname{Frame}_{i-1})/\Delta t$ 

FIG. 5. We show the 241st frame of a 4k, 60 fps grayscale video of a peacock walking across a road in Split, Croatia. Each frame is rank-500, and different  $\overline{\mathscr{F}}$  are used. The first-order retraction drifts, whereas the adaptive retraction is much more accurate.

**7.4.** Partial differential equations. Our perturbative retractions may also be applied to multistep integration schemes. Here, we apply our methodology to a deterministic PDE, where each solution slice for a fixed time is low-rank. We consider a two-dimensional diffusion equation with imaginary diffusivity.

$$\begin{split} \frac{\partial u}{\partial t} &= \frac{i}{2k} \nabla^2 u, \quad (x,y) \in \mathcal{D} = [0,250] \times [0,300], \quad t \geq 0, \\ u|_{\partial \mathcal{D}} &= 0, \quad u(x,y,0) = u_0(x,y). \end{split}$$

We set  $k(x, y) = 240\pi/(1500+500 \exp\{-(x-125)^2/62.5^2\} \exp\{-(y-150)^2/75^2\})$ . On different length and mass scales, this is Schrödinger's equation with spatially varying mass, which has relevance in crystal impurities, semiconductor heterostructure, and more [13, 76, 12, 53, 60, 4]. Alternatively, if we let k be constant and switch t to a range variable, this would correspond to the paraxial (or parabolic) scalar wave equation used extensively in optics [38, 28, 57, 16] and acoustics [17, 64, 27]. We use Dirichlet-zero boundary conditions which correspond to an infinite potential well in the Schrödinger interpretation and a pressure-release boundary in acoustics. Note that although we developed this methodology for real matrices and Riemannian manifolds, we may easily apply the same ideas to complex matrices and Hermitian matrices; in short, the transpose operator becomes the conjugate transpose operator.

For numerical integration, we employ the explicit and unconditionally stable Dufort-Frankel finite difference scheme [23], extended to two dimensions (see supplementary material section SM9, where the initial conditions  $u_0$  are also given). We use 250 points in x and 300 points in y, and we solve from t = 0 to 500 with onesecond time intervals. The results of our adaptive Algorithm 5.1, with perturbative retractions to solve the PDE at different ranks, are shown in Figure 6. We find that



FIG. 6. Solutions of the Schrödinger's or parabolic wave PDE,  $Im\{u(x, y = 74.7508, t)\}$ , for different ranks. For all but the full-rank solution (classic integration), we use our adaptive low-rank method with perturbative retractions. We see the low-rank approximations converge quickly to the full-rank solution.

as the rank increases, more of the solution energy is captured as is more detail. Furthermore, the low-rank solutions still seemingly respect the physics of the problem, and the integrity of the solution decays gracefully as the rank decreases.

**7.5.** Stochastic differential equations. Finally, we showcase our perturbative retractions on a stochastic variant of Burgers' equation [5, 7] with periodic boundary conditions using an implicit-explicit scheme.

$$\begin{aligned} \frac{\partial u}{\partial t} + \beta(\omega)u\frac{\partial u}{\partial x} &= \nu \frac{\partial^2 u}{\partial x^2} + f(x,t;\omega), \quad x \in \mathcal{D} = [-1,1], \quad t \ge 0, \quad \omega \in \Omega, \\ u(-1,t;\omega) &= u(1,t;\omega), \quad \frac{\partial u}{\partial x} \Big|_{x=-1} &= \frac{\partial u}{\partial x} \Big|_{x=1}, \quad u(x,0;\omega) = u_0(x;\omega). \end{aligned}$$

We set  $f(x,t;\omega) = \frac{1}{100}\gamma_1(\omega)\sin(t)e^{-100\left(x+\frac{1}{2}\right)^2} + \frac{1}{100}\gamma_2(\omega)\cos(t)e^{-64\left(x-\frac{1}{2}\right)^2}$  and note that the advection speed is stochastic as are the initial conditions and the forcing



FIG. 7. The first row compares 10,000 stochastic realizations of the solution at t = 10, where all but the Monte Carlo solutions use the adaptive Algorithm 5.1 with perturbative retractions. The rank-5 solution does relatively well but misses some variability, especially at the edges. This is recovered at rank-15. Figures 7(d) and 7(e) show non-Gaussian statistics at u(x = 0) via a comparison of lowrank and Monte Carlo histograms. Figures 7(f) and 7(i) show the spatial covariance of the Monte Carlo simulations at t = 0 and t = 10. Figures 7(g) and 7(h) show that the low-rank approximations capture the spatial covariance quite well.

function. We let  $\beta \sim \gamma_1 \sim \gamma_2 \sim U(-3/4, 3/4)$ , each sampled independently, and set  $\nu = 0.01$ . In the simulation shown, we use 10,000 Monte Carlo realizations, discretize  $\mathcal{D}$  into 150 points, and solve for  $t \in [0, 10]$  with 1,001 points. To construct the initial conditions, we define a new function  $g(x) = \frac{1}{2}\operatorname{sin}(4\pi x) + \frac{3}{4}\sin(15\pi x)$ , where  $\operatorname{sin}(x) = \sin(x)/x$ . Furthermore, let DST and IDST denote the discrete and inverse discrete sine transform. After spatial discretization, we write the initial conditions as  $u_0 = \operatorname{IDST}((1+z(\omega) \oslash \left[\sqrt{1} \quad \sqrt{2} \quad \cdots \quad \sqrt{150}\right]^T) \odot \operatorname{DST}(g))$ , where  $\odot$  denotes the Hadamard product,  $\oslash$  the elementwise division, and each  $z_i \sim \mathcal{N}(0, 1)$  are independent. With such initial conditions, each realization is smooth by reducing the stochasticity in the higher spatial frequencies. Further implementation details, including the semi-implicit time finite difference scheme, are given in supplementary material section SM10. We employ our perturbative retractions with the adaptive method (Algorithm 5.1) and hyperparameter  $\varepsilon = 0.025$ .

Figure 7 shows the fast convergence of low-rank solutions to the full-rank Monte Carlo solution. The reconstructed stochastic realizations of the low-rank solution approximate the full-rank solution well, suggesting convergence in probability (e.g., [67]). Furthermore, the histograms and spatial covariances indicate convergence in distribution (though this was expected since convergence in distribution is implied by convergence in probability). From the spatial covariances at t = 0 and t = 10, we see that this stochastic process is neither spatially nor temporally stationary. Hence, the dynamical low-rank approximation is an effective and general method for uncertainty quantification in systems with complex dynamics and non-Gaussian statistics.

8. Conclusions. The dynamical low-rank approximation has proven useful in approximating solutions to a plethora of problems. In this paper, we utilized the nonintrusive version of the spatially discrete dynamically orthogonal (DO) equations. In the continuous-time setting, projecting a system's dynamics onto the tangent space of the low-rank manifold is sufficient to exactly track the best instantaneous low-rank approximation. However, in the discrete time setting, numerical integration errors occur as we employ the projective DO equations and we introduce the projection-retraction error.

We classify errors that arise from numerically approximating a low-rank dynamical system into four categories: spatial discretization error, time-integration error, closure error, and manifold curvature error. The first two we combine into  $\overline{\mathscr{Q}}$ . The closure error, involving the normal  $\varepsilon_{\mathcal{N}}$ , dynamical  $\varepsilon_D$ , and Dirac–Frenkel  $\varepsilon_{DF}$  closure errors, is somewhat inevitable. The manifold curvature error is dynamic in the sense that local curvature depends on the present state of the low-rank approximation. It is realized in the form of projection-retraction error,  $\varepsilon_{pr}$ , and we demonstrate that it may be significantly reduced by projecting a system's dynamics onto higher-order polynomial approximations to the low-rank manifold via perturbative retractions. We show that these new retractions span a dense set of matrices in the low-rank manifold (proven in SM5) and that they converge with high order to the truncated SVD. Furthermore, they are explicit and relatively computationally inexpensive when the rank of the solution approximation is much smaller than the system's dimension. We give numerical examples that show high-order convergence to the best low-rank approximation in local error and global error.

By writing integration schemes as a perturbation series, we introduce and derive novel dynamically orthogonal Runge–Kutta (DORK) schemes that account for the evolution of the reduced-order integration space during the time-step. We show that DORK schemes are (i) highly accurate by incorporating knowledge of the dynamic, nonlinear manifold's high-order curvature from stage to stage and (ii) computationally efficient by limiting the growing rank needed to represent the evolving dynamics  $\overline{\mathscr{G}}$ . In our real-time data compression example, we show that these retractions may be used when the dynamics are data-driven (rather than model-driven), and the perturbative retractions with high-order corrections drift from the best low-rank approximation at a much slower rate than retractions with low-order corrections. Note that adjusting how the dynamics are calculated to correct the drift in the first place is important to obtain the most accurate scheme. Last, we show that the retractions may be used in deterministic and stochastic differential equations: the dynamical low-rank approximation and the retractions are agnostic to the nature of the mathematical spaces we choose to compress, whether deterministic or stochastic. In both cases, the low-rank approximation converges quickly not just to the best approximation, but to the full-rank solution. In the stochastic case, this offers an efficient methodology for uncertainty quantification in nonlinear problems with non-Gaussian statistics (e.g., [13, 42, 44, 41, 40]).

Several future research directions are possible. To increase the stability of these methods, a scheme with an adaptive rank [9, 14, 75] may be adopted using metrics given in [62, 21, 46], and a pseudoinverse may be used in instances where our system is ill-conditioned. Furthermore, a more rigorous stability analysis of the retractions is necessary to obtain sufficient stability criteria similar to the results on projector-splitting methods [51, 32] in [30]. One may also develop implicit schemes (see, e.g., [33]) perhaps using inexpensive low-rank inversion methods. Generalizing this methodology to PDEs with high-order time derivatives is also possible by using a phase space representation which has been explored in the context of the wave equation in [55, 56]. One could also generalize these retractions to low-rank tensors [31, 10, 32]. Our results can be useful to other fields such as constrained optimization [77, 35, 1] and autonomy with dynamic reduced-order prediction and control [58, 37, 43, 45, 25].

Acknowledgments. We thank our MSEAS group members for their collaboration. We also thank the anonymous reviewers for their constructive feedback.

### REFERENCES

- P.-A. ABSIL, L. AMODEI, AND G. MEYER, Two Newton methods on the manifold of fixed-rank matrices endowed with Riemannian quotient geometries, Comput. Statist., 29 (2014), pp. 569–590.
- P.-A. ABSIL AND J. MALICK, Projection-like retractions on matrix manifolds, SIAM J. Optim., 22 (2012), pp. 135–158, https://doi.org/10.1137/100802529.
- [3] P.-A. ABSIL AND I. V. OSELEDETS, Low-rank retractions: A survey and new results, Comput. Optim. Appl., 62 (2015), pp. 5–29.
- [4] G. BASTARD, Wave Mechanics Applied to Semiconductor Heterostructures, Les Éditions de Physique Les Ulis, Wiley, New York, 1988.
- [5] H. BATEMAN, Some recent research on the motion of fluids, Monthly Weather Rev., 43 (1915), pp. 163–170, https://doi.org/10.1175/1520-0493(1915)43(163:SRROTM)2.0.CO;2.
- [6] G. BERKOOZ, P. HOLMES, AND J. L. LUMLEY, The proper orthogonal decomposition in the analysis of turbulent flows, Annu. Rev. Fluid Mech., 25 (1993), pp. 539–575.
- J. BURGERS, A Mathematical Model Illustrating the Theory of Turbulence, Adv. Appl. Mech. 1, Elsevier, New York, 1948, pp. 171–199, https://doi.org/10.1016/S0065-2156(08)70100-5.
- [8] R. H. CAMERON AND W. T. MARTIN, The orthogonal development of non-linear functionals in series of Fourier-Hermite functionals, Ann. of Math., 48 (1947), pp. 385–392.
- G. CERUTI, J. KUSCH, AND C. LUBICH, A Rank-Adaptive Robust Integrator for Dynamical Low-Rank Approximation, https://arxiv.org/abs/2104.05247, 2021.
- [10] G. CERUTI AND C. LUBICH, Time integration of symmetric and anti-symmetric low-rank matrices and Tucker tensors, BIT, 60 (2020), pp. 1–24.

- [11] A. CHAROUS, High-Order Retractions for Reduced-Order Modeling and Uncertainty Quantification, Master's thesis, Massachusetts Institute of Technology, Center for Computational Science and Engineering, Cambridge, MA, 2021.
- [12] G. CHEN AND Z. DONG CHEN, Exact solutions of the position-dependent mass Schrödinger equation in D dimensions, Phys. Lett. A, 331 (2004), pp. 312–315, https://doi.org/10.1016/ j.physleta.2004.09.012.
- [13] A. DE SOUZA DUTRA AND C. ALMEIDA, Exact solvability of potentials with spatially dependent effective masses, Phys. Lett. A, 275 (2000), pp. 25–30, https://doi.org/10.1016/S0375-9601(00)00533-8.
- [14] A. DEKTOR, A. RODGERS, AND D. VENTURI, Rank-Adaptive Tensor Methods for High-Dimensional Nonlinear PDEs, https://arxiv.org/abs/2012.05962, 2021.
- [15] P. A. M. DIRAC, Note on exchange phenomena in the Thomas atom, Math. Proc. Cambridge Philos. Soc., 26 (1930), pp. 376–385, https://doi.org/10.1017/S0305004100016108.
- [16] G. D. DOCKERY, Modeling electromagnetic wave propagation in the troposphere using the parabolic equation, IEEE Trans. Antennas and Propagation, 36 (1988), pp. 1464–1470.
- [17] D. J. DONOHUE AND J. KUTTLER, Propagation modeling over terrain using the parabolic wave equation, IEEE Trans. Antennas and Propagation, 48 (2000), pp. 260–277.
- [18] C. ECKART AND G. YOUNG, The approximation of one matrix by another of lower rank, Psychometrika, 1 (1936), pp. 211–218, https://doi.org/10.1007/BF02288367.
- [19] E. FEHLBERG, Low-Order Classical Runge-Kutta Formulas with Stepsize Control and Their Application to Some Heat Transfer Problems, Technical report, National Aeronautics and Space Administration, Washington, DC, 1969.
- [20] F. FEPPON AND P. F. J. LERMUSIAUX, Dynamically orthogonal numerical schemes for efficient stochastic advection and Lagrangian transport, SIAM Rev., 60 (2018), pp. 595–625, https://doi.org/10.1137/16M1109394.
- [21] F. FEPPON AND P. F. J. LERMUSIAUX, A geometric approach to dynamical model-order reduction, SIAM J. Matrix Anal. Appl., 39 (2018), pp. 510–538, https://doi.org/10.1137/ 16M1095202.
- [22] M. GU, Subspace iteration randomization and singular value problems, SIAM J. Sci. Comput., 37 (2015), pp. A1139–A1173.
- [23] B. GUSTAFSSON, H. KREISS, AND J. OLIGER, Model Equations, Wiley, New York, 2013, pp. 1–46, https://doi.org/10.1002/9781118548448.ch1.
- [24] N. HALKO, P.-G. MARTINSSON, AND J. A. TROPP, Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions, SIAM Rev., 53 (2011), pp. 217–288.
- [25] J. P. HEUSS, P. J. HALEY, JR., C. MIRABITO, E. COELHO, M. C. SCHÖNAU, K. HEANEY, AND P. F. J.LERMUSIAUX, Reduced order modeling for stochastic prediction onboard autonomous platforms at sea, in Proceedings of OCEANS 2020 IEEE/MTS, IEEE, 2020, pp. 1–10, https://doi.org/10.1109/IEEECONF38699.2020.9389149.
- [26] L. O. JAY, A note on q-order of convergence, BIT, 41 (2001), pp. 422–429.
- [27] F. B. JENSEN, W. A. KUPERMAN, M. B. PORTER, AND H. SCHMIDT, Computational Ocean Acoustics, Springer, New York, 2011.
- [28] R. H. JORDAN AND D. G. HALL, Free-space azimuthal paraxial wave equation: The azimuthal Bessel-Gauss beam solution, Optim. Lett., 19 (1994), pp. 427–429.
- [29] K. KARHUNEN, Über lineare Methoden in der Wahrscheinlichkeitsrechnung, Ann. Acad. Sci. Fenn. Math. 37, Sana, 1947.
- [30] Y. KAZASHI, F. NOBILE, AND E. VIDLIČKOVÁ, Stability Properties of a Projector-Splitting Scheme for Dynamical Low Rank Approximation of Random Parabolic Equations, https:// arxiv.org/abs/2006.05211, 2020.
- [31] B. N. KHOROMSKIJ, I. V. OSELEDETS, AND R. SCHNEIDER, Efficient Time-Stepping Scheme for Dynamics on TT-Manifolds, Max-Planck-Institut f
  ür Mathematik in den Naturwissenschaften, Leipzig, 2012.
- [32] 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.
- [33] E. KIERI AND B. VANDEREYCKEN, Projection methods for dynamical low-rank approximation of high-dimensional problems, Comput. Methods Appl. Math., 19 (2019), pp. 73–92.
- [34] 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.
- [35] D. A. KOLESNIKOV AND I. V. OSELEDETS, Convergence analysis of projected fixed-point iteration on a low-rank matrix manifold, Numer. Linear Algebra Appl., 25 (2018), e2140.
- [36] D. D. KOSAMBI, Statistics in function space, in D. D. Kosambi: Selected Works in Mathematics and Statistics, Springer, New York, 2016, pp. 115–123.

- [37] J. N. KUTZ, S. L. BRUNTON, B. W. BRUNTON, AND J. L. PROCTOR, Dynamic mode decomposition: Data-driven modeling of complex systems, SIAM, Philadelphia, 2016
- [38] M. LAX, W. H. LOUISELL, AND W. B. MCKNIGHT, From Maxwell to paraxial wave optics, Phys. Rev. A, 11 (1975), 1365.
- [39] P. F. J. LERMUSIAUX, Evolving the subspace of the three-dimensional multiscale ocean variability: Massachusetts Bay, J. Marine Syst., 29 (2001), pp. 385–422, https://doi.org/10.1016/ S0924-7963(01)00025-2.
- [40] P. F. J. LERMUSIAUX, Uncertainty estimation and prediction for interdisciplinary ocean dynamics, J. Comput. Phys., 217 (2006), pp. 176–199, https://doi.org/10.1016/j.jcp.2006.02.010.
- [41] P. F. J. LERMUSIAUX, C.-S. CHIU, G. G. GAWARKIEWICZ, P. ABBOT, A. R. ROBINSON, R. N. MILLER, P. J. HALEY, JR., W. G. LESLIE, S. J. MAJUMDAR, A. PANG, AND F. LEKIEN, *Quantifying uncertainties in ocean predictions*, Oceanography, 19 (2006), pp. 92–105, https://doi.org/10.5670/oceanog.2006.93.
- [42] P. F. J. LERMUSIAUX, C.-S. CHIU, AND A. R. ROBINSON, Modeling uncertainties in the prediction of the acoustic wavefield in a shelfbreak environment, in Proceedings of the 5th International Conference on Theoretical and Computational Acoustics, E.-C. Shang, Q. Li, and T. F. Gao, eds., World Scientific, River Edge, NJ, 2002, pp. 191–200, https://doi.org/10.1142/9789812777362\_0020.
- [43] P. F. J. LERMUSIAUX, T. LOLLA, K. YIGIT, P. J. HALEY, JR., M. P. UECKERMANN, T. SON-DERGAARD, AND W. G. LESLIE, Science of autonomy: Time-optimal path planning and adaptive sampling for swarms of ocean vehicles, in Springer Handbook of Ocean Engineering: Autonomous Ocean Vehicles, Subsystems and Control, T. Curtin, ed., Springer, New York, 2016, pp. 481–498, https://doi.org/10.1007/978-3-319-16649-0\_21.
- [44] P. F. J. LERMUSIAUX, A. R. ROBINSON, P. J. HALEY, AND W. G. LESLIE, Advanced interdisciplinary data assimilation: Filtering and smoothing via error subspace statistical estimation, in Proceedings of The OCEANS 2002 MTS/IEEE Conference, Holland Publications, 2002, pp. 795–802, https://doi.org/10.1109/oceans.2002.1192071.
- [45] P. F. J. LERMUSIAUX, D. N. SUBRAMANI, J. LIN, C. S. KULKARNI, A. GUPTA, A. DUTT, T. LOLLA, W. H. ALI, P. J. HALEY, JR., C. MIRABITO, AND S. JANA, A future for intelligent autonomous ocean observing systems, J. Marine Res., 75 (2017), pp. 765–813, https://doi.org/10.1357/002224017823524035.
- [46] J. LIN, Bayesian Learning for High-Dimensional Nonlinear Dynamical Systems: Methodologies, Numerics and Applications to Fluid Flows, Ph.D. thesis, Massachusetts Institute of Technology, Department of Mechanical Engineering, Cambridge, MA, 2020.
- [47] J. LIN AND P. F. J. LERMUSIAUX, Minimum-correction second-moment matching: Theory, algorithms and applications, Numer. Math., 147 (2021), pp. 611–650, https://doi.org/ 10.1007/s00211-021-01178-8.
- [48] J. D. LOGAN, Applied Partial Differential Equations, Springer, New York, 2014.
- [49] M. LOÉVE, Probability Theory II, 4th ed., Grad. Texts in Math. 2, Springer, New York, 1978.
  [50] C. LUBICH, From Quantum to Classical Molecular Dynamics: Reduced Models and Numerical Analysis, European Mathematical Society, Helskini, 2008.
- [51] C. LUBICH AND I. V. OSELEDETS, A projector-splitting integrator for dynamical low-rank approximation, BIT, 54 (2014), pp. 171–188.
- [52] Z. LUO AND G. CHEN, Proper Orthogonal Decomposition Methods for Partial Differential Equations, Academic Press, New York, 2018.
- [53] J. LUTTINGER AND W. KOHN, Motion of electrons and holes in perturbed periodic fields, Phys. Rev., 97 (1955), pp. 869–883, https://doi.org/10.1103/PhysRev.97.869.
- [54] L. MIRSKY, Symmetric gauge functions and unitarily invariant norms, Q. J. Math., 11 (1960), pp. 50–59, https://doi.org/10.1093/qmath/11.1.50.
- [55] E. MUSHARBASH, Dynamical Low Rank Approximation of PDEs with Random Parameters, Technical report, EPFL, 2017.
- [56] E. MUSHARBASH, F. NOBILE, AND E. VIDLIČKOVÁ, Symplectic dynamical low rank approximation of wave equations with random parameters, BIT, 60 (2020), pp. 1153–1201.
- [57] G. NIENHUIS AND L. ALLEN, Paraxial wave optics and harmonic oscillators, Phys. Rev. A, 48 (1993), 656.
- [58] B. R. NOACK, M. MORZYNSKI, AND G. TADMOR, Reduced-Order Modelling for Flow Control, CISM Internat. Centre Mech. Sci. 528, Springer, New York, 2011.
- [59] F. POTRA, On Q-order and R-order of convergence, J. Optim. Theory Appl., 63 (1989), pp. 415–431.
- [60] O. ROJO AND J. LEVINGER, Integrated cross section for a velocity-dependent potential, Phys. Rev., 123 (1961), pp. 2177–2179, https://doi.org/10.1103/PhysRev.123.2177.

Downloaded 05/06/23 to 24.61.23.148 by Pierre Lermusiaux (pierrel@mit.edu). Redistribution subject to SIAM license or copyright; see https://epubs.siam.org/terms-privacy

- [61] T. P. SAPSIS AND P. F. J. LERMUSIAUX, Dynamically orthogonal field equations for continuous stochastic dynamical systems, Phys. D, 238 (2009), pp. 2347–2360, https://doi.org/ 10.1016/j.physd.2009.09.017.
- [62] T. P. SAPSIS AND P. F. J. LERMUSIAUX, Dynamical criteria for the evolution of the stochastic dimensionality in flows with uncertainty, Phys. D, 241 (2012), pp. 60–76, https://doi.org/10.1016/j.physd.2011.10.001.
- [63] E. SCHMIDT, Zur theorie der linearen und nichtlinearen integralgleichungen, Math. Ann., 63 (1907), pp. 433–476, https://doi.org/10.1007/BF01449770.
- [64] F. D. TAPPERT, The parabolic approximation method, in Wave Propagation and Underwater Acoustics, Springer, New York, 1977, pp. 224–287.
- [65] J. THOMPSON, J. C. MCWHORTER, S. SIDDIQI, AND S. SHANKS, A Study of Numerical Methods of Solution of the Equations of Motion of a Controlled Satellite Under the Influence of Gravity Gradient Torque, Technical report, National Aeronautics and Space Administration, Washington, DC, 1973.
- [66] L. TREFETHEN AND D. BAU, Numerical Linear Algebra, SIAM, Phildelphia, 1997.
- [67] M. P. UECKERMANN, P. F. J. LERMUSIAUX, AND T. P. SAPSIS, Numerical schemes for dynamically orthogonal equations of stochastic fluid and ocean flows, J. Comput. Phys., 233 (2013), pp. 272–294, https://doi.org/10.1016/j.jcp.2012.08.041.
- [68] A. USCHMAJEW AND B. VANDEREYCKEN, Geometric Methods on Low-Rank Matrix and Tensor Manifolds, Springer, Cham, 2020, pp. 261–313, https://doi.org/10.1007/978-3-030-31351-7\_9.
- [69] B. VANDEREYCKEN, Low-rank matrix completion by Riemannian optimization, SIAM J. Optim., 23 (2013), pp. 1214–1236, https://doi.org/10.1137/110845768.
- [70] T. VU, E. CHUNIKHINA, AND R. RAICH, Perturbation Expansions and Error Bounds for the Truncated Singular Value Decomposition, https://arxiv.org/abs/2009.07542, 2020.
- [71] N. WIENER, The homogeneous chaos, Amer. J. Math., 60 (1938), pp. 897–936.
- [72] N. WIENER, Nonlinear Problems in Random Theory, M.I.T. Paperback Ser., Technology Press of Massachusetts Institute of Technology, Cambridge, MA, 1958.
- [73] D. XIU AND G. E. KARNIADAKIS, The Wiener-Askey polynomial chaos for stochastic differential equations, SIAM J. Sci. Comput., 24 (2002), pp. 619–644.
- [74] D. XIU AND G. E. KARNIADAKIS, Modeling uncertainty in flow simulations via generalized polynomial chaos, J. Comput. Phys., 187 (2003), pp. 137–167.
- [75] M. YANG AND S. R. WHITE, Time-dependent variational principle with ancillary Krylov subspace, Phys. Rev. B, 102 (2020), 094315, https://doi.org/10.1103/PhysRevB.102.094315.
- [76] J. YU AND S.-H. DONG, Exactly solvable potentials for the Schrödinger equation with spatially dependent mass, Phys. Lett. A, 325 (2004), pp. 194–198, https://doi.org/10. 1016/j.physleta.2004.03.056.
- [77] J. ZHANG AND S. ZHANG, A Cubic Regularized Newton's Method over Riemannian Manifolds, preprint, arXiv:1805.05565, 2018.
- [78] J. ŠIMŠA, The best L2-approximation by finite sums of functions with separable variables, Aequationes Math., 43 (1992), pp. 248–263.