

# Galerkin v. discrete-optimal projection in nonlinear model reduction

Kevin Carlberg<sup>a,\*</sup>, Matthew Barone<sup>a,\*</sup>, Harbir Antil<sup>b,\*\*</sup>

<sup>a</sup>*Sandia National Laboratories*

<sup>b</sup>*George Mason University*

## Abstract

Discrete-optimal model-reduction techniques such as the Gauss–Newton with Approximated Tensors (GNAT) method have shown promise, as they have generated stable, accurate solutions for large-scale turbulent, compressible flow problems where standard Galerkin techniques have failed. However, there has been limited comparative analysis of the two approaches. This is due in part to difficulties arising from the fact that Galerkin techniques perform projection at the time-continuous level, while discrete-optimal techniques do so at the time-discrete level.

This work provides a detailed theoretical and experimental comparison of the two techniques for two common classes of time integrators: linear multistep schemes and Runge–Kutta schemes. We present a number of new findings, including conditions under which the discrete-optimal ROM has a time-continuous representation, conditions under which the two techniques are equivalent, and time-discrete error bounds for the two approaches. Perhaps most surprisingly, we demonstrate both theoretically and experimentally that decreasing the time step does not necessarily decrease the error for the discrete-optimal ROM; instead, the time step should be ‘matched’ to the spectral content of the reduced basis. In numerical experiments carried out on a turbulent compressible-flow problem with over one million unknowns, we show that increasing the time step to an intermediate value decreases both the error and the simulation time of the discrete-optimal reduced-order model by an order of magnitude.

*Keywords:* model reduction, GNAT, discrete optimality, Galerkin projection, CFD

## 1. Introduction

While modeling and simulation of parameterized systems has become an essential tool in many industries, the computational cost of executing high-fidelity simulations is infeasibly high for many time-critical applications. For example, real-time scenarios (e.g., model predictive control) require simulations to execute in seconds or minutes, while many-query scenarios (e.g., sampling statistical inversion) can require thousands of simulations corresponding to different input-parameter instances of the system.

Reduced-order models (ROMs) have been developed to mitigate this computational bottleneck. First, they execute an *offline* stage during which computationally expensive training tasks (e.g., evaluating the high-fidelity model at several points in the input-parameter space) compute a representative low-dimensional ‘trial’ basis for the system state. Then, during the inexpensive *online* stage, these methods quickly compute approximate solutions for arbitrary points in the input space via projection: they compute solutions in the span of the trial basis while enforcing the high-fidelity-model residual to be orthogonal to a low-dimensional ‘test’ basis. They also introduce other approximations in the presence of general nonlinearities (i.e., nonlinear terms that are not necessarily low-order polynomials) or non-affine parameter dependence. See Ref. [1] and references within for a survey of current methods.

\*7011 East Ave, MS 9159, Livermore, CA 94550. Sandia is a multiprogram laboratory operated by Sandia Corporation, a Lockheed Martin Company, for the United States Department of Energy under contract DE-AC04-94-AL85000.

\*\*4400 University Drive, MS: 3F2, Exploratory Hall, Room 4201, Fairfax, Virginia 22030.

*Email addresses:* [ktcarlb@sandia.gov](mailto:ktcarlb@sandia.gov) (Kevin Carlberg), [mbarone@sandia.gov](mailto:mbarone@sandia.gov) (Matthew Barone), [hantil@gmu.edu](mailto:hantil@gmu.edu) (Harbir Antil)

*URL:* [sandia.gov/~ktcarlb](http://sandia.gov/~ktcarlb) (Kevin Carlberg)

By far the most popular model-reduction technique for nonlinear ordinary differential equations (ODEs) is Galerkin projection [2], wherein the test basis is set to be equal to the trial basis, which is often computed via proper orthogonal decomposition (POD) [3]. This approach can be considered *continuous optimal*, as an orthogonal projection process is performed on the (time-continuous) ODE such that the approximated velocity vector is optimal in the  $\ell^2$  sense. In addition, for specialized dynamical systems (e.g., Lagrangian dynamical systems), performing Galerkin projection is necessary to preserve problem structure [4, 5, 6]. However, theoretical analysis—in the form of time-continuous error bounds [7] and stability analysis [8]—as well as numerical experiments have shown that Galerkin projection can lead to significant problems when applied to general nonlinear ODEs: instability [9], inaccurate long-time responses [10, 11], and no guarantee of *a priori* convergence (i.e., adding basis vectors can degrade the solution) [12]. In turbulent fluid flows, some of this poor performance can be attributed to the trial basis ‘filtering out’ small-scale modes essential for energy dissipation.

To address these shortcomings, alternative projection techniques have been developed, particularly in fluid dynamics. These include stabilizing inner products [13, 14, 15]; introducing dissipation via closure models [16, 10, 17, 18, 19] or numerical dissipation [20]; performing nonlinear Galerkin projection based on approximate inertial manifolds [21, 22, 23]; including a pressure-term representation [11, 24]; modifying the POD basis by including many modes (such that dissipative modes are captured), changing the norm [20], by enabling adaptivity [17], or by including basis functions that resolve a range of scales [25] or respect the attractor’s power balance [26]; and performing Petrov–Galerkin projection [27].

Alternatively, a promising new model-reduction methodology eschews Galerkin projection in favor of performing projection at the *fully discrete level*, i.e., after the ODE has been discretized in time [28]. This *discrete-optimal* method computes the solution that minimizes the  $\ell^2$  norm of the nonlinear residual arising at each time step; this leads to a notion of *a priori* convergence, as adding basis vectors guarantees monotonic decrease in the least-squares objective function. When equipped with a gappy POD [29] approximation of the discrete residual as a complexity-reduction mechanism, this approach is known as the Gauss–Newton with Approximated Tensors (GNAT) method [30]. It has been demonstrated to significantly outperform Galerkin projection on large-scale problems in turbulent, compressible fluid dynamics [31, 30].

In spite of its promise, theoretical analysis has been limited to developing consistency conditions for snapshot collection [28, 30] and discrete-time error bounds for simple time integrators [30, 32]. In particular, major outstanding questions include: (1) What are time-continuous and time-discrete representations of the Galerkin and discrete-optimal ROMs for broad classes of time integrators? (2) Are there conditions under which the two techniques are equivalent? (3) What discrete-time error bounds are available for the two techniques for broad classes time integrators? Related to the third issue is how parameters (e.g., time step or basis dimension) for the discrete-optimal ROM should be chosen to optimize performance. This work aims to fill this gap by performing a number of detailed theoretical and experimental studies that compare Galerkin and discrete-optimal ROMs for the two most important classes of time integrators: linear multistep methods and Runge–Kutta schemes. We summarize the most important theoretical results (which map to the three questions above) as follows:

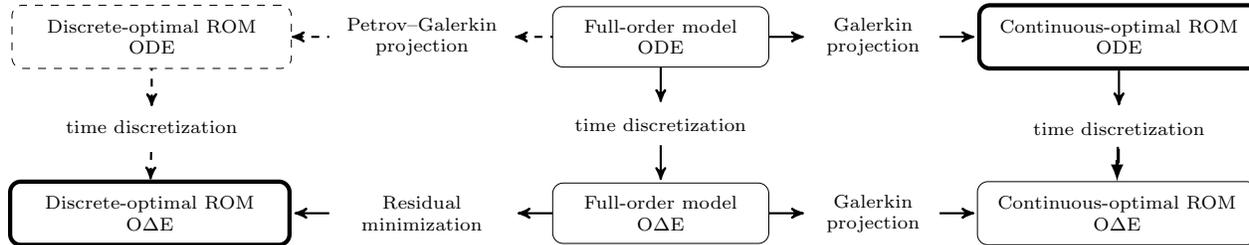


Figure 1: Relationship between Galerkin and discrete-optimal ROMs at the time-continuous and time-discrete levels. Bolded outlines imply an optimal ROM. Dashed lines imply the relationships are valid under certain conditions (see Theorems 4.2 and 4.3).

1. Continuous and discrete representations

- Projection and time discretization are commutative for Galerkin projection (Theorem 3.4, Figure 1).
- Discrete-optimal ROMs can be derived for Runge–Kutta schemes (Section 4.1).
- The discrete-optimal ROM has a time-continuous (i.e., ODE) representation under certain conditions (Theorems 4.2 and 4.3, Figure 1). This ODE depends on the time step  $\Delta t$ .

## 2. Equivalence conditions

- Galerkin ROMs are discrete optimal for explicit time integrators (Corollaries 5.1 and 5.2).
- Galerkin ROMs are discrete optimal in the limit of  $\Delta t \rightarrow 0$  for implicit time integrators (Theorem 5.3).
- Galerkin ROMs are discrete optimal for positive-definite residual Jacobians (Theorems 5.4 and 5.5).

## 3. Error analysis

- We provide discrete-time error bounds for both Galerkin and discrete-optimal ROMs for implicit linear multistep schemes (Theorem 6.1).
- We provide discrete-time error bounds for the Galerkin ROM for Runge–Kutta schemes (Theorem 6.2).
- For the backward Euler time integrator, we show that the discrete-optimal ROM yields a lower global state-space error bound than the Galerkin ROM because it solves a time-global optimization problem (over a time window) rather than a time-local optimization problem (Corollary 6.4).
- For the backward Euler time integrator, we show that *an intermediate* time step should yield the lowest error bound (Corollary 6.5).
- For the backward Euler time integrator, we show that a larger basis size leads to a smaller optimal time step for the discrete-optimal ROM (Corollary 6.5).

Figure 1 summarizes time-continuous and time-discrete representations of the two techniques.

In addition to the above theoretical results, we present numerical results for a large-scale compressible fluid-dynamics problem with turbulence model characterized by over one million degrees of freedom. These results illustrate the practical significance of the above theoretical results. Critically, we show that employing an intermediate time step for the discrete-optimal ROM can decrease both the error and the simulation time by an order of magnitude, which is a highly non-intuitive result that is of immense practical significance.

The remainder of the paper is organized as follows. Section 2 formulates the full-order model, including its representation at the time-continuous and time-discrete levels. Section 3 presents the Galerkin ROM at the continuous and discrete levels, and Section 4 does so for the discrete-optimal ROM. In particular, we provide conditions under which the discrete-optimal ROM has a time-continuous representation. Section 5 provides conditions under which the Galerkin and discrete-optimal ROMs are equivalent; in particular, equivalence holds for explicit integrators (Section 5.1), in the limit of the time step going to zero for implicit integrators (Section 5.2), and for symmetric-positive-definite residual Jacobians (Section 5.3). Section 6 provides error analysis for Galerkin and discrete-optimal ROMs for linear multistep schemes (Section 6.1), Runge–Kutta schemes (Section 6.2), and a detailed analysis in the case of backward Euler (Section 6.3). Section 7 provides detailed numerical examples that illustrate the practical importance of the analysis. Finally, Section 8 provides conclusions.

In the remainder of this paper, matrices are denoted by capitalized bold letters, vectors by lowercase bold letters, scalars by unbolded letters. The columns of a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  are denoted by  $\mathbf{a}_i \in \mathbb{R}^m$ ,  $i \in \mathbb{N}(n)$  with  $\mathbb{N}(a) := \{1, \dots, a\}$  such that  $\mathbf{A} := [\mathbf{a}_1 \cdots \mathbf{a}_n]$ . The scalar-valued matrix elements are denoted by  $a_{ij} \in \mathbb{R}$  such that  $\mathbf{a}_j := [a_{1j} \cdots a_{mj}]^T$ ,  $j \in \mathbb{N}(n)$ . A superscript denotes the value of a variable at that time instance, e.g.,  $\mathbf{x}^n$  is the value of  $\mathbf{x}$  at time  $n\Delta t$ , where  $\Delta t$  is the time step.

## 2. Full-order model

We begin by formulating both the time-continuous (ODE) and time-discrete (ODE) representations of the full-order model (FOM).

### 2.1. Continuous representation

In this work, we consider the full-order model (FOM) to be an initial value problem characterized by a system of nonlinear ODEs

$$\frac{d\mathbf{x}}{dt} = \mathbf{f}(\mathbf{x}, t), \quad \mathbf{x}(0) = \mathbf{x}_0, \quad (1)$$

where  $\mathbf{x} : [0, T] \rightarrow \mathbb{R}^N$  denotes the (time-dependent) state,  $\mathbf{x}_0 \in \mathbb{R}^N$  denotes the initial condition, and  $\mathbf{f} : \mathbb{R}^N \times [0, T] \rightarrow \mathbb{R}^N$  with  $(\boldsymbol{\xi}, \tau) \mapsto \mathbf{f}(\boldsymbol{\xi}, \tau)$ . This ODE can arise, for example, from applying spatial discretization (e.g., finite element, finite volume, or finite difference) to a partial differential equation with time dependence. We note that most model-reduction techniques are applied to parameterized systems wherein the velocity  $\mathbf{f}$  also depends on parametric inputs. However, we limit our presentation to unparameterized systems for notational simplicity, as we are interested comparing Galerkin and discrete-optimal ROMs for a given instance of the ODE.

### 2.2. Discrete representation

A time-discretization method is required to solve Eq. (1) numerically. We now characterize the full-order-model O $\Delta$ E, which is the time-discrete representation of the model, for two classes of time integrators: linear multistep schemes and Runge–Kutta schemes.

#### 2.2.1. Linear multistep schemes

A linear  $k$ -step method applied to numerically solve Eq. (1) can be expressed as

$$\sum_{j=0}^k \alpha_j \mathbf{x}^{n-j} = \Delta t \sum_{j=0}^k \beta_j \mathbf{f}(\mathbf{x}^{n-j}, t^{n-j}), \quad (2)$$

where  $\Delta t$  is the time step,  $\alpha_0 \neq 0$ , and  $\sum_{j=0}^k \alpha_j = 0$  is necessary for consistency. In this case, the O $\Delta$ E is characterized by the following algebraic system of equations to be solved at each time instance  $n \in \mathbb{N}(T/\Delta t)$ :

$$\mathbf{r}^n(\mathbf{w}^n) = 0, \quad (3)$$

where  $\mathbf{w}^n \in \mathbb{R}^N$  is the unknown variable and  $\mathbf{r}^n : \mathbb{R}^N \rightarrow \mathbb{R}^N$  denotes the linear multistep residual defined as

$$\mathbf{r}^n(\mathbf{w}) := \alpha_0 \mathbf{w} - \Delta t \beta_0 \mathbf{f}(\mathbf{w}, t^n) + \sum_{j=1}^k \alpha_j \mathbf{x}^{n-j} - \Delta t \sum_{j=1}^k \beta_j \mathbf{f}(\mathbf{x}^{n-j}, t^{n-j}). \quad (4)$$

Then, the state can be updated explicitly as

$$\mathbf{x}^n = \mathbf{w}^n.$$

Hence, the unknown is equal to the state. These methods are implicit if  $\beta_0 \neq 0$ .

#### 2.2.2. Runge–Kutta schemes

For an  $s$ -stage Runge–Kutta scheme, the O $\Delta$ E is characterized by the following algebraic system of equations to be solved at each time step:

$$\mathbf{r}_i^n(\mathbf{w}_1^n, \dots, \mathbf{w}_s^n) = 0, \quad i \in \mathbb{N}(s). \quad (5)$$

Here, the Runge–Kutta residual is defined as

$$\mathbf{r}_i^n(\mathbf{w}_1, \dots, \mathbf{w}_s) := \mathbf{w}_i - \mathbf{f}(\mathbf{x}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \mathbf{w}_j, t^{n-1} + c_i \Delta t), \quad i \in \mathbb{N}(s) \quad (6)$$

and the state is explicitly updated as

$$\mathbf{x}^n = \mathbf{x}^{n-1} + \Delta t \sum_{i=1}^s b_i \mathbf{w}_i^n. \quad (7)$$

Here, the unknowns  $\mathbf{w}_i^n$  correspond to the velocity  $d\mathbf{x}/dt$  at times  $t^{n-1} + c_i \Delta t$ ,  $i \in \mathbb{N}(s)$ . These methods are implicit if  $a_{ij} \neq 0$  for some  $j \geq i$ .

### 3. Galerkin ROM

This section provides the time-continuous and time-discrete representations of the Galerkin ROM, as well as key results related to optimality and commutativity of projection and time discretization.

#### 3.1. Continuous representation

Galerkin-projection reduced-order models compute an approximate solution  $\tilde{\mathbf{x}} \approx \mathbf{x}$  with  $\tilde{\mathbf{x}} \in \mathbb{R}^N$  to Eq. (1) by introducing two approximations. First, they restrict the approximate solution to lie in a low-dimensional affine trial subspace  $\mathbf{x}_0 + \text{range}(\Phi)$ , where  $\Phi \in \mathbb{R}^{N \times p}$  with  $\Phi^T \Phi = I$  denotes the given reduced basis (in matrix form) of dimension  $p \ll N$ . This basis can be computed by a variety of techniques, e.g., eigenmode analysis, POD [3], or the reduced-basis method [33, 34, 35, 36, 37]. Then, the approximate solution can be expressed as

$$\tilde{\mathbf{x}}(t) = \mathbf{x}_0 + \Phi \hat{\mathbf{x}}(t), \quad (8)$$

where  $\hat{\mathbf{x}} : [0, T] \rightarrow \mathbb{R}^p$  denotes the (time-dependent) generalized coordinates of the approximate solution. Second, these methods substitute  $\mathbf{x} \leftarrow \tilde{\mathbf{x}}$  into Eq. (1) and enforce the ODE residual to be orthogonal to  $\text{range}(\Phi)$ , which results in a low-dimensional system of nonlinear ODEs

$$\frac{d\hat{\mathbf{x}}}{dt} = \Phi^T \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t), \quad \hat{\mathbf{x}}(0) = 0. \quad (9)$$

**Remark 3.1.** *In order to obtain computational efficiency, it is necessary to reduce the computational complexity of repeatedly computing matrix–vector products of the form  $\Phi^T \mathbf{f}$ . This can be done using a variety of methods, e.g., collocation [38, 39, 40], gappy POD [29, 41, 38, 28, 30], the discrete empirical interpolation method (DEIM) [42, 43, 44, 45, 46], reduced-order quadrature [47], finite-element subassembly methods [48, 49], or reduced-basis-sparsification techniques [6]. However, in this work we limit ourselves to comparatively analyzing different projection techniques. For this reason, we do not perform additional analysis for such complexity-reduction mechanisms; this is the subject of follow-on work.*

We now restate the well-known result that Galerkin projection leads to a notion of optimality at the continuous level. This is reflected in the top-right box of Figure 1, where the bolded outline indicates optimality.

**Theorem 3.2 (Galerkin: continuous optimality).** *The Galerkin ROM (8)–(9) is continuous optimal in the sense that the approximated velocity minimizes the error in the velocity  $\mathbf{f}$  over  $\text{range}(\Phi)$ , i.e.,*

$$\frac{d\tilde{\mathbf{x}}}{dt}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t) = \arg \min_{\mathbf{v} \in \text{range}(\Phi)} \|\mathbf{v} - \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t)\|_2^2. \quad (10)$$

PROOF. Because  $\frac{d\tilde{\mathbf{x}}}{dt} = \Phi \frac{d\hat{\mathbf{x}}}{dt}$ , problem (10) can be written as

$$\frac{d\hat{\mathbf{x}}}{dt}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t) = \arg \min_{\hat{\mathbf{v}} \in \mathbb{R}^p} g(\hat{\mathbf{v}}) \quad (11)$$

where  $g(\hat{\mathbf{v}}) := \|\Phi \hat{\mathbf{v}} - \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t)\|_2^2$ . We now assess whether Eq. (11) holds, i.e., whether  $\frac{d\hat{\mathbf{x}}}{dt}$  as defined by Eq. (9) is the minimizer of  $g$ .

The function  $g$  can be expressed as  $g(\hat{\mathbf{v}}) = \hat{\mathbf{v}}^T \Phi^T \Phi \hat{\mathbf{v}} - 2\hat{\mathbf{v}}^T \Phi^T \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t) + \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t)^T \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t)$ . Due to the strict convexity of the function  $g$ , the global minimizer  $\hat{\mathbf{v}}^*$  is equal to the stationary point of  $g$ , i.e.,  $\hat{\mathbf{v}}^*$  satisfies

$$0 = \frac{dg}{d\hat{\mathbf{v}}}(\hat{\mathbf{v}}^*) = 2\Phi^T \Phi \hat{\mathbf{v}}^* - 2\Phi^T \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t) \quad (12)$$

$$\hat{\mathbf{v}}^* = \Phi^T \frac{d\mathbf{x}}{dt} = \Phi^T \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t), \quad (13)$$

where orthogonality of  $\Phi$  has been used. Comparing Eqs. (13) and (9) shows  $\frac{d\hat{\mathbf{x}}}{dt}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t) = \hat{\mathbf{v}}^*$ , which is the desired result.

**Remark 3.3 (Continuous a priori convergence of the Galerkin ROM).** *Due to optimality property (10), the Galerkin ROM can be considered a priori convergent at the continuous level, as adding vectors to the trial basis—which expands the trial subspace  $\text{range}(\Phi)$ —results in a monotonic decrease of the objective function in problem (10).*

### 3.2. Discrete representation

As before, a time-discretization method is needed to numerically solve Eq. (9). We now characterize the O $\Delta$ E for the Galerkin ROM.

#### 3.2.1. Linear multistep schemes

A linear  $k$ -step method applied to numerically solve Eq. (9) can be expressed as

$$\sum_{j=0}^k \alpha_j \hat{\mathbf{x}}^{n-j} = \Delta t \sum_{j=0}^k \beta_j \Phi^T \mathbf{f} \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}^{n-j}, t^{n-j} \right).$$

Here, the O $\Delta$ E is characterized by the following algebraic system of nonlinear equations to be solved at each time step:

$$\hat{\mathbf{r}}^n (\hat{\mathbf{w}}^n) = 0. \quad (14)$$

Here, the discrete residual corresponds to

$$\hat{\mathbf{r}}^n (\hat{\mathbf{w}}) := \alpha_0 \hat{\mathbf{w}} - \Delta t \beta_0 \Phi^T \mathbf{f} (\mathbf{x}_0 + \Phi \hat{\mathbf{w}}, t^n) + \sum_{j=1}^k \alpha_j \hat{\mathbf{x}}^{n-j} - \Delta t \sum_{j=1}^k \beta_j \Phi^T \mathbf{f} \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}^{n-j}, t^{n-j} \right) \quad (15)$$

and the generalized state is explicitly updated as

$$\hat{\mathbf{x}}^n = \hat{\mathbf{w}}^n.$$

#### 3.2.2. Runge–Kutta schemes

Applying an  $s$ -stage Runge–Kutta method to solve Eq. (9) leads to an O $\Delta$ E characterized by the following algebraic system of equations to be solved at each time step:

$$\hat{\mathbf{r}}_i^n (\hat{\mathbf{w}}_1^n, \dots, \hat{\mathbf{w}}_s^n) = 0, \quad i \in \mathbb{N}(s). \quad (16)$$

Here, the residual is defined as

$$\hat{\mathbf{r}}_i^n (\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_s) := \hat{\mathbf{w}}_i - \Phi^T \mathbf{f} (\mathbf{x}_0 + \Phi \hat{\mathbf{x}}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \Phi \hat{\mathbf{w}}_j, t^{n-1} + c_i \Delta t), \quad i \in \mathbb{N}(s) \quad (17)$$

and the generalized state is computed explicitly as

$$\hat{\mathbf{x}}^n = \hat{\mathbf{x}}^{n-1} + \Delta t \sum_{i=1}^s b_i \hat{\mathbf{w}}_i^n. \quad (18)$$

We now show that projection and time discretization are commutative for Galerkin projection. This corresponds to the rightmost part of Figure 1.

#### **Theorem 3.4 (Galerkin: commutativity of projection and time discretization).**

*Performing a Galerkin projection on the governing ODE and subsequently applying time discretization yields the same model as first applying time discretization on the governing ODE and subsequently performing Galerkin projection.*

**PROOF.** Linear multistep schemes. Eq. (14) was derived by performing Galerkin projection on the continuous representation of the FOM and subsequently applying time discretization. If instead we apply Galerkin projection to the discrete representation of the FOM in Eq. (3), set  $\mathbf{w} = \mathbf{x}_0 + \Phi \hat{\mathbf{w}}$  and  $\mathbf{x}^i = \mathbf{x}_0 + \Phi \hat{\mathbf{x}}^i$ ,  $i \in \mathbb{N}(n)$ , and use  $\sum_{j=1}^k \alpha_j = 0$  and  $\Phi^T \Phi = I$ , we obtain the following O $\Delta$ E to be solved at each time step:  $\Phi^T \mathbf{r}^n (\mathbf{x}_0 + \Phi \hat{\mathbf{w}}) = 0$ . Comparing the definition of the linear multistep residual (4) with Eq. (15) reveals

$$\hat{\mathbf{r}}^n (\hat{\mathbf{w}}) = \Phi^T \mathbf{r}^n (\mathbf{x}_0 + \Phi \hat{\mathbf{w}}), \quad (19)$$

which shows that the same discrete equations  $\hat{\mathbf{r}}^n(\hat{\mathbf{w}}) = 0$  are obtained at each time step regardless of the ordering of time discretization and Galerkin projection.

Runge–Kutta schemes. Eq. (16) was derived by performing Galerkin projection on the continuous FOM representation and then applying time discretization. If instead we apply Galerkin projection to the discrete FOM representation in Eq. (5), set  $\mathbf{x}^{n-1} = \mathbf{x}_0 + \Phi \hat{\mathbf{x}}^{n-1}$ ,  $\mathbf{w}_i = \Phi \hat{\mathbf{w}}_i$ ,  $i \in \mathbb{N}(s)$ , and use  $\Phi^T \Phi = I$ , we obtain the following OΔE to be solved at each time step:  $\Phi^T \mathbf{r}_i^n(\Phi \hat{\mathbf{w}}_1, \dots, \Phi \hat{\mathbf{w}}_s) = 0$ ,  $i \in \mathbb{N}(s)$ . Comparing the definition of the Runge–Kutta residual (6) with Eq. (17) reveals

$$\hat{\mathbf{r}}_i^n(\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_s) = \Phi^T \mathbf{r}_i^n(\Phi \hat{\mathbf{w}}_1, \dots, \Phi \hat{\mathbf{w}}_s), \quad i \in \mathbb{N}(s), \quad (20)$$

which shows that the same discrete equations  $\hat{\mathbf{r}}_i^n(\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_s) = 0$ ,  $i \in \mathbb{N}(s)$  are obtained at each time step regardless of the ordering of time discretization and Galerkin projection.

## 4. Discrete-optimal ROM

Rather than perform the projection on the full-order model ODE (i.e., at the continuous level), projection can be performed on the full-order model OΔE (i.e., at the discrete level). Doing so is useful if the associated projection associates with an optimization problem at the discrete level. In particular, we consider *residual-minimizing* projections that minimize the discrete residual(s) (in some norm) arising at each time instance.

We note that other residual-minimizing approaches have been developed in the case of steady-state problems [40], linear dynamical systems [50], and space–time solutions [51]. In addition, a recently developed approach [52] has suggested  $\ell^1$  minimization of the residual arising at each time instance for hyperbolic problems.

### 4.1. Discrete representation

We begin by developing the time-discrete representation for the discrete-optimal ROM for both linear multistep schemes and Runge–Kutta schemes. The latter is a novel contribution, as previous work has derived discrete-optimal ROMs only for linear multistep schemes [28, 30]. Optimality of this approach corresponds to the bolded bottom-left box of Figure 1.

#### 4.1.1. Linear multistep schemes

As before with Galerkin projection, discrete-optimal ROMs compute solutions using two approximations. First, they restrict the approximate solution to lie in the same low-dimensional affine trial subspace  $\tilde{\mathbf{x}} \in \mathbf{x}_0 + \text{range}(\Phi)$  as Galerkin; thus, the approximate solution can be written according to Eq. (8). In the case of linear multistep schemes, the unknown at time step  $n$  is simply the state, i.e.,  $\mathbf{w}^n = \mathbf{x}^n$ , which implies that  $\tilde{\mathbf{w}}^n = \mathbf{x}_0 + \Phi \hat{\mathbf{w}}^n$ . Second, the discrete-optimal ROM substitutes  $\mathbf{w}^n \leftarrow \tilde{\mathbf{w}}^n$  into the OΔE (3) and solves a minimization problem to ensure the approximate solution is optimal in some sense at the discrete level:

$$\tilde{\mathbf{w}}^n = \arg \min_{z \in \mathbf{x}_0 + \text{range}(\Phi)} \|\mathbf{A}(z) \mathbf{r}^n(z)\|_2^2 \quad (21)$$

or equivalently

$$\hat{\mathbf{w}}^n = \arg \min_{\hat{\mathbf{z}} \in \mathbb{R}^p} \|\mathbf{A}(\mathbf{x}_0 + \Phi \hat{\mathbf{z}}) \mathbf{r}^n(\mathbf{x}_0 + \Phi \hat{\mathbf{z}})\|_2^2. \quad (22)$$

Here,  $\mathbf{A} \in \mathbb{R}^{z \times N}$  with  $z \leq N$  is a weighting matrix that enables the definition of a weighted (semi)norm. Examples of such reduced-order models include the least-squares Petrov–Galerkin method [28, 30, 40] ( $\mathbf{A} = \mathbf{I}$ ) and the related GNAT method [28, 30] ( $\mathbf{A} = (\mathbf{P}\Phi_r)^+ \mathbf{P}$  with  $\Phi_r$  a basis for the residual, the superscript  $+$  denoting the Moore–Penrose pseudoinverse, and  $\mathbf{P}$  consisting of selected rows of the identity matrix).

Note that the solution to Eq. (22) corresponds to a stationary point of the objective function in Eq. (22), i.e., it satisfies

$$\Psi^n(\hat{\mathbf{w}}^n)^T \mathbf{r}^n(\mathbf{x}_0 + \Phi \hat{\mathbf{w}}^n) = 0 \quad (23)$$

where the entries of  $\Psi^n \in \mathbb{R}^{N \times p}$  are

$$\begin{aligned} \psi_{ij}^n(\hat{\mathbf{w}}) &= a_{mi}(\mathbf{x}_0 + \Phi \hat{\mathbf{w}}) \frac{\partial a_{ml}(\mathbf{x}_0 + \Phi \hat{\mathbf{w}})}{\partial w_k} \phi_{kj} r_\ell^n(\mathbf{x}_0 + \Phi \hat{\mathbf{w}}) + \\ & a_{mi}(\mathbf{x}_0 + \Phi \hat{\mathbf{w}}) a_{ml}(\mathbf{x}_0 + \Phi \hat{\mathbf{w}}) \frac{\partial r_\ell^n}{\partial w_k}(\mathbf{x}_0 + \Phi \hat{\mathbf{w}}) \phi_{kj}, \quad i \in \mathbb{N}(N), j \in \mathbb{N}(p), \end{aligned} \quad (24)$$

where a repeated index implies summation. Because Eq. (23) corresponds to a Petrov–Galerkin projection with trial subspace  $\text{range}(\Phi)$  and test subspace  $\text{range}(\Psi)$ , the discrete-optimal projection is sometimes referred to as a least-squares Petrov–Galerkin projection [30, 28].

#### 4.1.2. Runge–Kutta schemes

Discrete-optimal ROMs for Runge–Kutta schemes also approximate the solution according to Eq. (8). However, because the unknown at time step  $n$  and stage  $i$  is the velocity at an intermediate time point, i.e.,  $\mathbf{w}_i^n = \dot{\mathbf{x}}(t^{n-1} + c_i \Delta t)$  for  $i \in \mathbb{N}(s)$ , we have  $\tilde{\mathbf{w}}_i^n = \Phi \dot{\mathbf{x}}(t^{n-1} + c_i \Delta t)$  for this case. Then, these techniques substitute  $\mathbf{w}^n \leftarrow \tilde{\mathbf{w}}^n$  into the ODE (5) and solve the following minimization problem:

$$(\tilde{\mathbf{w}}_1^n, \dots, \tilde{\mathbf{w}}_s^n) = \arg \min_{(\mathbf{z}_1, \dots, \mathbf{z}_s) \in \text{range}(\Phi)^s} \sum_{i=1}^s \|\mathbf{A}_i(\mathbf{z}_1, \dots, \mathbf{z}_s) \mathbf{r}_i^n(\mathbf{z}_1, \dots, \mathbf{z}_s)\|_2^2 \quad (25)$$

or equivalently

$$(\hat{\mathbf{w}}_1^n, \dots, \hat{\mathbf{w}}_s^n) = \arg \min_{(\hat{\mathbf{z}}_1, \dots, \hat{\mathbf{z}}_s) \in \mathbb{R}^{p \times s}} \sum_{i=1}^s \|\mathbf{A}_i(\Phi \hat{\mathbf{z}}_1, \dots, \Phi \hat{\mathbf{z}}_s) \mathbf{r}_i^n(\Phi \hat{\mathbf{z}}_1, \dots, \Phi \hat{\mathbf{z}}_s)\|_2^2. \quad (26)$$

Here,  $\mathbf{A}_i \in \mathbb{R}^{z \times N}$ ,  $i \in \mathbb{N}(s)$  with  $z \leq N$  are weighting matrices. As before, the solution to Eq. (26) corresponds to a stationary point of the objective function, i.e., it satisfies

$$\sum_{j=1}^s \Psi_{ij}^n(\hat{\mathbf{w}}_1^n, \dots, \hat{\mathbf{w}}_s^n)^T \mathbf{r}_j^n(\Phi \hat{\mathbf{w}}_1^n, \dots, \Phi \hat{\mathbf{w}}_s^n) = 0, \quad i = 1, \dots, s, \quad (27)$$

where entries of the test bases  $\Psi_{ij}^n \in \mathbb{R}^{N \times p}$ ,  $i, j \in \mathbb{N}(s)$  are

$$\begin{aligned} [\Psi_{ij}^n]_{k\ell}(\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_s) &= [\mathbf{A}_i]_{uk}(\Phi \hat{\mathbf{w}}_1, \dots, \Phi \hat{\mathbf{w}}_s) \frac{\partial [\mathbf{A}_i]_{um}(\Phi \hat{\mathbf{w}}_1, \dots, \Phi \hat{\mathbf{w}}_s)}{\partial [\mathbf{w}_j]_n} \phi_{n\ell}[\mathbf{r}_i^n]_m(\Phi \hat{\mathbf{w}}_1, \dots, \Phi \hat{\mathbf{w}}_s) + \\ &[\mathbf{A}_i]_{uk}(\Phi \hat{\mathbf{w}}_1, \dots, \Phi \hat{\mathbf{w}}_s) [\mathbf{A}_i]_{um}(\Phi \hat{\mathbf{w}}_1, \dots, \Phi \hat{\mathbf{w}}_s) \frac{\partial [\mathbf{r}_i^n]_m}{\partial [\mathbf{w}_j]_n}(\Phi \hat{\mathbf{w}}_1, \dots, \Phi \hat{\mathbf{w}}_s) \phi_{n\ell}, \end{aligned} \quad (28)$$

where  $[\cdot]_{ij}$  denotes entry  $(i, j)$  of the argument. This again leads to a least-squares Petrov–Galerkin interpretation for the discrete-optimal ROM.

In the explicit case, we can consider another notion of discrete optimality. Explicit Runge–Kutta schemes are characterized by  $a_{ij} = 0$ ,  $\forall j \geq i$ . In this case, solutions  $\mathbf{w}_i^n$ ,  $i \in \mathbb{N}(s)$  can be computed sequentially, i.e.,

$$\mathbf{q}_i^n(\mathbf{w}_i^n) = 0, \quad i \in \mathbb{N}(s)$$

with

$$\mathbf{q}_i^n(\mathbf{w}) := \mathbf{w} - \mathbf{f}(\mathbf{x}^{n-1} + \Delta t \sum_{j=1}^{i-1} a_{ij} \mathbf{w}_j^n, t^{n-1} + c_i \Delta t), \quad i \in \mathbb{N}(s). \quad (29)$$

We can then formulate the following sequence of optimization problems to compute discrete-optimal approximations:

$$\tilde{\mathbf{w}}_i^n = \arg \min_{\mathbf{z} \in \text{range}(\Phi)} \|\mathbf{A}_i(\mathbf{z}) \mathbf{q}_i^n(\mathbf{z})\|_2^2, \quad i \in \mathbb{N}(s), \quad (30)$$

or equivalently

$$\hat{\mathbf{w}}_i^n = \arg \min_{\hat{\mathbf{z}} \in \mathbb{R}^p} \|\mathbf{A}_i(\Phi \hat{\mathbf{z}}) \mathbf{q}_i^n(\Phi \hat{\mathbf{z}})\|_2^2, \quad i \in \mathbb{N}(s). \quad (31)$$

Here, the associated Petrov–Galerkin projection is

$$\Psi_i^n(\hat{\mathbf{w}}_i^n)^T \mathbf{q}_i^n(\Phi \hat{\mathbf{w}}_i^n) = 0, \quad i \in \mathbb{N}(s), \quad (32)$$

with test-basis entries of

$$[\Psi_i^n]_{jk}(\Phi \hat{\mathbf{w}}) = [\mathbf{A}_i]_{uj}(\Phi \hat{\mathbf{w}}) \frac{\partial [\mathbf{A}_i]_{u\ell}(\Phi \hat{\mathbf{w}})}{\partial w_m} \phi_{mk}[\mathbf{q}_i^n]_\ell(\Phi \hat{\mathbf{w}}) + [\mathbf{A}_i]_{uj}(\Phi \hat{\mathbf{w}}) [\mathbf{A}_i]_{um}(\Phi \hat{\mathbf{w}}) \phi_{mk}, \quad (33)$$

where we have used  $\partial \mathbf{q}_i^n / \partial \mathbf{w} = \mathbf{I}$ .

**Remark 4.1 (Discrete *a priori* convergence of the discrete-optimal ROM).** *Due to optimality property (21), the discrete-optimal ROM can be considered a priori convergent at the discrete level for linear multistep schemes, as adding vectors to the trial basis—which expands the trial subspace range ( $\Phi$ )—results in a monotonic decrease in the objective function in problem (21). This result also holds for discrete-optimal ROMs applied to Runge–Kutta schemes, as the computed solutions satisfy alternative optimality properties in the implicit (25) and explicit (30) cases.*

#### 4.2. Continuous representation

Because the discrete-optimal ROM introduces approximations at the discrete level, it is unclear whether it can be interpreted at the continuous level. We now show that an ODE representation of the discrete-optimal ROM exists for both linear multistep schemes and Runge–Kutta schemes under certain conditions; however, the ODE depends on the time step used to define the discrete-optimal ROM. This associates with the top-left section of the relationship diagram in Figure 1.

**Theorem 4.2 (Discrete-optimal ROM continuous representation: linear multistep schemes).**

*The discrete-optimal ROM for linear multistep integrators is equivalent to applying a Petrov–Galerkin projection to the ODE with test basis (in matrix form)*

$$\Psi(\hat{\mathbf{x}}, t) = \mathbf{A}^T \mathbf{A} \left( \alpha_0 \mathbf{I} - \Delta t \beta_0 \frac{\partial \mathbf{f}}{\partial \xi}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t) \right) \Phi \quad (34)$$

and subsequently applying time integration with a linear multistep scheme with time step  $\Delta t$  if  $\mathbf{A}$  is a constant matrix and (at least) one of the following conditions holds:

1.  $\beta_j = 0$ ,  $j \geq 1$  (e.g., a single-step method),
2. the velocity  $\mathbf{f}$  is linear in the state, or
3.  $\beta_0 = 0$  (i.e., explicit schemes).

**PROOF.** Applying Petrov–Galerkin projection to the full-order model ODE (1) using a trial subspace  $\mathbf{x}_0 +$  range( $\Phi$ ) and test subspace range( $\Psi$ ) yields the following ODE

$$\Psi(\hat{\mathbf{x}}, t)^T \Phi \frac{d\hat{\mathbf{x}}}{dt} = \Psi(\hat{\mathbf{x}}, t)^T \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t), \quad \hat{\mathbf{x}}(0) = 0, \quad (35)$$

which can be written in standard form as

$$\frac{d\hat{\mathbf{x}}}{dt} = \left( \Psi(\hat{\mathbf{x}}, t)^T \Phi \right)^{-1} \Psi(\hat{\mathbf{x}}, t)^T \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t), \quad \hat{\mathbf{x}}(0) = 0. \quad (36)$$

Case 1 Applying a linear multistep time integrator with the stated assumption of  $\beta_j = 0$ ,  $j \geq 1$  to numerically solve Eq. (36) results in the following discrete equations to be solved at each time instance:

$$\alpha_0 \hat{\mathbf{y}}^n - \Delta t \beta_0 \left( \Psi(\hat{\mathbf{y}}^n, t^n)^T \Phi \right)^{-1} \Psi(\hat{\mathbf{y}}^n, t^n)^T \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{y}}^n, t^n) + \sum_{j=1}^k \alpha_j \hat{\mathbf{x}}^{n-j} = 0. \quad (37)$$

Pre-multiplying by  $\Psi(\hat{\mathbf{y}}^n, t^n)^T \Phi$  yields discrete equations  $\hat{\mathbf{r}}^n(\hat{\mathbf{y}}^n) = 0$  with residual

$$\hat{\mathbf{r}}^n(\hat{\mathbf{w}}) := \alpha_0 \Psi(\hat{\mathbf{w}}, t^n)^T \Phi \hat{\mathbf{w}} - \Delta t \beta_0 \Psi(\hat{\mathbf{w}}, t^n)^T \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{w}}, t^n) + \sum_{j=1}^k \alpha_j \Psi(\hat{\mathbf{w}}, t^n)^T \Phi \hat{\mathbf{x}}^{n-j}. \quad (38)$$

Comparing Eqs. (38) and (4) reveals  $\hat{\mathbf{r}}^n(\hat{\mathbf{w}}) = \Psi(\hat{\mathbf{w}}, t^n)^T \mathbf{r}^n(\mathbf{x}_0 + \Phi \hat{\mathbf{w}})$  and so the solution  $\hat{\mathbf{y}}^n$  satisfies

$$\Psi(\hat{\mathbf{y}}^n, t^n)^T \mathbf{r}^n(\mathbf{x}_0 + \Phi \hat{\mathbf{y}}^n) = 0. \quad (39)$$

Under the stated assumptions, we have  $\partial \mathbf{r}^n / \partial \mathbf{w}(\mathbf{x}) = \alpha_0 \mathbf{I} - \Delta t \beta_0 \frac{\partial \mathbf{f}}{\partial \boldsymbol{\xi}}(\mathbf{x}, t^n)$  and so the discrete-optimal test basis  $\boldsymbol{\Psi}^n$  defined in Eq. (24) is equal to the test basis in Eq. (34) evaluated at time instance  $n$ , i.e.,  $\boldsymbol{\Psi}^n(\hat{\mathbf{w}}) = \boldsymbol{\Psi}(\hat{\mathbf{w}}, t^n)$ . Therefore, the solution  $\hat{\mathbf{w}}^n$  to the discrete-optimal O $\Delta$ E (23) satisfies

$$\boldsymbol{\Psi}(\hat{\mathbf{w}}^n, t^n)^T \mathbf{r}^n(\mathbf{x}_0 + \boldsymbol{\Phi} \hat{\mathbf{w}}^n) = 0. \quad (40)$$

This shows that  $\hat{\mathbf{w}}^n = \hat{\mathbf{y}}^n$ , i.e., the solutions to the discrete-optimal O $\Delta$ E and the O $\Delta$ E obtained after applying Petrov–Galerkin projection with test basis  $\boldsymbol{\Psi}(\mathbf{x}, t)$  defined by Eq. (34) to the full-order model ODE and subsequently applying time integration are equivalent under the stated assumptions, which is the desired result.

Case 2 In this case, the test basis is independent of the state, i.e.,

$$\boldsymbol{\Psi}(t) = \mathbf{A}^T \mathbf{A} \left( \alpha_0 \mathbf{I} - \Delta t \beta_0 \frac{\partial \mathbf{f}}{\partial \boldsymbol{\xi}}(\cdot, t) \right) \boldsymbol{\Phi}. \quad (41)$$

Applying a linear multistep time integrator to solve Eq. (36) and subsequently pre-multiplying by the constant matrix  $\boldsymbol{\Psi}(t^n)^T \boldsymbol{\Phi}$  yields the following discrete equations arising at each time step

$$\hat{\mathbf{r}}^n(\hat{\mathbf{y}}^n) = 0, \quad (42)$$

where the residual is defined as

$$\begin{aligned} \hat{\mathbf{r}}^n(\hat{\mathbf{w}}) := & \alpha_0 \boldsymbol{\Psi}(t^n)^T \boldsymbol{\Phi} \hat{\mathbf{w}} - \Delta t \beta_0 \boldsymbol{\Psi}(t^n)^T \mathbf{f}(\mathbf{x}_0 + \boldsymbol{\Phi} \hat{\mathbf{w}}, t^n) + \sum_{j=1}^k \alpha_j \boldsymbol{\Psi}(t^n)^T \boldsymbol{\Phi} \hat{\mathbf{x}}^{n-j} - \\ & \Delta t \sum_{j=1}^k \beta_j \boldsymbol{\Psi}(t^n)^T \mathbf{f}(\mathbf{x}_0 + \boldsymbol{\Phi} \hat{\mathbf{x}}^{n-j}, t^{n-j}). \end{aligned} \quad (43)$$

Comparing Eqs. (43) and (4) reveals  $\hat{\mathbf{r}}^n(\hat{\mathbf{w}}) = \boldsymbol{\Psi}(t^n)^T \mathbf{r}^n(\mathbf{x}_0 + \boldsymbol{\Phi} \hat{\mathbf{w}})$  and so the solution  $\hat{\mathbf{y}}^n$  satisfies

$$\boldsymbol{\Psi}(t^n)^T \mathbf{r}^n(\mathbf{x}_0 + \boldsymbol{\Phi} \hat{\mathbf{y}}^n) = 0. \quad (44)$$

Under these assumptions, we have  $\partial \mathbf{r}^n / \partial \mathbf{w} = \alpha_0 \mathbf{I} - \Delta t \beta_0 \partial \mathbf{f} / \partial \boldsymbol{\xi}(\cdot, t^n)$  and so the discrete-optimal test basis  $\boldsymbol{\Psi}^n$  defined in Eq. (24) is equal to the test basis in Eq. (41) at time instance  $n$ , i.e.,  $\boldsymbol{\Psi}^n(\hat{\mathbf{w}}) = \boldsymbol{\Psi}(t^n)$ . Therefore, the discrete-optimal O $\Delta$ E (23) can be expressed as

$$\boldsymbol{\Psi}(t^n)^T \mathbf{r}^n(\mathbf{x}_0 + \boldsymbol{\Phi} \hat{\mathbf{w}}^n) = 0. \quad (45)$$

This shows that  $\hat{\mathbf{w}}^n = \hat{\mathbf{y}}^n$ , i.e., the solutions to the discrete-optimal O $\Delta$ E and the O $\Delta$ E obtained after applying Petrov–Galerkin projection with test basis  $\boldsymbol{\Psi}(t)$  defined by Eq. (34) to the full-order model ODE and subsequently applying time integration are equivalent under the stated assumptions.

Case 3 The assumption  $\beta_0 = 0$  results in a constant test basis

$$\boldsymbol{\Psi} = \alpha_0 \mathbf{A}^T \mathbf{A} \boldsymbol{\Phi}. \quad (46)$$

Applying a linear multistep time integrator to solve Eq. (36) and subsequently pre-multiplying by the constant matrix  $\boldsymbol{\Psi}^T \boldsymbol{\Phi}$  yields

$$\hat{\mathbf{r}}^n(\hat{\mathbf{y}}^n) = 0, \quad (47)$$

which is to be solved at each time step with a residual defined as

$$\hat{\mathbf{r}}^n(\hat{\mathbf{w}}) := \alpha_0 \boldsymbol{\Psi}^T \boldsymbol{\Phi} \hat{\mathbf{w}} - \Delta t \beta_0 \boldsymbol{\Psi}^T \mathbf{f}(\mathbf{x}_0 + \boldsymbol{\Phi} \hat{\mathbf{w}}, t^n) + \sum_{j=1}^k \alpha_j \boldsymbol{\Psi}^T \boldsymbol{\Phi} \hat{\mathbf{x}}^{n-j} - \Delta t \sum_{j=1}^k \beta_j \boldsymbol{\Psi}^T \mathbf{f}(\mathbf{x}_0 + \boldsymbol{\Phi} \hat{\mathbf{x}}^{n-j}, t^{n-j}). \quad (48)$$

As in Case 2, this leads to  $\hat{\mathbf{r}}^n(\hat{\mathbf{w}}) = \boldsymbol{\Psi}^T \mathbf{r}^n(\mathbf{x}_0 + \boldsymbol{\Phi} \hat{\mathbf{w}})$ . Because  $\frac{\partial \mathbf{r}^n}{\partial \mathbf{w}}(\mathbf{x}) = \alpha_0 \mathbf{I}$ , we also again have  $\boldsymbol{\Psi}^n(\hat{\mathbf{w}}) = \boldsymbol{\Psi}$ . This leads to the desired result, as the O $\Delta$ Es for the discrete-optimal ROM and the ROM obtained after applying Petrov–Galerkin projection with test basis  $\boldsymbol{\Psi}$  to the full-order model ODE and subsequently applying time integration both satisfy  $\boldsymbol{\Psi}^T \mathbf{r}^n(\mathbf{x}_0 + \boldsymbol{\Phi} \hat{\mathbf{w}}^n) = 0$  under the stated assumptions.

We now provide conditions under which the discrete-optimal ROM for Runge–Kutta schemes can be expressed as an ODE.

**Theorem 4.3 (Discrete-optimal ROM continuous representation: Runge–Kutta schemes).** *The discrete-optimal ROM for linear multistep integrators is equivalent to applying a Petrov–Galerkin projection to the ODE with test basis (in matrix form)*

$$\Psi(\hat{\mathbf{x}}, t) = \mathbf{A}^T \mathbf{A} \left( \mathbf{I} - \Delta t a_{11} \frac{\partial \mathbf{f}}{\partial \xi}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t) \right) \Phi \quad (49)$$

and subsequently applying time integration if  $\mathbf{A}_i = \mathbf{A} \forall i$  are constant matrices and either

1.  $a_{ij} = 0 \forall i \neq j$  and  $a_{ii} = a_{jj} \forall i, j$ , or
2. the scheme is explicit, i.e.,  $a_{ij} = 0, \forall j \geq i$ .

PROOF. Case 1 Applying Petrov–Galerkin projection to Eq. (1) using a trial subspace  $\mathbf{x}_0 + \text{range}(\Phi)$  and test subspace  $\text{range}(\Psi)$  yields the following ODE (in standard form)

$$\frac{d\hat{\mathbf{x}}}{dt} = \left( \Psi(\hat{\mathbf{x}}, t)^T \Phi \right)^{-1} \Psi(\hat{\mathbf{x}}, t)^T \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t), \quad \hat{\mathbf{x}}(0) = 0. \quad (50)$$

Applying a Runge–Kutta time integrator with  $a_{ij} = 0 \forall i \neq j$  and  $a_{ii} = a_{jj} \forall i, j$  to numerically solve Eq. (50) results in the following discrete equations to be solved at each time step

$$\begin{aligned} \hat{\mathbf{y}}_i^n - \left( \Psi(\hat{\mathbf{x}}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \hat{\mathbf{y}}_j^n, t^{n-1} + c_i \Delta t)^T \Phi \right)^{-1} \Psi(\hat{\mathbf{x}}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \hat{\mathbf{y}}_j^n, t^{n-1} + c_i \Delta t)^T \\ \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \Phi \hat{\mathbf{y}}_j^n, t^{n-1} + c_i \Delta t) = 0, \quad i \in \mathbb{N}(s) \end{aligned} \quad (51)$$

Pre-multiplying by  $\Psi(\hat{\mathbf{x}}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \hat{\mathbf{y}}_j^n, t^{n-1} + c_i \Delta t)^T \Phi$  yields the following discrete equations

$$\hat{\mathbf{r}}_i^n(\hat{\mathbf{y}}_1^n, \dots, \hat{\mathbf{y}}_s^n) = 0, \quad (52)$$

with residual

$$\begin{aligned} \hat{\mathbf{r}}_i^n(\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_s) := \Psi(\hat{\mathbf{x}}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \hat{\mathbf{w}}_j, t^{n-1} + c_i \Delta t)^T \Phi \hat{\mathbf{w}}_i - \\ \Psi(\hat{\mathbf{x}}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \hat{\mathbf{w}}_j, t^{n-1} + c_i \Delta t)^T \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \Phi \hat{\mathbf{w}}_j, t^{n-1} + c_i \Delta t) = 0, \quad i \in \mathbb{N}(s) \end{aligned} \quad (53)$$

Comparing Eqs. (53) and (6) reveals

$$\hat{\mathbf{r}}_i^n(\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_s) = \Psi(\hat{\mathbf{x}}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \hat{\mathbf{w}}_j, t^{n-1} + c_i \Delta t)^T \mathbf{r}_i^n(\Phi \hat{\mathbf{w}}_1, \dots, \Phi \hat{\mathbf{w}}_s), \quad i \in \mathbb{N}(s)$$

such that the solution  $(\hat{\mathbf{y}}_1^n, \dots, \hat{\mathbf{y}}_s^n)$  satisfies

$$\Psi(\hat{\mathbf{x}}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \hat{\mathbf{y}}_j^n, t^{n-1} + c_i \Delta t)^T \mathbf{r}_i^n(\Phi \hat{\mathbf{y}}_1^n, \dots, \Phi \hat{\mathbf{y}}_s^n) = 0 \quad (54)$$

Under the stated assumptions, we have

$$\frac{\partial \mathbf{r}_i^n}{\partial \mathbf{w}_j}(\mathbf{u}_1, \dots, \mathbf{u}_s) = \begin{cases} \mathbf{I} - \Delta t a_{ii} \frac{\partial \mathbf{f}}{\partial \xi}(\mathbf{x}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \mathbf{u}_j, t^{n-1} + c_i \Delta t), & i = j \\ 0, & \text{otherwise} \end{cases}$$

such that the discrete-optimal test basis  $\Psi_{ij}^n$  defined in Eq. (28) is relates to the test basis in Eq. (49) as follows:

$$\Psi_{ij}^n(\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_s) = \begin{cases} \Psi(\hat{\mathbf{x}}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \hat{\mathbf{w}}_j, t^{n-1} + c_i \Delta t), & i = j \\ 0, & \text{otherwise.} \end{cases}$$

Therefore, the solution  $(\hat{\mathbf{w}}_1^n, \dots, \hat{\mathbf{w}}_s^n)$  to the discrete-optimal O $\Delta$ E (27) satisfies

$$\Psi(\hat{\mathbf{x}}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \hat{\mathbf{w}}_j^n, t^{n-1} + c_i \Delta t)^T \mathbf{r}_i^n(\Phi \hat{\mathbf{w}}_1^n, \dots, \Phi \hat{\mathbf{w}}_s^n) = 0. \quad (55)$$

This shows that the  $\hat{\mathbf{w}}_i^n = \hat{\mathbf{y}}_i^n$ ,  $i \in \mathbb{N}(s)$ , i.e., the solutions to the discrete-optimal O $\Delta$ E and the O $\Delta$ E obtained after applying Petrov–Galerkin projection with test basis  $\Psi(\mathbf{x}, t)$  defined by Eq. (49) to the full-order model ODE and subsequently applying time integration are equivalent under the stated assumptions, which is the desired result.

Case 2 Explicit schemes are characterized by  $a_{ii} = 0$  and therefore result in a constant test basis of

$$\Psi = \mathbf{A}^T \mathbf{A} \Phi.$$

Applying Petrov–Galerkin projection to Eq. (1) using a trial subspace  $\mathbf{x}_0 + \text{range}(\Phi)$  and test subspace  $\text{range}(\Psi)$  yields the following ODE (in standard form)

$$\frac{d\hat{\mathbf{x}}}{dt} = \left( \Psi^T \Phi \right)^{-1} \Psi^T \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t), \quad \hat{\mathbf{x}}(0) = 0. \quad (56)$$

Applying an explicit Runge–Kutta time integrator and pre-multiplying the residual by the constant matrix  $\Psi^T \Phi$  results in the following sequence of discrete equations to be solved at each time step

$$\hat{\mathbf{q}}_i^n(\hat{\mathbf{y}}_i^n) = 0, \quad i \in \mathbb{N}(s)$$

with residual

$$\hat{\mathbf{q}}_i^n(\hat{\mathbf{w}}_i) := \Psi^T \Phi \hat{\mathbf{w}}_i - \Psi^T \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}^{n-1} + \Delta t \sum_{j=1}^{i-1} a_{ij} \Phi \hat{\mathbf{w}}_j^n, t^{n-1} + c_i \Delta t), \quad i \in \mathbb{N}(s) \quad (57)$$

Comparing Eqs. (57) and (29) reveals  $\hat{\mathbf{q}}_i^n(\hat{\mathbf{w}}_i) = \Psi^T \mathbf{q}_i^n(\Phi \hat{\mathbf{w}}_i)$ ,  $i \in \mathbb{N}(s)$ . Therefore, the solutions  $\hat{\mathbf{y}}_i^n$  satisfies

$$\Psi^T \mathbf{q}_i^n(\Phi \hat{\mathbf{y}}_i^n) = 0, \quad i \in \mathbb{N}(s).$$

Under the stated assumptions, the weighting matrices are equal and constant  $\mathbf{A}_i(\Phi \hat{\mathbf{z}}) = \mathbf{A}$ ,  $\forall i$  and such that the discrete-optimal test basis defined in Eq. (33) is equal to the constant test basis above, i.e.,  $\Psi_i^n = \Psi = \mathbf{A}^T \mathbf{A} \Phi$ . Therefore, the solution  $\hat{\mathbf{w}}_i^n$  to the discrete-optimal O $\Delta$ E (32) satisfies

$$\Psi^T \mathbf{q}_i^n(\Phi \hat{\mathbf{w}}_i^n) = 0, \quad i \in \mathbb{N}(s).$$

This shows that  $\hat{\mathbf{w}}_i^n = \hat{\mathbf{y}}_i^n$ ,  $i \in \mathbb{N}(s)$ , i.e., the solutions to the discrete-optimal O $\Delta$ E and the O $\Delta$ E obtained after applying Petrov–Galerkin projection with test basis  $\Psi(\mathbf{x}, t)$  defined by Eq. (49) to the full-order model ODE and subsequently applying time integration are equivalent under the stated assumptions, which is the desired result.

We now show that the discrete-optimal ROM has a time-continuous representation for all single-state Runge–Kutta schemes.

**Corollary 4.4 (Discrete-optimal ROM continuous representation: single-stage Runge–Kutta).**

*The discrete-optimal ROM for linear multistep integrators is equivalent to applying a Petrov–Galerkin projection to the ODE with test basis (in matrix form)*

$$\Psi(\hat{\mathbf{x}}, t) = \mathbf{A}^T \mathbf{A} \left( \mathbf{I} - \Delta t a_{11} \frac{\partial \mathbf{f}}{\partial \xi}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}, t) \right) \Phi$$

*and subsequently applying time integration if  $\mathbf{A}_i = \mathbf{A} \forall i$  are constant matrices and a single-stage Runge–Kutta scheme is employed.*

PROOF. Single-stage Runge–Kutta schemes are characterized by  $s = 1$  and so they satisfy the conditions of case 1 of Theorem 4.3.

## 5. Equivalence conditions

This section performs theoretical analysis that highlights cases in which Galerkin and discrete-optimal ROMs are equivalent. Section 5.1 shows that equivalence holds for explicit time integrators, Section 5.2 demonstrates equivalence in the limit of  $\Delta t \rightarrow 0$ , and Section 5.3 shows equivalence in the case of symmetric-positive-definite residual Jacobians.

### 5.1. Equivalence for explicit integrators

**Corollary 5.1 (Galerkin discrete optimality: explicit linear multistep scheme).**

*Galerkin projection is discrete optimal for explicit linear multistep schemes.*

PROOF. In the case of explicit linear multistep schemes,  $\beta_0 = 0$  and so Galerkin projection corresponds to Case 3 of Theorem 4.2 with  $\mathbf{A} = \frac{1}{\sqrt{\alpha_0}}\mathbf{I}$ , as  $\Psi = \Phi$  in this case.

**Corollary 5.2 (Galerkin discrete optimality: explicit Runge–Kutta scheme).** *Galerkin projection is discrete optimal for explicit Runge–Kutta schemes.*

PROOF. In the case of explicit Runge–Kutta schemes,  $a_{11} = 0$  and so Galerkin projection corresponds to a special of Case 2 of Theorem 4.3 with  $\mathbf{A} = \mathbf{I}$ , as  $\Psi = \Phi$  in this case.

### 5.2. Equivalence in the limit of $\Delta t \rightarrow 0$

**Theorem 5.3 (Limiting equivalence of Galerkin and discrete-optimal ROMs).**

*In the limit of  $\Delta t \rightarrow 0$ , continuous-optimal Galerkin ROMs are also discrete optimal.*

PROOF. Linear multistep schemes. Consider solving the discrete-optimal OΔE (23) with  $\mathbf{A} = \frac{1}{\sqrt{\alpha_0}}\mathbf{I}$ . Then, the test basis defined in Eq. (24) is simply

$$\Psi^n(\hat{\mathbf{w}}) = \frac{1}{\alpha_0} \frac{\partial \mathbf{r}^n}{\partial \mathbf{w}} (\mathbf{x}_0 + \Phi \hat{\mathbf{w}}) \Phi.$$

From Eq. (4), we can write the residual Jacobian as

$$\frac{\partial \mathbf{r}^n}{\partial \mathbf{w}} (\mathbf{u}) = \alpha_0 \mathbf{I} - \Delta t \beta_0 \frac{\partial \mathbf{f}}{\partial \boldsymbol{\xi}} (\mathbf{u}, t^n).$$

Therefore, we have

$$\lim_{\Delta t \rightarrow 0} \Psi^n(\hat{\mathbf{w}}) = \lim_{\Delta t \rightarrow 0} \frac{1}{\alpha_0} \left( \alpha_0 \mathbf{I} - \Delta t \beta_0 \frac{\partial \mathbf{f}}{\partial \boldsymbol{\xi}} (\mathbf{x}_0 + \Phi \hat{\mathbf{w}}, t^n) \right) \Phi = \Phi$$

and so in the limit of  $\Delta t \rightarrow 0$ , the discrete-optimal ROM solution satisfies

$$\lim_{\Delta t \rightarrow 0} \Psi^n(\hat{\mathbf{w}})^T \mathbf{r}^n (\mathbf{x}_0 + \Phi \hat{\mathbf{w}}^n) = \Phi^T \mathbf{r}^n (\mathbf{x}_0 + \Phi \hat{\mathbf{w}}^n) = 0. \quad (58)$$

Because the Galerkin ROM solution also satisfies Eq. (58) (see Eq. (19) of Theorem 3.4), the two techniques are equivalent in this limit, which is the desired result.

Runge–Kutta schemes. Consider solving the discrete-optimal OΔE (27) with  $\mathbf{A}_i = \mathbf{I}$ ,  $i \in \mathbb{N}(s)$ . Then, the test basis defined in Eq. (28) is simply

$$\Psi_{ij}^n(\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_s) = \frac{\partial \mathbf{r}_i^n}{\partial \mathbf{w}_j} (\Phi \hat{\mathbf{w}}_1, \dots, \Phi \hat{\mathbf{w}}_s) \Phi.$$

Now, from Eq. (6) the Jacobian can be expressed as

$$\frac{\partial \mathbf{r}_i^n}{\partial \mathbf{w}_j} (\mathbf{u}_1, \dots, \mathbf{u}_s) = \mathbf{I} \delta_{ij} - \Delta t a_{ij} \frac{\partial \mathbf{f}}{\partial \boldsymbol{\xi}} (\mathbf{x}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \mathbf{u}_j, t^{n-1} + c_i \Delta t).$$

Therefore, we have

$$\lim_{\Delta t \rightarrow 0} \Psi_{ij}^n(\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_s) = \lim_{\Delta t \rightarrow 0} \left( \mathbf{I} \delta_{ij} - \Delta t a_{ij} \frac{\partial \mathbf{f}}{\partial \boldsymbol{\xi}}(\mathbf{x}^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \mathbf{u}_j, t^{n-1} + c_i \Delta t) \right) \boldsymbol{\Phi} = \boldsymbol{\Phi} \delta_{ij}$$

and so in the limit of  $\Delta t \rightarrow 0$ , the discrete-optimal ROM solution satisfies

$$\lim_{\Delta t \rightarrow 0} \sum_{j=1}^s \Psi_{ij}^n(\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_s)^T \mathbf{r}_j^n(\boldsymbol{\Phi} \hat{\mathbf{w}}_1^n, \dots, \boldsymbol{\Phi} \hat{\mathbf{w}}_s^n) = \boldsymbol{\Phi}^T \mathbf{r}_i^n(\mathbf{x}_0 + \boldsymbol{\Phi} \hat{\mathbf{w}}^n) = 0, \quad i \in \mathbb{N}(s). \quad (59)$$

Because the Galerkin ROM solution also satisfies Eq. (59) (see Eq. (20) of Theorem 3.4), the two techniques are equivalent in this limit, which is the desired result.

### 5.3. Equivalence for symmetric-positive-definite residual Jacobians

**Theorem 5.4 (Galerkin discrete optimality: linear multistep schemes).** *In the case of linear multistep schemes, Galerkin projection satisfies Eq. (21) (i.e., exhibits discrete optimality) with  $\mathbf{A}(\mathbf{z}) = \mathbf{U}(\mathbf{z})$ , where  $\mathbf{U}$  is the Cholesky factor<sup>1</sup> of the residual-Jacobian inverse*

$$\left[ \frac{\partial \mathbf{r}^n}{\partial \mathbf{w}} \right]^{-1} = \mathbf{U}^T \mathbf{U}, \quad (60)$$

if  $\partial \mathbf{r}^n / \partial \mathbf{w}(\mathbf{w}^n, t^n) = \alpha_0 \mathbf{I} - \Delta t \beta_0 \frac{\partial \mathbf{f}}{\partial \boldsymbol{\xi}}(\mathbf{w}^n, t^n)$  is symmetric positive definite and if

$$\frac{\partial u_{i\ell}}{\partial w_k} \phi_{kj} r_\ell^n = 0, \quad \forall i, k. \quad (61)$$

Here, index notation has been used.

PROOF. Under the stated assumptions, the discrete-optimal test basis defined in Eq. (24) is equal to the trial basis, i.e.,  $\Psi^n(\hat{\mathbf{w}}^n) = \boldsymbol{\Phi}$ . By invoking Eq. (19), we can see that the ODEs for the the discrete-optimal ROM (23) and Galerkin ROM (14) both satisfy  $\boldsymbol{\Phi}^T \mathbf{r}^n(\mathbf{x}_0 + \boldsymbol{\Phi} \hat{\mathbf{w}}^n) = 0$ , which is the desired result.

**Theorem 5.5 (Galerkin discrete optimality: Runge–Kutta schemes).** *In the case of Runge–Kutta schemes, Galerkin projection exhibits discrete optimality if  $\partial \bar{\mathbf{r}}^n / \partial \bar{\mathbf{w}}(\bar{\mathbf{w}}^n, t^n)$  is symmetric positive definite and if*

$$\frac{\partial \bar{u}_{i\ell}}{\partial \bar{w}_k} \bar{\phi}_{kj} \bar{r}_\ell^n = 0, \quad \forall i, k. \quad (62)$$

Here, index notation has been used and  $\bar{\mathbf{U}}$  is the Cholesky factor of the residual-Jacobian inverse, i.e.,

$$\left[ \frac{\partial \bar{\mathbf{r}}^n}{\partial \bar{\mathbf{w}}} \right]^{-1} = \bar{\mathbf{U}}^T \bar{\mathbf{U}}. \quad (63)$$

Here,

$$\bar{\mathbf{w}} := \begin{bmatrix} \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_s \end{bmatrix} \in \mathbb{R}^{sN}, \quad \bar{\mathbf{r}}^n : \bar{\mathbf{w}} \mapsto \begin{bmatrix} \mathbf{r}_1^n(\mathbf{w}_1, \dots, \mathbf{w}_s) \\ \vdots \\ \mathbf{r}_s^n(\mathbf{w}_1, \dots, \mathbf{w}_s) \end{bmatrix} \in \mathbb{R}^{sN}, \quad \bar{\boldsymbol{\Phi}} := \begin{bmatrix} \boldsymbol{\Phi} & & \\ & \ddots & \\ & & \boldsymbol{\Phi} \end{bmatrix} \in \mathbb{R}^{sN \times sp}.$$

<sup>1</sup> Its derivative can be computed by solving the Lyapunov equation  $\frac{\partial \mathbf{U}}{\partial w_k}^T \mathbf{U} + \mathbf{U} \frac{\partial \mathbf{U}}{\partial w_k} = - \left[ \frac{\partial \mathbf{r}^n}{\partial \mathbf{w}} \right]^{-1} \frac{\partial^2 \mathbf{r}^n}{\partial \mathbf{w} \partial w_k} \left[ \frac{\partial \mathbf{r}^n}{\partial \mathbf{w}} \right]^{-1}$ .

PROOF. First, note that solution  $(\hat{\boldsymbol{w}}_1^n, \dots, \hat{\boldsymbol{w}}_s^n)$  to the Galerkin OΔE (20) equivalently satisfies

$$\bar{\boldsymbol{\Phi}}^T \bar{\boldsymbol{r}}^n (\bar{\boldsymbol{\Phi}} \hat{\boldsymbol{w}}^n) = 0,$$

where

$$\hat{\boldsymbol{w}} := \begin{bmatrix} \hat{\boldsymbol{w}}_1 \\ \vdots \\ \hat{\boldsymbol{w}}_s \end{bmatrix} \in \mathbb{R}^{sp}.$$

We are now precisely in the situation of Theorem 5.4: the Galerkin solution is the solution to the (discrete) optimization problem

$$\underset{\boldsymbol{z} \in \text{range}(\bar{\boldsymbol{\Phi}})}{\text{minimize}} \|\bar{\boldsymbol{U}}(\boldsymbol{z}) \bar{\boldsymbol{r}}^n(\boldsymbol{z})\|_2^2 \quad (64)$$

under the assumed conditions.

## 6. Error analysis

This section performs time-discrete state-space error analyses for Galerkin and discrete-optimal ROMs applied to different time integrators.

### 6.1. Linear multistep schemes

Here, we perform error analysis for implicit linear multistep schemes. We will use subscripts  $*$ ,  $G$  and  $D$  to denote the solution to full-order model OΔE (3), Galerkin ROM OΔE (14), and the discrete-optimal ROM OΔE (23), respectively. We also define  $\boldsymbol{\Psi}^n := \boldsymbol{\Psi}^n(\hat{\boldsymbol{x}}_D^n)$  whose entries are defined by Eq. (24).

$$\alpha_0 \boldsymbol{x}_*^n = \beta_0 \Delta t \boldsymbol{f}(\boldsymbol{x}_0 + \boldsymbol{x}_*^n, t^n) + \boldsymbol{r}_* [\boldsymbol{x}_*^{n-k}, \dots, \boldsymbol{x}_*^{n-1}], \quad \boldsymbol{x}_*^0 = \mathbf{0} \quad (65)$$

$$\alpha_0 \hat{\boldsymbol{x}}_G^n = \beta_0 \Delta t \bar{\boldsymbol{\Phi}}^T \boldsymbol{f}(\boldsymbol{x}_0 + \bar{\boldsymbol{\Phi}} \hat{\boldsymbol{x}}_G^n, t^n) + \hat{\boldsymbol{r}}_G [\hat{\boldsymbol{x}}_G^{n-k}, \dots, \hat{\boldsymbol{x}}_G^{n-1}], \quad \hat{\boldsymbol{x}}_G^0 = \mathbf{0} \quad (66)$$

$$\alpha_0 \hat{\boldsymbol{x}}_D^n = \beta_0 \Delta t \left( (\boldsymbol{\Psi}^n)^T \bar{\boldsymbol{\Phi}} \right)^{-1} (\boldsymbol{\Psi}^n)^T \boldsymbol{f}(\boldsymbol{x}_0 + \bar{\boldsymbol{\Phi}} \hat{\boldsymbol{x}}_D^n, t^n) + \hat{\boldsymbol{r}}_D^n [\hat{\boldsymbol{x}}_D^{n-k}, \dots, \hat{\boldsymbol{x}}_D^{n-1}], \quad \hat{\boldsymbol{x}}_D^0 = \mathbf{0}, \quad (67)$$

where

$$\begin{aligned} \boldsymbol{r}_* [\boldsymbol{x}^{n-k}, \dots, \boldsymbol{x}^{n-1}] &:= \sum_{\ell=1}^k \left( \beta_\ell \Delta t \boldsymbol{f}(\boldsymbol{x}_0 + \boldsymbol{x}^{n-\ell}, t^{n-\ell}) - \alpha_\ell \boldsymbol{x}^{n-\ell} \right) \\ \hat{\boldsymbol{r}}_G [\hat{\boldsymbol{x}}^{n-k}, \dots, \hat{\boldsymbol{x}}^{n-1}] &:= \sum_{\ell=1}^k \left( \beta_\ell \Delta t \bar{\boldsymbol{\Phi}}^T \boldsymbol{f}(\boldsymbol{x}_0 + \bar{\boldsymbol{\Phi}} \hat{\boldsymbol{x}}^{n-\ell}, t^{n-\ell}) - \alpha_\ell \hat{\boldsymbol{x}}^{n-\ell} \right) \\ \hat{\boldsymbol{r}}_D^n [\hat{\boldsymbol{x}}^{n-k}, \dots, \hat{\boldsymbol{x}}^{n-1}] &:= \sum_{\ell=1}^k \left( \beta_\ell \Delta t \left( (\boldsymbol{\Psi}^n)^T \bar{\boldsymbol{\Phi}} \right)^{-1} (\boldsymbol{\Psi}^n)^T \boldsymbol{f}(\boldsymbol{x}_0 + \bar{\boldsymbol{\Phi}} \hat{\boldsymbol{x}}^{n-\ell}, t^{n-\ell}) - \alpha_\ell \hat{\boldsymbol{x}}^{n-\ell} \right). \end{aligned} \quad (68)$$

We define the Galerkin and least-squares Petrov–Galerkin operators as

$$\mathbb{V} := \bar{\boldsymbol{\Phi}} \bar{\boldsymbol{\Phi}}^T, \quad \text{and} \quad \mathbb{P}^n := \bar{\boldsymbol{\Phi}} \left( (\boldsymbol{\Psi}^n)^T \bar{\boldsymbol{\Phi}} \right)^{-1} (\boldsymbol{\Psi}^n)^T,$$

respectively, and Galerkin and discrete-optimal state-space errors at time instance  $n$  as

$$\delta \boldsymbol{x}_G^n := \boldsymbol{x}_*^n - \bar{\boldsymbol{\Phi}} \hat{\boldsymbol{x}}_G^n, \quad \text{and} \quad \delta \boldsymbol{x}_D^n := \boldsymbol{x}_*^n - \bar{\boldsymbol{\Phi}} \hat{\boldsymbol{x}}_D^n,$$

respectively. As the second argument in  $\boldsymbol{f}$  does not play any role in this section, will drop it for notational convenience. Moreover, we assume Lipschitz continuity of  $\boldsymbol{f}$  in the first argument:

(A<sub>1</sub>) There exist a constant  $\kappa > 0$  such that for  $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^N$

$$\|\boldsymbol{f}(\boldsymbol{x}) - \boldsymbol{f}(\boldsymbol{y})\| \leq \kappa \|\boldsymbol{x} - \boldsymbol{y}\|.$$

**Theorem 6.1.** *If  $(\mathbf{A}_1)$  holds and  $\Delta t$  is such that  $0 < |\alpha_0| - \Delta t|\beta_0|\kappa$ , then*

$$\|\delta \mathbf{x}_G^n\| \leq \frac{\Delta t}{h} \sum_{\ell=0}^k |\beta_\ell| \left\| (\mathbf{I} - \mathbb{V}) \mathbf{f} \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_G^{n-\ell} \right) \right\| + \frac{1}{h} \sum_{\ell=1}^k (|\beta_\ell| \kappa \Delta t + |\alpha_\ell|) \|\delta \mathbf{x}_G^{n-\ell}\|, \quad (69)$$

$$\|\delta \mathbf{x}_D^n\| \leq \frac{\Delta t}{h} \sum_{\ell=0}^k |\beta_\ell| \left\| (\mathbf{I} - \mathbb{P}^n) \mathbf{f} \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^{n-\ell} \right) \right\| + \frac{1}{h} \sum_{\ell=1}^k (|\beta_\ell| \kappa \Delta t + |\alpha_\ell|) \|\delta \mathbf{x}_D^{n-\ell}\|, \quad (70)$$

where  $h := |\alpha_0| - |\beta_0|\kappa\Delta t$ .

PROOF. It is enough to show bound (70), as the arguments for (69) are similar. Let  $n$  be fixed but arbitrary, then subtracting Eq. (67) from Eq. (65) yields

$$|\alpha_0| \|\delta \mathbf{x}_D^n\| \leq |\beta_0| \Delta t \left\| \mathbf{f} \left( \mathbf{x}_0 + \mathbf{x}_*^n \right) - \mathbb{P}^n \mathbf{f} \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^n \right) \right\| + \left\| \delta \mathbf{r}_D^{n-1} \right\|,$$

where  $\delta \mathbf{r}_D^{n-1} := \mathbf{r}_* [\mathbf{x}_*^{n-k}, \dots, \mathbf{x}_*^{n-1}] - \Phi \hat{\mathbf{r}}_D^n [\hat{\mathbf{x}}_D^{n-k}, \dots, \hat{\mathbf{x}}_D^{n-1}]$ . Adding and subtracting  $\mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^n)$  and applying the triangle inequality leads to

$$|\alpha_0| \|\delta \mathbf{x}_D^n\| \leq |\beta_0| \Delta t \left( \left\| (\mathbf{I} - \mathbb{P}^n) \mathbf{f} \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^n \right) \right\| + \left\| \mathbf{f} \left( \mathbf{x}_0 + \mathbf{x}_*^n \right) - \mathbf{f} \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^n \right) \right\| \right) + \left\| \delta \mathbf{r}_D^{n-1} \right\|.$$

Invoking  $(\mathbf{A}_1)$ , and using  $\Delta t < \frac{|\alpha_0|}{|\beta_0|\kappa}$ , we deduce

$$\|\delta \mathbf{x}_D^n\| \leq \frac{|\beta_0| \Delta t}{h} \left\| (\mathbf{I} - \mathbb{P}^n) \mathbf{f} \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^n \right) \right\| + \frac{1}{h} \left\| \delta \mathbf{r}_D^{n-1} \right\|. \quad (71)$$

Next, we will estimate  $\left\| \delta \mathbf{r}_D^{n-1} \right\|$ . Using the definition of  $\mathbf{r}_*$ ,  $\hat{\mathbf{r}}_D^n$  from (68) we derive

$$\left\| \delta \mathbf{r}_D^{n-1} \right\| \leq \sum_{\ell=1}^k \left( |\beta_\ell| \Delta t \left\| \mathbf{f} \left( \mathbf{x}_0 + \mathbf{x}_*^{n-\ell} \right) - \mathbb{P}^n \mathbf{f} \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^{n-\ell} \right) \right\| + |\alpha_\ell| \left\| \delta \mathbf{x}_D^{n-\ell} \right\| \right).$$

Adding and subtracting  $\mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^{n-\ell})$ , applying the triangle inequality in conjunction with  $(\mathbf{A}_1)$  yields

$$\left\| \delta \mathbf{r}_D^{n-1} \right\| \leq \sum_{\ell=1}^k |\beta_\ell| \Delta t \left\| (\mathbf{I} - \mathbb{P}^n) \mathbf{f} \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^{n-\ell} \right) \right\| + \sum_{\ell=1}^k (|\beta_\ell| \kappa \Delta t + |\alpha_\ell|) \left\| \delta \mathbf{x}_D^{n-\ell} \right\|. \quad (72)$$

Then (71) and (72) implies (70).

## 6.2. Runge-Kutta schemes

For simplicity, we state the error estimate only for the Galerkin ROM ODE (17). We will use subscript  $i$  on  $\mathbf{f}$  to indicate dependence of  $\mathbf{f}$  on  $i$  in the second argument. Since the second argument in  $\mathbf{f}$  does not play another further role, in this section, we will suppress it for notational simplicity.

We rewrite Eqs. (5) and (16) as

$$\mathbf{w}_{*,i}^n = \mathbf{f}_i \left( \mathbf{x}_0 + \mathbf{x}_*^{n-1} + \Delta t \sum_j a_{ij} \mathbf{w}_{*,j}^n \right), \quad i \in \mathbb{N}(s) \quad \mathbf{x}_*^0 = \mathbf{0} \quad (73)$$

$$\hat{\mathbf{w}}_{G,i}^n = \Phi^T \mathbf{f}_i \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_G^{n-1} + \Delta t \sum_j a_{ij} \Phi \hat{\mathbf{w}}_{G,j}^n \right), \quad i \in \mathbb{N}(s) \quad \mathbf{x}_G^0 = \mathbf{0}. \quad (74)$$

**Theorem 6.2.** *If  $(\mathbf{A}_1)$  holds and  $\Delta t$  is such that*

(a) *the matrix  $\mathbf{D} \in \mathbb{R}^{s \times s}$  with entries  $d_{ij} := \delta_{ij} - \kappa \Delta t |a_{ij}|$  is invertible, and*

(b) for every  $\mathbf{x}, \mathbf{y} \geq 0$ , if  $\mathbf{D}\mathbf{x} \leq \mathbf{y}$  then  $\mathbf{x} \leq \mathbf{D}^{-1}\mathbf{y}$ ,

then

$$\|\delta \mathbf{x}_G^n\| \leq \Delta t \sum_{\ell=0}^{n-1} \left( 1 + \kappa \Delta t \sum_{k=1}^s |b_k| \sum_{i=1}^s [\mathbf{D}^{-1}]_{ki} \right)^\ell. \quad (75)$$

$$\left( \sum_{k=1}^s |b_k| \sum_{i=1}^s [\mathbf{D}^{-1}]_{ki} \left\| (\mathbf{I} - \mathbb{V}) \mathbf{f}_i \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_G^{n-\ell-1} + \Delta t \sum_{j=1}^s a_{ij} \Phi \hat{\mathbf{w}}_{G,j}^{n-\ell} \right) \right\| \right). \quad (76)$$

PROOF. Subtracting (74) from (73) and applying the triangle inequality yields

$$\|\delta \mathbf{w}_{G,i}^n\| \leq \left\| \mathbf{f}_i \left( \mathbf{x}_0 + \mathbf{x}_*^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \mathbf{w}_{*,j}^n \right) - \mathbb{V} \mathbf{f}_i \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_G^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \Phi \hat{\mathbf{w}}_{G,j}^n \right) \right\|, \quad i \in \mathbb{N}(s),$$

where  $\delta \mathbf{w}_{G,i}^n := \mathbf{w}_{*,i}^n - \Phi \hat{\mathbf{w}}_{G,i}^n$ . Adding and subtracting  $\mathbf{f}_i \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_G^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \Phi \hat{\mathbf{w}}_{G,j}^n \right)$  and invoking assumption  $(\mathbf{A}_1)$ , we deduce

$$\|\delta \mathbf{w}_{G,i}^n\| - \kappa \Delta t \sum_{j=1}^s |a_{ij}| \|\delta \mathbf{w}_{G,j}^n\| \leq \left\| (\mathbf{I} - \mathbb{V}) \mathbf{f}_i \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_G^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \Phi \hat{\mathbf{w}}_{G,j}^n \right) \right\| + \kappa \|\delta \mathbf{x}_G^{n-1}\|, \quad i \in \mathbb{N}(s).$$

Selecting  $\Delta t$  small enough such that (a) and (b) hold yields

$$\eta := \|\delta \mathbf{w}_{G,k}^n\| \leq \sum_{i=1}^s [\mathbf{D}^{-1}]_{ki} \left\| (\mathbf{I} - \mathbb{V}) \mathbf{f}_i \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_G^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \Phi \hat{\mathbf{w}}_{G,j}^n \right) \right\| + \kappa \|\delta \mathbf{x}_G^{n-1}\| \sum_{i=1}^s [\mathbf{D}^{-1}]_{ki},$$

where  $[\cdot]_{ij}$  denotes entry  $(i, j)$  of the argument. From explicit state updates (7) and (18), we obtain

$$\|\delta \mathbf{x}_G^n\| \leq \|\delta \mathbf{x}_G^{n-1}\| + \Delta t \sum_{k=1}^s |b_k| \|\delta \mathbf{w}_{G,k}^n\|.$$

Using the upper bound for  $\eta$  yields

$$\begin{aligned} \|\delta \mathbf{x}_G^n\| &\leq \Delta t \sum_{k=1}^s |b_k| \sum_{i=1}^s [\mathbf{D}^{-1}]_{ki} \left\| (\mathbf{I} - \mathbb{V}) \mathbf{f}_i \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_G^{n-1} + \Delta t \sum_{j=1}^s a_{ij} \Phi \hat{\mathbf{w}}_{G,j}^n \right) \right\| \\ &\quad + \left( 1 + \kappa \Delta t \sum_{k=1}^s |b_k| \sum_{i=1}^s [\mathbf{D}^{-1}]_{ki} \right) \|\delta \mathbf{x}_G^{n-1}\|. \end{aligned}$$

Finally, an induction argument produces the desired result.

### 6.3. Backward Euler

We now derive error bounds and comparative results for the backward Euler scheme.

**Corollary 6.3 (Backward Euler).** *Under the assumptions of Theorem 6.1, for Backward Euler we obtain*

$$\|\delta \mathbf{x}_G^n\| \leq \Delta t \sum_{j=0}^{n-1} \frac{1}{(h)^{j+1}} \left\| (\mathbf{I} - \mathbb{V}) \mathbf{f} \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_G^{n-j} \right) \right\| \quad (77)$$

$$\|\delta \mathbf{x}_D^n\| \leq \Delta t \sum_{j=0}^{n-1} \frac{1}{(h)^{j+1}} \left\| (\mathbf{I} - \mathbb{P}^{n-j}) \mathbf{f} \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^{n-j} \right) \right\|. \quad (78)$$

where  $h := 1 - \kappa \Delta t$ .

PROOF. Backward Euler is a single-step method that can be characterized by Eq. (2) with  $k = 1$ ,  $\alpha_0 = 1$ ,  $\alpha_1 = -1$ ,  $\beta_0 = 1$ , and  $\beta_1 = 0$ . Substituting these values into error bound (70) yields

$$\|\delta \mathbf{x}_D^n\| \leq \frac{\Delta t}{h} \left\| (\mathbf{I} - \mathbb{P}^n) \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^n) \right\| + \frac{1}{h} \|\delta \mathbf{x}_D^{n-1}\| \quad (79)$$

$$\leq \Delta t \sum_{j=0}^{n-1} \frac{1}{(h)^{j+1}} \left\| (\mathbf{I} - \mathbb{P}^{n-j}) \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^{n-j}) \right\|, \quad (80)$$

where we have used  $\delta \mathbf{x}_D^0 = \mathbf{0}$ . Thus, we obtain bound (78). Derivation of bound (77) is identical and is thus omitted.

It is not clear how to directly compare the Galerkin and the discrete optimal error bounds (69) and (70). However, our numerical experiments in Section 7, which use the three-point backward-difference scheme, suggest that discrete optimal ROM uniformly outperforms the Galerkin ROM. We will provide further theoretical justification for these numerical observations for the backward Euler scheme, i.e., we will compare bounds (77) and (78) for  $\|\delta \mathbf{x}_G^n\|$  and  $\|\delta \mathbf{x}_D^n\|$ . Similar arguments can be applied to the more general schemes.

Towards this end, for  $j = 0, \dots, n-1$ , it is sufficient to compare

$$\Delta t \left\| (\mathbf{I} - \mathbb{V}) \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}_G^{n-j}) \right\| \quad \text{and} \quad \Delta t \left\| (\mathbf{I} - \mathbb{P}^n) \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^{n-j}) \right\|.$$

Invoking Eq. (66), we can rewrite the first term as

$$\Delta t \left\| (\mathbf{I} - \mathbb{V}) \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}_G^{n-j}) \right\| = \left\| \Phi \hat{\mathbf{x}}_G^{n-j} - \Delta t \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}_G^{n-j-1}) - \Phi \hat{\mathbf{x}}_G^{n-j-1} \right\|. \quad (81)$$

Similarly, using Eq. (67) and the optimality property of  $\hat{\mathbf{x}}_D^{n-j}$ , we deduce

$$\begin{aligned} \Delta t \left\| (\mathbf{I} - \mathbb{P}^{n-j}) \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^{n-j}) \right\| &= \left\| \Phi \hat{\mathbf{x}}_D^{n-j} - \Delta t \mathbf{f}(\mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^{n-j}) - \Phi \hat{\mathbf{x}}_D^{n-j-1} \right\| \\ &= \min_{\mathbf{y}} \left\| \Phi \mathbf{y} - \Delta t \mathbf{f}(\mathbf{x}_0 + \Phi \mathbf{y}) - \Phi \hat{\mathbf{x}}_D^{n-j-1} \right\|. \end{aligned} \quad (82)$$

For the same previous state  $\hat{\mathbf{x}}_D^{n-j-1} = \hat{\mathbf{x}}_G^{n-j-1}$ , a direct comparison of Eqs. (81) and (82) yields that (82) will always be less than (81). We state this result below.

**Corollary 6.4.** *If  $\hat{\mathbf{x}}_D^j = \hat{\mathbf{x}}_G^j$ ,  $j \in \mathbb{N}(n-1)$ , then under the assumptions of Theorem 6.1, the upper bound for  $\|\delta \hat{\mathbf{x}}_D^k\|$  in Eq. (78) will be less than the upper bound for  $\|\delta \hat{\mathbf{x}}_G^k\|$  in Eq. (77) for  $k \in \mathbb{N}(n)$ .*

**Corollary 6.5.** *If  $\bar{\mathbf{x}}$  solves the auxiliary problem centered on the discrete-optimal ROM trajectory*

$$\bar{\mathbf{x}}^j = \Delta t \mathbf{f}(\mathbf{x}_0 + \bar{\mathbf{x}}^j) + \Phi \hat{\mathbf{x}}_D^{j-1}, \quad j \in \mathbb{N}(n), \quad (83)$$

then the following holds:

$$\|\delta \mathbf{x}_D^n\| \leq (1 + \kappa \Delta t) \sum_{j=0}^{n-1} \frac{\mu^{n-j}}{(h)^{j+1}} \quad (84)$$

$$= \Delta t (1 + \kappa \Delta t) \sum_{j=0}^{n-1} \frac{\bar{\mu}^{n-j}}{(h)^{j+1}} \|\mathbf{f}(\bar{\mathbf{x}}^{n-j})\|. \quad (85)$$

Here,  $\mu^j := \left\| \Phi \Delta \hat{\mathbf{x}}_D^j - \Delta \bar{\mathbf{x}}^j \right\|$  denotes the difference in solution increments at time instance  $j$ , where  $\Delta \hat{\mathbf{x}}_D^j := \hat{\mathbf{x}}_D^j - \hat{\mathbf{x}}_D^{j-1}$  and  $\Delta \bar{\mathbf{x}}^j := \bar{\mathbf{x}}^j - \Phi \hat{\mathbf{x}}_D^{j-1}$ . We denote the relative solution increment at time instance  $j$  by  $\bar{\mu}^j := \mu^j / \|\Delta \bar{\mathbf{x}}^j\|$ .

PROOF. Eq. (78) in conjunction with (82) implies

$$\|\delta \mathbf{x}_D^n\| \leq \sum_{j=0}^{n-1} \frac{1}{(h)^{j+1}} \left\| \Phi \Delta \hat{\mathbf{x}}_D^{n-j} - \Delta t \mathbf{f} \left( \mathbf{x}_0 + \Phi \hat{\mathbf{x}}_D^{n-j} \right) \right\|. \quad (86)$$

We can also write the auxiliary equation (83) as  $\Delta \bar{\mathbf{x}}^j = \Delta t \mathbf{f}(\bar{\mathbf{x}}^j)$ ,  $j \in \mathbb{N}(n)$ , which allows us to rewrite bound (86) as

$$\|\delta \mathbf{x}_D^n\| \leq \sum_{j=0}^{n-1} \frac{1}{(h)^{j+1}} \cdot \left\| \left( \Phi \Delta \hat{\mathbf{x}}_D^{n-j} - \Delta \bar{\mathbf{x}}^{n-j} \right) - \Delta t \left( \mathbf{f} \left( \mathbf{x}_0 + \Phi \Delta \hat{\mathbf{x}}_D^{n-j} + \Phi \hat{\mathbf{x}}_D^{n-j-1} \right) - \mathbf{f} \left( \mathbf{x}_0 + \Delta \bar{\mathbf{x}}^{n-j} + \Phi \hat{\mathbf{x}}_D^{n-j-1} \right) \right) \right\|.$$

Lipschitz continuity of  $\mathbf{f}$  leads to the bound (84). To obtain Eq. (85), we multiply and divide by  $\|\Delta \bar{\mathbf{x}}^{n-j}\|$  for each term in the summation and use  $\Delta \bar{\mathbf{x}}^{n-j} = \Delta t \mathbf{f}(\bar{\mathbf{x}}^{n-j})$ .

**Remark 6.6.** *The time step  $\Delta t$  in the error bound (85) for the discrete-optimal ROM solution plays an important role. In particular, decreasing the time step produces both beneficial effects (bound decrease) and deleterious effects (bound increase), which we denote by ‘+’ and ‘-’, respectively as follows:*

- + *The time-discretization error decreases (this does not appear in the time-discrete error analysis above).*
- *The number of overall time steps  $n$  increases, so there are more terms in the summation.*
- + *The terms  $\Delta t(1 + \kappa \Delta t)$  and  $1/(h)^{j+1}$  decrease.*
- ? *The term  $\bar{\mu}^{n-j}$  may increase or decrease, depending on the spectral content of the basis  $\Phi$ .*

We now discuss this final ambiguous effect. The term  $\bar{\mu}^n$  can be interpreted as the relative error in solution increment over  $[(n-1)\Delta t, n\Delta t]$ . Clearly, the ability of the discrete-optimal ROM to make  $\bar{\mu}^n$  small depends on the spectral content of the basis  $\Phi$ : if the basis only captures modes that evolve over long time scales, then  $\bar{\mu}^n$  will be large (i.e., close to one), as the basis does not contain the ‘fast evolving’ solution components that change over a single time step. This suggests that the time step should be ‘matched’ to the spectral content of the reduced basis  $\Phi$ . In Section 7.5 of the experiments, we explore this issue numerically, and demonstrate that the error bound is minimized for an intermediate value of the time step  $\Delta t$ .

We note that the above arguments do not hold for the Galerkin ROM, which is simply an ODE that does not depend on the time step. Instead, decreasing the time step should increase accuracy, as it has the effect of reducing the time-discretization error.

## 7. Numerical experiments

This section compares the performance of Galerkin and discrete-optimal ROMs on a CFD application using a basis constructed by proper orthogonal decomposition. These experiments highlight the importance of the previous analyses, in particular the limiting equivalence of Galerkin and discrete-optimal ROMs (Theorem 5.3), superior accuracy of the discrete-optimal ROM compared with the Galerkin ROM (Corollary 6.4), and performance improvement of the discrete-optimal ROM when an intermediate time step is selected (Corollary 6.5 and Remark 6.6).

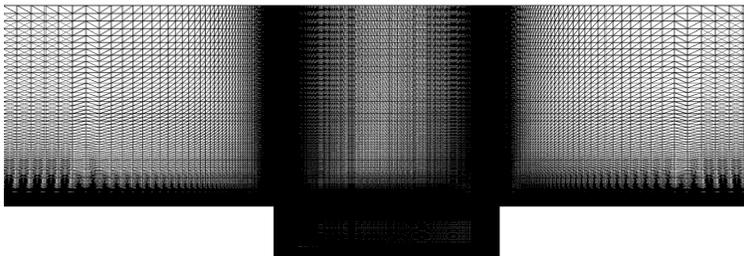
### 7.1. Problem description

The Galerkin and discrete-optimal ROMs are implemented in AERO-F [53, 54], a massively parallel compressible-flow solver. AERO-F solves the steady or unsteady compressible Navier–Stokes equations with various closure models available for turbulent flow, and employs a second-order node-centered finite-volume scheme. For model-reduction algorithms, all linear least-squares problems and singular value decompositions are computed in parallel using the ScaLAPACK library [55].

The full-order model corresponds to an unsteady Navier–Stokes simulation of a two-dimensional open cavity using AERO-F’s DES turbulence model (based on the Spalart–Almaras one-equation model) and a wall-function boundary condition applied on solid surface boundaries. The fluid domain is discretized by a mesh with 192,816 nodes and 573,840 tetrahedra (Figure 2). The two-dimensional geometry is discretized in three dimensions by considering a slab of thin, but finite thickness, in the  $z$ -direction. The viscosity is assumed to be constant, and the Reynolds number based on cavity length is  $6.30 \times 10^6$ , while the free-stream Mach number is 0.6. Due to the turbulence model and three-dimensional domain, the number of conservation equations per node is 6, and therefore the dimension of the CFD model is  $N = 1,156,896$ . Roe’s scheme is employed to discretize the convective fluxes, and a linear variation of the solution is assumed within each control volume, which leads to a second-order space-accurate scheme. We employ a low-numerical-dissipation scheme that gives fifth-order formal order of accuracy on inviscid, one-dimensional problems.



(a) Full domain



(b) Detail around cavity

Figure 2: Computational mesh:  $x - y$  plane cut.

Flow simulations are performed within a time interval  $t \in [0, T]$  with  $T = 12.5$  seconds. We employ the second-order accurate implicit three-point backward difference scheme, which is a linear multistep scheme characterized by  $k = 2$ ,  $\alpha_0 = 1$ ,  $\alpha_1 = -4/3$ ,  $\alpha_2 = 1/3$ ,  $\beta_0 = 2/3$ ,  $\beta_1 = \beta_2 = 0$ , for time integration. The ODE (3) arising at each time step is solved by a Newton–Krylov method, where GMRES is employed as the iterative linear solver with a restrictive additive Schwarz preconditioner (with no fill in) and the previous 50 Krylov vectors are employed for orthogonalization. Convergence is declared when the residual norm is reduced to a factor of  $10^{-3}$  of its starting value. All flow computations are performed in a non-dimensional setting.

The initial condition  $\mathbf{x}_0$  is provided by first computing a steady-state solution, and using that solution as an initial guess for an unsteady ‘transient’ simulation (which captures the initial transient before the flow reaches a quasi-periodic state) of 7.5 seconds. The state at the end of the unsteady transient simulation is then used as the initial condition for the subsequent simulations. The steady-state calculation is characterized by the same parameters as above, except that it employs local time stepping with a maximum CFL number of 100, it uses the first-order implicit backward Euler time integration scheme, it assumes a linear variation of the solution within each control volume, it employs a Spalart–Allmaras turbulence model, and it employs only one Newton iteration per (pseudo) time step.

The output of interest is the pressure at location  $(0.0001, -0.0508, 0.0025)$ , which is shown in the bottom row of Figure 4. All errors are reported as the  $\ell^2$  relative error in this quantity, i.e.,

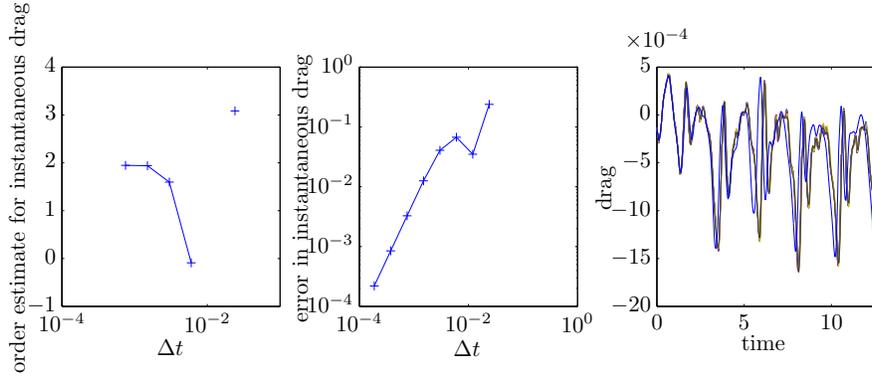
$$\varepsilon(p, p_\star) = \frac{\sqrt{\sum_{n=1}^{T/\Delta t_\star} (P_\star(p)(n\Delta t_\star) - p_\star(n\Delta t_\star))^2}}{\sqrt{\sum_{n=1}^{T/\Delta t} p_\star(n\Delta t_\star)^2}}$$

where  $p : \mathbb{N}(T/\Delta t) \rightarrow \mathbb{R}$  is the pressure for the model of interest,  $p_* : \mathbb{N}(T/\Delta t_*) \rightarrow \mathbb{R}$  is this pressure response of the designated ‘truth’ model (typically the full-order model), and  $P_*$  is a linear interpolation of the pressure response onto the grid based on the truth-model time step  $\Delta t_*$ .

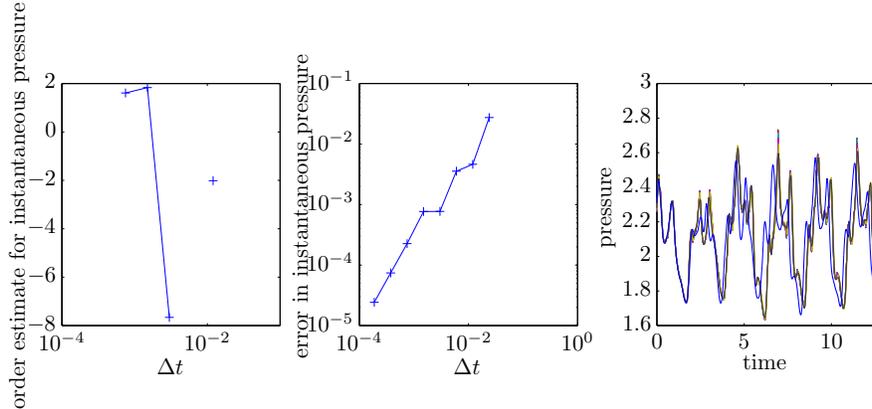
All computations are performed in double-precision arithmetic on a parallel Linux cluster<sup>2</sup> using 48 cores across 6 nodes.

### 7.2. Time-step verification

Because this paper considers the time step to be an important parameter in model reduction, we first perform a time-step verification study to ensure we employ an appropriate ‘nominal’ time step. Figure 3 reports these results using a time-step refinement factor of two. A time step of  $\Delta t_* = 0.0015$  seconds yields observed convergence rates in both the instantaneous drag force on the lower wall and instantaneous pressure at  $t = T$  that are close to the asymptotic rate of convergence (2.0) of three-point BDF2 scheme. Further, this value also leads to sub-2% errors in both quantities, which we deem to be sufficient for this set of experiments.



(a) Drag:  $\Delta t_* = 0.0015$  yields an approximate rate of convergence of 1.94 and an estimated error in the output quantity (computed via Richardson extrapolation) of  $1.26 \times 10^{-2}$ . The rightmost plot shows the time-dependent response for all tested time steps.



(b) Pressure:  $\Delta t_* = 0.0015$  yields an approximate rate of convergence of 1.83 and an estimated error in the output quantity (computed via Richardson extrapolation) of  $7.68 \times 10^{-4}$ . The rightmost plot shows the time-dependent response for all tested time steps.

Figure 3: Time-step verification study. Note that the approximated convergence rates are close to the asymptotic value of 2.0 for the BDF2 scheme.

Figure 4 shows several instantaneous snapshots of the vorticity field and corresponding pressure field

<sup>2</sup>The cluster contains 8-core compute nodes that each contain a 2.93 GHz dual socket/quad core Nehalem X5570 processor with 12 GB of memory. The interconnect is a 3D torus InfiniBand.

generated by the high-fidelity CFD model. The flow within the cavity is quasi-periodic; during one cycle, vorticity is shed from the leading edge of the cavity, convects downstream, and impinges on the aft edge of the cavity. Upon impingement, an acoustic disturbance is generated which propagates upstream and scatters on the leading edge of the cavity, generating a new vortical disturbance to initiate the next oscillation cycle. The pressure fields in the bottom row of Figure 4 reveal regions of low pressure (blue contours) associated with vortices, as well as acoustic disturbances both within the cavity and radiating outside the cavity. This complex flow is governed by the interactions of several nonlinear processes, including roll-up of the shear layer vortices, impingement of the vortices on the aft wall resulting in sound generation, propagation of nonlinear acoustic waves, and interaction of these waves with the shear layer vorticity.

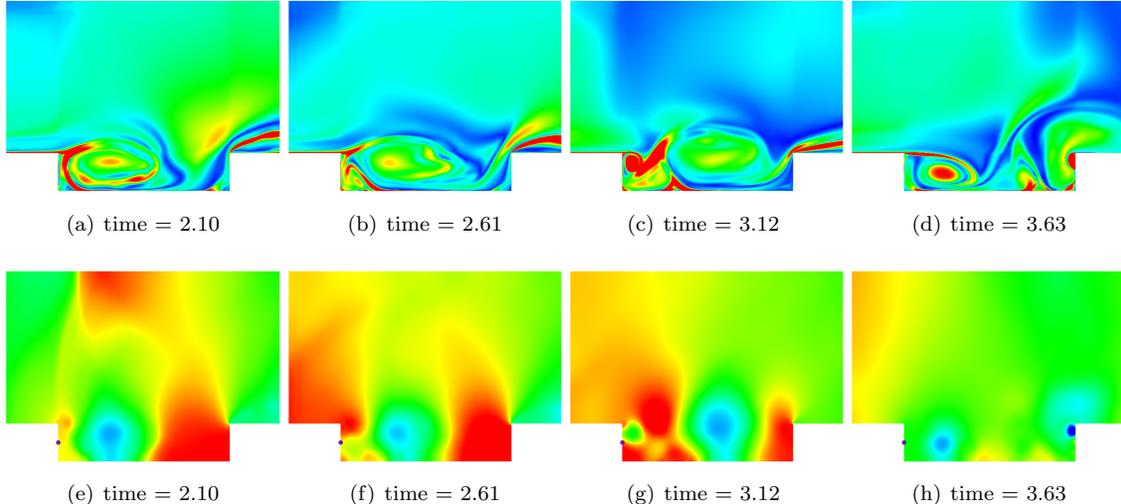


Figure 4: Instantaneous CFD vorticity field (top) and pressure field (bottom) during one oscillation cycle. The dot on the forward wall of the cavity indicates the location of the pressure signal output.

### 7.3. Reduced-order models

To construct both the Galerkin and discrete-optimal ROMs, we employ the proper orthogonal decomposition (POD) technique. In particular, we set  $\Phi \leftarrow \Phi(\mathcal{X}, \nu)$ , where  $\Phi$  is computed via Algorithm 1 of the appendix with snapshots consisting of the initial-condition-centered full-order model states  $\mathcal{X} = \{\mathbf{x}_*(k\Delta t_*) - \mathbf{x}_0\}_{k=1}^{8334}$ , where  $\mathbf{x}_*$  denotes the FOM response computed for a time step of  $\Delta t_*$ . Three values of the energy criterion  $\nu \in [0, 1]$  are used during the experiments:  $\nu = 1 - 10^{-4}$  ( $p = 204$ ),  $\nu = 1 - 10^{-5}$  ( $p = 368$ ), and  $\nu = 1 - 10^{-6}$  ( $p = 564$ ). Figure 5 shows a selection of the energy component of the computed POD modes. Note that as the mode number increases, the modes capture finer spatial-scale behavior, which we expect to be associated with finer time-scale behavior; this will be verified in Section 7.5.1.

We first repeat the time-step verification study, but we do so for the reduced-order models in the time interval  $0 \leq t \leq 0.55$ , as all Galerkin ROMs remain stable in this time interval. Figure 6 reports these results. First, we note that the Galerkin ROM converges an approximated rate of 2.0, which is what we expect given that the Galerkin ROM simply associates with a time-step-independent ODE (9). However, the discrete-optimal ROM does not exhibit this behavior; in fact the error convergence is not even monotonic. This is due to the fact that the method does not associate with a time-step-independent ODE.

We next perform simulations for both reduced-order models for all tested basis dimensions and time steps; Figure 7 reports the time-dependent responses. When a response stops before the end of the time interval, this indicates that a negative pressure was encountered, which causes AERO-F to exit the simulation. We interpret this phenomenon as a non-physical instability.

First, note that the Galerkin ROMs become unstable (i.e., generate a negative pressure) for all time steps and all basis dimensions. This is consistent with previously reported results [31, 30, 28, 56] that indicate Galerkin projection almost always leads to inaccurate responses for compressible fluid-dynamics problems. In contrast, the discrete-optimal ROM results in many stable, accurate responses for all basis dimensions.

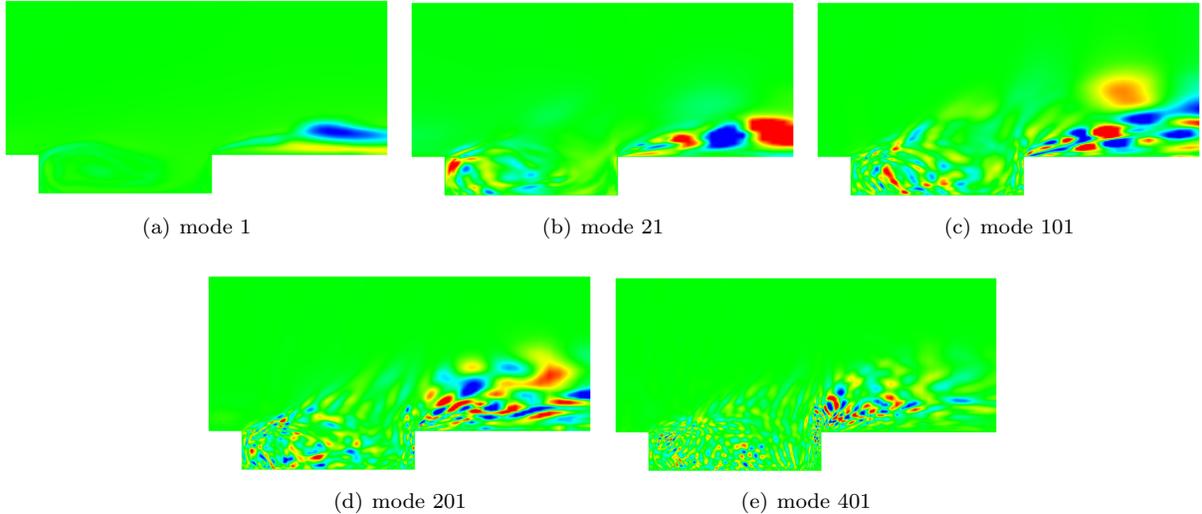


Figure 5: Visualization of the energy component of the POD modes.

Further, discrete-optimal responses exhibit a clear dependence on the time step  $\Delta t$ . Subsequent sections provide a deeper analysis of this dependence.

#### 7.4. Limiting case: comparison

We next compare the responses of the Galerkin and discrete-optimal ROMs for small time windows (when the Galerkin responses remain stable) and small time steps. Figure 8 reports  $\varepsilon(p_{\text{discrete opt.}}, p_{\text{Gal.}})$ —which is the difference between the discrete-optimal ROM pressure response and the Galerkin pressure response for  $\Delta t = 1.875 \times 10^{-4}$  (the smallest tested time step)—as a function of the time step for two different time windows. These responses support an important conclusion (see Theorem 5.3): the Galerkin and discrete-optimal ROMs are equal in the limit of  $\Delta t \rightarrow 0$ . This has significant consequences for the discrete-optimal ROM, as decreasing the time step leads to the same *unstable* response as Galerkin; larger time steps are needed to ensure the discrete-optimal ROM is stable for the entire time interval.

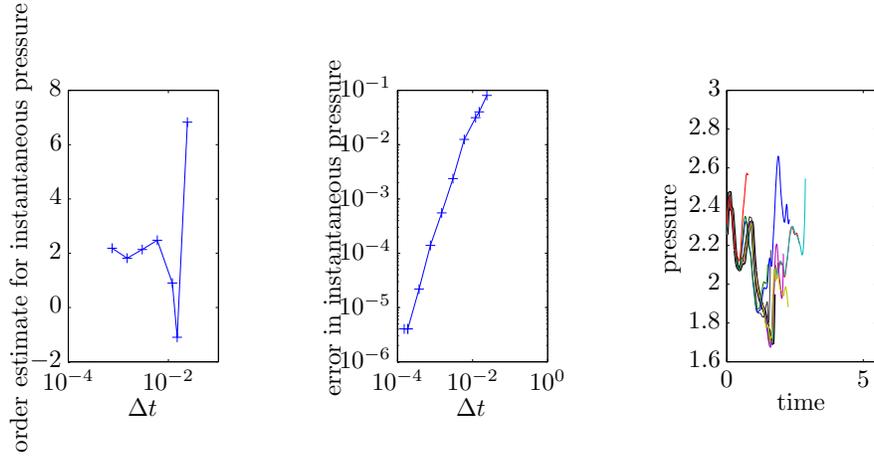
Figure 9 reports  $\varepsilon(p_{\text{discrete opt.}}, p_{\text{FOM}_*})$  and  $\varepsilon(p_{\text{Gal.}}, p_{\text{FOM}_*})$ —which are the differences between the two ROM-generated pressure responses and the full-order model pressure response for  $\Delta t = 1.875 \times 10^{-4}$ — as a function of the time step for all three basis dimensions and three time intervals. These results highlight a critical observation: the discrete-optimal ROM is *more accurate* for an intermediate time step. This not only supports the result of Corollary 6.5, but provides an interesting insight: taking a larger time step not only leads to better speedups (i.e., the end of the time interval is reached in fewer time steps), but it also decreases the error, sometimes significantly. This is further explored in the next section.

#### 7.5. Time-step selection

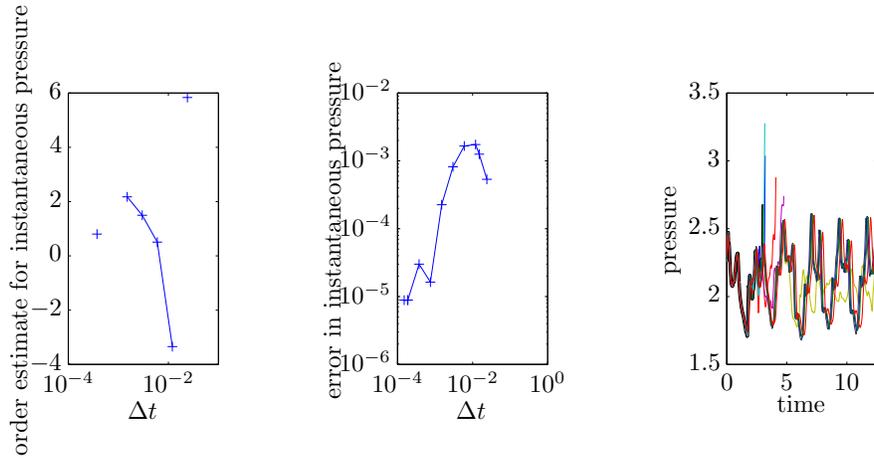
Recall from Corollary 6.5 and Remark 6.6 that decreasing the time step  $\Delta t$  has a non-obvious effect on the error bound for the discrete-optimal ROM. We now assess these effects for the current problem.

##### 7.5.1. Spectral content of POD basis

In our interpretation of the error bound (85) for the discrete-optimal ROM applied to the backward Euler scheme, we noted that the time step should be ‘matched’ to the spectral content of the trial basis  $\Phi$ . This is of practical importance, as selecting an appropriate time step for the ROM should take into account the relevant temporal dynamics associated with the basis. For example, a time step may be too small if the basis has filtered out modes with a time scale matching that of the time step. If we assume that the basis  $\Phi$  is computed via POD, then we would expect the vectors to be naturally ordered such that lower mode numbers are associated with lower temporal frequencies. Then, including additional modes has the effect of

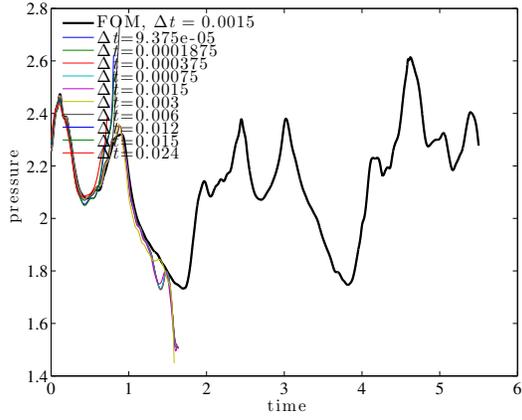


(a) Galerkin reduced-order model

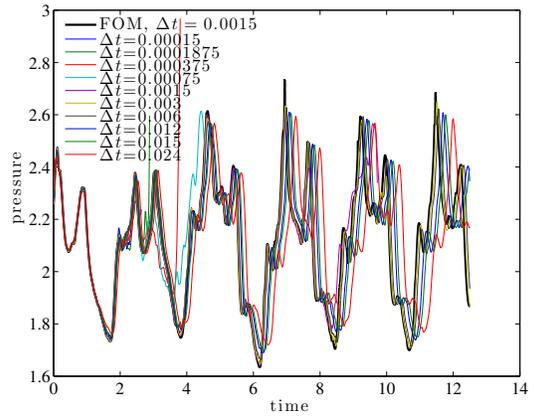


(b) Discrete-optimal reduced-order model

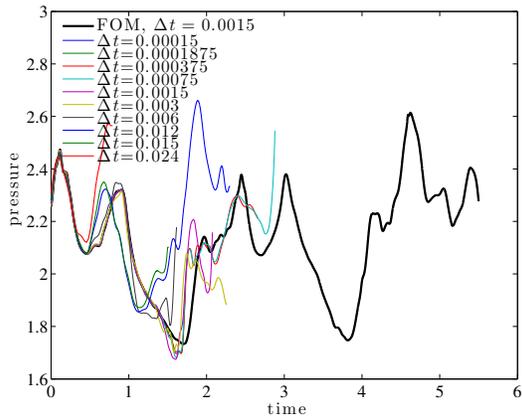
Figure 6: Time-step verification study for Galerkin and discrete-optimal reduced-order models for  $p = 368$  and  $0 \leq t \leq 0.55$ . While the approximated convergence rate for the Galerkin reduced-order model is close to the asymptotic value of 2.0 for the BDF2 scheme, this is not observed for the discrete-optimal reduced-order model.



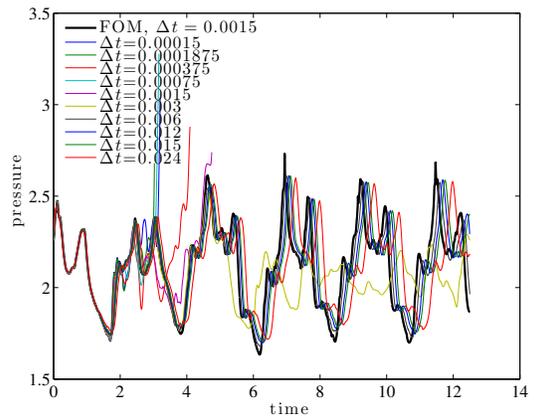
(a) Galerkin,  $p = 204$



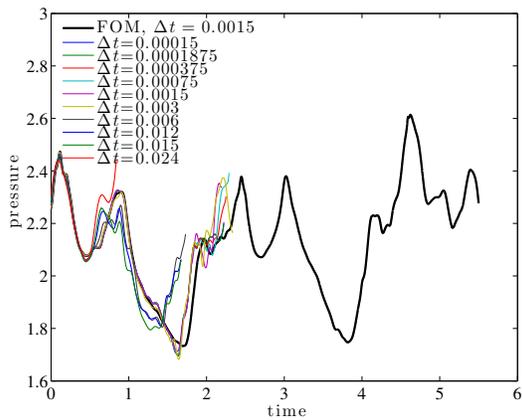
(b) Discrete optimal,  $p = 204$



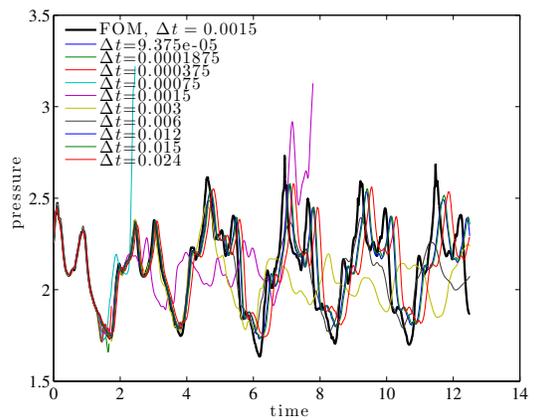
(c) Galerkin,  $p = 368$



(d) Discrete optimal,  $p = 368$



(e) Galerkin,  $p = 564$



(f) Discrete optimal,  $p = 564$

Figure 7: Responses generated by Galerkin and discrete optimal reduced-order models for different basis sizes  $p$  and timesteps  $\Delta t$

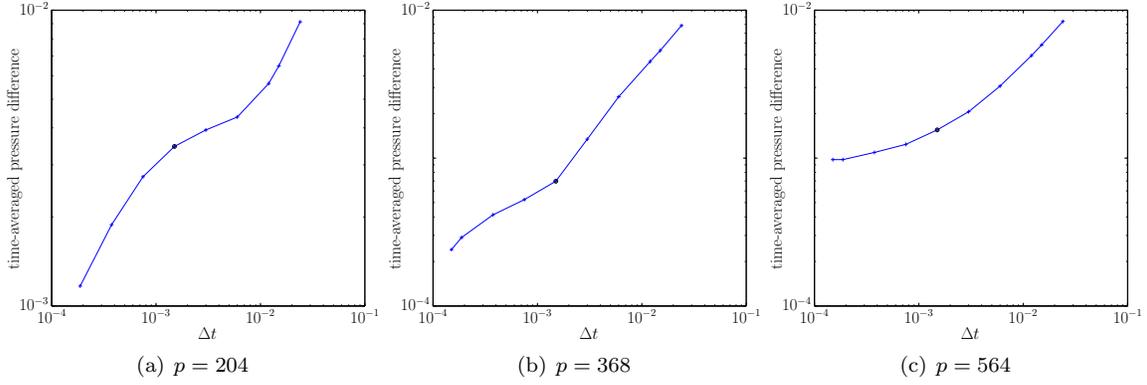


Figure 8: Error between Galerkin and discrete-optimal reduced-order models  $\varepsilon(p_{\text{discrete-opt.}}, p_{\text{Gal.}})$  for different timesteps in  $0 \leq t \leq 1.1$ . This demonstrates convergence of the discrete-optimal ROM to Galerkin as  $\Delta t \rightarrow 0$ .

encoding information at higher frequencies. It follows that the time step should be decreased as additional modes are retained in construction of the ROM.

Here we investigate the validity of this assumption by examining the spectral content of the POD basis vectors for the current cavity-flow problem. We compute the time histories of the generalized coordinates by projecting the FOM solution onto the POD basis as  $\hat{\mathbf{x}}_*(k\Delta t_*) := \Phi^T(\mathbf{x}_*(k\Delta t_*) - \mathbf{x}_0)$ ,  $k \in \mathbb{N}(8334)$ . We then compute power spectral densities of the generalized coordinates  $\hat{\mathbf{x}}_*(t)$ . Figure 10(a) shows sample spectra, normalized by the total energy in each signal,<sup>3</sup> for several of the POD modes. The figure shows that energy shifts to higher frequencies as the POD mode number increases, confirming our assumption for this example. This is further quantified by calculating a characteristic time-scale  $\tau_{95}$  associated with each mode; we define this time scale as the inverse of the frequency below which 95 percent of the energy is captured for that mode. Figure 10(b) plots this time scale versus the mode number, showing a clear trend of decreasing time scale with increasing mode number.

Thus, at least for the present application problem, we expect the optimal time step for the discrete-optimal ROM to decrease as modes are added to the POD basis (this will be verified by Figure 12). Note that systematic calibration could be performed to attempt to automate selection of the ROM time step as a function of basis dimension. We do not attempt this exercise here, but note that nonlinear interactions between modes may complicate such an effort.

### 7.5.2. Error bound behavior

Having verified that higher POD mode numbers correspond to smaller wavelengths, we now numerically assess quantities related to the error bound (85). First, Figure 11(a) reports the dependence of the maximum relative projection error  $\max_k \bar{\mu}_*^k(\Phi, \Delta t)$  on the time step  $\Delta t$  and the basis dimension, where

$$\bar{\mu}_*^k(\Phi, \Delta t) := \frac{\|(\mathbf{I} - \Phi\Phi^T)(\mathbf{x}_*(k\Delta t) - \mathbf{x}_*((k-1)\Delta t))\|}{\|\mathbf{x}_*(k\Delta t) - \mathbf{x}_*((k-1)\Delta t)\|}$$

Note that  $\bar{\mu}_*^k$  is closely related to  $\bar{\mu}^k$  from error bound (85), as they are equal if  $\mathbf{x}_0 + \Phi\hat{\mathbf{x}}_D(t) = \mathbf{x}(t)$  and the discrete-optimal ROM computes  $\hat{\mathbf{x}}_D^k$  such that  $\bar{\mu}^k$  is minimized.

These results confirm that adding basis vectors—which we know has the effect of encoding higher frequency content—significantly reduces the projection error for small time steps  $\Delta t$ , but has less of an effect on larger time steps, as retaining the first POD vectors already enables dynamics at that scale to be captured.

Next, Figure 11(b) plots the error bound (85) for a value of  $\kappa = 1$  and with  $\bar{\mu}^k = \bar{\mu}_*^k$ . This highlights an important result: *selecting an intermediate time step  $\Delta t$  leads to the lowest error bound, regardless of the basis dimension.* Even though this result corresponds to the backward Euler integrator, we expect a

<sup>3</sup>The energy in a time series within some frequency range is obtained by integrating the power spectral density over that range.

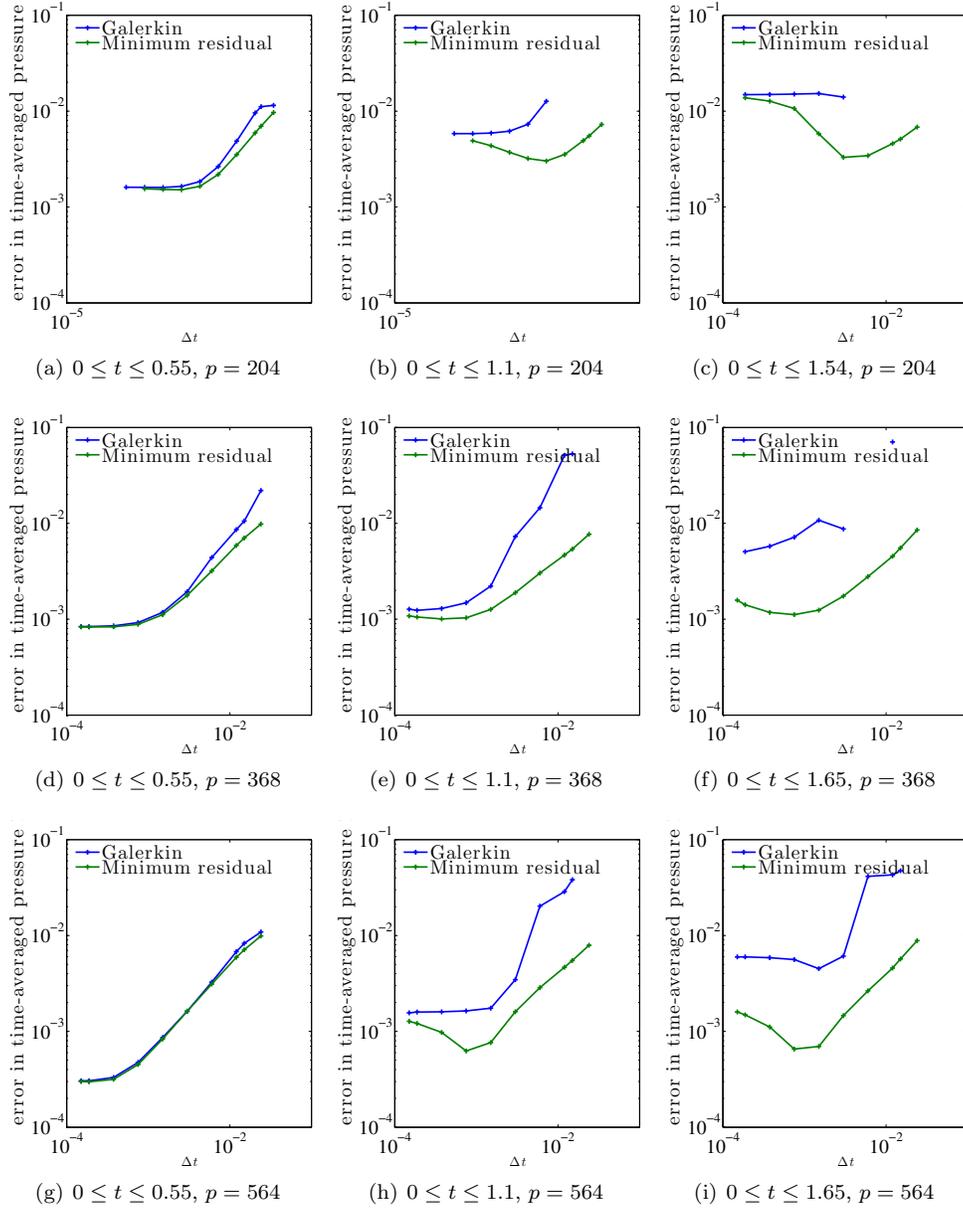


Figure 9: Galerkin errors  $\varepsilon(p_{\text{discrete opt.}}, p_{\text{FOM}_*})$  and Petrov–Galerkin errors  $\varepsilon(p_{\text{Gal.}}, p_{\text{FOM}_*})$  over different time intervals, time steps, and basis dimensions.

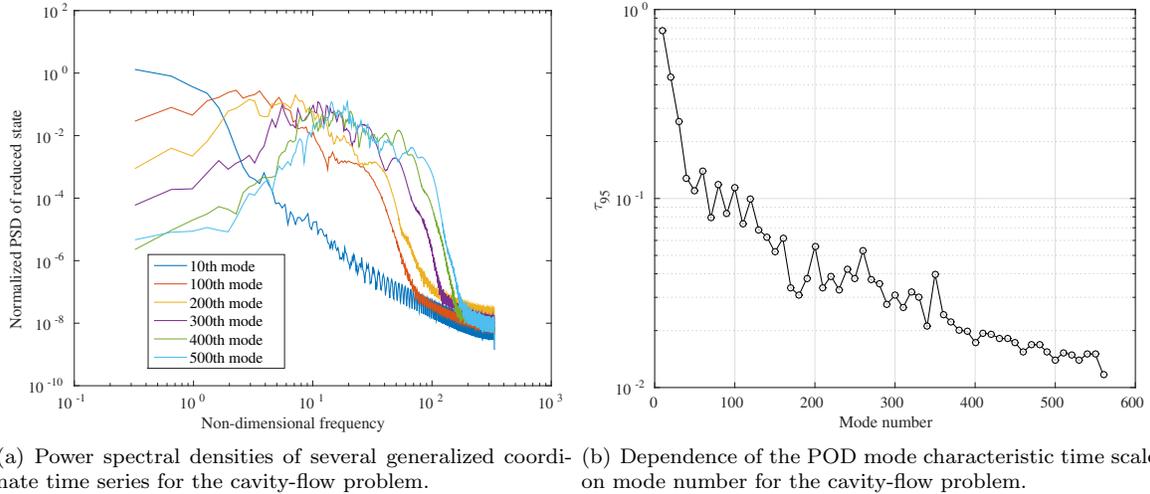


Figure 10: Spectral content of the POD basis.

similar trend to hold for the present experiment, which uses the BDF2 scheme. The next section assesses the performance of the discrete-optimal ROM, including its dependence on the time step.

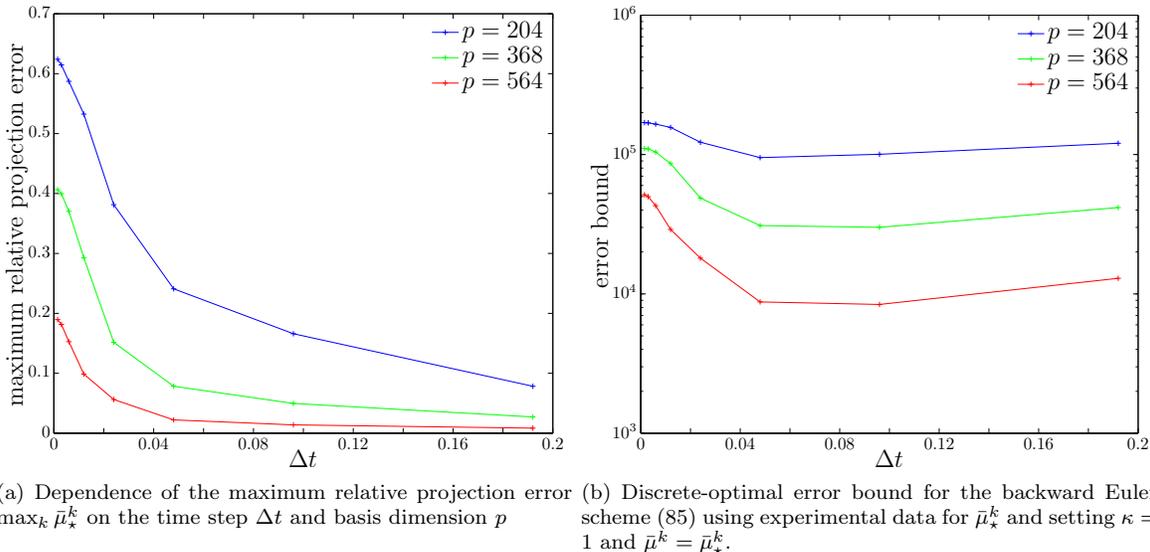


Figure 11: Assessment of quantities appearing in error bound (85). This analysis suggests that an intermediate time step  $\Delta t$  can reduce errors for the discrete-optimal ROM.

### 7.6. Discrete optimality ROM performance

We now compare the accuracy and walltime performance of the discrete-optimal ROM as the dimension of the basis, time step, and time interval change. The most salient result from Figure 12 is that choosing an intermediate time step leads to both better accuracy and faster simulation times. This shows that our theoretical analysis of the error bound performed in Section 7.5.2 leads to an actual observed performance improvement. For example, consider the  $p = 564$  case over the time interval  $0 \leq t \leq 2.5$ . In this case, a time step of  $\Delta t = 1.875 \times 10^{-4}$  leads to a relative error of 0.0140 and a simulation time of 289 hours; increasing this value to  $\Delta t = 1.5 \times 10^{-3}$  reduces the relative error to  $9.46 \times 10^{-4}$  and the simulation time to 35.8 hours, which constitutes roughly an order of magnitude improvement in both quantities. Again, this supports the

theoretical results of Corollary 6.5 and highlights the critical importance of the time step for discrete-optimal reduced-order models.

In addition, Figure 12 shows that as the basis dimension increases, the optimal time step decreases; this was anticipated from the spectral analysis performed in Section 7.5.1. In addition, adding POD basis vectors does not improve accuracy for large time steps. We interpret this effect as follows: for larger time steps, the first few POD modes accurately capture ‘coarse’ phenomena on the scale of the time step. Therefore, accuracy improvement is not achieved by adding modes that encode dynamics that evolve on a time scale finer than the time step itself.

Further, Figure 12(g) highlights that as the basis dimension increases, the error generally decreases, which is an artifact of *a priori* convergence achieved by the discrete-optimal ROM (Remark 4.1). Finally, the figure shows that as the time interval grows, the optimal time step generally increases.

### 7.7. GNAT: ROM with complexity reduction

In this section, we perform a similar study, but equip the discrete-optimal ROM with complexity reduction in order to achieve computational savings. In particular, we employ the GNAT method [28, 31, 30], which solves Eq. (21) with  $\mathbf{A} = (\mathbf{P}\Phi_r)^+ \mathbf{P}$ , where  $\Phi_r$  is a basis for the residual and  $\mathbf{P}$  consisting of selected rows of the identity matrix.

The problem is identical to that described in Section 7.1 except that we take  $T = 5.5$  seconds and employ a second-order space-accurate dissipation scheme wherein a linear variation of the solution is assumed within each control volume.<sup>4</sup> For this simulation, the full-order model consumes 5.0 hours on 48 cores across six compute nodes.

To construct the trial basis  $\Phi$  and basis for the residual  $\Phi_r$  for the GNAT models, we again employ POD. In particular, we set  $\Phi \leftarrow \Phi(\mathcal{X}, \nu)$ , where  $\Phi$  is computed via Algorithm 1 with snapshots consisting of the centered full-order model states  $\mathcal{X} = \{\mathbf{x}_*(k\Delta t_*) - \mathbf{x}_0\}_{k=1}^{3668}$ . An energy criterion of  $\nu = 1 - 10^{-5}$  ( $p = 179$ ) is used during the experiments. For the residual, we employ  $\Phi_r \leftarrow \Phi(\mathcal{X}_r, \nu_r)$  via Algorithm 1 with snapshots  $\mathcal{X}_r = \{\mathbf{r}^n(\mathbf{x}_0 + \Phi \hat{\mathbf{w}}^{n(k)}), k \in \mathbb{N}(K(n)), n \in \mathbb{N}(2228)\}$  and  $\hat{\mathbf{w}}^{n(k)}$  corresponding to the discrete-optimal ROM solution at Gauss–Newton iteration  $k$  within time step  $n$  using a time step of  $\Delta t = 6 \times 10^{-3}$ . Here,  $K(n)$  denotes the number of Newton iterations required for convergence of at time instance  $n$ . An energy criterion of  $\nu_r = 1.0$  is employed. In addition, the GNAT model sets the Jacobian basis equal to residual basis  $\Phi_J = \Phi_r$  and employs  $n_s = 743$  sample nodes that define  $\mathbf{P}$ , which leads to 4458 rows in  $\mathbf{P}$  as there are six conservation equations per node due to the turbulence model (see Ref. [30] for definitions).

The GNAT implementation in AERO-F is characterized by the sample-mesh concept [30]. Figure 13 depicts the sample mesh for this problem, which was constructed using  $n_c = 2228$  working columns [30, Algorithm 3], and includes two layers of nodes around the sample nodes (to enable the residual to be computed at the sample nodes). It is characterized by 7,974 total nodes (4.1% of the original mesh) and 17,070 total volumes (3.0% of the original mesh). Due to the small footprint of the sample mesh, the GNAT simulations are run using only 2 cores on a single compute node.

Figure 14 reports the results obtained with the GNAT ROM using different time steps. Critically, note that the GNAT ROM also exhibits a ‘dip’ in the optimal time step, with a time step of  $6.0 \times 10^{-3}$  yielding the lowest error. In fact, increasing the time step from  $1.5 \times 10^{-3}$  to  $6.0 \times 10^{-3}$  decreases the error from 3.32% to 2.25% and also significantly increases the computational savings relative to the full-order model (as measured in core–hours) from 14.9 to 55.7. This highlights that the analysis is also relevant to ROMs equipped with complexity reduction.

### 7.8. Summary of experimental results

We now briefly summarize the main experimental results:

- Galerkin ROMs are unstable for long time intervals (Figure 7).
- Discrete-optimal ROMs are only unstable for small time steps (Figure 7).
- Galerkin and discrete-optimal ROMs are equivalent as  $\Delta t \rightarrow 0$  (Figure 8).

<sup>4</sup>This is done to ensure the sample mesh requires two layers of neighboring nodes for each sample node.

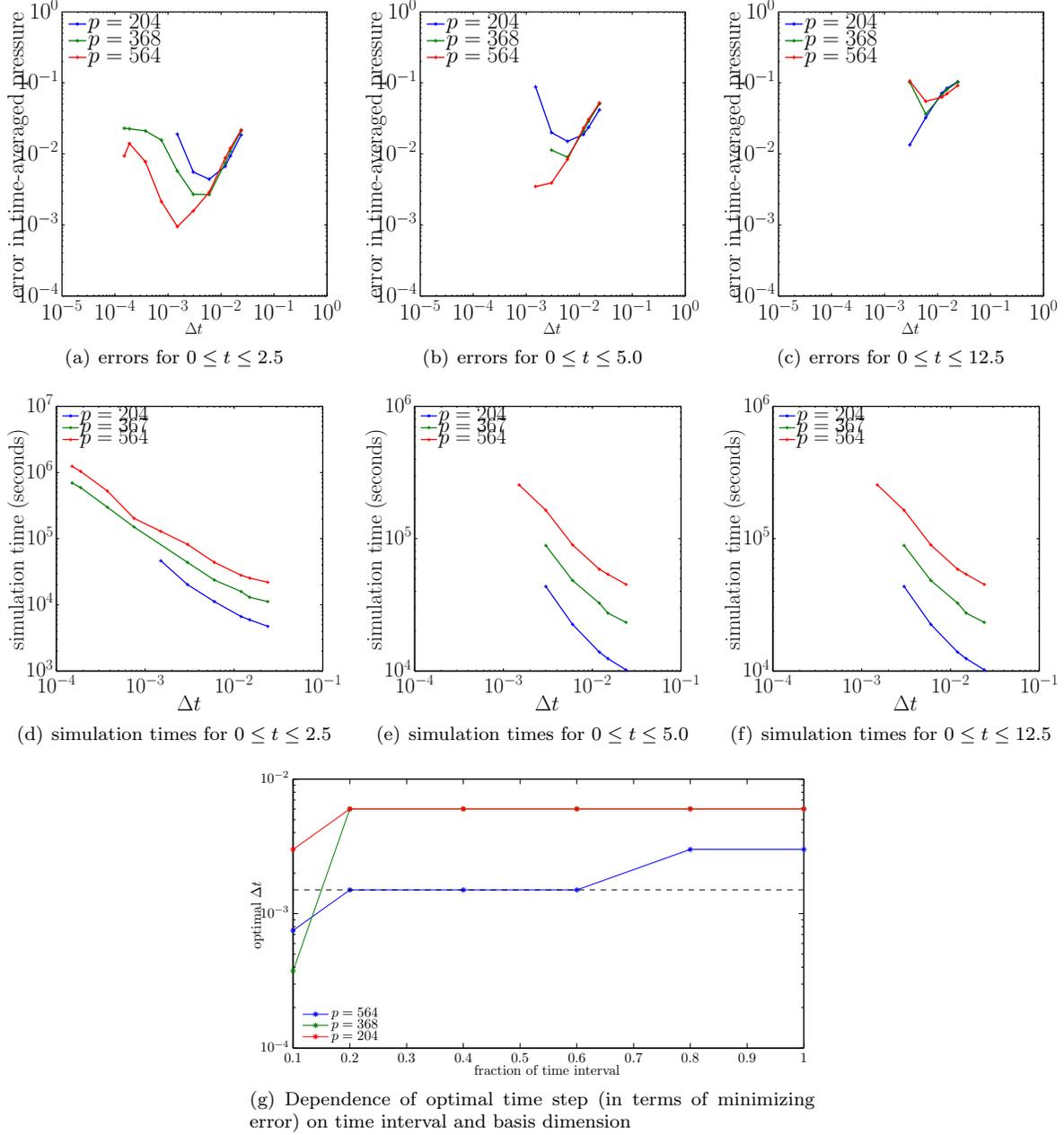
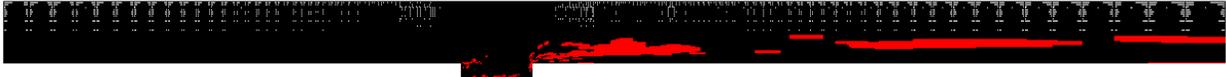
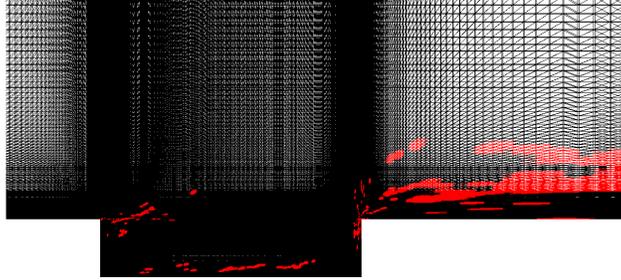


Figure 12: Dependence of error and simulation time for the discrete-optimal reduced-order model on the time step  $\Delta t$ , basis dimension, and time interval

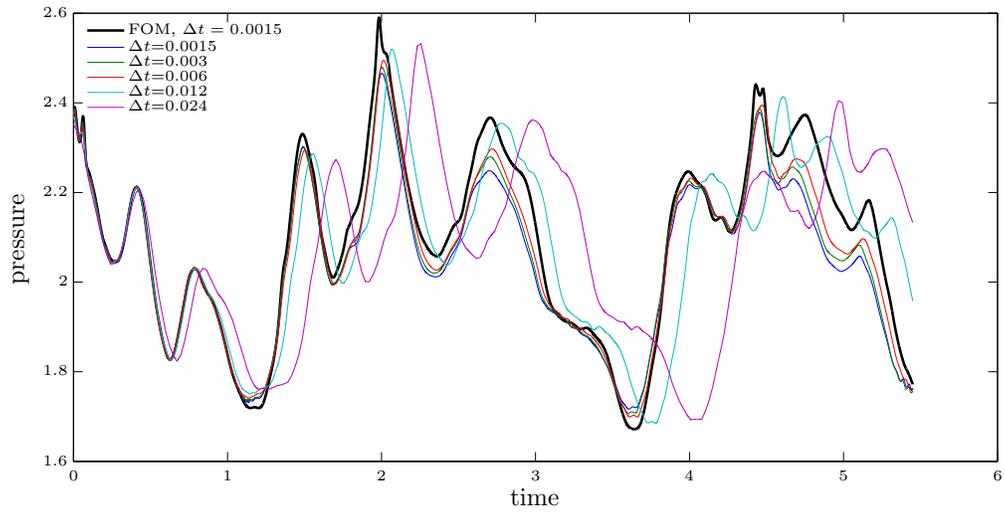


(a) Full domain

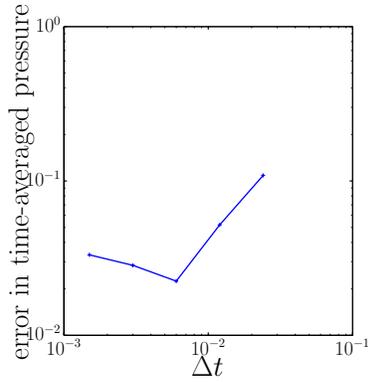


(b) Zoom on cavity

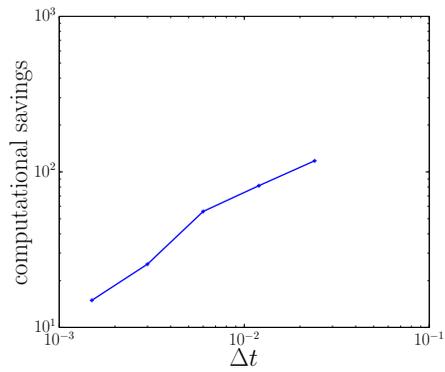
Figure 13: Sample mesh (red) embedded within original mesh.



(a) Responses



(b) Error



(c) Speedup

Figure 14:

- Discrete-optimal ROMs are more accurate than Galerkin ROMs over small time windows where Galerkin is stable (Figure 9).
- Discrete-optimal ROMs are most accurate for an intermediate time step (Figure 9).
- Adding POD modes has the effect of including higher-frequency response components (Figure 10).
- The theoretical error bound for the discrete-optimal ROM exhibits the same time step ‘dip’ as the experimentally observed error (Figure 11).
- The optimal time step for the discrete-optimal ROM decreases as modes are added to the POD basis (Figure 12).
- Adding modes to the POD basis has little effect on discrete-optimal ROM accuracy for large time steps (Figure 12).
- The optimal time step for the discrete-optimal ROM tends to increase as the time interval increases (Figure 12(g)).
- The GNAT ROM, which is discrete optimal and is equipped with complexity reduction, also produces minimal error for an intermediate time step (Figure 14).

## 8. Conclusions

This work has performed a comparative theoretical and experimental analysis of Galerkin and discrete-optimal reduced-order models for linear multistep schemes and Runge–Kutta schemes. We have demonstrated a number of new findings that have important practical implications, including conditions under which the discrete-optimal ROM has a time-continuous representation, conditions under which the two techniques are equivalent, and time-discrete error bounds for the two approaches.

Perhaps most surprisingly, we demonstrated that decreasing the time step does not necessarily decrease the error for the discrete-optimal ROM. This phenomenon arose in both the theoretical analysis and in numerical experiments. In particular, our results suggest that the time step should be ‘matched’ to the spectral content of the reduced basis. In the experiments, we showed that increasing the time step to an intermediate value decreased both the error and the simulation time by an order of magnitude in certain cases. Alternatively, decreasing the time step cause the discrete-optimal ROM to become unstable for longer time intervals. This highlights the critical importance of time-step selection for discrete-optimal ROMs.

## Acknowledgments

We thank Prof. Stephen Pope for insightful conversations related to comparing Galerkin and discrete-optimal reduced-order models; these conversations inspired this work. We also thank Prof. Charbel Farhat for permitting us the use of AERO-F, as well as Julien Cortial, Charbel Bou-Mosleh, and David Amsalem for their previous contributions in implementing nonlinear reduced-order models in AERO-F. The first author acknowledges an appointment to the Sandia National Laboratories Truman Fellowship in National Security Science and Engineering. The Truman Fellowship is sponsored by Sandia National Laboratories. Sandia National Laboratories is a multi-program laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy’s National Nuclear Security Administration under contract DE-AC04-94AL85000. The content of this publication does not necessarily reflect the position or policy of any of these institutions, and no official endorsement should be inferred.

Sandia National Laboratories is a multi-program laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy’s National Nuclear Security Administration under contract DE-AC04-94AL85000.

H. Antil would like to thank George Mason University for the appointment and startup funds.

**Input:** Set of snapshots  $\mathcal{X} \equiv \{\mathbf{w}_i\}_{i=1}^{n_w} \subset \mathbb{R}^N$ , energy criterion  $\nu \in [0, 1]$

**Output:**  $\Phi(\mathcal{X}, \nu)$

- 1: Compute thin singular value decomposition  $\mathbf{W} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ , where  $\mathbf{W} \equiv [\mathbf{w}_1/\|\mathbf{w}_1\| \cdots \mathbf{w}_{n_w}/\|\mathbf{w}_{n_w}\|]$ .
- 2: Choose dimension of truncated basis  $p = n_e(\nu)$ , where

$$n_e(\nu) \equiv \arg \min_{i \in \mathcal{V}(\nu)} i$$

$$\mathcal{V}(\nu) \equiv \{n \in \{1, \dots, n_w\} \mid \sum_{i=1}^n \sigma_i^2 / \sum_{i=1}^{n_w} \sigma_i^2 \geq \nu\},$$

and  $\mathbf{\Sigma} \equiv \text{diag}(\sigma_i)$ .

- 3:  $\Phi(\mathcal{X}, \nu) = [\mathbf{u}^1 \cdots \mathbf{u}^p]$ , where  $\mathbf{U} \equiv [\mathbf{u}^1 \cdots \mathbf{u}^{n_w}]$ .

**Algorithm 1:** Proper-orthogonal-decomposition basis computation (normalized snapshots)

## Appendix

Algorithm 1 reports the algorithm for computing a POD basis using normalized snapshots.

## References

- [1] Benner, P., Gugercin, S., and Willcox, K., “A survey of model reduction methods for parametric systems,” *Max Planck Institute Magdeburg Preprints*, Vol. MPIMD/13–14, 2013.
- [2] Sirovich, L., “Turbulence and the dynamics of coherent structures. III: dynamics and scaling,” *Quarterly of Applied Mathematics*, Vol. 45, No. 3, October 1987, pp. 583–590.
- [3] Holmes, P., Lumley, J., and Berkooz, G., *Turbulence, Coherent Structures, Dynamical Systems and Symmetry*, Cambridge University Press, 1996.
- [4] Lall, S., Krysl, P., and Marsden, J., “Structure-preserving model reduction for mechanical systems,” *Physica D: Nonlinear Phenomena*, Vol. 184, No. 1-4, 2003, pp. 304–318.
- [5] Carlberg, K., Tuminaro, R., and Boggs, P., “Efficient structure-preserving model reduction for nonlinear mechanical systems with application to structural dynamics,” *AIAA Paper 2012-1969, 53rd AIAA/ASME/ASCE/AHS/ASC Structures, Structural Dynamics and Materials Conference, Honolulu, Hawaii*, April 23–26 2012.
- [6] Carlberg, K., Tuminaro, R., and Boggs, P., “Preserving Lagrangian structure in nonlinear model reduction with application to structural dynamics,” *SIAM J. Sci. Comput.*, Vol. 37, No. 2, 2015, pp. B153–B184.
- [7] Rathinam, M. and Petzold, L. R., “A New Look at Proper Orthogonal Decomposition,” *SIAM Journal on Numerical Analysis*, Vol. 41, No. 5, 2003, pp. 1893–1925.
- [8] Foias, C., Jolly, M., Kevrekidis, I., and Titi, E., “Dissipativity of numerical schemes,” *Nonlinearity*, Vol. 4, No. 3, 1991, pp. 591.
- [9] Rempfer, D., “On low-dimensional Galerkin models for fluid flow,” *Theoretical and Computational Fluid Dynamics*, Vol. 14, No. 2, 2000, pp. 75–88.
- [10] Sirisup, S. and Karniadakis, G., “A spectral viscosity method for correcting the long-term behavior of POD models,” *Journal of Computational Physics*, Vol. 194, No. 1, 2004, pp. 92–116.
- [11] Noack, B. R., Papas, P., and Monkewitz, P. A., “The need for a pressure-term representation in empirical Galerkin models of incompressible shear flows,” *Journal of Fluid Mechanics*, Vol. 523, 2005, pp. 339–365.

- [12] Bui-Thanh, T., Willcox, K., Ghattas, O., and van Bloemen Waanders, B., “Goal-oriented, model-constrained optimization for reduction of large-scale systems,” *Journal of Computational Physics*, Vol. 224, No. 2, 2007, pp. 880–896.
- [13] Rowley, C. W., Colonius, T., and Murray, R. M., “Model reduction for compressible flows using POD and Galerkin projection,” *Physica D: Nonlinear Phenomena*, Vol. 189, No. 1–2, 2004, pp. 115–129.
- [14] Barone, M. F., Kalashnikova, I., Segalman, D. J., and Thornquist, H. K., “Stable Galerkin reduced order models for linearized compressible flow,” *Journal of Computational Physics*, Vol. 228, No. 6, 2009, pp. 1932–1946.
- [15] Kalashnikova, I. and Barone, M., “On the stability and convergence of a Galerkin reduced order model (ROM) of compressible flow with solid wall and far-field boundary treatment,” *International journal for numerical methods in engineering*, Vol. 83, No. 10, 2010, pp. 1345–1375.
- [16] Aubry, N., Holmes, P., Lumley, J. L., and Stone, E., “The dynamics of coherent structures in the wall region of a turbulent boundary layer,” *Journal of Fluid Mechanics*, Vol. 192, 1988, pp. 115–173.
- [17] Bergmann, M., Bruneau, C.-H., and Iollo, A., “Enablers for robust POD models,” *Journal of Computational Physics*, Vol. 228, No. 2, 2009, pp. 516 – 538.
- [18] Wang, Z., Akhtar, I., Borggaard, J., and Iliescu, T., “Proper orthogonal decomposition closure models for turbulent flows: a numerical comparison,” *Computer Methods in Applied Mechanics and Engineering*, Vol. 237, 2012, pp. 10–26.
- [19] San, O. and Iliescu, T., “Proper orthogonal decomposition closure models for fluid flows: Burgers equation,” *arXiv preprint arXiv:1308.3276*, 2013.
- [20] Iollo, A., Lanteri, S., and Desideri, J. A., “Stability Properties of POD–Galerkin Approximations for the Compressible Navier–Stokes Equations,” *Theoretical and Computational Fluid Dynamics*, Vol. 13, No. 6, 2000, pp. 377–396.
- [21] Marion, M. and Temam, R., “Nonlinear Galerkin methods,” *SIAM Journal on Numerical Analysis*, Vol. 26, No. 5, 1989, pp. 1139–1157.
- [22] Shen, J., “Long time stability and convergence for fully discrete nonlinear Galerkin methods,” *Applicable Analysis*, Vol. 38, No. 4, 1990, pp. 201–229.
- [23] Jolly, M., Kevrekidis, I., and Titi, E., “Preserving dissipation in approximate inertial forms for the Kuramoto–Sivashinsky equation,” *Journal of Dynamics and Differential Equations*, Vol. 3, No. 2, 1991, pp. 179–197.
- [24] Galletti, B., Bruneau, C., Zannetti, L., and Iollo, A., “Low-order modelling of laminar flow regimes past a confined square cylinder,” *Journal of Fluid Mechanics*, Vol. 503, 2004, pp. 161–170.
- [25] Balajewicz, M. and Dowell, E., “Stabilization of projection-based reduced order models of the Navier–Stokes equations,” *Nonlinear Dynamics*, Vol. 70, No. 2, 2012, pp. 1619–1632.
- [26] Balajewicz, M., Dowell, E., and Noack, B., “Low-dimensional modelling of high-Reynolds-number shear flows incorporating constraints from the Navier–Stokes equation,” *Journal of Fluid Mechanics*, Vol. 729, 2013, pp. 285–308.
- [27] Fang, F., Pain, C., Navon, I., Elsheikh, A., Du, J., and Xiao, D., “Non-linear Petrov–Galerkin methods for reduced order hyperbolic equations and discontinuous finite element methods,” *Journal of Computational Physics*, Vol. 234, No. 0, 2013, pp. 540 – 559.
- [28] Carlberg, K., Bou-Mosleh, C., and Farhat, C., “Efficient non-linear model reduction via a least-squares Petrov–Galerkin projection and compressive tensor approximations,” *International Journal for Numerical Methods in Engineering*, Vol. 86, No. 2, April 2011, pp. 155–181.

- [29] Everson, R. and Sirovich, L., “Karhunen–Loève procedure for gappy data,” *Journal of the Optical Society of America A*, Vol. 12, No. 8, 1995, pp. 1657–1664.
- [30] Carlberg, K., Farhat, C., Cortial, J., and Amsallem, D., “The GNAT method for nonlinear model reduction: effective implementation and application to computational fluid dynamics and turbulent flows,” *Journal of Computational Physics*, Vol. 242, 2013, pp. 623–647.
- [31] Carlberg, K., Cortial, J., Amsallem, D., Zahr, M., and Farhat, C., “The GNAT nonlinear model reduction method and its application to fluid dynamics problems,” *AIAA Paper 2011-3112, 6th AIAA Theoretical Fluid Mechanics Conference, Honolulu, HI*, June 27–30, 2011.
- [32] Amsallem, D., Zahr, M., and Washabaugh, K., “Fast Local Reduced Basis Updates for the Efficient Reduction of Nonlinear Systems with Hyper-Reduction,” *Advances in Computational Mathematics*, Vol. Special Issue on Model Reduction of Parameterized Systems, January 2015.
- [33] Prud’homme, C., Rovas, D. V., Veroy, K., Machiels, L., Maday, Y., Patera, A. T., and Turinici, G., “Reliable real-time solution of parameterized partial differential equations: Reduced-basis output bound methods,” *Journal of Fluids Engineering*, Vol. 124, No. 1, 2002, pp. 70–80.
- [34] Rozza, G., Huynh, D. B. P., and Patera, A. T., “Reduced basis approximation and a posteriori error estimation for affinely parametrized elliptic coercive partial differential equations,” *Archives of Computational Methods in Engineering*, Vol. 15, No. 3, 2008, pp. 229–275.
- [35] Veroy, K., Prud’homme, C., Rovas, D. V., and Patera, A. T., “A posteriori error bounds for reduced-basis approximation of parametrized noncoercive and nonlinear elliptic partial differential equations,” *AIAA Paper 2003-3847, 16th AIAA Computational Fluid Dynamics Conference, Orlando, FL*, June 23–26, 2003.
- [36] Ngoc Cuong, N., Veroy, K., and Patera, A. T., “Certified Real-Time Solution of Parametrized Partial Differential Equations,” *Handbook of Materials Modeling*, edited by S. Yip, Springer Netherlands, 2005, pp. 1529–1564.
- [37] Veroy, K. and Patera, A. T., “Certified real-time solution of the parametrized steady incompressible Navier-Stokes equations: Rigorous reduced-basis a posteriori error bounds,” *International Journal for Numerical Methods in Fluids*, Vol. 47, No. 8, 2005, pp. 773–788.
- [38] Astrid, P., Weiland, S., Willcox, K., and Backx, T., “Missing point estimation in models described by proper orthogonal decomposition,” *IEEE Transactions on Automatic Control*, Vol. 53, No. 10, 2008, pp. 2237–2251.
- [39] Ryckelynck, D., “A priori hyperreduction method: an adaptive approach,” *Journal of Computational Physics*, Vol. 202, No. 1, 2005, pp. 346–366.
- [40] LeGresley, P. A., *Application of Proper Orthogonal Decomposition (POD) to Design Decomposition Methods*, Ph.D. thesis, Stanford University, 2006.
- [41] Bos, R., Bombois, X., and Van den Hof, P., “Accelerating large-scale non-linear models for monitoring and control using spatial and temporal correlations,” *Proceedings of the American Control Conference*, Vol. 4, 2004, pp. 3705–3710.
- [42] Barrault, M., Maday, Y., Nguyen, N. C., and Patera, A. T., “An ‘empirical interpolation’ method: application to efficient reduced-basis discretization of partial differential equations,” *Comptes Rendus Mathématique Académie des Sciences*, Vol. 339, No. 9, 2004, pp. 667–672.
- [43] Chaturantabut, S. and Sorensen, D. C., “Nonlinear model reduction via discrete empirical interpolation,” *SIAM Journal on Scientific Computing*, Vol. 32, No. 5, 2010, pp. 2737–2764.
- [44] Galbally, D., Fidkowski, K., Willcox, K., and Ghattas, O., “Non-linear model reduction for uncertainty quantification in large-scale inverse problems,” *International Journal for Numerical Methods in Engineering*, Vol. 81, No. 12, published online September 2009, pp. 1581–1608.

- [45] Drohmann, M., Haasdonk, B., and Ohlberger, M., “Reduced Basis Approximation for Nonlinear Parametrized Evolution Equations based on Empirical Operator Interpolation,” *SIAM Journal on Scientific Computing*, Vol. 34, No. 2, 2012, pp. A937–A969.
- [46] Antil, H., Heinkenschloss, M., and Sorensen, D. C., “Application of the Discrete Empirical Interpolation method to Reduced Order Modeling of Nonlinear and Parametric systems,” Vol. 8 of *Springer MS&A series: Reduced Order Methods for modeling and computational r G. Rozza, Eds*, Springer-Verlag Italia, Milano, 2013.
- [47] Antil, H., Field, S., Herrmann, F., Nochetto, R., and Tiglio, M., “Two-Step Greedy Algorithm for Reduced Order Quadratures,” *Journal of Scientific Computing*, Vol. 57, 2013, pp. 604–637.
- [48] An, S., Kim, T., and James, D., “Optimizing cubature for efficient integration of subspace deformations,” *ACM Transactions on Graphics (TOG)*, Vol. 27, No. 5, 2008, pp. 165.
- [49] Farhat, C., Avery, P., Chapman, T., and Cortial, J., “Dimensional reduction of nonlinear finite element dynamic models with finite rotations and energy-based mesh sampling and weighting for computational efficiency,” *International Journal for Numerical Methods in Engineering*, Vol. 98, No. 9, 2014, pp. 625–662.
- [50] Bui-Thanh, T., Willcox, K., and Ghattas, O., “Model reduction for large-scale systems with high-dimensional parametric input space,” *SIAM Journal on Scientific Computing*, Vol. 30, No. 6, 2008, pp. 3270–3288.
- [51] Constantine, P. and Wang, Q., “Residual Minimizing Model Reduction for Parameterized Nonlinear Dynamical Systems,” *SIAM J. Sci. Comput.*, Vol. 34, No. 4, December 2012, pp. A2118–A2144.
- [52] Abgrall, R. and Amsallem, D., “Robust Model Reduction by  $L^1$ -norm Minimization and Approximation via Dictionaries: Application to Linear and Nonlinear Hyperbolic Problems,” *Stanford University Preprint*, 2015.
- [53] Geuzaine, P., Brown, G., Harris, C., and Farhat, C., “Aeroelastic dynamic analysis of a full F-16 configuration for various flight conditions,” *AIAA Journal*, Vol. 41, No. 3, 2003, pp. 363–371.
- [54] Farhat, C., Geuzaine, P., and Brown, G., “Application of a three-field nonlinear fluid-structure formulation to the prediction of the aeroelastic parameters of an F-16 fighter,” *Computers & Fluids*, Vol. 32, No. 1, 2003, pp. 3–29.
- [55] Blackford, L., Cleary, A., Choi, J., d’Azevedo, E., Demmel, J., Dhillon, I., Dongarra, J., Hammarling, S., Henry, G., Petitet, A., et al., *ScaLAPACK Users’ Guide*, Society for Industrial and Applied Mathematics, 1997.
- [56] Carlberg, K., *Model Reduction of Nonlinear Mechanical Systems via Optimal Projection and Tensor Approximation*, Ph.D. thesis, Stanford University, August 2011.