
Scalable nonlinear manifold reduced order model for dynamical systems

Ivan Zanardi¹
zanardi3@illinois.edu

Alejandro N. Diaz²
andiaz@sandia.gov

Seung Whan Chung³
chung28@llnl.gov

Marco Panesi¹
mpanesi@illinois.edu

Youngsoo Choi³
choi15@llnl.gov

¹University of Illinois Urbana-Champaign, Urbana, IL 61801

²Sandia National Laboratories, Albuquerque, NM 87123

³Lawrence Livermore National Laboratory, Livermore, CA 94550

Abstract

The domain decomposition (DD) nonlinear-manifold reduced-order model (NM-ROM) represents a computationally efficient method for integrating underlying physics principles into a neural network-based, data-driven approach. Compared to linear subspace methods, NM-ROMs offer superior expressivity and enhanced reconstruction capabilities, while DD enables cost-effective, parallel training of autoencoders by partitioning the domain into algebraic subdomains. In this work, we investigate the scalability of this approach by implementing a “bottom-up” strategy: training NM-ROMs on smaller domains and subsequently deploying them on larger, composable ones. The application of this method to the two-dimensional time-dependent Burgers’ equation shows that extrapolating from smaller to larger domains is both stable and effective. This approach achieves an accuracy of 1% in relative error and provides a remarkable speedup of nearly 700 times.

1 Introduction

Complex tasks such as design optimization and uncertainty quantification often require repeated simulations of a large-scale, parameterized, nonlinear system, commonly referred to as the full-order model (FOM). This approach can become impractical for large-scale problems due to the computational demands involved. Model reduction addresses this challenge by substituting the FOM with a more computationally efficient, lower-dimensional model known as a reduced-order model (ROM). Despite its advantages, constructing accurate and efficient ROMs presents its own set of challenges. In this work we explore the framework proposed by Diaz *et al.* [1, 2], which combines the nonlinear-manifold ROM (NM-ROM) approach with an algebraic domain-decomposition (DD) framework.

Various model reduction techniques have been integrated with DD, including reduced basis elements (RBE) [3–9] or the alternating Schwarz method [10–13], which are often tailored to specific problems and address the physical domain at the PDE level. An alternative approach is the algebraic method proposed by Hoang *et al.* [14], which involves decomposing the FOM at the discrete level and computing linear-subspace reduced-order models (LS-ROMs) for each subdomain. Although LS-ROMs perform effectively in many cases [15–40], they are known to struggle with advection-dominated problems and those with sharp gradients, which are characterized by slowly decaying Kolmogorov n -width [41]. Recent approaches, such as nonlinear-manifold reduced-order models

(NM-ROMs), address these challenges by approximating the FOM within a low-dimensional nonlinear manifold. This is typically achieved by training an autoencoder on FOM snapshot data [42–46]. However, training NM-ROMs is computationally expensive due to the high dimensionality of the FOM training data, leading to many neural network (NN) parameters. To mitigate this, Barnett *et al.* [47] first computed a low-dimensional proper orthogonal decomposition (POD) model and then trained the NN on the POD coefficients. Instead, we adopt the approach by Diaz *et al.* [1, 2], which integrates an autoencoder framework with DD. This method allows for the computation of FOM training data on subdomains, thereby reducing the dimensionality of the subdomain NM-ROM training data. Consequently, fewer parameters are required for training in each subdomain.

In this work, we adopt the DD NM-ROM approach for its proven effectiveness in tackling large-scale problems with slowly decaying Kolmogorov n -width. This method showed superior accuracy and performance compared to both LS-ROMs and monolithic NN-ROMs [1]. We specifically examine the scalability of the framework by implementing a “bottom-up” training strategy. This involves using snapshots from subdomains of a small-sized domain for training, and subsequently deploying the trained autoencoders to a larger, composable domain. The architecture employed is a wide, shallow, and sparse autoencoder with a sparsity mask applied to both the encoder input layer and decoder output layer, as described in [45]. The DD NM-ROM approach is applied to the two-dimensional time-dependent Burgers’ equation.

2 DD full order model

First consider the monolithic FOM as a system of ODEs

$$\frac{d}{dt}\mathbf{x}(t) = \mathbf{f}(\mathbf{x}(t); \boldsymbol{\mu}), \quad t \in [0, T], \quad \mathbf{x}(0) = \mathbf{x}_0(\boldsymbol{\mu}), \quad (1)$$

where $\mathbf{x} : [0, T] \rightarrow \mathbb{R}^{N_x}$ is the high-dimensional state, $\boldsymbol{\mu} \in \mathcal{D} \subset \mathbb{R}^{N_u}$ is a parameter, and $\mathbf{f} : \mathbb{R}^{N_x} \times \mathcal{D} \rightarrow \mathbb{R}^{N_x}$. Applying Backward Euler (BE) time stepping to (1) results in the following (nonlinear) system of equations, which must be solved at each time step $k = 1, \dots, N_t$:

$$\mathbf{r}(\mathbf{x}^{(k)}, \mathbf{x}^{(k-1)}; \boldsymbol{\mu}) = \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)} - \tau \mathbf{f}(\mathbf{x}^{(k)}; \boldsymbol{\mu}) = \mathbf{0}, \quad \mathbf{x}^{(0)} = \mathbf{x}_0(\boldsymbol{\mu}), \quad (2)$$

where $\tau = T/N_t$ denotes the time step, and $\mathbf{r} : \mathbb{R}^{N_x} \times \mathbb{R}^{N_x} \times \mathcal{D} \rightarrow \mathbb{R}^{N_x}$ is the residual function. FOMs of the form (2) typically arise from discretizations of partial differential equations (PDEs). One can reformulate (2) into a DD formulation by partitioning the residual equation into n_Ω systems of equations (so-called *algebraic* subdomains), coupling them via *compatibility constraints*, and converting the systems of equations into a least-squares problem, resulting in

$$\min_{(\mathbf{x}_i^{\Omega(k)}, \mathbf{x}_i^{\Gamma(k)})} \frac{h}{2} \sum_{i=1}^{n_\Omega} \left\| \mathbf{r}_i(\mathbf{x}_i^{\Omega(k)}, \mathbf{x}_i^{\Gamma(k)}, \mathbf{x}_i^{\Omega(k-1)}, \mathbf{x}_i^{\Gamma(k-1)}; \boldsymbol{\mu}) \right\|_2^2, \quad \text{s.t.} \quad \sum_{i=1}^{n_\Omega} \mathbf{A}_i \mathbf{x}_i^{\Gamma(k)} = \mathbf{0}, \quad (3)$$

accompanied by the proper initial conditions $\mathbf{x}_i^{\Omega(0)}$ and $\mathbf{x}_i^{\Gamma(0)}$, with $h > 0$ being a scaling factor and $\mathbf{x}_i^{\Omega(k)} \in \mathbb{R}^{N_i^\Omega}$, $\mathbf{x}_i^{\Gamma(k)} \in \mathbb{R}^{N_i^\Gamma}$, $\mathbf{r}_i : \mathbb{R}^{N_i^\Omega} \times \mathbb{R}^{N_i^\Gamma} \times \mathbb{R}^{N_i^\Omega} \times \mathbb{R}^{N_i^\Gamma} \times \mathcal{D} \rightarrow \mathbb{R}^{N_i^r}$, and $\mathbf{A}_i \in \{-1, 0, 1\}^{N_A \times N_i^\Gamma}$ being the i -th subdomain interior-state, interface-state, residual function, and compatibility constraint matrix, respectively. The structure of the subdomain residual functions \mathbf{r}_i and the division of the state \mathbf{x} into subdomain states $(\mathbf{x}_i^{\Omega}, \mathbf{x}_i^{\Gamma})$ are influenced by the sparsity pattern of the overall residual function \mathbf{r} . Specifically, the interior states \mathbf{x}_i^{Ω} are used solely for computing the residual \mathbf{r}_i within the i -th subdomain, while the interface states \mathbf{x}_i^{Γ} are also involved in the residual computations for adjacent subdomains. Additionally, we can define port states \mathbf{x}_j^p as a set of nodes uniquely shared by multiple subdomains. By combining multiple ports, we can generate the interface states $\mathbf{x}_i^{\Gamma} = \bigcup_{j \in i} \mathbf{x}_j^p$. The equality constraint determined by \mathbf{A}_i enforces equality on the overlapping interface states. For further details, see [1, Sec. 2] or [14, Sec. 2].

The dependence of the residual function on the previous time step and the parameter $\boldsymbol{\mu}$ will be omitted until required.

3 DD nonlinear-manifold reduced order model

In the general formulation, for each subdomain $i \in \mathcal{S}^\Omega = \{1, \dots, n_\Omega\}$, let $\mathbf{g}_i^\Omega : \mathbb{R}^{n_i^\Omega} \rightarrow \mathbb{R}^{N_i^\Omega}$, $n_i^\Omega \ll N_i^\Omega$, and $\mathbf{g}_i^\Gamma : \mathbb{R}^{n_i^\Gamma} \rightarrow \mathbb{R}^{N_i^\Gamma}$, $n_i^\Gamma \ll N_i^\Gamma$, be decoders such that $\mathbf{x}_i^\Omega \approx \mathbf{g}_i^\Omega(\widehat{\mathbf{x}}_i^\Omega)$ and

$\mathbf{x}_i^\Gamma \approx \mathbf{g}_i^\Gamma(\widehat{\mathbf{x}}_i^\Gamma)$. The corresponding encoders are denoted by \mathbf{h}_i^Ω and \mathbf{h}_i^Γ . Also let $\mathbf{B}_i \in \{0, 1\}^{N_i^B \times N_i^\Gamma}$, $N_i^B \leq N_i^\Gamma$, denote a row-sampling matrix for collocation hyper-reduction (HR). The DD NM-ROM is evaluated by solving

$$\min_{(\widehat{\mathbf{x}}_i^{\Omega(k)}, \widehat{\mathbf{x}}_i^{\Gamma(k)})} \frac{h}{2} \sum_{i=1}^{n_\Omega} \left\| \mathbf{B}_i \mathbf{r}_i \left(\mathbf{g}_i^\Omega \left(\widehat{\mathbf{x}}_i^{\Omega(k)} \right), \mathbf{g}_i^\Gamma \left(\widehat{\mathbf{x}}_i^{\Gamma(k)} \right) \right) \right\|_2^2, \quad \text{s.t.} \quad \sum_{i=1}^{n_\Omega} \widehat{\mathbf{A}}_i \widehat{\mathbf{x}}_i^{\Gamma(k)} = \mathbf{0}. \quad (4)$$

with $\widehat{\mathbf{x}}_i^{\Omega(0)} = \mathbf{h}_i^\Omega(\mathbf{x}_i^{\Omega(0)})$ and $\widehat{\mathbf{x}}_i^{\Gamma(0)} = \mathbf{h}_i^\Gamma(\mathbf{x}_i^{\Gamma(0)})$. In this work, the following set up is used:

- HR is not applied, meaning $\mathbf{B}_i = \mathbf{I}$. However, in most cases, applying HR is crucial for achieving greater speedups, as it avoids the need to evaluate the full residual of the FOM.
- The interface states encoder/decoder are formulated as a combination of port encoders/decoders, e.g., $\mathbf{g}_i^\Gamma = \bigcup_{j \in i} \mathbf{g}_j^p$. Similarly, the latent interface states are expressed as $\widehat{\mathbf{x}}_i^\Gamma = \bigcup_{j \in i} \widehat{\mathbf{x}}_j^p$.
- The compatibility constraints in (4) are expressed at the ROM latent interface states, similar to the approach used for the FOM, with $\widehat{\mathbf{A}}_i \in \{-1, 0, 1\}^{n_A \times n_i^\Gamma}$. This approach is known as the strong ROM-port constraint (SRPC) formulation [1].

The DD NM-ROM formulation (4) provides several advantages. Its training process, which involves computing \mathbf{g}_i^Ω and \mathbf{g}_j^p , is localized, requires few parameters, and can be executed in parallel. This allows for the adaptation of ROMs to specific features of the problem, potentially resulting in more compact ROMs. Additionally, the framework supports parallelization to accelerate both ROM computation/training and execution, and it enables the implementation of a ‘‘bottom-up’’ training strategy, as detailed in the introductory section.

Solver The DD FOM (3) and DD NM-ROM (4) are solved using an inexact Lagrange-Newton sequential quadratic programming (SQP) solver. In this approach, the Hessian of the Lagrangian is approximated using a Gauss-Newton method, thus eliminating the need for computing second-order derivatives of residuals and constraints in (4). Despite this approximation, the method still ensures effective convergence for both (3) and (4). For more details, refer to [1].

Autoencoder We use single-layer, wide, and sparse decoders with smooth activation functions to represent the maps \mathbf{g}_i^Ω and \mathbf{g}_j^p . The corresponding encoders, denoted \mathbf{h}_i^Ω and \mathbf{h}_j^p , are also single-layer, wide, and sparse. Shallow networks are used for computational efficiency; fewer layers correspond to fewer repeated matrix-vector multiplications when evaluating the decoders. The shallow depth necessitates a wide network to maintain enough expressiveness for use in NM-ROM. Smooth activations (i.e., *swish*) are used to ensure that \mathbf{g}_i^Ω and \mathbf{g}_j^p are continuously differentiable. Sparsity is imposed on both the input and output layers of the autoencoder to maintain symmetry across the latent layer. This sparsity pattern adopts a tri-banded structure, inspired by 2D finite difference stencils, with the number of nonzero elements per band and the spacing between bands serving as hyperparameters [1]. Normalization and de-normalization layers are also applied at the encoder input and decoder output layers, respectively, and Gaussian noise is added during training. To train the autoencoders, we first gather N_t snapshots of the interior and port states in an *offline* stage by solving (2) for each parameter μ_ℓ for $\ell = 1, \dots, M$. The snapshots are then randomly divided into an 80-20 split for training and validation. The Adam optimizer is used to minimize the MSE loss over 1000 epochs, with a batch size of 1024. We also incorporate early stopping and reduce the learning rate on plateau, starting with an initial value of 10^{-3} . The implementation is carried out using PyTorch, leveraging the PyTorch Sparse and SparseLinear packages.

4 Numerical experiment

We consider the 2D Burgers’ equation

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} + v \frac{\partial u}{\partial y} = \nu \left(\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right), \quad \frac{\partial v}{\partial t} + u \frac{\partial v}{\partial x} + v \frac{\partial v}{\partial y} = \nu \left(\frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial y^2} \right), \quad (5)$$

with viscosity $\nu = 10^{-3}$. The PDE is discretized using centered finite differences over a uniform structured mesh, with grid sizes h_x and h_y . The integration is performed for $T = 2$ s with a time step of $\tau = 0.02$ s. To evaluate the performances of the proposed DD NM-ROM, the following discrete $L^\infty - L^2$ error is used:

$$e = \max_k \left[\frac{h_x h_y}{n_\Omega} \sum_{i=1}^{n_\Omega} \left(\left\| \mathbf{x}_i^{\Omega(k)} - \mathbf{g}_i^\Omega \left(\widehat{\mathbf{x}}_i^{\Omega(k)} \right) \right\|_2^2 + \left\| \mathbf{x}_i^{\Gamma(k)} - \mathbf{g}_i^\Gamma \left(\widehat{\mathbf{x}}_i^{\Gamma(k)} \right) \right\|_2^2 \right) \right]^{1/2}. \quad (6)$$

The implementation was performed sequentially. However, to highlight the potential benefits of a parallel implementation, the reported wall clock time for computing subdomain-specific quantities in the SQP solver is based on the maximum wall clock time recorded among all subdomains. The wall clock time for the remaining steps of the SQP solver is measured as the total wall clock time.

Training was conducted using the Lassen machine, while testing and computations were carried out on the Dane machine at Lawrence Livermore National Laboratory. For more details, refer to [48]. The code can be found at <https://github.com/LLNL/DD-NM-ROM>.

Model construction For training, we employ a DD problem with periodic boundary conditions, divided into four uniformly sized subdomains ($n_\Omega = 4$). These subdomains are arranged in a 2×2 configuration within the domain $(x, y) \in [0, 1] \times [0, 1]$, with each direction containing 100 grid nodes. We explore the effects of the scalar parameter μ_i , which represents the peak value of the initial sinusoidal function in the i -th subdomain, as defined below:

$$u_i^{(0)} = v_i^{(0)} = |\mu_i \sin(2\pi x) \sin(2\pi y)| \quad \forall i \in \mathcal{S}^\Omega, \quad (7)$$

with $\mu_i = \gamma_i \xi_i$, $\gamma_i \sim \mathcal{U}(0.5, 1.5)$, $\xi_i \sim \mathcal{B}(1, 0.5)$,

where \mathcal{U} and \mathcal{B} denote uniform and binomial distributions, respectively. We sample a total of 2000 different initial 2×2 configurations from (7) using a Latin hypercube sampling (LHS) strategy, which ensures thorough exploration of all possible configurations. The sole constraint is that the bottom-left subdomain always has $\xi_0 = 1$, ensuring the model can observe waves traveling across the entire domain.

In our study, we constructed three distinct autoencoders: one for the interior states ($\mathbf{g}_i^\Omega = \mathbf{g}^\Omega \forall i \in \mathcal{S}^\Omega$), and two others for the vertical (\mathbf{g}_v^p) and horizontal (\mathbf{g}_h^p) ports. We explored various latent space dimensions, ranging from 12 to 36 for the interior states and from 6 to 14 for both the vertical and horizontal port states. The results shown in Fig. 1 are based on averaged speedups and errors from 50 distinct 2×2 test cases, evaluated across 25 different combinations of latent dimensions for vertical/horizontal port nodes (\mathcal{P}) and interior nodes (\mathcal{I}). These results indicate that the dimensions of the latent space for the port states have a notably smaller impact on performance compared to those for the interior states, likely due to differences in compression rates.

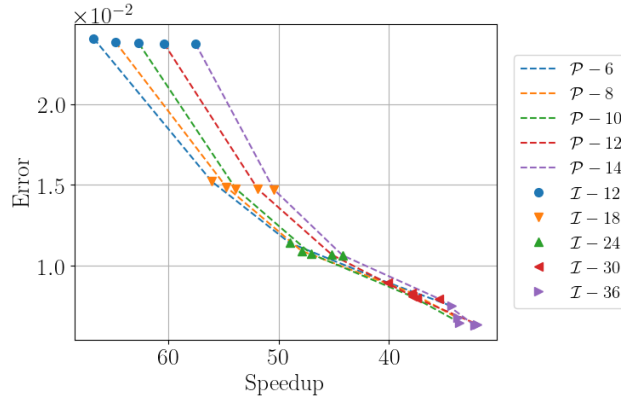


Figure 1: Averaged speedups and errors from 50 distinct 2×2 test cases, evaluated across 25 different combinations of latent dimensions for vertical/horizontal port nodes (\mathcal{P}) and interior nodes (\mathcal{I}). The results illustrate the impact of varying latent space dimensions on performance metrics.

Model deployment To assess our “bottom-up” strategy, we applied the Pareto optimal NM-ROM from the previous section (with $\mathcal{I} = 24$ and $\mathcal{P} = 10$) to a DD problem with homogeneous Neumann boundary conditions and $n_\Omega = 100$ subdomains, each uniformly sized and arranged in a 10×10 configuration within the spatial range $(x, y) \in [0, 5] \times [0, 5]$. Each direction was discretized using 500 grid nodes and the initial condition was randomly sampled from (7). At each time step, the decoder \mathbf{g}^Ω is evaluated for each subdomain to update the latent representation of the interior states. Similarly, the decoders \mathbf{g}_v^p and \mathbf{g}_h^p are evaluated for each vertical and horizontal port, respectively, to update the latent representations associated with those ports. As illustrated in Fig. 2, which shows the u velocity component at three different time instants, the predicted solution from the DD NM-ROM approximates well the true DD FOM solution. The error, quantified at 1.21×10^{-2} , is consistent

with the average errors shown in Fig. 1 for the same latent dimensions, while achieving an excellent speedup of 662.62 when compared to the DD FOM. However, increased error is observed at the shock wave fronts, attributed to accumulation error phenomena. This issue can be addressed by employing dynamic training approaches, such as adjoint-based optimization, which allows the model to consider the evolution of the approximated dynamics. The notable speedup, achieved without using HR, is attributed to the inexact SQP solver, where the size of the linear system solved at each iteration is significantly smaller for the ROM.

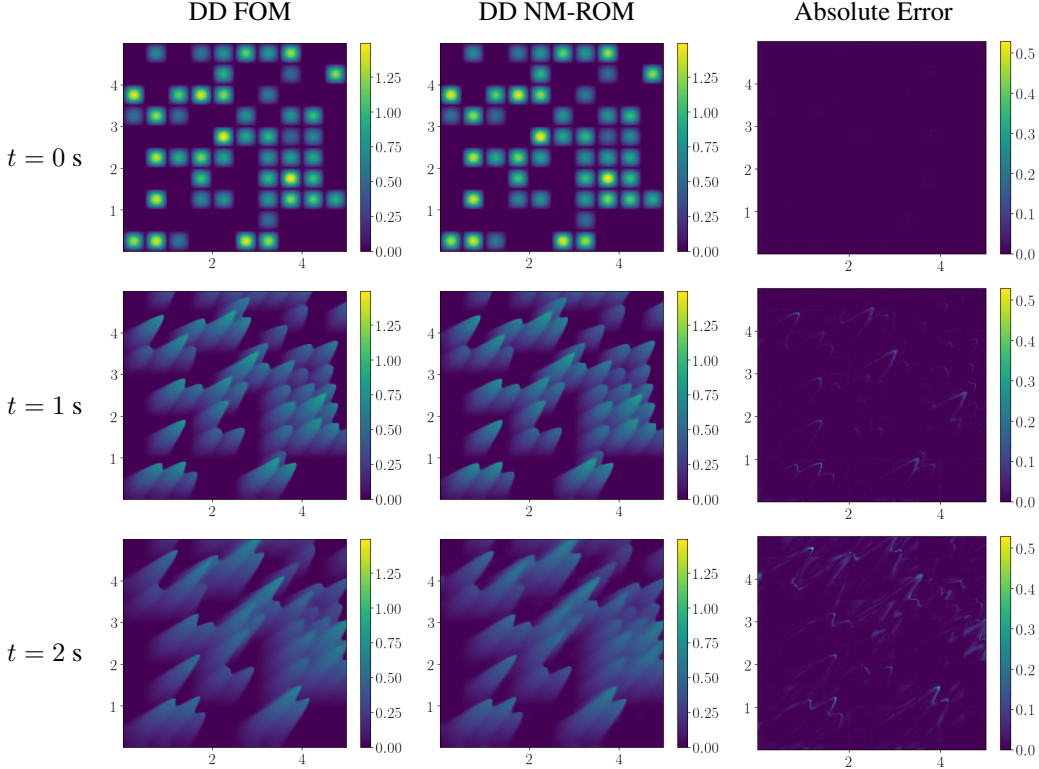


Figure 2: The u velocity predicted by the DD FOM and DD NM-ROM models, along with the absolute error between them, at various time instants. The initial condition was randomly sampled using (7) on a 10×10 configuration.

We want to emphasize that DD NM-ROM approach offers a significant advantage by enabling faster training on a smaller 2×2 domain, avoiding the high computational cost and time required to train a monolithic model on a larger 10×10 domain with random initial conditions.

5 Conclusion

In this study, we assessed the effectiveness of the “bottom-up” strategy for the DD NM-ROM framework proposed in [1, 2]. Our approach began with the development of three distinct autoencoders tailored for interior states and vertical and horizontal ports, using a small training domain with a 2×2 subdomain configuration. These models were then applied to a larger composable domain featuring a 10×10 subdomain configuration. The results indicate that extrapolating from the smaller to the larger domain is both stable and effective, achieving accuracy and speedup improvements, with the DD NM-ROM operating approximately 700 times faster than the DD FOM. Future research will aim to integrate adjoint-based optimization techniques to mitigate potential accumulation errors. Additionally, the framework’s applicability will be extended to more complex physical systems, including the Kuramoto–Sivashinsky equation, the Korteweg–De Vries equation, and the Navier–Stokes equations.

Acknowledgments and Disclosure of Funding

This work was performed at Lawrence Livermore National Laboratory. **I. Zanardi** was supported by the Data Science Summer Institute (DSSI) and LDRD (22-SI-006) at Lawrence Livermore National Laboratory and the Vannevar Bush Faculty Fellowship OUSD(RE) Grant N00014-21-1-295. **A. N. Diaz** was supported in part by a 2021 DoD National Defense Science and Engineering Graduate Fellowship and the S. Scott Collis Fellowship at Sandia National Laboratory. **S. W. Chung** was supported by LDRD (22-SI-006). **M. Panesi** was supported by the Vannevar Bush Faculty Fellowship OUSD(RE) Grant N00014-21-1-295. **Y. Choi** was supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, as part of the CHARMNET Mathematical Multifaceted Integrated Capability Center (MMICC) program, under Award Number DE-SC0023164 and partially by LDRD (22-SI-006). Lawrence Livermore National Laboratory is operated by Lawrence Livermore National Security, LLC, for the U.S. Department of Energy, National Nuclear Security Administration under Contract DE-AC52-07NA27344. IM release number: LLNL-CONF-869013.

References

- [1] Alejandro N. Diaz, Youngsoo Choi, and Matthias Heinkenschloss. A fast and accurate domain decomposition nonlinear manifold reduced order model. *Computer Methods in Applied Mechanics and Engineering*, 425:116943, 2024.
- [2] Alejandro N. Diaz. *Domain decomposition-based reduced-order models using nonlinear-manifolds and interpolatory projections*. Phd thesis, Rice University, April 2024.
- [3] Y. Maday and E. M. Rønquist. A reduced-basis element method. *J. Sci. Comput.*, 17(1-4):447–459, 2002.
- [4] Y. Maday and E. M. Rønquist. The reduced basis element method: application to a thermal fin problem. *SIAM J. Sci. Comput.*, 26(1):240–258, 2004.
- [5] L. Iapichino, A. Quarteroni, and G. Rozza. A reduced basis hybrid method for the coupling of parametrized domains represented by fluidic networks. *Comput. Methods Appl. Mech. Engrg.*, 221/222:63–82, 2012.
- [6] P. F. Antonietti, P. Pacciarini, and A. Quarteroni. A discontinuous Galerkin reduced basis element method for elliptic problems. *ESAIM Math. Model. Numer. Anal.*, 50(2):337–360, 2016.
- [7] J. L. Eftang, D. B. P. Huynh, D. J. Knezevic, E. M. Ronquist, and A. T. Patera. Adaptive port reduction in static condensation. *IFAC Proceedings Volumes*, 45(2):695–699, 2012. 7th Vienna International Conference on Mathematical Modelling.
- [8] D. B. P. Huynh, D. J. Knezevic, and A. T. Patera. A static condensation reduced basis element method: approximation and *a posteriori* error estimation. *ESAIM Math. Model. Numer. Anal.*, 47(1):213–251, 2013.
- [9] J. L. Eftang and A. T. Patera. Port reduction in parametrized component static condensation: approximation and *a posteriori* error estimation. *Internat. J. Numer. Methods Engrg.*, 96(5):269–302, 2013.
- [10] M. Buffoni, H. Telib, and A. Iollo. Iterative methods for model reduction by domain decomposition. *Comput. & Fluids*, 38(6):1160–1167, 2009.
- [11] J. Barnett, I. Tezaur, and A. Mota. The Schwarz alternating method for the seamless coupling of nonlinear reduced order models and full order models. *arXiv:2210.12551*, 2022.
- [12] K. Smetana and T. Taddei. Localized model reduction for nonlinear elliptic partial differential equations: localized training, partition of unity, and adaptive enrichment. *arXiv:2202.09872v1*, 2022.

- [13] A. Iollo, G. Sambataro, and T. Taddei. A one-shot overlapping Schwarz method for component-based model reduction: application to nonlinear elasticity. *Comput. Methods Appl. Mech. Engrg.*, 404:Paper No. 115786, 32, 2023.
- [14] C. Hoang, Y. Choi, and K. Carlberg. Domain-decomposition least-squares Petrov-Galerkin (DD-LSPG) nonlinear model reduction. *Comput. Methods Appl. Mech. Engrg.*, 384:Paper No. 113997, 41, 2021.
- [15] B. Haasdonk. Chapter 2: Reduced basis methods for parametrized PDEs - a tutorial introduction for stationary and instationary problems. In P. Benner, A. Cohen, M. Ohlberger, and K. Willcox, editors, *Model Reduction and Approximation: Theory and Algorithms*, Computational Science and Engineering, pages 65–136. SIAM, Philadelphia, 2017.
- [16] A. Quarteroni, A. Manzoni, and F. Negri. *Reduced Basis Methods for Partial Differential Equations. An Introduction*, volume 92 of *Unitext*. Springer, Cham, 2016.
- [17] M. Hinze and S. Volkwein. Proper orthogonal decomposition surrogate models for nonlinear dynamical systems: Error estimates and suboptimal control. In P. Benner, V. Mehrmann, and D. C. Sorensen, editors, *Dimension Reduction of Large-Scale Systems*, Lecture Notes in Computational Science and Engineering, Vol. 45, pages 261–306, Heidelberg, 2005. Springer-Verlag.
- [18] M. Gubisch and S. Volkwein. Chapter 1: Proper Orthogonal Decomposition for linear-quadratic optimal control. In P. Benner, A. Cohen, M. Ohlberger, and K. Willcox, editors, *Model Reduction and Approximation: Theory and Algorithms*, Computational Science and Engineering, pages 3–64, Philadelphia, 2017. SIAM.
- [19] Siu Wun Cheung, Youngsoo Choi, Dylan Matthew Copeland, and Kevin Huynh. Local lagrangian reduced-order modeling for the rayleigh-taylor instability by solution manifold decomposition. *Journal of Computational Physics*, 472:111655, 2023.
- [20] Dylan Matthew Copeland, Siu Wun Cheung, Kevin Huynh, and Youngsoo Choi. Reduced order models for lagrangian hydrodynamics. *Computer Methods in Applied Mechanics and Engineering*, 388:114259, 2022.
- [21] Kevin Carlberg, Youngsoo Choi, and Syuzanna Sargsyan. Conservative model reduction for finite-volume models. *Journal of Computational Physics*, 371:280–314, 2018.
- [22] A. C. Antoulas. *Approximation of Large-Scale Dynamical Systems*, volume 6 of *Advances in Design and Control*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2005.
- [23] P. Benner and T. Breiten. Chapter 6: Model order reduction based on system balancing. In P. Benner, A. Cohen, M. Ohlberger, and K. Willcox, editors, *Model Reduction and Approximation: Theory and Algorithms*, Computational Science and Engineering, pages 261–295, Philadelphia, 2017. SIAM.
- [24] A. C. Antoulas, C. A. Beattie, and S. Gugercin. *Interpolatory Model Reduction*, volume 21 of *Computational Science & Engineering*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2020.
- [25] C. Gu. QLMOR: A projection-based nonlinear model order reduction approach using quadratic-linear representation of nonlinear systems. *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, 30(9):1307–1320, sept. 2011.
- [26] P. Benner and T. Breiten. Two-sided projection methods for nonlinear model order reduction. *SIAM J. Sci. Comput.*, 37(2):B239–B260, 2015.
- [27] A. J. Mayo and A. C. Antoulas. A framework for the solution of the generalized realization problem. *Linear Algebra Appl.*, 425(2-3):634–662, 2007.
- [28] A. C. Antoulas, I. V. Gosea, and A. C. Ionita. Model reduction of bilinear systems in the Loewner framework. *SIAM J. Sci. Comput.*, 38(5):B889–B916, 2016.

- [29] I. V. Gosea and A. C. Antoulas. Data-driven model order reduction of quadratic-bilinear systems. *Numer. Linear Algebra Appl.*, 25(6):e2200, 2018.
- [30] Youngsoo Choi, Peter Brown, William Arrighi, Robert Anderson, and Kevin Huynh. Space-time reduced order model for large-scale linear dynamical systems with application to boltzmann transport problems. *Journal of Computational Physics*, 424:109845, 2021.
- [31] Youngkyu Kim, Karen Wang, and Youngsoo Choi. Efficient space-time reduced order model for linear dynamical systems in python using less than 120 lines of code. *Mathematics*, 9(14):1690, 2021.
- [32] Youngsoo Choi and Kevin Carlberg. Space-time least-squares petrov-galerkin projection for nonlinear model reduction. *SIAM Journal on Scientific Computing*, 41(1):A26–A58, 2019.
- [33] S. W. Cheung, Y. Choi, H. K. Springer, and T. Kadeethum. Data-scarce surrogate modeling of shock-induced pore collapse process. *Shock Waves*, 34(3):237–256, Jun 2024.
- [34] Ping-Hsuan Tsai, Seung Whan Chung, Debojyoti Ghosh, John Loffeld, Youngsoo Choi, and Jonathan L. Belof. Accelerating kinetic simulations of electrostatic plasmas with reduced-order modeling, 2023.
- [35] Quincy A Huhn, Mauricio E Tano, Jean C Ragusa, and Youngsoo Choi. Parametric dynamic mode decomposition for reduced order modeling. *Journal of Computational Physics*, 475:111852, 2023.
- [36] Sean McBane, Youngsoo Choi, and Karen Willcox. Stress-constrained topology optimization of lattice-like structures using component-wise reduced order models. *Computer Methods in Applied Mechanics and Engineering*, 400:115525, 2022.
- [37] Sean McBane and Youngsoo Choi. Component-wise reduced order model lattice-type structure design. *Computer methods in applied mechanics and engineering*, 381:113813, 2021.
- [38] Youngsoo Choi, Gabriele Boncoraglio, Spenser Anderson, David Amsallem, and Charbel Farhat. Gradient-based constrained optimization using a database of linear reduced-order models. *Journal of Computational Physics*, 423:109787, 2020.
- [39] Seung Whan Chung, Youngsoo Choi, Pratanu Roy, Thomas Moore, Thomas Roy, Tiras Y. Lin, Du T. Nguyen, Christopher Hahn, Eric B. Duoss, and Sarah E. Baker. Train small, model big: Scalable physics simulators via reduced order modeling and domain decomposition. *Computer Methods in Applied Mechanics and Engineering*, 427:117041, 2024.
- [40] Seung Whan Chung, Youngsoo Choi, Pratanu Roy, Thomas Roy, Tiras Y Lin, Du T Nguyen, Christopher Hahn, Eric B Duoss, and Sarah E Baker. Scaled-up prediction of steady Navier–Stokes equation with component reduced order modeling. *arXiv preprint arXiv:2410.21534*, 2024.
- [41] M. Ohlberger and S. Rave. Reduced basis methods: Success, limitations and future challenges. *Proceedings of the Conference Algoritmy*, pages 1–12, 2016.
- [42] K. Kashima. Nonlinear model reduction by deep autoencoder of noise response data. In *2016 IEEE 55th Conference on Decision and Control (CDC)*, pages 5750–5755, 2016.
- [43] D. Hartman and L. K. Mestha. A deep learning framework for model reduction of dynamical systems. In *2017 IEEE Conference on Control Technology and Applications (CCTA)*, pages 1917–1922, 2017.
- [44] K. Lee and K. T. Carlberg. Model reduction of dynamical systems on nonlinear manifolds using deep convolutional autoencoders. *J. Comput. Phys.*, 404:108973, 32, 2020.
- [45] Y. Kim, Y. Choi, D. Widemann, and T. Zohdi. A fast and accurate physics-informed neural network reduced order model with shallow masked autoencoder. *J. Comput. Phys.*, 451:Paper No. 110841, 29, 2022.
- [46] Youngkyu Kim, Youngsoo Choi, David Widemann, and Tarek Zohdi. Efficient nonlinear manifold reduced order model. *arXiv preprint arXiv:2011.07727*, 2020.

- [47] J. Barnett, C. Farhat, and Y. Maday. Neural-network-augmented projection-based model order reduction for mitigating the Kolmogorov barrier to reducibility. *J. Comput. Phys.*, 492:Paper No. 112420, 20, 2023.
- [48] LLNL. Compute platforms. Available at <https://hpc.llnl.gov/hardware/compute-platforms>. Accessed on 2024-08-09.