

Contents lists available at ScienceDirect

Comput. Methods Appl. Mech. Engrg.

journal homepage: www.elsevier.com/locate/cma





A fast and accurate domain decomposition nonlinear manifold reduced order model

Alejandro N. Diaz a,*, Youngsoo Choi b, Matthias Heinkenschloss a

- ^a Department of Computational Applied Mathematics and Operations Research, Rice University, Houston, 77005, TX, United States of America
- b Center for Applied Scientific Computing, Lawrence Livermore National Laboratory, Livermore, 94550, CA, United States of America

ARTICLE INFO

Keywords: Reduced order model Domain decomposition Nonlinear manifold Sparse autoencoders Neural networks Least-squares Petrov–Galerkin

ABSTRACT

This paper integrates nonlinear-manifold reduced order models (NM-ROMs) with domain decomposition (DD). NM-ROMs approximate the full order model (FOM) state in a nonlinearmanifold by training a shallow, sparse autoencoder using FOM snapshot data. These NM-ROMs can be advantageous over linear-subspace ROMs (LS-ROMs) for problems with slowly decaying Kolmogorov n-width. However, the number of NM-ROM parameters that need to be trained scales with the size of the FOM. Moreover, for "extreme-scale" problems, the storage of high-dimensional FOM snapshots alone can make ROM training expensive. To alleviate the training cost, this paper applies DD to the FOM, computes NM-ROMs on each subdomain, and couples them to obtain a global NM-ROM. This approach has several advantages: Subdomain NM-ROMs can be trained in parallel, involve fewer parameters to be trained than global NM-ROMs, require smaller subdomain FOM dimensional training data, and can be tailored to subdomain-specific features of the FOM. The shallow, sparse architecture of the autoencoder used in each subdomain NM-ROM allows application of hyper-reduction (HR), reducing the complexity caused by nonlinearity and yielding computational speedup of the NM-ROM. This paper provides the first application of NM-ROM (with HR) to a DD problem. In particular, this paper details an algebraic DD reformulation of the FOM, training a NM-ROM with HR for each subdomain, and a sequential quadratic programming (SQP) solver to evaluate the coupled global NM-ROM. Theoretical convergence results for the SQP method and a priori and a posteriori error estimates for the DD NM-ROM with HR are provided. The proposed DD NM-ROM with HR approach is numerically compared to a DD LS-ROM with HR on the 2D steady-state Burgers' equation, showing an order of magnitude improvement in accuracy of the proposed DD NM-ROM over the DD LS-ROM.

1. Introduction

Many applications in science and engineering require the high-fidelity numerical simulation of a parameterized, large-scale, nonlinear system, referred to as the full order model (FOM). For example, in the design of the airfoil of an aircraft, one repeatedly simulates the airflow around the wing to compute the lift and drag for a number of shapes to determine the optimal shape. Alternatively, in the case of digital twins, one simulates the high-fidelity FOM in real-time for given system inputs. To guarantee a high-fidelity simulation, a high-dimensional numerical model is required, resulting in high computational expense when simulating the FOM. Consequently, both many-query and real-time applications are infeasible for large-scale problems. Model

E-mail addresses: and5@rice.edu (A.N. Diaz), choi15@llnl.gov (Y. Choi), heinken@rice.edu (M. Heinkenschloss).

^{*} Corresponding author.

reduction alleviates the computational burden of simulating the high-dimensional FOM by replacing it with a low-dimensional, computationally inexpensive model, referred to as a reduced order model (ROM), that approximates the dynamics of the FOM within a tunable accuracy. This ROM can then be used in place of the FOM in real-time and many-query applications. In this work, we integrate model reduction, specifically the nonlinear-manifold ROM (NM-ROM) approach, with an algebraic domain-decomposition (DD) framework.

There are a large number of works that consider the integration of DD with model reduction. One family of approaches is based on the reduced basis element (RBE) method [1,2], in which reduced bases are computed locally for each subdomain. In the RBE method, continuity of the reduced basis solution across subdomains can be enforced via Lagrange multipliers as in [3], while others consider a discontinuous Galerkin approach [4]. Several modifications to the RBE method have been proposed, including the so-called static condensation RBE method [5-7], which computes a reduced basis (RB) approximation of the Schur complement and provides rigorous a posteriori error estimators. The reduced basis hybrid method (RBHM) [3] is another modification of the RBE method, in which a global coarse-grid solution is included in the reduced basis to ensure continuity of normal fluxes at subdomain interfaces. For RBHM, this continuity is enforced using Lagrange multipliers. Another well-studied approach uses the alternating Schwarz method, which decomposes the physical domain into two or more subdomains with or without overlap, and produces a global solution by iteratively solving the PDE on separate subdomains with boundary conditions coming from the state of neighboring subdomains at the previous iteration. The Schwarz method has been developed for both FOM-ROM and ROM-ROM couplings in, e.g., [8,9], where the ROM is projection-based using Proper Orthogonal Decomposition (POD). The approach in [10] also considers FOM-ROM and ROM-ROM couplings, but couple subdomain solutions using Lagrange multipliers, and compute bases such that the Schur complement system required for recovering interface solutions is nonsingular. The authors in [11] also consider a Schwarz approach, but use an optimization-based coupling that minimizes the jump between PDE state solutions on the interface of neighboring subdomains. The authors in [12] compute component-based ROMs based on a partition-of-unity to couple local solutions. Others have considered using DD to compute ROMs for problems with spatially localized nonlinearities [13], and for use in design optimization [14,15]. While these approaches have been successful, they are often problem-specific. That is, both RBEand Schwarz-based methods typically formulate the DD problem at the PDE level and decompose the physical domain into separate subdomains. In contrast, the authors in [16] integrate DD and ROM for a general nonlinear FOM at the fully discrete level rather than the PDE level, and algebraically decompose the FOM rather than considering a decomposition of the physical domain. The authors then use POD to compute ROMs for each subdomain, and use an optimization-based coupling to minimize the discrete PDE residual while enforcing compatibility constraints at the interfaces. In this paper, we extend the DD ROM framework of [16] to incorporate the NM-ROM approach.

We integrate NM-ROM with DD to reduce the *offline* computational cost required for training an NM-ROM, and to allow NM-ROMs to scale with increasingly large FOMs. Indeed, in the monolithic single-domain case, training NM-ROMs is expensive due to the high-dimensionality of the FOM training data, which results in a large number of neural network (NN) parameters requiring training. By coupling NM-ROM with DD, one can compute FOM training data on subdomains, thus reducing the dimensionality of subdomain NM-ROM training data, resulting in fewer parameters that need to be trained per subdomain NM-ROM. Furthermore, the subdomain NM-ROMs can be trained in parallel and adapted to subdomain-specific features of the FOM. We also note that couplings of NNs and DD for solutions of partial differential equations (PDEs) have been considered in previous work (e.g., [17–20]). However, these approaches use deep learning to solve a PDE by representing its solution as a NN and minimizing a corresponding physics-informed loss function. In contrast, our work uses autoencoders to reduce the dimensionality of an existing numerical model. The autoencoders are pretrained in an *offline* stage to find low-dimensional representations of FOM snapshot data, and used in an *online* stage to significantly reduce the computational cost and runtime of numerical simulations. Our work is the first to couple autoencoders with DD in the reduced-order modeling context.

A number of current model reduction approaches approximate the FOM solution in a low-dimensional linear subspace. In this paper, we collectively refer to this class of approaches as linear subspace ROM (LS-ROM). The LS-ROM approach supposes that the state solutions of the FOM are contained in a low-dimensional linear subspace. A basis for the linear subspace is then computed, resulting in a ROM whose state consists of the generalized coordinates of the approximate state solution in the reduced subspace. ROM approaches that follow LS-ROM include the reduced basis (RB) method [21,22], proper orthogonal decomposition (POD) [23–27], balanced truncation and balanced POD [28,29], interpolation and moment-matching based approaches [30–32], the Loewner framework [33–35], and the space–time POD [36–38] that expands the POD modes to temporal domain. Although LS-ROM approaches have been successful for a number of applications, it is well known that for advection-dominated problems and problems with sharp gradients, LS-ROM based approaches cannot produce low-dimensional subspaces where the state is well-approximated. More precisely, LS-ROM struggles when applied to problems with slowly decaying Kolmogorov *n*-width [39].

In recent years, a number of model reduction approaches have been developed to address the Kolmogorov *n*-width barrier. For example, one class of approaches leverages knowledge of the advection behavior of the given problem to enhance the approximation capabilities of linear subspaces. These approaches include composing transport maps with the reduced bases [40–42], shifting the POD basis [43], transforming the physical domain of the snapshots [44], and computing a reduced basis for a Lagrangian formulation of the PDE [45]. Other approaches consider the use of multiple linear subspaces, where instead of using a global reduced basis, one constructs multiple subspaces for separate regions in the time domain [25,26,46,47], physical domain [48], or state space [49,50]. However, each of these approaches relies upon a substantial amount of *a priori* knowledge of the governing PDE in order to improve the local approximation capabilities of linear subspaces. In contrast, another class of methods circumvents these drawbacks by approximating the FOM solution in a low-dimensional nonlinear manifold rather than a low-dimensional linear subspace. While LS-ROM approaches map the low-dimensional ROM state space to the high-dimensional FOM state space via an affine mapping, the

approaches in [51,52] consider the use of quadratic manifolds, where the ROM state space is mapped to the FOM state space via a quadratic mapping. As a further generalization of this mapping, researchers have investigated the use of neural networks to represent general nonlinear mappings from the ROM state space to the FOM state space. In particular, the use of autoencoders in the context of model reduction was first considered in the papers [53,54]. Autoencoders are a type of neural network that aims to learn the identity mapping by first encoding the inputs to some latent representation via the encoder, then decoding the latent representation to the original input space via the decoder. In [55], the authors consider the use of deep convolutional autoencoders, which augment the autoencoder architecture with convolutional layers. While their approach was successful in addressing the Kolmogorov n-width issue, the computational speedup was limited because hyper-reduction (HR) was not incorporated into their framework to properly reduce the complexity caused by nonlinear terms. The authors in [56,57] successfully apply HR in the context of NM-ROM and achieve a considerable speed-up, and do so by choosing a shallow, wide, and sparse architecture for the autoencoder. The approach in [58] also incorporates HR into an NM-ROM approach, but do so by employing a teacher-student training approach, where an autoencoder is first trained to reduce the entire state, and a second decoder is trained to only reproduce the HR nodes. This approach also permits the use of more general autoencoder architectures than the shallow, wide, and sparse architectures of [56,57]. However, an advantage of the approach in [56,57] is that autoencoder training only happens once rather than requiring a teacher-student training approach. Furthermore, this approach allows for different choices of HR nodes after NM-ROM training, whereas the approach in [58] requires fixed HR nodes.

In this paper, we extend the work of [16] on DD LS-ROM and integrate the NM-ROM approach with HR discussed in [56]. We incorporate the NM-ROM approach into this framework because of its success when applied to problems with slowly decaying Kolmogorov *n*-width. Specifically, to build ROMs on each subdomain of the DD problem, we apply NM-ROM with HR by using wide, shallow, sparse-masked autoencoders. The wide, shallow, and sparse architecture allows for hyper-reduction to be efficiently applied, thus reducing the complexity caused by nonlinearity and yielding computational speedup. Additionally, we modify the wide, shallow, and sparse architecture used in [56] to also include a sparsity mask for the encoder input layer as well as the decoder output layer. The sparsity mask at the encoder input layer results in an architecture that is symmetric across the latent layer of the autoencoder. Using *sparse* linear layers also allows one to make the encoders and decoders very wide while keeping memory costs low. Integrating NM-ROM with DD allows one to compute the FOM training snapshots on subdomains, thus significantly reducing the number of NN parameters requiring training for each subdomain.

A summary of the key contributions from this paper are as follows.

- We develop the first application of NM-ROM with HR to a DD problem.
- We modify the autoencoder architecture discussed in [56] to also include sparsity in the encoder input layer as well as the decoder output layer.
- We develop an inexact Lagrange–Newton sequential quadratic programming (SQP) method for the DD NM-ROM, and provide a theoretical convergence result for the SQP solver.
- We provide a priori and a posteriori error estimates for the DD ROM which are valid for both LS-ROM and NM-ROM.
- We numerically compare DD LS-ROM with DD NM-ROM, both with and without HR, for a number of different problem configurations using the 2D steady-state Burgers' equation.

This paper is structured as follows. Section 2 discusses the algebraic DD FOM formulation that we consider. Section 3 discusses the constrained least-squares Petrov–Galerkin (LSPG) formulation for the ROM, which respects the DD FOM formulation. We then review the LS-ROM approach based on POD in Section 3.3, and detail the NM-ROM approach in Section 3.4. We develop an inexact Lagrange–Newton sequential quadratic programming (SQP) method for the constrained LSPG-ROM in Section 4, followed by standard theoretical convergence results for the SQP solver in Section 4.2. We then discuss the autoencoder architecture used in Section 5, the application of hyper-reduction in Section 5.3, and the construction of a HR subnet in Section 5.4. In Section 6, we provide both a posteriori and a priori error bounds for the ROM solution in Theorems 5 and 6, respectively. We numerically compare the DD LS-ROM and DD NM-ROM performance, both with and without HR, on the 2D steady-state Burgers' equation in Section 7 for a number of different problem configurations. Lastly, we conclude the paper and discuss future directions in Section 8.

2. Domain-decomposition FOM formulation

This section presents the algebraic domain-decomposition formulation [16]. We consider a FOM parameterized by $\mu \in \mathcal{D} \subset \mathbb{R}^{N_{\mu}}$. Given $\mu \in \mathcal{D}$, the FOM is expressed as a parameterized system of nonlinear algebraic equations

$$r(x(\mu);\mu) = 0, \tag{1}$$

where $r: \mathbb{R}^{N_x} \times \mathcal{D} \to \mathbb{R}^{N_x}$ denotes the residual and $x(\mu) \in \mathbb{R}^{N_x}$ denotes the state. For notational simplicity, the dependence on μ is suppressed until needed. Typically r corresponds to a discretized PDE (e.g., using finite differences or finite elements) and in our target applications r is nonlinear in x.

Next we decompose the system (1) into n_{Ω} algebraic subdomains. Before giving the technical details, we describe the decomposition using a simple example illustrated in Fig. 1. In this example, suppose the system (1) is obtained from a finite difference discretization with a 5-point stencil of a scalar PDE in a rectangular domain in \mathbb{R}^2 with Dirichlet boundary conditions. The situation would be similar if the PDE was discretized using linear finite elements on a regular grid. Moreover, the decomposition into algebraic subdomains is not limited to the finite difference discretization with a 5-point stencil; this discretization is simply used for illustration. In the left plot in Fig. 1, the finite difference discretization uses a 14×5 grid of $N_x = 70$ nodes. Each node corresponds to

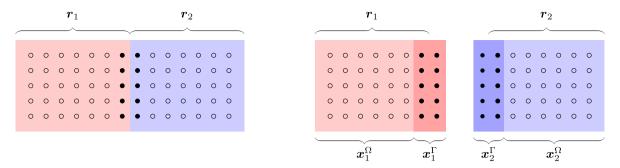


Fig. 1. Left plot: Each node in the domain corresponds to an unknown x and an equation in the system (1). The domain is subdivided into $n_{\Omega} = 2$ subdomains. Residuals corresponding to nodes marked by filled circles near the boundary in subdomain 1 depend on nodes marked by filled circles in subdomain 2, and residuals corresponding to nodes marked by filled circles near the boundary in subdomain 2 depend on nodes marked by filled circles in subdomain 1. Right plot: Variables that enter computations of residuals in one or more subdomains are duplicated as interface state variables x_1^T and x_2^T . Variables that only enter the computations of the residuals in one subdomains are the interior state variables x_1^{Ω} and x_2^{Ω} , respectively. Equality $x_1^T = x_2^T$ of interface state variables will be enforced via constraints.

a component in the vector of unknown states x and an equation in the system (1). The domain is subdivided into $n_{\Omega}=2$ subdomains. Nodes near the interface between the two subdomains are marked by filled circles. Because the PDE is discretized using a 5-point stencil, residuals corresponding to nodes marked by filled circles in subdomain 1 depend on state variables corresponding to nodes marked by filled circles in subdomain 2. Similarly, residuals corresponding to nodes marked by filled circles in subdomain 2 depend on state variables corresponding to nodes marked by filled circles in subdomain 1. In the right plot in Fig. 1, these variables are duplicated as components of the vectors of the *interface states* x_1^{Γ} and x_2^{Γ} . These have to satisfy $x_1^{\Gamma} = x_2^{\Gamma}$ and this compatibility condition will later be enforced via constraints. The other state variables are the *interior states* x_i^{Ω} , i = 1, 2. These are the state variables that only enter the residuals corresponding to subdomain i. Next we provide a detailed description of the general case.

We decompose the system (1) into $n_{\Omega} \leq N_x$ algebraic subdomains by defining so-called residual sampling matrices $P_i^v \in \{0,1\}^{N_i^v \times N_x}$ and computing subdomain residuals as

$$\mathbf{P}_{i}^{r}\mathbf{r}(\mathbf{x}) \in \mathbb{R}^{N_{i}^{r}}, \quad i = 1, \dots, n_{O}.$$

The residual sampling matrices are assumed to be algebraically non-overlapping, i.e.

$$\mathbf{P}_{:}^{r}(\mathbf{P}_{:}^{r})^{T} = \mathbf{0}, \quad \forall \ i \neq j,$$

and $\sum_{i=1}^{n_{\Omega}} N_i^r = N_x$. For problems (1) arising from a PDE discretization, the sparsity structure of the monolithic residual function r implies that subdomain residuals $P_i^r r(x)$ only depend on a subset of the full state x. Furthermore, the residual corresponding to points at the boundary of subdomain i depend on the state x at points within subdomain i and at points that belong to neighboring subdomains. Therefore, for subdomain i, we decompose the state components into *interior states*

$$\mathbf{x}_{i}^{\Omega} := \mathbf{P}_{i}^{\Omega} \mathbf{x} \in \mathbb{R}^{N_{i}^{\Omega}} \tag{2a}$$

and interface states

$$\mathbf{x}_{i}^{\Gamma} := \mathbf{P}_{i}^{\Gamma} \mathbf{x} \in \mathbb{R}^{N_{i}^{\Gamma}}, \tag{2b}$$

where $P_i^{\Omega} \in \{0,1\}^{N_i^{\Omega} \times N_x}$ denotes the *i*th interior-state sampling matrix and $P_i^{\Gamma} \in \{0,1\}^{N_i^{\Gamma} \times N_x}$ denotes the *i*th interface-state sampling matrix. The interior states $\mathbf{x}_i^{\Omega} := P_i^{\Omega} \mathbf{x}$ only enter the evaluation of the *i*th subdomain residual $P_j^{\Gamma} \mathbf{r}(\mathbf{x})$. The interface states $\mathbf{x}_i^{\Gamma} := P_i^{\Gamma} \mathbf{x}$ also enter the evaluation of another subdomain residual $P_j^{\Gamma} \mathbf{r}(\mathbf{x})$, $j \neq i$. Since the *i*th interior states only enter the evaluation of the *i*th subdomain, the interior-state sampling matrices are algebraically non-overlapping,

$$\mathbf{P}_{i}^{\Omega}(\mathbf{P}_{i}^{\Omega})^{T} = \mathbf{0}, \quad \forall i \neq j.$$

The interface state variables are duplicated across one or more subdomains, and we will describe later how to enforce equality among duplicated interface state variables.

With these specifications we can now define subdomain residual functions $\mathbf{r}_i : \mathbb{R}^{N_i^{\Omega}} \times \mathbb{R}^{N_i^{\Gamma}} \to \mathbb{R}^{N_i^{\Gamma}}$ as

$$r_i(\mathbf{x}_i^{\Omega}, \mathbf{x}_i^{\Gamma}) = \mathbf{P}_i^r \mathbf{r} \left(\left(\mathbf{P}_i^{\Omega} \right)^T \mathbf{x}_i^{\Omega} + \left(\mathbf{P}_i^{\Gamma} \right)^T \mathbf{x}_i^{\Gamma} \right). \tag{3}$$

Furthermore, the monolithic residual function (1) can be decomposed as

$$\mathbf{r}(\mathbf{x}) = \sum_{i=1}^{n_{\Omega}} \left(\mathbf{P}_{i}^{r} \right)^{T} \mathbf{r}_{i} (\mathbf{P}_{i}^{\Omega} \mathbf{x}, \mathbf{P}_{i}^{\Gamma} \mathbf{x}), \qquad \forall \ \mathbf{x} \in \mathbb{R}^{N_{X}}.$$

$$(4)$$

Eqs. (1), (3), and (4) imply that the solution (2) of (1) restricted to the ith subdomain satisfies

$$r_i(\mathbf{x}_i^{\Omega}, \mathbf{x}_i^{\Gamma}) = \mathbf{0}, \quad i = 1, \dots, n_O.$$
 (5)

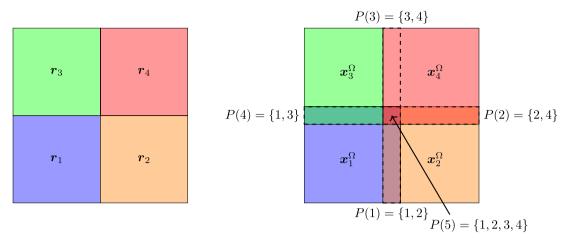


Fig. 2. Left: Residual decomposition using 4 subdomains. Notice that the residuals do not overlap. Right: State decomposition. The regions without overlap correspond to interior states \mathbf{x}_i^Q while regions with overlap correspond to interface states \mathbf{x}_i^Γ . The overlapping regions enclosed by black dashed lines represent the ports $P(j) \subset \{1, \dots, n_Q\}$.

In addition to (5), compatibility conditions must be imposed that enforce equality between overlapping interface states for neighboring subdomains. These compatibility conditions are enforced by defining n_p non-overlapping ports. Geometrically, the jth port is a subset of subdomains. The jth port has $N_j^p \leq N_x$ overlapping interface-state variables. The indices of subdomains that intersect with the jth port are $P(j) \subseteq \left\{1,\dots,n_{\Omega}\right\}$. Fig. 2 displays the ports for a 4-subdomain example configuration.

Using the ports, the compatibility conditions can be expressed as

$$\mathbf{P}_{i}^{j}\mathbf{x}_{i}^{\Gamma} = \mathbf{P}_{\ell}^{j}\mathbf{x}_{\ell}^{\Gamma}, \quad i, \ell \in P(j), \ j = 1, \dots, n_{p}, \tag{6}$$

where $P_i^j \in \{0,1\}^{N_j^p \times N_i^T}$ denotes the *j*th port sampling matrix for subdomain *i*. Because the ports are non-overlapping, if $Q(i) := \{j \mid i \in P(j)\}$ is the set of ports associated with subdomain *i*, then

$$P_i^j(P_i^j)^T = \mathbf{0}, \quad \forall j, \ell \in Q(i), j \neq \ell, \tag{7}$$

and the sum of numbers of variables in the ports associated with subdomain i is equal to the number of interface variables in the ith subdomain, $\sum_{i \in O(i)} N_i^p = N_i^\Gamma$.

As written, most conditions in (6) are redundant. Instead, for port j one needs $(|P(j)| - 1)N_j^p$ conditions, where |P(j)| denotes cardinality of P(j). For example, in Fig. 2 one needs the conditions $P_1^1 x_1^\Gamma = P_2^1 x_2^\Gamma$ for the first port $P(1) = \{1, 2\}$, and the conditions $P_1^5 x_1^\Gamma = P_2^5 x_2^\Gamma$, $P_2^5 x_2^\Gamma = P_3^5 x_3^\Gamma$, $P_3^5 x_3^\Gamma = P_4^5 x_4^\Gamma$ for the fifth port $P(5) = \{1, 2, 3, 4\}$. Removing redundant conditions (6), the port compatibility conditions (6) can be written as

$$\sum_{i=1}^{n_{\Omega}} \mathbf{A}_i \mathbf{x}_i^{\Gamma} = \mathbf{0},\tag{8}$$

where the $\mathbf{A}_i \in \{-1,0,1\}^{N_A \times N_i^{\Gamma}}$ denote the constraint matrices associated with the port compatibility conditions (6) and the total number of compatibility conditions is $N_A = \sum_{i=1}^{n_p} (|P(j)| - 1)N_j^p$. The matrix $(\mathbf{A}_1, \dots, \mathbf{A}_{n_Q})$ has full row rank.

In summary, the algebraic DD formulation of the FOM (1) is given by

$$\mathbf{r}_i(\mathbf{x}_i^{\Omega}, \mathbf{x}_i^{\Gamma}) = \mathbf{0}, \qquad i = 1, \dots, n_{\Omega},$$
 (9a)

$$\sum_{i=1}^{n_{\Omega}} \mathbf{A}_{i} \mathbf{x}_{i}^{\Gamma} = \mathbf{0}. \tag{9b}$$

Associated with (9) we also consider the nonlinear least-squares problem with equality constraints,

$$\min_{\left(\boldsymbol{x}_{i}^{\Omega}, \boldsymbol{x}_{i}^{\Gamma}\right), i=1, \dots, n_{\Omega}} \quad \frac{1}{2} \sum_{i=1}^{n_{\Omega}} \left\| \boldsymbol{r}_{i} \left(\boldsymbol{x}_{i}^{\Omega}, \boldsymbol{x}_{i}^{\Gamma}\right) \right\|_{2}^{2} \tag{10a}$$

s.t.
$$\sum_{i=1}^{n_{\Omega}} \mathbf{A}_i \mathbf{x}_i^{\Gamma} = \mathbf{0}. \tag{10b}$$

The connections between the formulations (1), (9), and (10) are summarized in the following theorem.

Theorem 1. If x solves the FOM (1) then $(\mathbf{x}_i^{\Omega}, \mathbf{x}_i^{\Gamma})$ with $\mathbf{x}_i^{\Omega} := \mathbf{P}_i^{\Omega} \mathbf{x}$ and $\mathbf{x}_i^{\Gamma} := \mathbf{P}_i^{\Gamma} \mathbf{x}$, $i = 1, \dots, n_{\Omega}$, solves the algebraic DD formulation of the FOM (9) and vice versa. A solution of (9) also solves (10), and a solution of (10) with objective function value equal to zero solves (9).

The proof of Theorem 1 follows immediately from the construction of (9).

In the FOM context, the constrained nonlinear least-squares problem formulation (10) is not needed, but we include it here because it will become important for the model reduction derivation in Section 3, where we use a least squares formulation for the subdomain ROMs. The constrained nonlinear least-squares formulation (10) of the FOM could be solved using a Lagrange–Newton sequential quadratic programming (SQP) method with Gauss–Newton Hessian approximation and the constrained nonlinear least-squares formulation of the NM-ROM corresponding to (10) will be solved using the Lagrange–Newton SQP method discussed in Section 4.

In principle, an alternative DD formulation of the FOM is possible, which reverses the role of satisfying the subdomain equations and of the compatibility conditions. Instead of (9), one can impose the subdomain equations $\mathbf{r}_i(\mathbf{x}_i^{\Omega}, \mathbf{x}_i^{\Gamma}) = \mathbf{0}$, $i = 1, ..., n_{\Omega}$, as constraints and use a least squares formulation of the compatibility conditions as the objective. This is used, e.g., in [11,59]. However, since we use LSPG-ROMs, which in general do not have zero residual, these cannot be incorporated as equality constraints. In contrast, the formulation (10) can be used with subdomain LSPG-ROMs, as we will describe in the next section.

3. Domain-decomposition ROM

The ROM construction is built on the assumption that the high-dimensional subdomain state variables $\mathbf{x}_i^{\Omega} \in \mathbb{R}^{N_i^{\Omega}}$ and $\mathbf{x}_i^{\Gamma} \in \mathbb{R}^{N_i^{\Gamma}}$, $\mathbf{x}_i^{\Omega} \in \mathbb{R}^{N_i^{\Gamma}}$, $\mathbf{x}_i^{\Omega} \in \mathbb{R}^{N_i^{\Gamma}}$, $\mathbf{x}_i^{\Omega} \in \mathbb{R}^{N_i^{\Gamma}}$, $\mathbf{x}_i^{\Gamma} \in \mathbb{R}^{$

$$\mathbf{x}_{i}^{\Omega} \approx \mathbf{g}_{i}^{\Omega}(\widehat{\mathbf{x}}_{i}^{\Omega}), \quad \mathbf{x}_{i}^{\Gamma} \approx \mathbf{g}_{i}^{\Gamma}(\widehat{\mathbf{x}}_{i}^{\Gamma}), \quad i = 1, \dots, n_{O}.$$
 (11)

In the traditional LS-ROM the maps \mathbf{g}_i^Ω and \mathbf{g}_i^Γ are linear, whereas in our NM-ROM these maps are computed via autoencoders/decoders. Assuming that we have maps \mathbf{g}_i^Ω and \mathbf{g}_i^Γ such that (11) holds, we discuss how to construct the ROM in Section 3.1. Specifically, our ROM is based on the constrained nonlinear least-squares formulation (10) of the FOM. One issue in the ROM construction based on (10) is the formulation of compatibility constraints for the ROM. In Section 3.1 we use a formulation following [16] and in Section 3.2 we provide an alternative formulation of the ROM compatibility constraints by constructing the maps $\mathbf{g}_i^\Gamma(\hat{\mathbf{x}}_i^\Gamma)$ in a suitable way. The detailed construction of maps \mathbf{g}_i^Ω and \mathbf{g}_i^Γ such that (11) holds is discussed in Sections 3.3, 3.4, and 5. Specifically, we will review the traditional LS-ROM in Section 3.3. Sections 3.4 and 5 discuss how to compute these maps via the NM-ROM approach.

3.1. Least-squares formulation

Given maps \mathbf{g}_i^{Ω} and \mathbf{g}_i^{Γ} such that (11) holds, a naive way of computing the ROM is to simply replace \mathbf{x}_i^{Ω} and \mathbf{x}_i^{Γ} in the constrained nonlinear least-squares formulation (10) of the FOM by $\mathbf{g}_i^{\Omega}(\hat{\mathbf{x}}_i^{\Omega})$ and $\mathbf{g}_i^{\Gamma}(\hat{\mathbf{x}}_i^{\Gamma})$. An evaluation of this ROM requires the solution of

$$\min_{(\widehat{\mathbf{x}}_{i}^{\Omega}, \widehat{\mathbf{x}}_{i}^{\Gamma}), i=1, \dots, n_{\Omega}} \frac{1}{2} \sum_{i=1}^{n_{\Omega}} \left\| \mathbf{r}_{i} \left(\mathbf{g}_{i}^{\Omega} (\widehat{\mathbf{x}}_{i}^{\Omega}), \mathbf{g}_{i}^{\Gamma} (\widehat{\mathbf{x}}_{i}^{\Gamma}) \right) \right\|_{2}^{2}$$
(12a)

s.t.
$$\sum_{i=1}^{n_{\Omega}} \mathbf{A}_i \mathbf{g}_i^{\Gamma}(\hat{\mathbf{x}}_i^{\Gamma}) = \mathbf{0}.$$
 (12b)

This corresponds to a (naive) LSPG-ROM. There are two issues with this formulation.

The first issue is that, just as in the case of the classical LSPG-ROM, the complexity of the evaluation of the subdomain residuals, i.e., $(\hat{x}_i^\Omega, \hat{x}_i^\Gamma) \to (g_i^\Omega(\hat{x}_i^\Omega), g_i^\Gamma(\hat{x}_i^\Gamma)) \to r_i(g_i^\Omega(\hat{x}_i^\Omega), g_i^\Gamma(\hat{x}_i^\Gamma))$ scales with the size N_i^Ω and N_i^Γ of the FOM. This issue is addressed using so-called hyper-reduction (HR). See, e.g., [60] for an overview. HR replaces the residual $r_i(g_i^\Omega(\hat{x}_i^\Omega), g_i^\Gamma(\hat{x}_i^\Gamma))$ in ((12)a) by $B_i r_i(g_i^\Omega(\hat{x}_i^\Omega), g_i^\Gamma(\hat{x}_i^\Gamma))$, where $B_i \in \mathbb{R}^{N_i^B \times N_i^F}$, $N_i^B \leq N_i^F$, is determined by the HR approach. For example, $B_i = I$ corresponds to vanilla LSPG, $B_i = Z_i$, where $Z_i \in \{0,1\}^{N_i^Z \times N_i^F}$, $N_i^Z < N_i^F$, denotes a row-sampling matrix, corresponds to collocation HR, and $B_i = (Z_i \Phi_i^F)^\dagger Z_i$, where Z_i is as before, $\Phi_i^F \in \mathbb{R}^{N_i^T \times n_i^F}$, $i = 1, \dots, n_\Omega$, denotes a reduced subspace for the corresponding subdomain residual and the superscript \dagger denotes the Moore–Penrose pseudoinverse, corresponds to gappy POD HR [61–63]. Further details on HR for our DD NM-ROM are discussed in Section 5.3. For the application of HR to DD LS-ROM, we refer the reader to [16].

The second issue with (12) is that it involves the same number of constraints ((12)b) as the FOM (10), but fewer degrees of freedom to satisfy them. In the extreme case, it may be impossible to satisfy the constraints ((12)b). One approach, following [16], is to replace A_i in ((12)b) by CA_i , where $C \in \mathbb{R}^{n_C \times N_A}$, $n_C \ll N_A$, is a test matrix that converts ((12)b) into a so-called "weak compatibility constraint". We will choose C to be a Gaussian matrix, but in principle other choices of C can be used.

To summarize, given maps $\mathbf{g}_i^{\Omega}: \mathbb{R}^{n_i^{\Omega}} \to \mathbb{R}^{N_i^{\Omega}}$ and $\mathbf{g}_i^{\Gamma}: \mathbb{R}^{n_i^{\Gamma}} \to \mathbb{R}^{N_i^{\Gamma}}$ such that (11) holds, given HR matrices $\mathbf{B}_i \in \mathbb{R}^{N_i^B \times N_i^r}$, $N_i^B \leq N_i^r$, $i = 1, \dots, n_{\Omega}$, and given $C \in \mathbb{R}^{n_C \times N_A}$, $n_C \ll N_A$, our DD-LSPG-ROM is evaluated by solving

$$\min_{(\widehat{\mathbf{x}}_{i}^{\Omega}, \widehat{\mathbf{x}}_{i}^{\Gamma}), i=1, \dots, n_{\Omega}} \frac{1}{2} \sum_{i=1}^{n_{\Omega}} \left\| \boldsymbol{B}_{i} \boldsymbol{r}_{i} \left(\boldsymbol{g}_{i}^{\Omega} \left(\widehat{\mathbf{x}}_{i}^{\Omega} \right), \boldsymbol{g}_{i}^{\Gamma} \left(\widehat{\mathbf{x}}_{i}^{\Gamma} \right) \right) \right\|_{2}^{2}$$

$$(13a)$$

s.t.
$$\sum_{i=1}^{n_{\Omega}} C A_i g_i^{\Gamma}(\hat{\mathbf{x}}_i^{\Gamma}) = \mathbf{0}.$$
 (13b)

The DD-LSPG-ROM formulation (13) will be referred to as the weak FOM-port constraint (WFPC) formulation.

While the FOM (9) or (10) has linear constraints, the WFPC formulation has nonlinear constraints in general. Corresponding to $\mathbf{g}_i^\Gamma: \mathbb{R}^{n_i^\Gamma} \to \mathbb{R}^{N_i^\Gamma}$ is a function $\mathbf{h}_i^\Gamma: \mathbb{R}^{N_i^\Gamma} \to \mathbb{R}^{n_i^\Gamma}$ such that $\left\|\mathbf{g}_i^\Gamma(\mathbf{h}_i^\Gamma(\mathbf{x}_i^{\Gamma,\text{train}})) - \mathbf{x}_i^{\Gamma,\text{train}}\right\|$ is small for some training/snapshot data $\mathbf{x}_i^{\Gamma,\text{train}}$, $i=1,\ldots,n_\Omega$, that satisfy the linear FOM constraints ((9)b). See Sections 3.3 and 3.4. Thus $\sum_{i=1}^{n_\Omega} C\mathbf{A}_i \mathbf{g}_i^\Gamma(\hat{\mathbf{x}}_i^\Gamma)$ is guaranteed to be small at these training/snapshot data. Existence of points that satisfy ((13)b) in the nonlinear case is still an open issue. However, in our numerical examples we have not observed any issues related to existence of feasible points for (13). The existence of solutions of (13) can be guaranteed under mild conditions that are typical for optimization problems.

Theorem 2. Let $\widetilde{\mathbf{x}}_i^{\Gamma}$, $i = 1, ..., n_{\Omega}$, satisfy the constraints ((13)b) and let $\widetilde{\mathbf{x}}_i^{\Omega}$, $i = 1, ..., n_{\Omega}$, be arbitrary. If the residual function \mathbf{r} and the maps \mathbf{g}_i^{Ω} , \mathbf{g}_i^{Γ} , $i = 1, ..., n_{\Omega}$, are continuous and if the level set

$$L = \left\{ (\widehat{\mathbf{x}}_{1}^{\Omega}, \widehat{\mathbf{x}}_{1}^{\Gamma}, \dots, \widehat{\mathbf{x}}_{n_{\Omega}}^{\Omega}, \widehat{\mathbf{x}}_{n_{\Omega}}^{\Gamma}) : \sum_{i=1}^{n_{\Omega}} C A_{i} \mathbf{g}_{i}^{\Gamma}(\widehat{\mathbf{x}}_{i}^{\Gamma}) = \mathbf{0}, \right.$$

$$\left. \sum_{i=1}^{n_{\Omega}} \left\| \mathbf{B}_{i} \mathbf{r}_{i} \left(\mathbf{g}_{i}^{\Omega} \left(\widehat{\mathbf{x}}_{i}^{\Omega} \right), \mathbf{g}_{i}^{\Gamma} \left(\widehat{\mathbf{x}}_{i}^{\Gamma} \right) \right) \right\|_{2}^{2} \le \sum_{i=1}^{n_{\Omega}} \left\| \mathbf{B}_{i} \mathbf{r}_{i} \left(\mathbf{g}_{i}^{\Omega} \left(\widehat{\mathbf{x}}_{i}^{\Omega} \right), \mathbf{g}_{i}^{\Gamma} \left(\widehat{\mathbf{x}}_{i}^{\Gamma} \right) \right) \right\|_{2}^{2} \right\}$$

is bounded, then (13) has a solution.

Proof. If $(\hat{\mathbf{x}}_i^{\Omega}, \hat{\mathbf{x}}_i^{\Gamma})$, $i = 1, ..., n_{\Omega}$, solves (13), then it also solves the minimization problem with the constraint $(\hat{\mathbf{x}}_1^{\Omega}, \hat{\mathbf{x}}_1^{\Gamma}, ..., \hat{\mathbf{x}}_{n_{\Omega}}^{\Omega}, \hat{\mathbf{x}}_{n_{\Omega}}^{\Gamma}) \in L$ added. The feasible set of this new minimization problem is compact, the objective function is continuous, and therefore this minimization problem has a solution, which is also a solution of (13).

Instead of the weak compatibility constraint ((13)b) one can also construct the maps \mathbf{g}_i^{Γ} , $i = 1, ..., n_{\Omega}$, such that compatibility is enforced strongly for appropriate components of $\mathbf{g}_i^{\Gamma}(\hat{\mathbf{x}}_i^{\Gamma})$, $i = 1, ..., n_{\Omega}$. This approach is introduced in the following Section 3.2.

3.2. Strong ROM-port constraints

In the general formulation, the maps (11) are computed separately for each subdomain i. However, since the interface variables $\mathbf{x}_{\ell}^{\Gamma}$, $\mathbf{x}_{\ell}^{\Gamma}$ are identical on each port j associated with the subdomains i, ℓ , i.e., on each port j with $i, \ell \in P(j)$ (see (6)), one can instead reduce the interface variables on each port and then combine the reduced port interface variables to a reduced interface variable. By (6) the interface variables for port j must satisfy

$$\mathbf{x}_{i}^{p} = \mathbf{P}_{i}^{j} \mathbf{x}_{i}^{\Gamma} \in \mathbb{R}^{N_{j}^{p}}, \qquad \forall i \in P(j).$$

$$(14)$$

For each port j we reduce the FOM port variables x_j^p , i.e., for each port j we compute a single map $g_j^p : \mathbb{R}^{n_j^p} \to \mathbb{R}^{N_j^p}$, where

$$g_i^p(\hat{\mathbf{x}}_i^p) \approx \mathbf{x}_i^p = \mathbf{P}_i^j \mathbf{x}_i^\Gamma, \quad \forall i \in P(j).$$
 (15)

The reduced interface variable \hat{x}_i^{Γ} is now computed by concatenating all port variables \hat{x}_j^{ρ} with $i \in P(j)$. This leads to the ROM port sampling matrices $\hat{P}_i^j \in \{0,1\}^{n_j^{\rho} \times n_i^{\Gamma}}$ which are defined through

$$\hat{\mathbf{x}}_{i}^{p} = \hat{\mathbf{P}}_{i}^{j} \hat{\mathbf{x}}_{i}^{T}, \qquad i \in P(j). \tag{16}$$

Equation (16) implies that on the *i*th port we have

$$\widehat{\boldsymbol{P}}_{i}^{j}\widehat{\boldsymbol{x}}_{i}^{\Gamma} = \widehat{\boldsymbol{P}}_{\ell}^{j}\widehat{\boldsymbol{x}}_{\ell}^{\Gamma}, \quad i, \ell \in P(j). \tag{17}$$

To introduce parallelism, ROM interface variables \hat{x}_i^{Γ} are introduced for each subdomain, and are coupled by enforcing (17). By construction, the ROM ports are non-overlapping, i.e.

$$\widehat{\boldsymbol{P}}_{i}^{j} \left(\widehat{\boldsymbol{P}}_{i}^{\ell}\right)^{T} = \boldsymbol{0}, \quad \forall j, \ell \in Q(i), j \neq \ell,$$

$$(18)$$

and $n_i^{\Gamma} = \sum_{j \in Q(i)} n_j^{p}$. As discussed in Section 2, after removing redundant conditions in (17), one can write the ROM port compatibility conditions (17) as

$$\sum_{i=1}^{n_{\Omega}} \hat{\mathbf{A}}_{i} \hat{\mathbf{x}}_{i}^{\Gamma} = \mathbf{0}, \tag{19}$$

where $\hat{A}_i \in \{-1,0,1\}^{n_A \times n_i^{\Gamma}}$, $n_A = \sum_{j=1}^{n_p} (|P(j)| - 1)n_j^p$, denote the constraint matrices associated with port compatibility conditions

$$(\widehat{\boldsymbol{A}}_1, \dots, \widehat{\boldsymbol{A}}_{n_Q}) \in \mathbb{R}^{n_A \times \sum_{i=1}^{n_Q} n_i^{\Gamma}}$$
(20)

has full row rank n_A , and $n_A < \sum_{i=1}^{n_\Omega} n_i^{\Gamma}$.

The map $\mathbf{g}_i^{\Gamma} : \mathbb{R}^{n_i^{\Gamma}} \to \mathbb{R}^{N_i^{\Gamma}}$ that approximates the interface state \mathbf{x}_i^{Γ} is implied by the port maps \mathbf{g}_j^{P} . To see this, note that the FOM compatibility conditions (6) and the non-overlapping condition (7) allow one to rewrite \mathbf{x}_i^{Γ} as

$$\boldsymbol{x}_i^{\Gamma} = \sum_{i \in O(i)} (\boldsymbol{P}_i^j)^T \boldsymbol{P}_i^j \boldsymbol{x}_i^{\Gamma}. \tag{21}$$

Thus, using (15), (16), and (21) the map $\mathbf{g}_i^{\Gamma}: \mathbb{R}^{n_i^{\Gamma}} \to \mathbb{R}^{N_i^{\Gamma}}$ that approximates the interface state \mathbf{x}_i^{Γ} is given by

$$\mathbf{g}_{i}^{\Gamma}(\hat{\mathbf{x}}_{i}^{\Gamma}) := \sum_{i \in O(i)} (\mathbf{P}_{i}^{j})^{T} \mathbf{g}_{j}^{p} \left(\hat{\mathbf{P}}_{i}^{j} \hat{\mathbf{x}}_{i}^{\Gamma}\right). \tag{22}$$

In particular, the definition (22) of g_i^{Γ} and the ROM compatibility conditions (17) imply that

$$\boldsymbol{P}_{i}^{j}\boldsymbol{g}_{i}^{\Gamma}(\widehat{\boldsymbol{x}}_{i}^{\Gamma}) = \boldsymbol{g}_{j}^{p}\left(\widehat{\boldsymbol{P}}_{i}^{j}\widehat{\boldsymbol{x}}_{i}^{\Gamma}\right) = \boldsymbol{g}_{j}^{p}\left(\widehat{\boldsymbol{P}}_{\ell}^{j}\widehat{\boldsymbol{x}}_{\ell}^{\Gamma}\right) = \boldsymbol{P}_{\ell}^{j}\boldsymbol{g}_{\ell}^{\Gamma}(\widehat{\boldsymbol{x}}_{\ell}^{\Gamma})$$

for all $i, \ell \in P(j)$ and for all ports P(j). This implies that strong compatibility holds for the FOM ports:

$$\sum_{i=1}^{n_{\Omega}} \mathbf{A}_i \mathbf{g}_i^{\Gamma}(\hat{\mathbf{x}}_i^{\Gamma}) = \mathbf{0}.$$

In summary, if port maps g_i^{Γ} are constructed such that (15) holds and the implied interface maps g_i^{Γ} are (22), then the DD-LSPG-ROM is evaluated by solving

$$\min_{(\widehat{\mathbf{x}}_{i}^{\Omega}, \widehat{\mathbf{x}}_{i}^{\Gamma}), i=1, \dots, n_{\Omega}} \frac{1}{2} \sum_{i=1}^{n_{\Omega}} \left\| \boldsymbol{B}_{i} \boldsymbol{r}_{i} \left(\boldsymbol{g}_{i}^{\Omega} \left(\widehat{\mathbf{x}}_{i}^{\Omega} \right), \boldsymbol{g}_{i}^{\Gamma} \left(\widehat{\mathbf{x}}_{i}^{\Gamma} \right) \right) \right\|_{2}^{2}$$
(23a)

s.t.
$$\sum_{i=1}^{n_{\Omega}} \widehat{A}_i \widehat{\mathbf{x}}_i^{\Gamma} = \mathbf{0}. \tag{23b}$$

The formulation (23) will be referred to as the strong ROM-port constraint (SRPC) formulation.

In contrast to (13) the constraints in (23) are linear and the set of feasible points for (23) is the null-space of the constraint matrix (23). Thus existence of feasible points for (23) is now trivial. Existence of solutions of (23) can be guaranteed analogously

Theorem 3. i. The null-space of the constraint matrix (23) has dimension $(\sum_{i=1}^{n_{\Omega}} n_i^{\Gamma}) - n_A \ge 1$. ii. Let $\widetilde{\mathbf{x}}_i^{\Gamma}$, $i = 1, \dots, n_{\Omega}$, satisfy the constraints ((23)b) and let $\widetilde{\mathbf{x}}_i^{\Omega}$, $i = 1, \dots, n_{\Omega}$, be arbitrary. If the residual function \mathbf{r} and the maps \mathbf{g}_i^{Ω} , \mathbf{g}_i^{Γ} , $i = 1, \dots, n_{\Omega}$, are continuous and if the level set

$$L = \left\{ (\widehat{\mathbf{x}}_{1}^{\Omega}, \widehat{\mathbf{x}}_{1}^{\Gamma}, \dots, \widehat{\mathbf{x}}_{n_{\Omega}}^{\Omega}, \widehat{\mathbf{x}}_{n_{\Omega}}^{\Gamma}) : \sum_{i=1}^{n_{\Omega}} \widehat{\mathbf{A}}_{i} \widehat{\mathbf{x}}_{i}^{\Gamma} = \mathbf{0}, \right.$$

$$\left. \sum_{i=1}^{n_{\Omega}} \left\| \mathbf{B}_{i} \mathbf{r}_{i} \left(\mathbf{g}_{i}^{\Omega} \left(\widehat{\mathbf{x}}_{i}^{\Omega} \right), \mathbf{g}_{i}^{\Gamma} \left(\widehat{\mathbf{x}}_{i}^{\Gamma} \right) \right) \right\|_{2}^{2} \le \sum_{i=1}^{n_{\Omega}} \left\| \mathbf{B}_{i} \mathbf{r}_{i} \left(\mathbf{g}_{i}^{\Omega} \left(\widehat{\mathbf{x}}_{i}^{\Omega} \right), \mathbf{g}_{i}^{\Gamma} \left(\widehat{\mathbf{x}}_{i}^{\Gamma} \right) \right) \right\|_{2}^{2} \right\}$$

is bounded, then (23) has a solution.

Proof. The first part follows immediately from the properties of the constraint matrix (23). The proof of ii. is analogous to the proof of Theorem 2ii.

So far, we have specified our DD-LSPG-ROM (13) or (23) given maps g_i^Ω and g_i^Γ such that (11) holds, or given maps g_i^Ω , g_i^p and implied interface maps (22) such that (11) holds. Next we discuss how these maps can be computed. In the following Section 3.3 we first review traditional approaches based on linear subspaces to compute g_i^{Ω} and g_i^{Γ} (or g_i^{Ω} and g_i^{ρ}). In Section 3.4 we will then introduce the nonlinear-manifold ROM.

3.3. Linear-subspace ROM

First we review linear subspace approximation to construct the maps g_i^{Ω} and g_i^{Γ} , or g_i^{Ω} and g_i^{ρ} . We will refer to resulting ROM as LS-ROM. The LS-ROM approach supposes that the state solutions of the FOM are contained in a low-dimensional linear subspace. A basis for the linear subspace is then computed, resulting in a ROM whose state consists of the generalized coordinates of the state solution in the reduced subspace. The use of LS-ROM for the DD problem (13) has already been considered in [16], where the LS-ROM bases are computed using POD, but in principle any choice of basis can be used. The numerics in Section 7 also use

POD for consistency with previous works. We briefly review POD here for completeness. A thorough treatment of POD can be found in [23].

As mentioned above, the LS-ROM approach approximates the FOM states \mathbf{x}_i^{Ω} , \mathbf{x}_i^{Γ} in a linear subspace. Hence $\mathbf{g}_i^{\Omega}: \mathbb{R}^{n_i^{\Omega}} \to \mathbb{R}^{N_i^{\Omega}}$ and $\mathbf{g}_i^{\Gamma}: \mathbb{R}^{n_i^{\Gamma}} \to \mathbb{R}^{N_i^{\Gamma}}$ are linear maps,

$$\mathbf{g}_{i}^{\Omega}: \quad \widehat{\mathbf{x}}_{i}^{\Omega} \mapsto \mathbf{\Phi}_{i}^{\Omega} \widehat{\mathbf{x}}_{i}^{\Omega}, \mathbf{g}_{i}^{\Gamma}: \quad \widehat{\mathbf{x}}_{i}^{\Gamma} \mapsto \mathbf{\Phi}_{i}^{\Gamma} \widehat{\mathbf{x}}_{i}^{\Gamma},$$

where $\Phi_i^{\Omega} \in \mathbb{R}^{N_i^{\Omega} \times n_i^{\Omega}}$ and $\Phi_i^{\Gamma} \in \mathbb{R}^{N_i^{\Gamma} \times n_i^{\Gamma}}$ are basis matrices corresponding to the reduced linear subspaces. Consequently, the Jacobians $\frac{d}{d\hat{\mathbf{x}}_i^{\Omega}} \mathbf{g}_i^{\Omega}(\hat{\mathbf{x}}_i^{\Omega}) = \Phi_i^{\Omega}$ and $\frac{d}{d\hat{\mathbf{x}}_i^{\Gamma}} \mathbf{g}_i^{\Gamma}(\hat{\mathbf{x}}_i^{\Gamma}) = \Phi_i^{\Gamma}$ are constant and do not need to be recomputed at each iteration of the SQP solver described in Section 4 that is used to solve (13) or (23).

The POD bases are computed by minimizing the reconstruction error for a set of snapshots. First we focus on constructing POD bases for the WFPC formulation. Recall that the residual functions r_i are parameterized with parameter space $\mathcal{D} \subset \mathbb{R}^{N_\mu}$. Let $\left\{\mu_\ell^{\text{train}}\right\}_{\ell=1}^{n_\mu} \subset \mathcal{D}$ be a set of training parameters, and solve the DD FOM (9) for each parameter μ_ℓ^{train} to obtain FOM solutions $(\mathbf{x}_i^\Omega(\mu_\ell^{\text{train}}), \mathbf{x}_i^\Gamma(\mu_\ell^{\text{train}}))$, $i=1,\ldots,n_\Omega$. Of course, one can solve the monolithic, single-domain FOM (1) at $\mu=\mu_\ell^{\text{train}}$, and restrict the solution $\mathbf{x}(\mu_\ell^{\text{train}})$ to the subdomain interior and interface states, $\mathbf{x}_i^\Omega(\mu_\ell^{\text{train}}) = P_i^\Omega \mathbf{x}(\mu_\ell^{\text{train}})$, $\mathbf{x}_i^\Gamma(\mu_\ell^{\text{train}}) = P_i^\Gamma \mathbf{x}(\mu_\ell^{\text{train}})$. One then computes bases Φ_i^Ω and Φ_i^Γ using the SVD applied to snapshot matrices for the interior and interface states

$$\boldsymbol{X}_{i}^{\Omega} = \begin{bmatrix} \boldsymbol{x}_{i}^{\Omega}(\boldsymbol{\mu}_{1}^{\text{train}}) & \dots & \boldsymbol{x}_{i}^{\Omega}(\boldsymbol{\mu}_{n_{\mu}}^{\text{train}}) \end{bmatrix} \in \mathbb{R}^{N_{i}^{\Omega} \times n_{\mu}}, \tag{24a}$$

$$\boldsymbol{X}_{i}^{\Gamma} = \begin{bmatrix} \boldsymbol{x}_{i}^{\Gamma}(\boldsymbol{\mu}_{1}^{\text{train}}) & \dots & \boldsymbol{x}_{i}^{\Gamma}(\boldsymbol{\mu}_{n_{\mu}}^{\text{train}}) \end{bmatrix} \in \mathbb{R}^{N_{i}^{\Gamma} \times n_{\mu}}.$$
 (24b)

The process is the same for Φ_i^{Ω} and Φ_i^{Γ} and therefore we drop the superscript Ω or Γ and describe the process to compute a basis Φ_i from a generic snapshot matrix $X_i \in \mathbb{R}^{N_i \times n_{\mu}}$.

One computes the 'thin' SVD $X_i = U_i \Sigma_i V_i^T$ of the snapshot matrix, where $U_i \in \mathbb{R}^{N_i \times m_i}$ is the matrix of left singular vectors, $\Sigma_i \in \mathbb{R}^{m_i \times m_i}$ is the diagonal matrix of singular values $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_{m_i} \geq 0$, $V_i \in \mathbb{R}^{n_\mu \times m_i}$ is the matrix of right singular vectors, and $m_i = \min\{N_i, n_\mu\}$. Given a tolerance $v_i \in (0, 1)$ one computes n_i as the smallest integer such that

$$\sum_{i=1}^{n_i} \sigma_j^2 \ge (1 - \nu_i) \sum_{i=1}^{m_i} \sigma_j^2, \tag{25}$$

and selects

$$\mathbf{\Phi}_i = \mathbf{U}_i(:,1:n_i).$$

The POD basis Φ_i minimize the snapshot reconstruction error $\|\boldsymbol{X}_i - \boldsymbol{\Phi}_i(\boldsymbol{\Phi}_i)^T \boldsymbol{X}_i\|_F^2$ among all possible orthogonal basis matrices of sizes $N_i \times n_i$. See, e.g., [23].

POD basis construction for the SRPC formulation from Section 3.2 is similar. In fact, the bases Φ_i^{Ω} , $i=1,\ldots,n_{\Omega}$, are computed as before, and the bases Φ_i^{Γ} , $i=1,\ldots,n_{\Omega}$, are computed from bases Φ_j^{ρ} on the ports. Because $\mathbf{x}_j^{\rho}(\boldsymbol{\mu}_{\ell}^{\text{train}}) = \mathbf{P}_i^j \mathbf{x}_i^{\Gamma}(\boldsymbol{\mu}_{\ell}^{\text{train}})$ for any $i \in P(j)$ and all ports P(j), the snapshot matrices restricted to port P(j) are

$$X_i^p = P_i^j X_i^\Gamma$$
 for any $i \in P(j)$. (26)

For each port P(j), the POD basis $\Phi_j^p = U_j^p(:, 1:n_j^p)$ is computed from the 'thin' SVD $X_j^p = U_j^p \Sigma_j^p (V_j^p)^T$ as described before. With the port basis matrices Φ_j^p the linear map \mathbf{g}_i^T is constructed following (22),

$$\mathbf{g}_{i}^{\Gamma}(\widehat{\mathbf{x}}_{i}^{\Gamma}) := \sum_{i \in O(i)} (\mathbf{P}_{i}^{j})^{T} \mathbf{\Phi}_{j}^{p} \widehat{\mathbf{p}}_{i}^{j} \widehat{\mathbf{x}}_{i}^{\Gamma} = \mathbf{\Phi}_{i}^{\Gamma} \widehat{\mathbf{x}}_{i}^{\Gamma}, \quad \text{where} \quad \mathbf{\Phi}_{i}^{\Gamma} = \sum_{i \in O(i)} (\mathbf{P}_{i}^{j})^{T} \mathbf{\Phi}_{j}^{p} \widehat{\mathbf{P}}_{i}^{j}.$$

3.4. Nonlinear-manifold ROM

The ROM approach that we focus on in this work is the nonlinear manifold approach, also referred to as NM-ROM. Rather than supposing that the state solutions of the FOM are contained in a low-dimensional linear subspace as in LS-ROM, one supposes that the FOM state solutions are contained in a low-dimensional nonlinear manifold. To build the NM-ROM, one must compute a suitable mapping from a low-dimensional coordinate space, often referred to as the *latent space*, to the manifold of candidate state solutions, or *trial manifold*. Solving the ROM then yields the generalized coordinates in the latent space for a solution in the trial manifold. The approach considered here for computing the nonlinear trial manifold is similar to the approach in [56, Sec. 3], which uses wide, shallow, and sparse autoencoders to compute suitable mappings from the latent space to the trial manifold. Further information regarding the autoencoder architecture we use is given in Section 5.

To compute a DD ROM using the NM-ROM approach, one must compute continuously differentiable nonlinear mappings from a suitably chosen latent space to the trial manifold. Hence the nonlinear functions $\mathbf{g}_i^\Omega: \mathbb{R}^{n_i^\Omega} \to \mathbb{R}^{N_i^\Omega}$ and $\mathbf{g}_i^\Gamma: \mathbb{R}^{n_i^\Gamma} \to \mathbb{R}^{N_i^\Gamma}$ defined in (11) are computed as the *decoders* of autoencoders $\mathbf{a}_i^\Omega: \mathbb{R}^{N_i^\Omega} \to \mathbb{R}^{N_i^\Omega}$ and $\mathbf{a}_i^\Gamma: \mathbb{R}^{N_i^\Gamma} \to \mathbb{R}^{N_i^\Gamma}$. The autoencoders \mathbf{a}_i^Ω and \mathbf{a}_i^Γ consist of two parts each: encoders $\mathbf{a}_i^\Omega: \mathbb{R}^{N_i^\Omega} \to \mathbb{R}^{n_i^\Omega}$ and $\mathbf{a}_i^\Gamma: \mathbb{R}^{N_i^\Gamma} \to \mathbb{R}^{n_i^\Gamma}$, and decoders $\mathbf{g}_i^\Omega: \mathbb{R}^{n_i^\Omega} \to \mathbb{R}^{N_i^\Omega}$ and $\mathbf{g}_i^\Gamma: \mathbb{R}^{n_i^\Gamma} \to \mathbb{R}^{N_i^\Gamma}$. The encoders map inputs from the high-dimensional state space to a low-dimensional latent space, while the decoders map elements

from the low-dimensional space to the high-dimensional state space. The autoencoders a_i^{Ω} and a_i^{Γ} are then defined via the function

$$a_i^{\Omega} = g_i^{\Omega} \circ h_i^{\Omega}, \qquad a_i^{\Gamma} = g_i^{\Gamma} \circ h_i^{\Gamma}.$$

In this work, the encoders and decoders are neural networks such that the autoencoder approximates its inputs:

$$\mathbf{x}_{i}^{\Omega} \approx \mathbf{a}_{i}^{\Omega}(\mathbf{x}_{i}^{\Omega}) = \mathbf{g}_{i}^{\Omega}(\mathbf{h}_{i}^{\Omega}(\mathbf{x}_{i}^{\Omega})), \qquad \mathbf{x}_{i}^{\Gamma} \approx \mathbf{a}_{i}^{\Gamma}(\mathbf{x}_{i}^{\Gamma}) = \mathbf{g}_{i}^{\Gamma}(\mathbf{h}_{i}^{\Gamma}(\mathbf{x}_{i}^{\Gamma})), \qquad i = 1, \dots, n_{\Omega}.$$

Further details on the neural network architecture used can be found in Section 5. The decoders \mathbf{g}_i^{Ω} and \mathbf{g}_i^{Γ} can be interpreted as approximate inverses of the encoders $\boldsymbol{h}_{i}^{\Omega}$ and $\boldsymbol{h}_{i}^{\Gamma}$. In the SRPC case, the autoencoder $\boldsymbol{a}_{i}^{\Gamma}$ is composed of autoencoders $\boldsymbol{a}_{j}^{p}: \mathbb{R}^{N_{j}^{p}} \to \mathbb{R}^{N_{j}^{p}}$ with encoder $\boldsymbol{h}_{j}^{p}: \mathbb{R}^{N_{j}^{p}} \to \mathbb{R}^{N_{j}^{p}}$ and decoder $\boldsymbol{g}_{j}^{p}: \mathbb{R}^{N_{j}^{p}} \to \mathbb{R}^{N_{j}^{p}}$ for each port P(j).

The mean-square-error (MSE) losses for the interior, interface, and port states are defined as

$$\mathcal{L}_{i}^{\Omega} := \frac{1}{n_{n}} \sum_{\ell=1}^{n_{\mu}} \left\| \mathbf{x}_{i}^{\Omega} (\boldsymbol{\mu}_{\ell}^{\text{train}}) - \mathbf{g}_{i}^{\Omega} (\boldsymbol{h}_{i}^{\Omega} (\mathbf{x}_{i}^{\Omega} (\boldsymbol{\mu}_{\ell}^{\text{train}}))) \right\|_{2}^{2}, \qquad i = 1, \dots, n_{\Omega},$$

$$(27a)$$

$$\mathcal{L}_{i}^{\Gamma} := \frac{1}{n_{\mu}} \sum_{\ell=1}^{n_{\mu}} \left\| \mathbf{x}_{i}^{\Gamma}(\boldsymbol{\mu}_{\ell}^{\text{train}}) - \mathbf{g}_{i}^{\Gamma}(\boldsymbol{h}_{i}^{\Gamma}(\mathbf{x}_{i}^{\Gamma}(\boldsymbol{\mu}_{\ell}^{\text{train}}))) \right\|_{2}^{2}, \qquad i = 1, \dots, n_{\Omega},$$

$$(27b)$$

$$\mathcal{L}_{j}^{p} := \frac{1}{n_{\mu}} \sum_{\ell=1}^{n_{\mu}} \left\| \mathbf{x}_{j}^{p}(\boldsymbol{\mu}_{\ell}^{\text{train}}) - \mathbf{g}_{j}^{p}(\boldsymbol{h}_{j}^{p}(\mathbf{x}_{j}^{p}(\boldsymbol{\mu}_{\ell}^{\text{train}}))) \right\|_{2}^{2}, \qquad j = 1, \dots, n_{p},$$
(27c)

where $\mathbf{x}_i^{\Omega}(\boldsymbol{\mu}_{\ell}^{\text{train}})$ and $\mathbf{x}_i^{\Gamma}(\boldsymbol{\mu}_{\ell}^{\text{train}})$ are snapshots of the interior and interface states on subdomain i at parameter $\boldsymbol{\mu}_{\ell}^{\text{train}}$ and $\mathbf{x}_i^{\rho}(\boldsymbol{\mu}_{\ell}^{\text{train}})$ is the state on port P(j), as discussed in Section 3.3. In the WPFC case, the interior state and interface state autoencoders a_i^{Ω} and a_i^{Γ} are trained by minimizing the interior and interface losses \mathcal{L}_i^{Ω} and \mathcal{L}_i^{Γ} , respectively. In the SPRC case, the interior state autoencoders a_i^{Ω} and the port autoencoders a_j^{ρ} are trained by minimizing the interior and interface losses \mathcal{L}_i^{Ω} and \mathcal{L}_j^{ρ} , respectively. The interface state autoencoders a_i^{Γ} are implied by the port autoencoders. Specifically, h_i^{Γ} is

$$\boldsymbol{h}_{i}^{\Gamma}(\boldsymbol{x}_{i}^{\Gamma}) = \sum_{i \in O(i)} (\widehat{\boldsymbol{P}}_{i}^{j})^{T} \boldsymbol{h}_{j}^{p} (\boldsymbol{P}_{i}^{j} \boldsymbol{x}_{i}^{\Gamma}), \tag{28}$$

and \mathbf{g}_{i}^{Γ} is defined using equation (22).

Notice that minimizing the MSE loss is equivalent to minimizing the snapshot reconstruction error, which is exactly how POD bases are constructed, as discussed in Section 3.3. Training the autoencoders can also be interpreted as "learning" the forward and inverse mappings from the latent space of generalized coordinates to the nonlinear trial manifold. After training, the decoders g_i^{Ω} and \mathbf{g}_{i}^{Γ} are used for the DD NM-ROM (13) or (23).

4. Sequential quadratic programming solver

4.1. Lagrange-Gauss-Newton SQP method

The problems (13) and (23) are nonlinear programs (NLPs) with equality constraints, and can be solved using sequential quadratic programming (SQP) [64], [65, Ch. 18]. The SQP solver detailed below amounts to applying a Newton-type method to the Karush-Kuhn-Tucker (KKT) necessary optimality conditions. Note that the SQP solver described in this section can also be applied to the FOM (10) by considering the case $B_i = I$ and g_i^{Ω} , g_i^{Γ} equal to the identity mapping.

To develop a solver that is applicable to either the WFPC (13) or the SPRC formulations (23), we define the constraint functions $\widetilde{\mathbf{A}}_i: \mathbb{R}^{n_i^{\Gamma}} \to \mathbb{R}^{n_A}$ and consider the general formulation

$$\min_{(\widehat{\mathbf{x}}_{i}^{\Omega}, \widehat{\mathbf{x}}_{i}^{\Gamma}), i=1, \dots, n_{\Omega}} \frac{1}{2} \sum_{i=1}^{n_{\Omega}} \left\| \boldsymbol{B}_{i} \boldsymbol{r}_{i} \left(\boldsymbol{g}_{i}^{\Omega} \left(\widehat{\mathbf{x}}_{i}^{\Omega} \right), \boldsymbol{g}_{i}^{\Gamma} \left(\widehat{\mathbf{x}}_{i}^{\Gamma} \right) \right) \right\|_{2}^{2}$$
(29a)

s.t.
$$\sum_{i=1}^{n_{\Omega}} \widetilde{\mathbf{A}}_{i}(\widehat{\mathbf{x}}_{i}^{\Gamma}) = \mathbf{0}.$$
 (29b)

In the WFPC case $\widetilde{A}_i(\widehat{x}_i^{\Omega}) = CA_i g_i^{\Gamma}(\widehat{x}_i^{\Gamma})$ and the constraints are nonlinear in the case of NM-ROMs. In the SRPC case $\widetilde{A}_i(\widehat{x}_i^{\Gamma}) = \widehat{A}_i \widehat{x}_i^{\Gamma}$ and the constraints are always linear.

To apply the SQP solver, one first writes the Lagrangian

$$\widehat{L}(\widehat{\mathbf{x}}_{1}^{\Omega}, \widehat{\mathbf{x}}_{1}^{\Gamma}, \dots, \widehat{\mathbf{x}}_{n_{\Omega}}^{\Omega}, \widehat{\mathbf{x}}_{n_{\Omega}}^{\Gamma}, \widehat{\lambda}) = \frac{1}{2} \sum_{i=1}^{n_{\Omega}} \left\| \mathbf{B}_{i} \mathbf{r}_{i} \left(\mathbf{g}_{i}^{\Omega} \left(\widehat{\mathbf{x}}_{i}^{\Omega} \right), \mathbf{g}_{i}^{\Gamma} \left(\widehat{\mathbf{x}}_{i}^{\Gamma} \right) \right) \right\|_{2}^{2} + \sum_{i=1}^{n_{\Omega}} \widehat{\lambda}^{T} \widetilde{\mathbf{A}}_{i}(\widehat{\mathbf{x}}_{i}^{\Gamma})$$

$$(30)$$

for the ROM NLP (29), where $\hat{\lambda} \in \mathbb{R}^{n_A}$ are the Lagrange multipliers associated with the DD-ROM constraints ((29)b). Let ∇_v and $\frac{\partial}{\partial v}$ denote the partial gradient and partial Jacobian with respect to v, respectively, and let $\frac{d}{dv}$ denote the Jacobian. The first order necessary optimality conditions are

$$\nabla_{\hat{\mathbf{x}}^{\Omega}} \hat{L}(\hat{\mathbf{x}}_{1}^{\Omega}, \hat{\mathbf{x}}_{1}^{\Gamma}, \dots, \hat{\mathbf{x}}_{n_{\Omega}}^{\Omega}, \hat{\mathbf{x}}_{n_{\Omega}}^{\Gamma}, \hat{\boldsymbol{\lambda}}) = \rho_{i}^{\Omega}(\hat{\mathbf{x}}_{i}^{\Omega}, \hat{\mathbf{x}}_{i}^{\Gamma}) = \mathbf{0}, \qquad i = 1, \dots, n_{\Omega},$$
(31a)

$$\nabla_{\hat{\mathbf{x}}^{\Gamma}} \hat{L}(\hat{\mathbf{x}}_{1}^{\Omega}, \hat{\mathbf{x}}_{1}^{\Gamma}, \dots, \hat{\mathbf{x}}_{n_{\Omega}}^{\Omega}, \hat{\mathbf{x}}_{n_{\Omega}}^{\Gamma}, \hat{\lambda}) = \rho_{i}^{\Gamma}(\hat{\mathbf{x}}_{i}^{\Omega}, \hat{\mathbf{x}}_{i}^{\Gamma}, \hat{\lambda}) = \mathbf{0}, \qquad i = 1, \dots, n_{\Omega},$$
(31b)

$$\nabla_{\widehat{\lambda}}\widehat{L}(\widehat{\mathbf{x}}_{1}^{\Omega},\widehat{\mathbf{x}}_{1}^{\Gamma},\ldots,\widehat{\mathbf{x}}_{n_{\Omega}}^{\Omega},\widehat{\mathbf{x}}_{n_{\Omega}}^{\Gamma},\widehat{\lambda}) = \sum_{i=1}^{n_{\Omega}} \widetilde{\mathbf{A}}_{i}(\widehat{\mathbf{x}}_{i}^{\Gamma}) = \mathbf{0}, \tag{31c}$$

where

$$\rho_i^{\Omega}(\hat{\mathbf{x}}_i^{\Omega}, \hat{\mathbf{x}}_i^{\Gamma}) = \frac{d\mathbf{g}_i^{\Omega}}{d\hat{\mathbf{x}}_i^{\Omega}}(\hat{\mathbf{x}}_i^{\Omega})^T \frac{\partial \mathbf{r}_i}{\partial \mathbf{x}_i^{\Omega}} \left(\mathbf{g}_i^{\Omega}(\hat{\mathbf{x}}_i^{\Omega}), \mathbf{g}_i^{\Gamma}(\hat{\mathbf{x}}_i^{\Gamma}) \right)^T \mathbf{B}_i^T \mathbf{B}_i \mathbf{r}_i \left(\mathbf{g}_i^{\Omega}(\hat{\mathbf{x}}_i^{\Omega}), \mathbf{g}_i^{\Gamma}(\hat{\mathbf{x}}_i^{\Gamma}) \right), \tag{32a}$$

$$\rho_{i}^{\Gamma}(\widehat{\mathbf{x}}_{i}^{\Omega},\widehat{\mathbf{x}}_{i}^{\Gamma},\widehat{\lambda}) = \frac{d\mathbf{g}_{i}^{\Gamma}}{d\widehat{\mathbf{x}}_{i}^{\Gamma}}(\widehat{\mathbf{x}}_{i}^{\Gamma})^{T} \frac{\partial \mathbf{r}_{i}}{\partial \mathbf{x}_{i}^{\Gamma}} \left(\mathbf{g}_{i}^{\Omega}(\widehat{\mathbf{x}}_{i}^{\Omega}), \mathbf{g}_{i}^{\Gamma}(\widehat{\mathbf{x}}_{i}^{\Gamma})\right)^{T} \mathbf{B}_{i}^{T} \mathbf{B}_{i} \mathbf{r}_{i} \left(\mathbf{g}_{i}^{\Omega}(\widehat{\mathbf{x}}_{i}^{\Omega}), \mathbf{g}_{i}^{\Gamma}(\widehat{\mathbf{x}}_{i}^{\Gamma})\right) + \frac{d\widetilde{\mathbf{A}}_{i}}{d\widehat{\mathbf{x}}_{i}^{\Gamma}}(\widehat{\mathbf{x}}_{i}^{\Gamma})^{T} \widehat{\lambda}$$
(32b)

are the gradients of the Lagrangian with respect to the subdomain variables $\hat{\mathbf{x}}_i^{\Omega}$ and $\hat{\mathbf{x}}_i^{\Gamma}$, respectively.

A Newton-type method applied to (31) yields the SQP iterations

$$\begin{bmatrix} \boldsymbol{H}_{1}(\widehat{\boldsymbol{x}}_{1}^{\Omega(k)}, \widehat{\boldsymbol{x}}_{1}^{\Gamma(k)}) & \dots & \boldsymbol{E}_{1}(\widehat{\boldsymbol{x}}_{1}^{\Gamma})^{T} \\ \vdots & \vdots & \vdots \\ \boldsymbol{H}_{n_{\Omega}}(\widehat{\boldsymbol{x}}_{n_{\Omega}}^{\Omega(k)}, \widehat{\boldsymbol{x}}_{n_{\Omega}}^{\Gamma(k)}) & \boldsymbol{E}_{n_{\Omega}}(\widehat{\boldsymbol{x}}_{1}^{\Gamma})^{T} \end{bmatrix} \begin{bmatrix} \boldsymbol{s}_{1}^{(k)} \\ \vdots \\ \boldsymbol{s}_{n_{\Omega}}^{(k)} \\ \boldsymbol{s}^{\hat{\lambda}(k)} \end{bmatrix} = - \begin{bmatrix} \boldsymbol{\rho}_{1}(\widehat{\boldsymbol{x}}_{1}^{\Omega(k)}, \widehat{\boldsymbol{x}}_{1}^{\Gamma(k)}, \widehat{\boldsymbol{\lambda}}^{(k)}) \\ \vdots \\ \boldsymbol{\rho}_{n_{\Omega}}(\widehat{\boldsymbol{x}}_{n_{\Omega}}^{\Omega(k)}, \widehat{\boldsymbol{x}}_{1}^{\Gamma(k)}, \widehat{\boldsymbol{\lambda}}^{(k)}) \\ \sum_{i=1}^{n_{\Omega}} \widehat{\boldsymbol{A}}_{i}(\widehat{\boldsymbol{x}}_{1}^{\Gamma(k)}) \end{bmatrix},$$

$$(33)$$

where k is the SQP iteration index, $\boldsymbol{H}_i(\widehat{\boldsymbol{x}}_i^{\varOmega(k)}, \widehat{\boldsymbol{x}}_i^{\varGamma(k)})$ is the Hessian of the Lagrangian with respect to the subdomain variables $(\widehat{\boldsymbol{x}}_i^\varOmega, \widehat{\boldsymbol{x}}_i^\varGamma)$ evaluated at $(\widehat{\boldsymbol{x}}_i^{\varOmega(k)}, \widehat{\boldsymbol{x}}_i^{\varGamma(k)})$ or an approximation of this Hessian, and where

$$E_{i}(\widehat{\mathbf{x}}_{i}^{\Gamma}) = \begin{bmatrix} \mathbf{0} & \frac{d\widetilde{\mathbf{A}}_{i}}{d\widehat{\mathbf{x}}_{i}^{\Gamma}}(\widehat{\mathbf{x}}_{i}^{\Gamma}) \end{bmatrix}, \mathbf{s}_{i}^{(k)} = \begin{bmatrix} \mathbf{s}_{i}^{\Omega(k)} \\ \mathbf{s}_{i}^{\Gamma(k)} \end{bmatrix}, \rho_{i}(\widehat{\mathbf{x}}_{i}^{\Omega}, \widehat{\mathbf{x}}_{i}^{\Gamma}, \widehat{\boldsymbol{\lambda}}) = \begin{bmatrix} \rho_{i}^{\Omega}(\widehat{\mathbf{x}}_{i}^{\Omega}, \widehat{\mathbf{x}}_{i}^{\Gamma}) \\ \rho_{i}^{\Gamma}(\widehat{\mathbf{x}}_{i}^{\Omega}, \widehat{\mathbf{x}}_{i}^{\Gamma}, \widehat{\boldsymbol{\lambda}}) \end{bmatrix}, \tag{34}$$

for $i = 1, ..., n_O$. The next result on the unique solvability of (33) follows from adapting standard results to the block structure of

Lemma 1. If for $i=1,\ldots,n_{\Omega}$ the matrices $H_i(\widehat{\mathbf{x}}_i^{\Omega(k)},\widehat{\mathbf{x}}_i^{\Gamma(k)})$ are positive definite on the null-space of $E_i(\widehat{\mathbf{x}}_i^{\Gamma})$, and if $(E_1(\widehat{\mathbf{x}}_1^{\Gamma}),\ldots,E_{n_{\Omega}}(\widehat{\mathbf{x}}_{n_{\Omega}}^{\Gamma}))$ has full row rank, then (33) has a unique solution.

For a proof see, e.g., [66, Thm. 3.2], [65, Lemma 16.12].

We use a Gauss-Newton approximation of the Hessian, which is motivated by the following consideration. If the residuals $B_i r_i (g_i^{\Omega}(\hat{x}_i^{\Omega}), g_i^{\Gamma}(\hat{x}_i^{\Omega}))$ are small at the solution of (29), then the first order optimality condition ((31)b) implies that $\hat{\lambda}$ is small. Thus all second derivative terms in the true Hessians $H_i(\hat{x}_i^{\Omega(k)}, \hat{x}_i^{\Gamma(k)})$ are multiplied by small residuals or small Lagrange multipliers. The Gauss-Newton Hessian approximation neglects these terms and approximates the Hessians by

$$H_{i}(\hat{\mathbf{x}}_{i}^{\Omega}, \hat{\mathbf{x}}_{i}^{\Gamma}) = R_{i}(\hat{\mathbf{x}}_{i}^{\Omega}, \hat{\mathbf{x}}_{i}^{\Gamma})^{T} B_{i}^{T} B_{i} R_{i}(\hat{\mathbf{x}}_{i}^{\Omega}, \hat{\mathbf{x}}_{i}^{\Gamma}), \tag{35a}$$

where

$$\boldsymbol{R}_{i}(\widehat{\boldsymbol{x}}_{i}^{\Omega},\widehat{\boldsymbol{x}}_{i}^{\Gamma}) = \begin{bmatrix} \frac{\partial \boldsymbol{r}_{i}}{\partial \boldsymbol{x}_{i}^{\Omega}}(\boldsymbol{g}_{i}^{\Omega}(\widehat{\boldsymbol{x}}_{i}^{\Omega}), \boldsymbol{g}_{i}^{\Gamma}(\widehat{\boldsymbol{x}}_{i}^{\Gamma})) \frac{d\boldsymbol{g}_{i}^{\Omega}}{d\widehat{\boldsymbol{x}}^{\Omega}}(\widehat{\boldsymbol{x}}_{i}^{\Omega}), & \frac{\partial \boldsymbol{r}_{i}}{\partial \boldsymbol{x}_{i}^{\Gamma}}(\boldsymbol{g}_{i}^{\Omega}(\widehat{\boldsymbol{x}}_{i}^{\Omega}), \boldsymbol{g}_{i}^{\Gamma}(\widehat{\boldsymbol{x}}_{i}^{\Gamma})) \frac{d\boldsymbol{g}_{i}^{\Gamma}}{d\widehat{\boldsymbol{x}}^{\Gamma}}(\widehat{\boldsymbol{x}}_{i}^{\Gamma}) \end{bmatrix}$$
(35b)

is the Jacobian of r_i with respect to $(\hat{x}_i^{\Omega}, \hat{x}_i^{\Gamma})$. The advantage is that (35) only requires first order derivatives. Note that the FOM solution satisfies (9), i.e., the residual in the least squares formulation (10) is zero. Thus, if the ROM well approximates the FOM (9) or, equivalently, its least squares formulation (10), then we expect the residuals $B_i r_i (g_i^{\Omega}(\widehat{\mathbf{x}}_i^{\Omega}), g_i^{\Gamma}(\widehat{\mathbf{x}}_i^{\Gamma}))$ to be small at the solution of (29) and the Gauss–Newton Hessian (35) to be good approximation of the true Hessian of the Lagrangian (30). Note that with the notation (35), the gradients $\rho_i(\hat{x}_i^{\Omega}, \hat{x}_i^{\Gamma}, \hat{\lambda})$ in (34) can be written as

$$\rho_{i}(\widehat{\mathbf{x}}_{i}^{\Omega}, \widehat{\mathbf{x}}_{i}^{\Gamma}, \widehat{\boldsymbol{\lambda}}) = \mathbf{R}_{i}(\widehat{\mathbf{x}}_{i}^{\Omega}, \widehat{\mathbf{x}}_{i}^{\Gamma})^{T} \mathbf{B}_{i}^{T} \mathbf{B}_{i} \mathbf{r}_{i}(\mathbf{g}_{i}^{\Omega}(\widehat{\mathbf{x}}_{i}^{\Omega}), \mathbf{g}_{i}^{\Gamma}(\widehat{\mathbf{x}}_{i}^{\Gamma})) + E_{i}(\widehat{\mathbf{x}}_{i}^{\Gamma})^{T} \widehat{\boldsymbol{\lambda}}, \quad i = 1, \dots, n_{\Omega}.$$

$$(36)$$

With the Gauss-Newton approximations (35), the SQP system (33) is essentially the optimality system for the quadratic program

$$\min_{\mathbf{s}_{i}=(\mathbf{s}_{i}^{\Omega},\mathbf{s}_{i}^{\Gamma}),i=1,\ldots,n_{\Omega}} \frac{1}{2} \sum_{i=1}^{n_{\Omega}} \left\| \mathbf{B}_{i} \mathbf{r}_{i} \left(\mathbf{g}_{i}^{\Omega} (\widehat{\mathbf{x}}_{i}^{\Omega(k)}), \mathbf{g}_{i}^{\Gamma} (\widehat{\mathbf{x}}_{i}^{\Gamma(k)}) \right) + \mathbf{B}_{i} \mathbf{R}_{i} (\widehat{\mathbf{x}}_{i}^{\Omega(k)}, \widehat{\mathbf{x}}_{i}^{\Gamma(k)}) \mathbf{s}_{i} \right\|_{2}^{2}$$
(37a)

s.t.
$$\sum_{i=1}^{n_{\Omega}} \widetilde{A}_{i}(\widehat{\mathbf{x}}_{i}^{\Gamma(k)}) + \frac{d}{d\widehat{\mathbf{x}}_{i}^{\Gamma}} \widetilde{A}_{i}(\widehat{\mathbf{x}}_{i}^{\Gamma(k)}) \mathbf{s}_{i} = \mathbf{0}.$$
 (37b)

More precisely, the following result holds.

Lemma 2. If the assumptions of Lemma 1 hold, then the quadratic program (37) has a unique solution $s_i^{(k)} = (s_i^{\Omega(k)}, s_i^{\Gamma(k)})$, $i = 1, ..., n_{\Omega}$, given by the solution of (33). The associated Lagrange multiplier for (37) is $\hat{\lambda}^{(k)} + s^{\hat{\lambda}(k)}$, where $\hat{\lambda}^{(k)}$ is the Lagrange multiplier estimate in (33) and $s^{\hat{\lambda}(k)}$ is the last component in the solution vector of (33).

The proof of Lemma 2 follows from the necessary and sufficient optimality conditions (e.g., [65, Sec 16.1]) for the quadratic program (37). The necessary and sufficient optimality conditions for (37) are given by (33) with the terms $E_i(\hat{x}_i^{\Gamma(k)})^T \hat{\lambda}^{(k)}$ (see (36)) moved from the right to the left hand side.

An advantage of the Gauss–Newton approximation is that no Lagrange multiplier estimate is needed in (37) or the associated optimality system. Of course, quantities like $\mathbf{g}_i^{\Omega}(\hat{\mathbf{x}}_i^{\Omega})$, $\mathbf{g}_i^{\Gamma}(\hat{\mathbf{x}}_i^{\Gamma})$, $\mathbf{B}_i\mathbf{r}_i(\mathbf{g}_i^{\Omega}(\hat{\mathbf{x}}_i^{\Omega}))$, $\mathbf{g}_i^{\Gamma}(\hat{\mathbf{x}}_i^{\Gamma})$, $\mathbf{B}_i\mathbf{R}_i(\hat{\mathbf{x}}_i^{\Omega},\hat{\mathbf{x}}_i^{\Gamma})$, $\widetilde{\mathbf{A}}_i(\hat{\mathbf{x}}_i^{\Gamma})$, and $\frac{d\widetilde{\mathbf{A}}_i}{d\hat{\mathbf{x}}_i^{\Gamma}}(\hat{\mathbf{x}}_i^{\Gamma})$ can be computed in parallel across the subdomains. Moreover, the block structure of the system (33) lends itself to a parallel solution strategy. However, since (33) corresponds to the ROM its size tends to be small and parallelism in its solution may yield less speedup than it would if applied to the DD formulation (9) of the FOM (1). The parallel implementation of the approach discussed in this paper is left to future work.

Given the solution of the SQP system (33), the new iterate, i.e., the new approximate solution of (29) is computed as

$$\widehat{\mathbf{x}}_{i}^{\Omega(k+1)} = \widehat{\mathbf{x}}_{i}^{\Omega(k)} + \alpha^{(k)} \mathbf{s}_{i}^{\Omega(k)}, \qquad i = 1, \dots, n_{\Omega}, \tag{38a}$$

$$\widehat{\mathbf{x}}_{i}^{\Gamma(k+1)} = \widehat{\mathbf{x}}_{i}^{\Gamma(k)} + \alpha^{(k)} \mathbf{s}_{i}^{\Gamma(k)}, \qquad i = 1, \dots, n_{\Omega}, \tag{38b}$$

with step size $\alpha^{(k)} \in (0,1]$. If one chooses to keep a Lagrange multiplier estimate, then $\hat{\lambda}^{(k+1)} = \hat{\lambda}^{(k)} + \alpha^{(k)} s^{\hat{\lambda}(k)}$, $i = 1, \dots, n_{\Omega}$, where $s^{\hat{\lambda}(k)}$ is the last component in the solution vector of (33). The step size $\alpha^{(k)}$ is computed via line-search using a merit function that coordinates progress of the iterates (and Lagrange multipliers) towards feasibility and optimality. In this work, we simply use the norm of the gradients (31). This is an appropriate criterion if one starts sufficiently close to a (local) minimizer of (29), and this criterion yielded good results in our examples. In our examples, the step size is computed using a backtracking line search with the Armijo rule.

4.2. Convergence of SQP solver

Convergence of the Lagrange–Gauss–Newton SQP method can be established using one of two related approaches. The iteration (38) with $s_i^{\Omega(k)}, s_i^{\Gamma(k)}, i=1,\ldots,n_\Omega$, computed as the solution of (33) with Gauss–Newton Hessian approximation (35) can be interpreted and analyzed as a generalized Gauss–Newton iteration. See [67]. Alternatively, this iteration can also be viewed as an inexact Newton method applied to the first-order optimality conditions (31). The local convergence result using either approach requires that the Lagrange–Gauss–Newton SQP method is started sufficiently close to a (local) minimizer of (29). We summarize the convergence theory for inexact Newton methods (see, e.g. [65, Thm. 11.3]) in Theorem 4. First we define the following notation to improve readability. We group the vectors $(\hat{\mathbf{x}}_1^\Omega, \hat{\mathbf{x}}_1^\Gamma, \dots, \hat{\mathbf{x}}_{n_\Omega}^\Omega, \hat{\mathbf{x}}_{n_\Omega}^\Gamma)$ and $(\mathbf{s}_1^{\Omega(k)}, \mathbf{s}_1^{\Gamma(k)}, \dots, \mathbf{s}_{n_\Omega}^{\Omega(k)}, \mathbf{s}_{n_\Omega}^{\Gamma(k)})$ as

$$\widehat{\mathbf{x}} = \begin{bmatrix} \widehat{\mathbf{x}}_{1}^{\Omega} \\ \widehat{\mathbf{x}}_{1}^{\Gamma} \\ \vdots \\ \widehat{\mathbf{x}}_{n_{\Omega}}^{\Omega} \\ \widehat{\mathbf{x}}_{n_{\Omega}}^{\Gamma} \end{bmatrix} \in \mathbb{R}^{n_{D}}, \mathbf{s}_{x}^{(k)} = \begin{bmatrix} \mathbf{s}_{1}^{\Omega(k)} \\ \mathbf{s}_{1}^{\Gamma(k)} \\ \vdots \\ \mathbf{s}_{n_{\Omega}}^{\Omega(k)} \\ \mathbf{s}_{n_{\Omega}}^{\Gamma(k)} \end{bmatrix} \in \mathbb{R}^{n_{D}}.$$

$$(39a)$$

where $n_D = \sum_{i=1}^{n_\Omega} (n_i^\Omega + n_i^\Gamma)$. Furthermore let $\hat{F}: \mathbb{R}^{n_D + n_A} \to \mathbb{R}^{n_D + n_A}$,

$$\widehat{F}(\widehat{\mathbf{x}},\widehat{\lambda}) = - \begin{bmatrix} \rho_1(\widehat{\mathbf{x}}_1^{\Omega(k)}, \widehat{\mathbf{x}}_1^{\Gamma(k)}, \widehat{\lambda}^{(k)}) \\ \vdots \\ \rho_{n_{\Omega}}(\widehat{\mathbf{x}}_{n_{\Omega}}^{\Omega(k)}, \widehat{\mathbf{x}}_{n_{\Omega}}^{\Gamma(k)}, \widehat{\lambda}^{(k)}) \\ \sum_{l=1}^{n_{\Omega}} \widetilde{\mathbf{A}}_l(\widehat{\mathbf{x}}_l^{\Gamma(k)}) \end{bmatrix}. \tag{39b}$$

denote the right hand side of the KKT system. Recall that if $\overline{\mathbf{x}} \in \mathbb{R}^{n_D}$ is a local minimizer of (29) with associated Lagrange multiplier $\overline{\lambda} \in \mathbb{R}^{n_A}$, then $\widehat{F}(\overline{\mathbf{x}}, \overline{\lambda}) = \mathbf{0}$. Next define the Hessian approximation $\mathbf{H} : \mathbb{R}^{n_D} \times \mathbb{R}^{n_A} \to \mathbb{R}^{n_D \times n_D}$,

$$\boldsymbol{H}(\widehat{\boldsymbol{x}},\widehat{\boldsymbol{\lambda}}) = \begin{bmatrix} \boldsymbol{H}_{1}(\widehat{\boldsymbol{x}}_{1}^{\Omega}, \widehat{\boldsymbol{x}}_{1}^{\Gamma}, \widehat{\boldsymbol{\lambda}}) & & & \\ & \ddots & & & \\ & & \boldsymbol{H}_{n_{\Omega}}(\widehat{\boldsymbol{x}}_{n_{\Omega}}^{\Omega}, \widehat{\boldsymbol{x}}_{n_{\Omega}}^{\Gamma}, \widehat{\boldsymbol{\lambda}}) \end{bmatrix}, \tag{39c}$$

and constraint Jacobian $\frac{d}{d\hat{x}^{\Gamma}}\tilde{A}: \mathbb{R}^{\sum_{i=1}^{n_{\Omega}}n_{i}^{\Gamma}} \to \mathbb{R}^{n_{A} \times n_{D}}$,

$$\frac{d}{d\widehat{\mathbf{x}}^{\Gamma}}\widetilde{\mathbf{A}}(\widehat{\mathbf{x}}_{1}^{\Gamma},\dots,\widehat{\mathbf{x}}_{n_{\Omega}}^{\Gamma}) = \begin{bmatrix} \mathbf{0} & \frac{d}{d\widehat{\mathbf{x}}_{1}^{\Gamma}}\widetilde{\mathbf{A}}_{1}(\widehat{\mathbf{x}}_{1}^{\Gamma}) & \dots & \mathbf{0} & \frac{d}{d\widehat{\mathbf{x}}_{n_{\Omega}}^{\Gamma}}\widetilde{\mathbf{A}}_{n_{\Omega}}(\widehat{\mathbf{x}}_{n_{\Omega}}^{\Gamma}) \end{bmatrix}. \tag{39d}$$

The convergence result can now be stated as follows.

Theorem 4. Let r_i , g_i^{Ω} , and g_i^{Γ} be continuously differentiable for all $i=1,\ldots,n_{\Omega}$. Let \overline{x} be a local minimizer of (29) such that the Jacobian $\frac{d}{d\overline{x}^{\Gamma}}\widetilde{\Lambda}(\overline{x}_1^{\Gamma},\ldots,\overline{x}_{n_{\Omega}}^{\Gamma})$ has full row rank, let $\overline{\lambda}$ denote the associated Lagrange multiplier and assume that $H(\overline{x},\overline{\lambda})$ is positive definite

on the null-space of $\frac{d}{d\hat{x}^\Gamma} \tilde{A}(\overline{x}_1^\Gamma, \dots, \overline{x}_{n_Q}^\Gamma)$. Furthermore, assume that $\hat{F}'(\hat{x}, \hat{\lambda})$ is Lipschitz continuous with Lipschitz constant K, and that the steps $s_x^{(k)}$ satisfy

$$\left\| \left(\nabla_{\widehat{\mathbf{x}}}^2 \widehat{L}(\widehat{\mathbf{x}}^{(k)}, \widehat{\boldsymbol{\lambda}}^{(k)}) - \boldsymbol{H}(\widehat{\mathbf{x}}^{(k)}, \widehat{\boldsymbol{\lambda}}^{(k)}) \right) \boldsymbol{s}_x^{(k)} \right\|_2 \leq \eta_k \left\| \widehat{\boldsymbol{F}}(\widehat{\mathbf{x}}^{(k)}, \widehat{\boldsymbol{\lambda}}^{(k)}) \right\|_2$$

for some sequence of forcing parameters η_k , and where \hat{L} is the Lagrangian defined in (30).

If $\{\eta_k\}$ satisfies $0 < \eta_k \le \eta$ where η is such that $4\eta\bar{\kappa} < 1$ with $\bar{\kappa} = \|\hat{F}'(\overline{x}, \overline{\lambda})^{-1}\|_2 \|\hat{F}'(\overline{x}, \overline{\lambda})\|_2$, then for all $\sigma \in (4\eta\bar{\kappa}, 1)$ there exists an $\epsilon > 0$ such that for any $(\hat{x}^{(0)}, \hat{\lambda}^{(0)})$ with $\|(\overline{x}, \overline{\lambda}) - (\hat{x}^{(0)}, \hat{\lambda}^{(0)})\|_2 < \epsilon$, the sequence of iterates $\{(\hat{x}^{(k)}, \hat{\lambda}^{(k)})\}$ generated by the SQP solver converges to $(\overline{x}, \overline{\lambda})$, and the iterates satisfy

$$\begin{split} \left\| (\widehat{\boldsymbol{x}}^{(k)}, \widehat{\boldsymbol{\lambda}}^{(k)}) - (\overline{\boldsymbol{x}}, \overline{\boldsymbol{\lambda}}) \right\|_2 &\leq K \left\| \widehat{\boldsymbol{F}}'(\overline{\boldsymbol{x}}, \overline{\boldsymbol{\lambda}})^{-1} \right\|_2 \left\| (\widehat{\boldsymbol{x}}^{(k)}, \widehat{\boldsymbol{\lambda}}^{(k)}) - (\overline{\boldsymbol{x}}, \overline{\boldsymbol{\lambda}}) \right\|_2^2 + 4 \eta_k \bar{\kappa} \left\| (\widehat{\boldsymbol{x}}^{(k)}, \widehat{\boldsymbol{\lambda}}^{(k)}) - (\overline{\boldsymbol{x}}, \overline{\boldsymbol{\lambda}}) \right\|_2 \\ &\leq \sigma \left\| (\widehat{\boldsymbol{x}}^{(k)}, \widehat{\boldsymbol{\lambda}}^{(k)}) - (\overline{\boldsymbol{x}}, \overline{\boldsymbol{\lambda}}) \right\|_2. \end{split}$$

See, e.g., [65, Thm. 11.3] for a proof of this theorem.

Remark 1. As stated, Theorem 4 only guarantees a solution to the KKT system (31), which are (first order) necessary optimality conditions. However one can give alternative inexactness conditions on Gauss-Newton Hessian approximations that ensure local convergence to a point at which the second order sufficient optimality conditions are satisfied. See [68, L. 2.5] for the unconstrained case and [67, Sec. 3.5] for the constrained case, but with ℓ_1 rather than ℓ_2 (=least squares) objective.

5. Autoencoder architecture

Following [56], we consider the use of shallow, wide, sparse-masked autoencoders with smooth activation functions for representing the maps, \mathbf{g}_i^{Ω} and \mathbf{g}_i^{Γ} . 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. Sparsity is applied at the decoder output layer so that hyper-reduction can be applied. Further details on hyper reduction are addressed in Section 5.3. Smooth activations are used to ensure that g_i^{Ω} and g_i^{Γ} are continuously differentiable. In contrast with [56], we also apply a sparsity mask to the encoder input layer so that the encoders and decoders are symmetric across the latent layer. We found that applying a sparsity mask to the encoder input layer permitted the use of a wider network for the encoder, resulting in improved performance over a dense input layer. See Section 7.4 for further details.

5.1. Weak FOM-port formulation

First we detail the architectures used for the weak FOM-port constraint formulation. We use a single-layer architecture for the encoders and decoders with a smooth, non-polynomial activation function. The encoders, \mathbf{h}_i^Ω and \mathbf{h}_i^Γ , and decoders, \mathbf{g}_i^Ω and \mathbf{g}_i^Γ , are of the form

$$\boldsymbol{h}_{i}^{\Omega}: \mathbb{R}^{N_{i}^{\Omega}} \to \mathbb{R}^{n_{i}^{\Omega}}, \qquad \boldsymbol{h}_{i}^{\Omega}(\boldsymbol{x}_{i}^{\Omega}) = \boldsymbol{W}_{i}^{h_{i}^{\Omega}} \boldsymbol{\sigma}_{i}^{\Omega}(\boldsymbol{W}_{i}^{h_{i}^{\Omega}} \boldsymbol{x}_{i}^{\Omega} + \boldsymbol{b}_{i}^{h_{i}^{\Omega}}), \tag{40a}$$

$$\mathbf{g}_{i}^{\Omega}: \mathbb{R}^{n_{i}^{\Omega}} \to \mathbb{R}^{N_{i}^{\Omega}}, \qquad \qquad \mathbf{g}_{i}^{\Omega}(\widehat{\mathbf{x}}_{i}^{\Omega}) = \mathbf{W}_{i}^{\mathbf{g}_{i}^{\Omega}} \widehat{\mathbf{\sigma}}_{i}^{\Omega}(\mathbf{W}_{i}^{\mathbf{g}_{i}^{\Omega}} \widehat{\mathbf{x}}_{i}^{\Omega} + \mathbf{b}_{i}^{\mathbf{g}_{i}^{\Omega}}), \tag{40b}$$

$$\boldsymbol{h}_{i}^{\Gamma}: \mathbb{R}^{N_{i}^{\Gamma}} \to \mathbb{R}^{n_{i}^{\Gamma}}, \qquad \boldsymbol{h}_{i}^{\Gamma}(\boldsymbol{x}_{i}^{\Gamma}) = \boldsymbol{W}_{i}^{h_{i}^{\Gamma}} \boldsymbol{\sigma}_{i}^{\Gamma}(\boldsymbol{W}_{i}^{h_{i}^{\Gamma}} \boldsymbol{x}_{i}^{\Gamma} + \boldsymbol{b}_{i}^{h_{i}^{\Gamma}}), \tag{40c}$$

$$\mathbf{g}_{i}^{\Gamma}: \mathbb{R}^{n_{i}^{\Gamma}} \to \mathbb{R}^{N_{i}^{\Gamma}}, \qquad \mathbf{g}_{i}^{\Gamma}(\hat{\mathbf{x}}_{i}^{\Gamma}) = \mathbf{W}_{2}^{\mathbf{g}_{i}^{\Gamma}} \sigma_{i}^{\Gamma}(\mathbf{W}_{i}^{\mathbf{g}_{i}^{\Gamma}} \hat{\mathbf{x}}_{i}^{\Gamma} + \mathbf{b}_{i}^{\mathbf{g}_{i}^{\Gamma}}), \tag{40d}$$

where

$$\boldsymbol{W}_{1}^{\boldsymbol{g}_{i}^{\Omega}}, \left(\boldsymbol{W}_{2}^{\boldsymbol{h}_{i}^{\Omega}}\right)^{T} \in \mathbb{R}^{w_{i}^{\Omega} \times \boldsymbol{n}_{i}^{\Omega}}, \qquad \qquad \boldsymbol{W}_{2}^{\boldsymbol{g}_{i}^{\Omega}}, \left(\boldsymbol{W}_{1}^{\boldsymbol{h}_{i}^{\Omega}}\right)^{T} \in \mathbb{R}^{N_{i}^{\Omega} \times w_{i}^{\Omega}}, \tag{41a}$$

$$\boldsymbol{b}_{1}^{\boldsymbol{g}_{i}^{\mathcal{U}}} \in \mathbb{R}^{\boldsymbol{u}_{i}^{\mathcal{Q}}}, \qquad \boldsymbol{b}_{1}^{\boldsymbol{g}_{i}^{\mathcal{U}}} \in \mathbb{R}^{\boldsymbol{u}_{i}^{\mathcal{Q}}}, \qquad (41b)$$

$$\boldsymbol{b}_{1}^{\boldsymbol{b}_{i}^{\Omega}} \in \mathbb{R}^{w_{i}^{\Omega}}, \qquad \qquad \boldsymbol{b}_{1}^{\boldsymbol{g}_{i}^{\Omega}} \in \mathbb{R}^{w_{i}^{\Omega}}, \qquad (41b)$$

$$\boldsymbol{W}_{1}^{\boldsymbol{g}_{i}^{\Gamma}}, \left(\boldsymbol{W}_{2}^{\boldsymbol{h}_{i}^{\Gamma}}\right)^{T} \in \mathbb{R}^{w_{i}^{\Gamma} \times \boldsymbol{n}_{i}^{\Gamma}}, \qquad \boldsymbol{W}_{2}^{\boldsymbol{g}_{i}^{\Gamma}}, \left(\boldsymbol{W}_{1}^{\boldsymbol{h}_{i}^{\Gamma}}\right)^{T} \in \mathbb{R}^{N_{i}^{\Gamma} \times w_{i}^{\Gamma}}, \qquad (41c)$$

$$\boldsymbol{b}_{1}^{l_{i}^{\Gamma}} \in \mathbb{R}^{w_{i}^{\Gamma}}, \qquad \boldsymbol{b}_{1}^{\boldsymbol{g}_{i}^{\Gamma}} \in \mathbb{R}^{w_{i}^{\Gamma}},$$
 (41d)

 $\sigma_i^{\Omega}, \sigma_i^{\Gamma}$ are smooth, non-polynomial activation functions (e.g. Sigmoid or Swish), and where $w_i^{\Omega}, w_i^{\Gamma}$ are the network widths for all subdomains $i = 1, \dots, n_{\Omega}$. The weight matrices $\boldsymbol{W}_2^{\mathbf{g}_i^{\Omega}}, \boldsymbol{W}_1^{h_i^{\Omega}}, \boldsymbol{W}_2^{\mathbf{g}_i^{\Gamma}}$ and $\boldsymbol{W}_1^{h_i^{\Gamma}}$ are all sparse, while the remaining weights and biases are

The widths, w_i^{Ω} and w_i^{Γ} , as well as the sparsity patterns of $\boldsymbol{W}_2^{\boldsymbol{g}_i^{\Omega}}$, $\boldsymbol{W}_1^{\boldsymbol{h}_i^{\Omega}}$, $\boldsymbol{W}_2^{\boldsymbol{g}_i^{\Gamma}}$ and $\boldsymbol{W}_1^{\boldsymbol{h}_i^{\Gamma}}$ are hyper-parameters that require tuning. The use of a single-layer architecture of arbitrary width and non-polynomial activation, as defined in ((40)a-d), is motivated by the well-known universal approximation theorem [69,70]. Furthermore, the use of a smooth activation function ensures that $\boldsymbol{g}_i^{\Omega}$ and g_i^T are continuously differentiable. This is important because the Jacobians $\frac{dg_i^\Omega}{d\hat{x}_i^\Omega}$ and $\frac{dg_i^\Gamma}{d\hat{x}_i^\Omega}$ are required by the SQP solver discussed in Section 4. The autoencoders are trained by minimizing the MSE loss defined in equation (27).

5.2. Strong ROM-port formulation

Next we detail the architectures used for the strong ROM-port constraint formulation. As before, we use a single-layer architecture for the encoders and decoders with a smooth, non-polynomial activation function. The interior state encoders h_i^{Ω} and decoders g_i^{Ω} have the same architecture as in the weak FOM-port constraint formulation. Thus we focus on the interface encoders h_i^{Γ} and decoders g_i^{Γ} . As stated in Section 3.4, in the strong ROM-port case, the interface encoders h_i^{Γ} and decoders g_i^{Γ} are composed of encoders h_j^{Γ} and decoders g_i^{Γ} for ports P(j). These encoders h_i^{Γ} and decoders g_i^{Γ} are of the form

$$h_{i}^{p}: \mathbb{R}^{N_{j}^{p}} \to \mathbb{R}^{n_{j}^{p}}, \qquad h_{i}^{p}(\mathbf{x}_{i}^{p}) = \mathbf{W}_{2}^{h_{j}^{p}} \sigma_{i}^{p}(\mathbf{W}_{1}^{h_{j}^{p}} \mathbf{x}_{i}^{p} + \mathbf{b}_{1}^{h_{j}^{p}}),$$
 (42a)

$$\boldsymbol{g}_{j}^{p}: \mathbb{R}^{n_{j}^{p}} \to \mathbb{R}^{N_{j}^{p}}, \qquad \qquad \boldsymbol{g}_{j}^{p}(\widehat{\boldsymbol{x}}_{1}^{p}) = \boldsymbol{W}_{2}^{\boldsymbol{g}_{j}^{p}} \boldsymbol{\sigma}_{j}^{p}(\boldsymbol{W}_{1}^{\boldsymbol{g}_{j}^{p}} \widehat{\boldsymbol{x}}_{j}^{p} + \boldsymbol{b}_{1}^{\boldsymbol{g}_{j}^{p}}), \tag{42b}$$

where

$$\boldsymbol{W}_{1}^{\boldsymbol{g}_{j}^{p}}, \left(\boldsymbol{W}_{2}^{\boldsymbol{h}_{j}^{p}}\right)^{T} \in \mathbb{R}^{\boldsymbol{w}_{j}^{p} \times \boldsymbol{n}_{j}^{p}}, \qquad \qquad \boldsymbol{W}_{2}^{\boldsymbol{g}_{j}^{p}}, \left(\boldsymbol{W}_{1}^{\boldsymbol{h}_{j}^{p}}\right)^{T} \in \mathbb{R}^{N_{j}^{p} \times \boldsymbol{w}_{j}^{p}}, \tag{43a}$$

$$\boldsymbol{b}_{1}^{\boldsymbol{b}_{j}^{p}} \in \mathbb{R}^{w_{j}^{p}}, \qquad \qquad \boldsymbol{b}_{1}^{\boldsymbol{g}_{j}^{p}} \in \mathbb{R}^{w_{j}^{p}}, \tag{43b}$$

 σ_j^p are smooth, non-polynomial activation functions (e.g., Sigmoid or Swish), and where w_j^p are the network widths for all ports P(j), $j=1,\ldots,n_p$. The weight matrices $\boldsymbol{W}_2^{\boldsymbol{g}_j^p}$ and $\boldsymbol{W}_1^{h_j^p}$ are sparse, while the remaining weights and biases are dense. As in the WFPC case, the width w_j^p and the sparsity patterns of $\boldsymbol{W}_2^{\boldsymbol{g}_j^p}$ and $\boldsymbol{W}_1^{h_j^p}$ are hyper-parameters that require tuning. The autoencoders are trained by minimizing the MSE loss defined in equation (27).

Recall that the interface encoders h_i^{Γ} and g_i^{Γ} are computed using equations (28) and (22), respectively. The encoders h_i^{Γ} and decoders g_i^{Γ} can be written in the form (40)c, d as follows. For a given subdomain i, let $j_1, \ldots, j_{|Q(i)|}$ denote the indices of the subdomains contained in Q(i). The weights and biases of h_i^{Γ} and g_i^{Γ} can then be assembled in block form as

$$\boldsymbol{W}_{1}^{\boldsymbol{h}_{i}^{\Gamma}} = \begin{bmatrix} \boldsymbol{W}_{1}^{\boldsymbol{h}_{j_{1}}^{P}} & & & \\ & \ddots & & & \\ & & \boldsymbol{W}_{1}^{\boldsymbol{h}_{j|Q(i)|}^{P}} \end{bmatrix} \begin{bmatrix} \boldsymbol{P}_{i}^{j_{1}} \\ \vdots \\ \boldsymbol{P}_{i}^{j_{|Q(i)|}} \end{bmatrix}, \qquad \boldsymbol{b}_{1}^{\boldsymbol{h}_{i}^{\Gamma}} = \begin{bmatrix} \boldsymbol{b}_{1}^{\boldsymbol{h}_{j_{1}}^{P}} \\ \boldsymbol{b}_{1}^{\vdots} \\ \boldsymbol{b}_{1}^{\boldsymbol{h}_{j|Q(i)|}} \end{bmatrix}, \tag{44a}$$

$$\boldsymbol{W}_{2}^{\boldsymbol{h}_{i}^{\Gamma}} = \left[(\widehat{\boldsymbol{P}}_{i}^{j_{1}})^{T} \quad \dots \quad (\widehat{\boldsymbol{P}}_{i}^{j|Q(i)})^{T} \right] \begin{bmatrix} \boldsymbol{W}_{2}^{\boldsymbol{h}_{j_{1}}^{p}} & & \\ & \ddots & \\ & & \boldsymbol{W}_{j|Q(i)} \end{bmatrix}, \tag{44b}$$

$$\boldsymbol{W}_{1}^{\boldsymbol{g}_{1}^{\Gamma}} = \begin{bmatrix} \boldsymbol{W}_{1}^{\boldsymbol{g}_{j_{1}}^{p}} & & \\ & \ddots & & \\ & & \boldsymbol{W}_{1}^{\boldsymbol{g}_{j|Q(i)|}^{p}} \end{bmatrix} \begin{bmatrix} \hat{\boldsymbol{p}}_{i}^{j_{1}} \\ \vdots \\ \hat{\boldsymbol{p}}_{i}^{j_{|Q(i)|}} \end{bmatrix}, \qquad \boldsymbol{b}_{1}^{\boldsymbol{g}_{1}^{\Gamma}} = \begin{bmatrix} \boldsymbol{b}_{1}^{\boldsymbol{g}_{j_{1}}^{p}} \\ \vdots \\ \boldsymbol{b}_{1}^{\boldsymbol{g}_{j|Q(i)|}^{p}} \end{bmatrix}, \tag{44c}$$

$$\boldsymbol{W}_{2}^{\mathbf{g}_{i}^{T}} = \left[(\boldsymbol{P}_{i}^{j_{1}})^{T} \dots (\boldsymbol{P}_{i}^{j_{|Q(i)|}})^{T} \right] \begin{bmatrix} \boldsymbol{W}_{2}^{\mathbf{g}_{j_{1}}^{p}} & & & & \\ & \ddots & & & \\ & & \boldsymbol{W}_{2}^{\mathbf{g}_{j_{|Q(i)|}}^{p}} \end{bmatrix}, \tag{44d}$$

with activation

$$\sigma_i^{\Gamma}(\cdot) = \begin{bmatrix} \sigma_{j_1}^{P}(\cdot) \\ \vdots \\ \sigma_{j_{|Q(j)|}}^{P}(\cdot) \end{bmatrix}. \tag{44e}$$

5.3. Hyper-reduction

If no hyper-reduction (HR) is applied (i.e. $B_i = I$) when solving DD ROM (13), the computational savings from the ROM is limited because the evaluation of the residuals r_i and their Jacobians still scales with the dimension of the FOM. Thus HR is applied to decrease the computational complexity caused by the nonlinearity of r_i , and increase the computational speedup. Possible HR approaches include collocation ($B_i = Z_i$) and gappy POD ($B_i = (Z_i \Phi_i^r)^{\dagger} Z_i$) [16,61]. In both cases, only a subsample of the residual components and their corresponding Jacobian components are computed. This subsample is determined by the row-sampling matrix Z_i , which is typically computed greedily (see Remark 2).

Now for both cases $B_i = Z_i$ and $B_i = (Z_i \Phi_i^r)^{\dagger} Z_i$, one must compute the products $Z_i r_i$, $Z_i \frac{\partial r_i}{\partial x_i^{iQ}}$, and $Z_i \frac{\partial r_i}{\partial x_i^{iP}}$. In implementation, instead of computing matrix-vector or matrix-matrix products, one only needs to compute the entries of r_i and rows of $\frac{\partial r_i}{\partial x_i^{iQ}}$ and

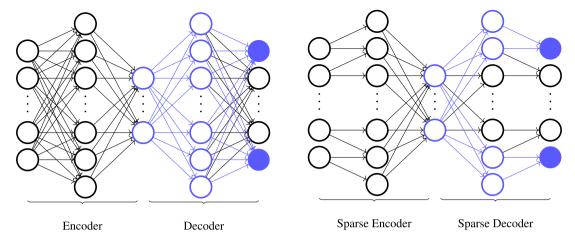


Fig. 3. Left: Dense autoencoder. The HR nodes are represented by solid blue neurons, and the nodes required to compute the HR nodes are outlined in blue. Notice that each node in the decoder hidden layer are required to compute the HR nodes in the output layer. Right: Sparse autoencoder. The encoder input layer and decoder output layer are sparsely connected, and only the blue-outlined hidden nodes are required to compute the HR nodes. The sparse output layer allows one to only keep track of the blue connections to evaluate g_i^{Ω} , g_i^{Γ} and their Jacobians, resulting in computational speedup.

 $\frac{\partial r_i}{\partial x_i^{\Gamma}}$ that are sampled by Z_i . Hence the application of HR is typically code-intrusive. Moreover, since only a subset of the entries of $g_i^{\Omega}(\hat{x}_i^{\Omega})$ and $g_i^{\Gamma}(\hat{x}_i^{\Gamma})$ are needed to evaluate $Z_i r_i$, $Z_i \frac{\partial r_i}{\partial x_i^{\Omega}}$, and $Z_i \frac{\partial r_i}{\partial x_i^{\Gamma}}$, only the corresponding outputs of the decoders g_i^{Ω} and g_i^{Γ} need to be kept track of. This motivates the use of a sparsity mask in the last layer of the decoders, which was first introduced in the context of NM-ROM in [56].

Indeed, in the case of a dense linear layer, each node in the hidden layer is needed to compute one node in the output layer, thus limiting the computational savings gained through HR. If instead a sparsity mask is applied to the layer, only a subset of the hidden nodes is required to compute a node in the output layer. This allows for the computation of a *subnet*, which only keeps track of the nodes used to compute the output nodes remaining after HR. We discuss the computation of a subnet in Section 5.4. Fig. 3 provides a visualization of a subnet. In this paper, we also apply the transpose of decoder sparsity mask to the input layer of the encoder, resulting in autoencoders whose architectures are (approximately) symmetric across the latent layer. We found that this choice gave improved performance (i.e. ROM accuracy) over the architectures used in [56], which use dense encoder input layers.

Remark 2. Following [16] and [56], we use [71, Algo. 3] to greedily compute a row sampling matrix Z_i . This approach relies upon the computation of a residual basis Φ_i^r for each subdomain. In practice, these residual bases are computed by applying POD to *residual snapshots*, which are collected from the iteration history of Newton's method when computing interior- and interface-state snapshots for ROM training.

Remark 3. For the WFPC case, the products $CA_ig_i^\Gamma(\hat{x}_i^\Gamma)$ and $\hat{\lambda}^TCA_i\frac{dg_i^\Gamma}{d\hat{x}_i^\Gamma}(\hat{x}_i^\Gamma)$ coming from the equality constraint appear in the KKT system (33) for the SQP solver. For LS-ROM, since the Jacobian of g_i^Γ is nothing but the POD basis matrix Φ_i^Γ , one can easily precompute $CA_i\Phi_i^\Gamma$, thus making HR unnecessary for these quantities. However for NM-ROM, $dg_i^\Gamma/d\hat{x}_i^\Gamma$ must be re-evaluated at each iteration of the SQP solver, thus introducing additional computation expense that is not present in LS-ROM. Currently, these quantities do not undergo HR in the NM-ROM case, and hence the decoders g_i^Γ for the entire interface states must be kept track of. While this limits the computational savings that can be obtained through HR, in practice the dimension of the FOM interface states N_i^Γ is much smaller than the dimension of the FOM interior states N_i^Ω , thus making the HR of g_i^Γ less critical than the HR of g_i^Ω . This issue is not present in the SRPC case because the constraints are purely linear.

Remark 4. In practice, the pattern of the sparsity mask is determined by a number of hyper parameters to be tuned by the user. Further details on the sparsity pattern used for our numerical results is discussed in Section 7.

5.4. Construction of a subnet

The authors in [56, Sec. 4.4.1] discuss the construction of a subnet in terms of gradients of the loss function with respect to the weights and biases of the sparse decoder. In this section, we present an alternative method for constructing the subnet solely by keeping track of indices of HR nodes. For simplicity, we consider a generic decoder $g: \mathbb{R}^n \to \mathbb{R}^N$ of the form ((40)b, d) with width w. Computing subnets for each decoder, e.g., g_i^Ω and g_i^Γ , follows the same procedure. Recall that in the architecture considered in this paper, $\mathbf{W}_2 \in \mathbb{R}^{N \times w}$ is sparse and $\mathbf{W}_1 \in \mathbb{R}^{w \times n}$ is dense.

Let $\mathcal{I}_o \subset \{1, \dots, N\}$ denote the indices of the outputs of g that are selected through HR. To find the indices of the hidden nodes required to compute the HR nodes, find the index set \mathcal{I}_b where

$$I_h = \{ j \in \{1, ..., w\} \mid \exists i \in I_o \text{ s.t. } (\mathbf{W}_2)_{ij} \neq 0 \}.$$

The index sets \mathcal{I}_o and \mathcal{I}_h contain the indices of all nonzero elements in \mathbf{W}_2 . Now let $i_1, \dots, i_{|\mathcal{I}_o|}$ and $j_1, \dots, j_{|\mathcal{I}_h|}$ denote the elements of \mathcal{I}_o and \mathcal{I}_h in ascending order, respectively. Define the matrix $\widetilde{\mathbf{W}}_2 \in \mathbb{R}^{|\mathcal{I}_o| \times |\mathcal{I}_h|}$ as

$$(\widetilde{\boldsymbol{W}}_2)_{\ell,k} = (\boldsymbol{W}_2)_{i_{\ell},i_{\ell}}, \quad \forall \ \ell = 1, \dots, |\mathcal{I}_o|, \ k = 1, \dots, |\mathcal{I}_h|.$$

The matrix \widetilde{W}_2 precisely consists of the connections in the subnet that remain after HR. Next, since the activation σ acts elementwise, the connections that remain in the first layer of the subnet can be represented by $\widetilde{W}_1 \in \mathbb{R}^{|\mathcal{I}_h| \times n}$, which consists of nothing but the rows in W_1 corresponding to the index set \mathcal{I}_h :

$$(\widetilde{\boldsymbol{W}}_1)_{k,:} = \boldsymbol{W}_{j_k,:}, \quad \forall \ k = 1, \dots, |\mathcal{I}_h|.$$

Lastly, the subnet bias $\widetilde{\boldsymbol{b}}_1 \in \mathbb{R}^{|\mathcal{I}_h|}$ is similarly defined as $(\widetilde{\boldsymbol{b}}_1)_k = (\boldsymbol{b}_1)_{j_k}$ for all $k = 1, \dots, |\mathcal{I}_h|$. The subnet $\widetilde{\boldsymbol{g}} : \mathbb{R}^n \to \mathbb{R}^{|\mathcal{I}_o|}$ is then defined as

$$\widetilde{g}(\widehat{x}) = \widetilde{W}_2 \sigma(\widetilde{W}_1 \widehat{x} + \widetilde{b}_1). \tag{45}$$

Remark 5. This framework for computing a subnet can easily be extended to neural networks with arbitrarily many sparse linear layers provided that the sparsity patterns for each layer's weight matrix is known. Thus one could construct deep sparse autoencoders with narrower width than the architectures considered here. However, for this paper we only consider single-layer, wide, sparse decoders.

6. Error analysis

We present *a priori* and *a posteriori* error bounds analogous to those found in [16]. To simplify notation, analogous to the notation in Section 4.2, we denote the optimal solutions to the FOM (10), to the ROM (29), and the ROM solution lifted to the FOM state space as

$$\boldsymbol{x}^{*} = \begin{bmatrix} \boldsymbol{x}_{1}^{\Omega*} \\ \boldsymbol{x}_{1}^{\Gamma*} \\ \vdots \\ \boldsymbol{x}_{n_{O}}^{\Omega*} \\ \boldsymbol{x}_{n_{O}}^{\Gamma*} \end{bmatrix} \in \mathbb{R}^{N_{D}}, \hat{\boldsymbol{x}}^{*} = \begin{bmatrix} \hat{\boldsymbol{x}}_{1}^{\Omega*} \\ \hat{\boldsymbol{x}}_{1}^{\Gamma*} \\ \vdots \\ \hat{\boldsymbol{x}}_{n_{O}}^{\Omega*} \\ \hat{\boldsymbol{x}}_{n_{O}}^{\Gamma*} \end{bmatrix} \in \mathbb{R}^{n_{D}}, \boldsymbol{g}(\hat{\boldsymbol{x}}^{*}) = \begin{bmatrix} \boldsymbol{g}_{1}^{\Omega}(\hat{\boldsymbol{x}}_{1}^{\Omega*}) \\ \boldsymbol{g}_{1}^{\Gamma}(\hat{\boldsymbol{x}}_{1}^{T*}) \\ \vdots \\ \boldsymbol{g}_{n_{O}}^{\Omega}(\hat{\boldsymbol{x}}_{n_{O}}^{\Omega*}) \\ \boldsymbol{g}_{n_{O}}^{T}(\hat{\boldsymbol{x}}_{n_{O}}^{T*}) \end{bmatrix} \in \mathbb{R}^{N_{D}},$$

$$(46)$$

respectively, where $N_D = \sum_{i=1}^{n_\Omega} (N_i^\Omega + N_i^\Gamma)$ and, as before, $n_D = \sum_{i=1}^{n_\Omega} (n_i^\Omega + n_i^\Gamma)$. We also define the FOM constraint matrix

$$\mathbf{A} = \begin{bmatrix} \mathbf{0} & \mathbf{A}_1 & \dots & \mathbf{0} & \mathbf{A}_{n_O} \end{bmatrix} \in \mathbb{R}^{N_A \times N_D} \tag{47}$$

so that the constraints ((10)b) can be written as $\mathbf{A}\mathbf{x}=\mathbf{0}$. As in Section 4, we define the constraint functions $\widetilde{\mathbf{A}}_i:\mathbb{R}^{n_i^\Gamma}\to\mathbb{R}^{n_A}$, where $\widetilde{\mathbf{A}}_i(\widehat{\mathbf{x}}_i^\Omega)=C\mathbf{A}_i\mathbf{g}_i^\Gamma(\widehat{\mathbf{x}}_i^\Gamma)$ in the WFPC case (13) and $\widetilde{\mathbf{A}}_i(\widehat{\mathbf{x}}_i^\Gamma)=\widehat{\mathbf{A}}_i\widehat{\mathbf{x}}_i^\Gamma$ in the SRPC case (23), so that the DD-LSPG-ROMs (13) and (23) can be written as (29). We define the ROM constraint function $\widetilde{\mathbf{A}}:\mathbb{R}^{n_D}\to\mathbb{R}^{n_A}$ as

$$\widetilde{A}(\widehat{\mathbf{x}}) = \sum_{i=1}^{n_{\Omega}} \widetilde{A}_i(\widehat{\mathbf{x}}_i^{\Gamma}),\tag{48}$$

so that the constraints ((29)b) can be written as $\widetilde{A}(\hat{x}) = 0$. Lastly we define the feasible set

$$S_{\text{ROM}} = \left\{ \hat{\mathbf{x}} \in \mathbb{R}^{n_D} : \widetilde{\mathbf{A}}(\hat{\mathbf{x}}) = \mathbf{0} \right\}$$
(49)

for (29).

The next two results provide basic error bounds between a solution to the FOM (9) and solutions to the DD-LSPG-ROM (13) or (23).

Theorem 5 (A Posteriori Error Bound). Let $\mathbf{x}^* \in \mathbb{R}^{N_D}$ be a solution to the FOM (9) and let $\hat{\mathbf{x}}^* \in \mathbb{R}^{n_D}$ be a (local) solution to the DD-LSPG-ROM (13) or (23). If the residual is inverse Lipschitz continuous, that is, if there exists $\kappa_{\ell} > 0$ such that

$$\left(\sum_{i=1}^{n_{\Omega}} \left\| \boldsymbol{r}_{i}(\boldsymbol{y}_{i}^{\Omega}, \boldsymbol{y}_{i}^{\Gamma}) - \boldsymbol{r}_{i}(\boldsymbol{z}_{i}^{\Omega}, \boldsymbol{z}_{i}^{\Gamma}) \right\|_{2}^{2} \right)^{1/2} \ge \kappa_{\ell} \|\boldsymbol{y} - \boldsymbol{z}\|_{2} \qquad \forall \, \boldsymbol{y}, \boldsymbol{z} \in \mathbb{R}^{N_{D}},$$
(50a)

and if there exists P > 0 such that

$$\left(\sum_{i=1}^{n_{\Omega}} \left\| \boldsymbol{B}_{i} \boldsymbol{r}_{i} (\boldsymbol{y}_{i}^{\Omega}, \boldsymbol{y}_{i}^{\Gamma}) \right\|_{2}^{2} \right)^{1/2} \geq P \left(\sum_{i=1}^{n_{\Omega}} \left\| \boldsymbol{r}_{i} (\boldsymbol{y}_{i}^{\Omega}, \boldsymbol{y}_{i}^{\Gamma}) \right\|_{2}^{2} \right)^{1/2} \qquad \forall \ \boldsymbol{y} \in \boldsymbol{g}(S_{\text{ROM}}),$$

$$(50b)$$

A.N. Diaz et al.

then

$$\|\mathbf{x}^* - \mathbf{g}(\hat{\mathbf{x}}^*)\|_2 \le \frac{1}{P\kappa_{\ell}} \left(\sum_{i=1}^{n_{\Omega}} \|\mathbf{B}_i \mathbf{r}_i \mathbf{g}_i^{\Omega}(\hat{\mathbf{x}}^{\Omega*}), \mathbf{g}_i^{\Gamma}(\hat{\mathbf{x}}^{\Gamma*}) \|_2^2 \right)^{1/2}.$$
(51)

Proof. Using ((50)a) and the fact that $x^* \in \mathbb{R}^{N_D}$ solves the FOM (9) gives

$$\left\|\boldsymbol{x}^* - \boldsymbol{g}(\widehat{\boldsymbol{x}}^*)\right\|_2^2 \leq \frac{1}{\kappa_\ell^2} \sum_{i=1}^{n_\Omega} \left\|\boldsymbol{r}_i(\boldsymbol{x}_i^{\Omega*}, \boldsymbol{x}_i^{\Gamma*}) - \boldsymbol{r}_i(\boldsymbol{g}_i^{\Omega}(\widehat{\boldsymbol{x}}_i^{\Omega*}), \boldsymbol{g}_i^{\Gamma}(\widehat{\boldsymbol{x}}_i^{\Gamma*}))\right\|_2^2 \leq \frac{1}{\kappa_\ell^2} \sum_{i=1}^{n_\Omega} \left\|\boldsymbol{r}_i(\boldsymbol{g}_i^{\Omega}(\widehat{\boldsymbol{x}}_i^{\Omega*}), \boldsymbol{g}_i^{\Gamma}(\widehat{\boldsymbol{x}}_i^{\Gamma*}))\right\|_2^2.$$

Applying ((50)b) with $(\mathbf{y}_i^{\Omega}, \mathbf{y}_i^{\Gamma}) = (\mathbf{g}_i^{\Omega}(\hat{\mathbf{x}}_i^{\Omega*}), \mathbf{g}_i^{\Gamma}(\hat{\mathbf{x}}_i^{\Gamma*}))$ gives the desired result. \square

Theorem 6 (A Priori Error Bound). Let $\mathbf{x}^* \in \mathbb{R}^{N_D}$ be a solution to the FOM (9) and let $\hat{\mathbf{x}}^* \in \mathbb{R}^{n_D}$ be a solution to the DD-LSPG-ROM (13) or (23). If the inequalities ((50)a, b) hold and the HR residual is Lipschitz continuous, i.e., there exists $\kappa_n > 0$ such that

$$\left(\sum_{i=1}^{n_{\Omega}} \left\| \boldsymbol{B}_{i} \boldsymbol{r}_{i} (\boldsymbol{y}_{i}^{\Omega}, \boldsymbol{y}_{i}^{\Gamma}) - \boldsymbol{B}_{i} \boldsymbol{r}_{i} (\boldsymbol{z}_{i}^{\Omega}, \boldsymbol{z}_{i}^{\Gamma}) \right\|_{2}^{2} \right)^{1/2} \leq \kappa_{u} \|\boldsymbol{y} - \boldsymbol{z}\|_{2} \qquad \forall \, \boldsymbol{y}, \boldsymbol{z} \in \mathbb{R}^{N_{D}},$$

$$(52)$$

then

$$\|\mathbf{x}^* - g(\hat{\mathbf{x}}^*)\|_2 \le \frac{\kappa_u}{P\kappa_\ell} \inf_{\hat{\mathbf{w}} \in S_{\text{DOM}}} \|\mathbf{x}^* - g(\hat{\mathbf{w}})\|_2.$$
 (53)

Proof. Since $\hat{x}^* \in \mathbb{R}^{n_D}$ is a solution to the DD-LSPG-ROM (13) or (23), any feasible \hat{w} , i.e., any $\hat{w} \in S_{\text{ROM}}$ satisfies

$$\sum_{i=1}^{n_{\Omega}} \left\| \boldsymbol{B}_{i} \boldsymbol{r}_{i}(\boldsymbol{g}_{i}^{\Omega}(\widehat{\boldsymbol{x}}_{i}^{\Omega*}), \boldsymbol{g}_{i}^{\Gamma}(\widehat{\boldsymbol{x}}_{i}^{\Gamma*})) \right\|_{2}^{2} \leq \sum_{i=1}^{n_{\Omega}} \left\| \boldsymbol{B}_{i} \boldsymbol{r}_{i}(\boldsymbol{g}_{i}^{\Omega}(\widehat{\boldsymbol{w}}_{i}^{\Omega*}), \boldsymbol{g}_{i}^{\Gamma}(\widehat{\boldsymbol{w}}_{i}^{\Gamma*})) \right\|_{2}^{2}. \tag{54}$$

Moreover, since $\mathbf{x}^* \in \mathbb{R}^{N_D}$ solves the FOM (9), $\mathbf{r}_i(\mathbf{x}_i^{\Omega*}, \mathbf{x}_i^{\Gamma*}) = \mathbf{0}$, for all $i = 1, ..., n_O$, (52) and (54) imply

$$\begin{split} &\sum_{i=1}^{n_{\Omega}} \left\| \boldsymbol{B}_{i} \boldsymbol{r}_{i}(\boldsymbol{g}_{i}^{\Omega}(\widehat{\boldsymbol{x}}_{i}^{\Omega*}), \boldsymbol{g}_{i}^{\Gamma}(\widehat{\boldsymbol{x}}_{i}^{\Gamma*})) \right\|_{2}^{2} \leq \sum_{i=1}^{n_{\Omega}} \left\| \boldsymbol{B}_{i} \boldsymbol{r}_{i}(\boldsymbol{x}_{i}^{\Omega*}, \boldsymbol{x}_{i}^{\Gamma*}) - \boldsymbol{B}_{i} \boldsymbol{r}_{i}(\boldsymbol{g}_{i}^{\Omega}(\widehat{\boldsymbol{w}}_{i}^{\Omega}), \boldsymbol{g}_{i}^{\Gamma}(\widehat{\boldsymbol{w}}_{i}^{\Gamma})) \right\|_{2}^{2} \\ &\leq \kappa_{u} \left\| \boldsymbol{x}^{*} - \boldsymbol{g}(\widehat{\boldsymbol{w}}) \right\|_{2} \quad \text{for all } \widehat{\boldsymbol{w}} \in S_{\text{ROM}}. \end{split}$$

Combining this result with the a-posterior bound (51) in Theorem 5 yields $\|x^* - g(\hat{x}^*)\|_2 \le \frac{\kappa_u}{P\kappa_{\mathcal{C}}} \|x^* - g(\hat{w})\|_2$ for all $\hat{w} \in \mathcal{S}_{ROM}$, which implies (53).

Remark 6. As a consequence of Theorem 6, if ((50)a, b) and (52) hold, and if x^* is in the image of the g over the feasible set S_{ROM} of (13), i.e. if $x^* \in g(S_{ROM})$, then $x^* = g(\hat{x}^*)$.

The error bounds in Theorems 5 and 6 only involve the FOM and ROM states, but not the Lagrange multipliers. However, the present error bounds require stronger assumptions such as ((50)a). Alternatively, one could try to extend the error analysis for ROMs applied to nonlinear systems. such as those in [22, Sec. 11.5], [72]. In the context of the FOM (9) and the DD-LSPG-ROM (13) or (23) the role of the nonlinear residual in [22, Sec. 11.5], [72] would now be played by the system of first order necessary optimality conditions, given by (31) for the general ROM formulation (29) and correspondingly for the FOM (9). However, these residuals involve the FOM states and Lagrange multipliers associated with ((9)b), and ROM states and Lagrange multipliers associated with ((29)b). Moreover, this analysis requires to relate the ROM states with the FOM states and the ROM Lagrange multipliers with the FOM Lagrange multipliers. The former is done via $g(\hat{x}) \approx x^*$. However, the connection between the Lagrange multipliers in the general case is still open. For linear PDEs and a PDE-based (as opposed to our algebraic) DD formulation, [10] construct appropriate, so called trace-compatible reduced bases for the Lagrange multipliers (see e.g., [10, Eq. (14)]). In the context of [10], the reduced bases for the Lagrange multipliers, but no explicit construction of a reduced bases for these ROM Lagrange multipliers. This is subject of future research.

7. Numerics

We apply LS-ROM and NM-ROM with and without HR to the DD ROM with WFPC (13) and with SRPC (23) for the 2D steady-state Burgers equation. We use the following formula for computing the relative error between the FOM and ROM solutions:

$$e = \left(\frac{1}{n_{\Omega}} \sum_{i=1}^{n_{\Omega}} \frac{\left\|\mathbf{x}_{i}^{\Omega} - \mathbf{g}_{i}^{\Omega}(\widehat{\mathbf{x}}_{i}^{\Omega})\right\|_{2}^{2} + \left\|\mathbf{x}_{i}^{\Gamma} - \mathbf{g}_{i}^{\Gamma}(\widehat{\mathbf{x}}_{i}^{\Gamma})\right\|_{2}^{2}}{\left\|\mathbf{x}_{i}^{\Omega}\right\|_{2}^{2} + \left\|\mathbf{x}_{i}^{\Gamma}\right\|_{2}^{2}}\right)^{1/2}.$$
(55)

The autoencoder training and subsequent computations in this section were performed on the Lassen machine at Lawrence Livermore National Laboratory, which consists of an IBM Power9 processor with NVIDIA V100 (Volta) GPUs, clock speed between 2.3–3.8 GHz, and 256 GB DDR4 memory.

The implementation of the DD FOM, DD LS-ROM, and DD NM-ROM is done sequentially. However, to highlight potential advantages of a parallel implementation, the recorded wall clock time for the computation of the subdomain-specific quantities required by the SQP solver is taken to be the largest wall clock time incurred among all subdomains. The wall clock time for the remaining steps of the SQP solver (e.g. assembling and solving the KKT system (33), updating the interior- and interface-states and Lagrange multipliers Eq. (38), etc.) is set to the overall wall clock time to execute the steps.

7.1. 2D Burgers' equation

In this experiment, we consider the 2D steady-state Burgers' equation on the domain $[-1,1] \times [0,0.05]$

$$u\frac{\partial u}{\partial x} + v\frac{\partial u}{\partial y} = v\left(\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2}\right), \qquad (x, y) \in [-1, 1] \times [0, 0.05],$$
(56a)

$$u\frac{\partial v}{\partial x} + v\frac{\partial v}{\partial y} = v\left(\frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial y^2}\right), \quad (x, y) \in [-1, 1] \times [0, 0.05],\tag{56b}$$

where $\nu > 0$ is the viscosity. As in [16], we consider the following exact solution, and its restriction to the boundary as Dirichlet boundary conditions:

$$u_{ax}(x, y; a, \lambda) = -2v \left[a + \lambda \left(e^{\lambda(x-1)} - e^{-\lambda(x-1)} \right) \cos(\lambda y) \right] / \psi(x, y; a, \lambda),$$
 (57a)

$$v_{ex}(x, y; a, \lambda) = 2\nu \left[\lambda \left(e^{\lambda(x-1)} + e^{-\lambda(x-1)}\right) \sin(\lambda y)\right] / \psi(x, y; a, \lambda), \tag{57b}$$

where

$$\psi(x, y; a, \lambda) = a(1+x) + \left(e^{\lambda(x-1)} + e^{-\lambda(x-1)}\right)\cos(\lambda y),$$
 (57c)

and where (a, λ) are parameters.

The PDE is discretized using centered finite differences with $n_x + 2$ uniformly spaced grid points in the x-direction and $n_y + 2$ uniformly spaced grid points in the y-direction, resulting in grid points (x_i, y_i) where

$$x_i = -1 + ih_x,$$
 $i = 0, ..., n_x + 1,$ $y_j = jh_y,$ $j = 0, ..., n_y + 1,$

where $h_x = 2/(n_x + 1)$ and $h_y = 0.05/(n_y + 1)$. The solutions u, v on the grid points are denoted $u_{ij} \approx u(x_i, y_j)$ and $v_{ij} \approx v(x_i, y_j)$. The PDE is then discretized using centered finite differences for the first and second derivative terms. The fully discretized system is given by

$$0 = r_{\mu}(u, v) = u \odot (B_{\nu}u - b_{\mu\nu}) + v \odot (B_{\nu}u - b_{\mu\nu}) + Cu + c_{\mu}, \tag{58a}$$

$$\mathbf{0} = r_{\nu}(u, v) = u \odot (B_{x}v - b_{\nu, x}) + v \odot (B_{y}v - b_{\nu, y}) + Cv + c_{\nu}, \tag{58b}$$

where o represents the Hadamard product, and where

$$\boldsymbol{u} = \begin{bmatrix} \boldsymbol{u}^{[1]} \\ \vdots \\ \boldsymbol{u}^{[n_y]} \end{bmatrix} \in \mathbb{R}^{n_x n_y}, \qquad \boldsymbol{u}^{[j]} = \begin{bmatrix} u_{1,j} \\ \vdots \\ u_{n_x,j} \end{bmatrix} \in \mathbb{R}^{n_x}, \quad j = 1, \dots, n_y,$$
 (59a)

$$\boldsymbol{v} = \begin{bmatrix} \boldsymbol{v}^{[1]} \\ \vdots \\ \boldsymbol{v}^{[n_y]} \end{bmatrix} \in \mathbb{R}^{n_x n_y}, \qquad \boldsymbol{v}^{[j]} = \begin{bmatrix} v_{1,j} \\ \vdots \\ v_{n...j} \end{bmatrix} \in \mathbb{R}^{n_x}, \quad j = 1, \dots, n_y,$$
 (59b)

$$\boldsymbol{B}_{x} = -\frac{1}{2h_{x}} \left(\boldsymbol{I}_{n_{y}} \otimes \widetilde{\boldsymbol{B}}_{x} \right) \in \mathbb{R}^{n_{x}n_{y} \times n_{x}n_{y}}, \qquad \widetilde{\boldsymbol{B}}_{x} = \begin{bmatrix} 0 & 1 \\ -1 & \ddots & 1 \\ -1 & 0 \end{bmatrix} \in \mathbb{R}^{n_{x} \times n_{x}}, \tag{59c}$$

$$\boldsymbol{B}_{y} = -\frac{1}{2h_{y}} \left(\widetilde{\boldsymbol{B}}_{y} \otimes \boldsymbol{I}_{n_{x}} \right) \in \mathbb{R}^{n_{x}n_{y} \times n_{x}n_{y}} \qquad \widetilde{\boldsymbol{B}}_{y} = \begin{bmatrix} 0 & 1 \\ -1 & \ddots & 1 \\ -1 & 0 \end{bmatrix} \in \mathbb{R}^{n_{y} \times n_{y}}, \tag{59d}$$

$$C = \frac{v}{h_x^2} \left(I_{n_y} \otimes \widetilde{C}_x \right) + \frac{v}{h_y^2} \left(\widetilde{C}_y \otimes I_{n_x} \right) \in \mathbb{R}^{n_x n_y \times n_x n_y}, \tag{59e}$$

$$\widetilde{\boldsymbol{C}}_{x} = \begin{bmatrix} -2 & 1 \\ 1 & \ddots \\ & 1 & -2 \end{bmatrix} \in \mathbb{R}^{n_{x} \times n_{x}}, \qquad \widetilde{\boldsymbol{C}}_{y} = \begin{bmatrix} -2 & 1 \\ 1 & \ddots \\ & 1 & -2 \end{bmatrix} \in \mathbb{R}^{n_{y} \times n_{y}}, \tag{59f}$$

$$\boldsymbol{b}_{u,x} = -\frac{1}{2h_x}(\boldsymbol{b}_{ux\ell} - \boldsymbol{b}_{uxr}), \qquad \boldsymbol{b}_{u,y} = -\frac{1}{2h_y}(\boldsymbol{b}_{uy\ell} - \boldsymbol{b}_{uyr}), \tag{59g}$$

$$c_{u} = \frac{v}{h_{x}^{2}}(\boldsymbol{b}_{ux\ell} + \boldsymbol{b}_{uxr}) + \frac{v}{h_{y}^{2}}(\boldsymbol{b}_{uy\ell} + \boldsymbol{b}_{uyr})$$
 (59h)

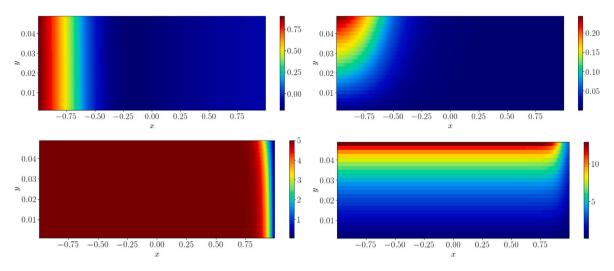


Fig. 4. Top left: *u*-component with $(a, \lambda) = (10^4, 5)$; Top right: *v*-component with $(a, \lambda) = (10^4, 5)$; Bottom left: *u*-component with $(a, \lambda) = (1, 25)$; Bottom right: *v*-component with $(a, \lambda) = (1, 25)$.

$$\boldsymbol{b}_{v,x} = -\frac{1}{2h_x}(\boldsymbol{b}_{vx\ell} - \boldsymbol{b}_{vxr}), \qquad \boldsymbol{b}_{v,y} = -\frac{1}{2h_y}(\boldsymbol{b}_{vy\ell} - \boldsymbol{b}_{vyr}), \tag{59i}$$

$$c_{v} = \frac{v}{h_{x}^{2}}(b_{vx\ell} + b_{vxr}) + \frac{v}{h_{y}^{2}}(b_{vy\ell} + b_{vyr})$$
 (59j)

$$\boldsymbol{b}_{ux\ell} = \begin{bmatrix} u_{ex}(x_0, y_1) \\ \vdots \\ u_{ex}(x_0, y_{n_y}) \end{bmatrix} \otimes \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}_{n_x \times 1} \in \mathbb{R}^{n_x n_y}, \qquad \boldsymbol{b}_{uxr} = \begin{bmatrix} u_{ex}(x_{n_x+1}, y_1) \\ \vdots \\ u_{ex}(x_{n_x+1}, y_{n_y}) \end{bmatrix} \otimes \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}_{n_x \times 1} \in \mathbb{R}^{n_x n_y}, \tag{59k}$$

$$\boldsymbol{b}_{uyb} = \begin{bmatrix} 1\\0\\\vdots\\0\\n_{x} \times 1 \end{bmatrix} \otimes \begin{bmatrix} u_{ex}(x_1, y_0)\\\vdots\\u_{ex}(x_{n_x}, y_0) \end{bmatrix} \in \mathbb{R}^{n_x n_y} \qquad \boldsymbol{b}_{uyt} = \begin{bmatrix} 0\\\vdots\\0\\1\\n_{x} \times 1 \end{bmatrix} \otimes \begin{bmatrix} u_{ex}(x_1, y_{n_y+1})\\\vdots\\u_{ex}(x_{n_x}, y_{n_y+1}) \end{bmatrix} \in \mathbb{R}^{n_x n_y}$$
(591)

$$\boldsymbol{b}_{vx\ell} = \begin{bmatrix} v_{ex}(x_0, y_1) \\ \vdots \\ v_{ex}(x_0, y_{n_y}) \end{bmatrix} \otimes \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}_{n \to 1} \in \mathbb{R}^{n_x n_y}, \qquad \boldsymbol{b}_{vxr} = \begin{bmatrix} v_{ex}(x_{n_x+1}, y_1) \\ \vdots \\ v_{ex}(x_{n_x+1}, y_{n_y}) \end{bmatrix} \otimes \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}_{n \to 1} \in \mathbb{R}^{n_x n_y}, \tag{59m}$$

$$\boldsymbol{b}_{vyb} = \begin{bmatrix} 1\\0\\\vdots\\0\\n_{x} \le \begin{bmatrix} v_{ex}(x_1, y_0)\\\vdots\\v_{ex}(x_{n_x}, y_0) \end{bmatrix} \in \mathbb{R}^{n_x n_y} \qquad \boldsymbol{b}_{vyt} = \begin{bmatrix} 0\\\vdots\\0\\1\\n_{x} \le \begin{bmatrix} v_{ex}(x_1, y_{n_y+1})\\\vdots\\v_{ex}(x_{n_x}, y_{n_y+1}) \end{bmatrix} \in \mathbb{R}^{n_x n_y}. \tag{59n}$$

For the monolithic (single domain) FOM, we take $n_x = 480$, $n_y = 24$, viscosity v = 0.1, and parameters $(a, \lambda) \in \mathcal{D} = [1, 10^4] \times [5, 25]$. The parameter a corresponds to the distance of the shock from the left boundary, whereas λ corresponds to the steepness of the shock, as illustrated in Fig. 4. We use the ROMs to predict the case where $(a, \lambda) = (7692.5384, 21.9230)$. The SQP solver for the DD FOM, DD LS-ROM, and DD NM-ROM terminates when the 2-norm of the right hand side of (33) is less than 10^{-4} , or after 15 iterations.

7.2. Snapshot data collection

To compute ROMs, we first collect 6400 snapshots for training with parameters (a, λ) uniformly sampled in a 80 × 80 grid for the full-domain problem. These full-domain snapshots are then restricted to the interior, interface, and port states, which are then used for training. This is the so-called "top-down" approach. The residual bases Φ_i^r for each subdomain are computed by taking the Newton iteration history for 400 state snapshots sampled on a 20 × 20 (a, λ) grid, and computing a POD basis with energy criterion $v = 10^{-10}$. These 400 state snapshots are then used to train RBF interpolator models (using Scipy's RBFInterpolator function) for each subdomain's interior and interface states, which are then used to compute an initial iterate for $(\hat{x}_i^\Omega, \hat{x}_i^\Gamma)$ for the SQP solver. The $(\hat{x}_i^\Omega, \hat{x}_i^\Gamma)$ initial iterates are then used to compute an initial iterate for the Lagrange multipliers λ by applying a least-squares

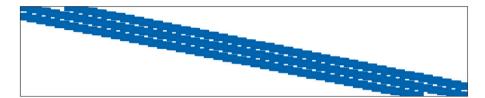


Fig. 5. Three-banded sparsity mask for decoder.

Table 1
Comparison of dense and sparse encoder architectures.

	Train loss	Test loss	Width	# encoder parameters
Dense encoder Sparse encoder	5.74×10^{-2} 7.21×10^{-5}	1.56×10^{-1} 7.48×10^{-5}	17 280 28 800	9.96×10^7 2.3×10^5

solver to equation ((32)b). The wall clock time to compute the initial iterates $(\hat{x}_i^{\Omega}, \hat{x}_i^{\Gamma})$ is taken to be the largest wall clock time incurred among all subdomains, while the wall clock time to compute the initial iterate for λ is the time required to sequentially solve the least-squares problem ((32)b). The wall clock time to compute the initial iterate for NM-ROM using the RBF interpolator is included in the computation times and speedups reported.

7.3. Autoencoder training

For NM-ROM, we randomly split the state snapshots into 5760 training snapshots and 640 testing snapshots. For training the autoencoders, we use the MSE loss, the Adam optimizer over 2000 epochs, and a batch size of 32. We also normalize the snapshots so that all snapshot components are in [-1,1], and apply a de-normalization layer to the output of the autoencoder. We apply early stopping with a stopping patience of 300, and reduce the learning rate on plateau with an initial learning rate of 10^{-3} and a patience of 50. The implementation was done using PyTorch, as well as the Pytorch Sparse and SparseLinear packages.

The sparsity masks used for the output layers of the decoders have a banded structure inspired by 2D finite difference stencils. Each row has three bands, where each band consists of contiguous nonzero entries, and where the band shifts to the right a specified amount from one row to the next. The number of nonzero entries per band and the number of columns the band shifts over are hyper-parameters. The separation between the bands in each row is equal to the product of the number of nonzeros per band and the column-shift per row. These parameters, as well as the dimension of the interior and interface states, determine the width of the decoders. Fig. 5 provides a visualization for the decoder mask used. The transpose of these masks is used at input layer of the encoders.

7.4. Dense vs sparse encoder

We first briefly compare the performance of a dense encoder with a sparse encoder. We train two autoencoders, one with a fully dense encoder and one with a sparse encoder, on a coarse, single domain problem with $n_x = 240$ and $n_y = 12$. Both decoders share the same sparse architecture. The data collection, training, and relevant sparsity masks are identical to the procedures discussed in Section 7.3. The discretization results in a FOM size of 5760 and we used a ROM of size 4.

Table 1 summarizes the key performance differences between the two encoders. Importantly, the sparse encoder architecture achieves losses 4 orders of magnitude smaller than the dense encoder with 2 orders of magnitude fewer parameters.

7.5. Single-domain NM-ROM vs DD NM-ROM

Next we examine the per-subdomain reduction in the required number of autoencoder parameters for different DD configurations compared to the monolithic single-domain NM-ROM. See Table 2. In the single domain case, we solve the LSPG problem

$$\min_{\widehat{\mathbf{x}}} \left\| \mathbf{Br} \left(\mathbf{g} \left(\widehat{\mathbf{x}} \right) \right) \right\|_{2}^{2} \tag{60}$$

using the Gauss–Newton method. The function $g: \mathbb{R}^{n_x} \to \mathbb{R}^{N_x}$ is the decoder of an autoencoder trained on snapshots of the monolithic single-domain FOM as discussed in Sections 3.4, 5, and 7.3, and $\mathbf{B} \in \{0,1\}^{N_B \times N_x}$ is a collocation HR matrix, as discussed in Section 5.3.

We use the notation 2×1 subdomains to indicate 2 subdomains in the *x*-direction and 1 subdomain in the *y*-direction. As expected, from Table 2, we see that the maximum number of NN parameters per subdomain decreases significantly as more subdomains are used. Furthermore, the total number of NN parameters in the DD cases also decreases relative to the single-domain case. We also note that the error increases as more subdomains are used. We kept $(n_i^\Omega, n_i^\Gamma) = (6, 3)$ constant for each subdomain configuration to isolate the effect of DD on the number of NN parameters, but this may cause overfitting in the 16 subdomain case. More careful hyper-parameter tuning is necessary to mitigate increases in error as the number of subdomains is increased.

Table 2Max number of NN parameters per subdomain, the per-subdomain reduction in number of NN parameters, the total number of parameters, and the corresponding error for different subdomain configurations. For the single-domain case, an NM-ROM of dimension n = 9 is used. For the DD cases, $(n_i^{\Omega}, n_i^{\Gamma}) = (6, 3)$, resulting in 9 DoF per subdomain. HR was not used to evaluate the NM-ROMs in these examples.

Max # subdomain params.	Reduction	Total # params.	Error
2.995 × 10 ⁶	0.0%	2.995×10^{6}	1.08×10^{-3}
1.147×10^6	61.7%	2.307×10^{6}	1.27×10^{-3}
5.257×10^5	82.4%	2.384×10^{6}	2.42×10^{-3}
2.617×10^5	91.3%	2.391×10^{6}	4.26×10^{-3}
1.297×10^5	95.7%	2.406×10^{6}	4.58×10^{-2}
	2.995×10^{6} 1.147×10^{6} 5.257×10^{5} 2.617×10^{5}	2.995×10^{6} 0.0% 1.147×10^{6} 61.7% 5.257×10^{5} 82.4% 2.617×10^{5} 91.3%	$\begin{array}{cccccccccccccccccccccccccccccccccccc$

7.6. LS-ROM vs NM-ROM comparison

Next we compare the DD LS-ROM and DD NM-ROM. We first focus on a DD configuration with 2 uniformly sized subdomains in the x-direction and 2 uniformly sized subdomains in the y-direction (4 subdomains total) using the WFPC formulation (13). The interior and interface states for the FOM were of dimension $N_i^{\Omega} = 5238$ and $N_i^{\Gamma} = 1006$, respectively, resulting in 25056 degrees of freedom (DoF) aggregated across all subdomains. For both the LS-ROM and NM-ROM, we use reduced state dimensions of $n_i^{\Omega} = 8$ for the interior states \hat{x}_i^{Ω} and $n_i^{\Gamma} = 4$ for the interface states \hat{x}_i^{Γ} for each subdomain, resulting in 48 DoF aggregated across all subdomains. In the HR case, $N_i^{B} = 100$ HR nodes are used for each subdomain, resulting in 400 total HR nodes aggregated all subdomains.

Each interior-state autoencoder has input/output dimension $N_i^{\Omega} = 5238$, width $w_i^{\Omega} = 26290$, latent dimension $n_i^{\Omega} = 8$, and Swish activation $\sigma_i^{\Omega}(z) = z/(1 + e^{-z})$. Each interface-state autoencoder has input/output dimension $N_i^{\Gamma} = 1006$, width $w_i^{\Gamma} = 5030$, latent dimension $n_i^{\Gamma} = 4$, and Swish activation. The number of nonzeros per row and column-shift were both set to 5 for both the interior and interface state decoders. The number of nonzeros for the interior-states masks is 78820, resulting in 99.94% sparsity, while the number of nonzeros for the interface-states masks is 15040, resulting in 99.70% sparsity.

Fig. 6 shows the FOM, LS-ROM, and NM-ROM solutions without HR, and Fig. 8 shows the solutions with collocation HR using 48 DoF in both cases. In both the HR and non-HR cases with the same DoF, NM-ROM achieves an order of magnitude lower relative error than LS-ROM, as evidenced in Figs. 7 and 9 and Table 3. Without HR, NM-ROM achieves a relative error of 1.28×10^{-3} while LS-ROM achieves a relative error 1.98×10^{-2} using the same number of DoF. LS-ROM also achieves a speedup of 30.0, whereas NM-ROM achieves a 21.7 times speedup. With HR, NM-ROM achieves a relative error of 1.64×10^{-3} while LS-ROM achieves a relative error 1.44×10^{-2} using the same number of DoF. In the HR case, LS-ROM achieves a speedup of 347.6, whereas NM-ROM achieves a 43.9 speedup.

Now we compare the performance of LS-ROM and NM-ROM for both the WFPC and SRPC formulations while varying the interior and interface ROM state dimensions n_i^{Ω} and n_i^{Γ} . For WFPC, the decoders \mathbf{g}_i^{Ω} and \mathbf{g}_i^{Γ} use Swish activations, and their sparsity masks have 5 nonzeros per band and a column-shift of 5 for each test. For SRPC, the interior-state decoders \mathbf{g}_i^{Ω} are reused from the WFPC formulation, and the port-state decoders \mathbf{g}_j^{p} were chosen to have Sigmoid activation and sparsity masks with 3 nonzero entries per band with column-shift of 3. We also define \widetilde{n}_i^{p} , which determines the port latent dimensions n_i^{p} via the relation

$$n_{j}^{p} = \max \left\{ \min \left\{ N_{j}^{p} - 1, \ \widetilde{n}_{j}^{p} \right\}, 1 \right\}, \qquad \forall \ j = 1, \dots, n_{p}.$$

This rule was chosen to ensure that $1 \le n_j^p < N_j^p$ because some DD configurations have ports with a very small number of nodes (e.g., $N_i^p = 2$). Recall that for SRPC, the ROM interface-state dimension is $n_i^\Gamma = \sum_{j \in O(i)} n_j^p$.

Table 3 shows the relative error and speedup for LS-ROM and NM-ROM for both WFPC and SRPC on the 2 \times 2 subdomain problem with and without HR while varying n_i^{Ω} , \widetilde{n}_{ij}^{p} , and n_i^{Γ} .

From Table 3, we see that NM-ROM consistently achieves an order of magnitude lower error than LS-ROM both with and without HR when comparing ROMs of the same dimensions and constraint formulations. In the non-HR case with WFPC, LS-ROM only achieves an order 10^{-3} error for a ROM with 96 total DoF (rel. error = 2.66×10^{-3}), while NM-ROM can achieve a similar error with only 36 DoF (rel. error = 2.42×10^{-3}) and a higher speedup (speedup = 26.2) compared to LS-ROM with similar accuracy (speedup = 18.3). For SRPC, LS-ROM was only able to achieve order 10^{-2} accuracy for all cases tested. We also see that LS-ROM achieves a much higher speedup in the HR cases while retaining similar errors from the non-HR cases. NM-ROM also retains high accuracy after HR, and gains an extra $15\text{-}20\times$ speedup after applying HR.

Next we examine the effect of subdomain configuration on the accuracy of LS-ROM and NM-ROM. Again let n_{Ω}^x and n_{Ω}^y denote the number of subdomains in the x- and y-directions, respectively. In each case, the subdomains are of uniform size. For the WFPC cases, we used $(n_i^{\Omega}, n_i^{\Gamma}) = (8, 4)$ for each subdomain, and for the SRPC cases, we used $(n_i^{\Omega}, \widetilde{n}_j^{\rho}) = (8, 2)$ for each subdomain and port, respectively. For the HR cases, we used $N_i^B = 100$ HR nodes. The results for the non-HR and HR cases are reported in Table 4.

Table 4 shows that LS-ROM is more sensitive to subdomain configuration than NM-ROM in the WFPC cases. Indeed, when using 2 subdomains in the *y*-direction, the relative error for LS-ROM increases to the order of 10^{-2} , compared to errors on the order of 10^{-3} when only 1 subdomain is used in the *y*-direction. In contrast, relative error for NM-ROM with WFPC is more stable with respect to subdomain configuration. The relative errors are on the order of 10^{-3} for all configuration considered. For SRPC, LS-ROM only achieves order 10^{-3} relative error for the $(n_{\Omega}^{x}, n_{\Omega}^{y}) = (2, 1)$ configuration, while the remaining configurations have order 10^{-2} relative error. For NM-ROM with SRPC, all configurations have order 10^{-3} relative error except for $(n_{\Omega}^{x}, n_{\Omega}^{y}) = (4, 2)$, which attains

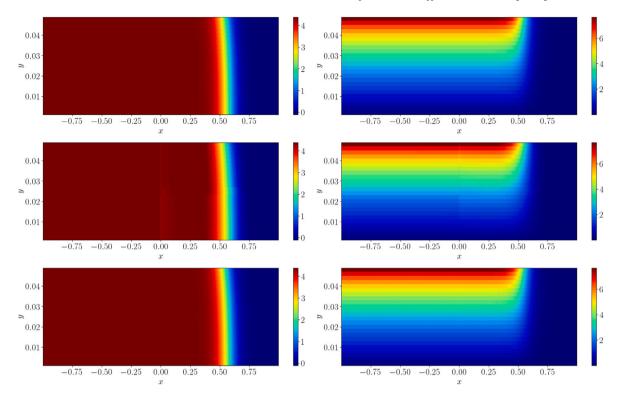


Fig. 6. Top left: u-component of FOM; Top right: v-component of FOM; Middle left: u-component of LS-ROM without HR, WFPC, 48 DoF; Middle right: v-component of LS-ROM without HR, WFPC, 48 DoF; Bottom left: u-component of NM-ROM without HR, WFPC, 48 DoF; Bottom right: v-component of NM-ROM without HR, WFPC, 48 DoF.

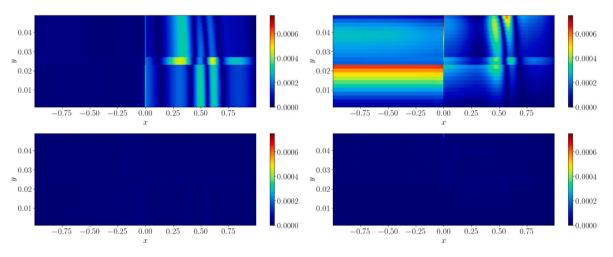


Fig. 7. Top left: u relative error for LS-ROM without HR, WFPC, 48 DoF; Top right: v relative error for LS-ROM without HR, WFPC, 48 DoF; Bottom left: u relative error for RM-ROM without HR, WFPC, 48 DoF; Bottom right: v relative error for NM-ROM without HR, WFPC, 48 DoF.

order 10^{-2} relative error. For both WFPC and SRPC, the NM-ROM error increases slightly as more subdomains are used, but we expect this error can be decreased by adjusting hyper-parameters during ROM training. Hyper-parameter tuning was only done for the 2×2 subdomain configuration.

Tables 3 and 4 show that SRPC has slightly worse performance than WFPC. In particular, Table 3 shows that the relative errors for both LS-ROM and NM-ROM with SRPC do not decrease monotonically as n_i^{Ω} and \tilde{n}_j^{P} increase. In contrast, the relative errors do decrease monotonically as n_i^{Ω} and n_i^{Γ} increase for both LS-ROM and NM-ROM with WFPC. Furthermore, Table 4 shows that LS-ROM with SRPC consistently has larger errors than with WFPC for each subdomain configuration. For NM-ROM, the relative errors and speedups are similar between WFPC and SRPC for each subdomain configuration except for $(n_{\Omega}^{x}, n_{\Omega}^{y}) = (4, 2)$, which has an order

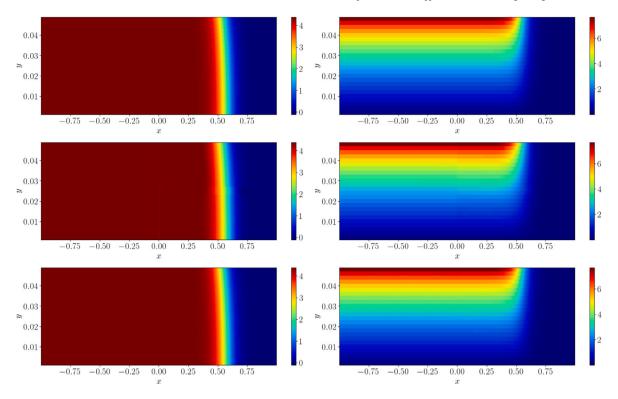


Fig. 8. Top left: *u*-component of FOM; Top right: *v*-component of FOM; Middle left: *u*-component of LS-ROM with collocation HR, WFPC, 48 DoF; Middle right: *v*-component of LS-ROM with collocation HR, WFPC, 48 DoF; Bottom left: *u*-component of NM-ROM with collocation HR, WFPC, 48 DoF; Bottom right: *v*-component of NM-ROM with collocation HR, WFPC, 48 DoF.

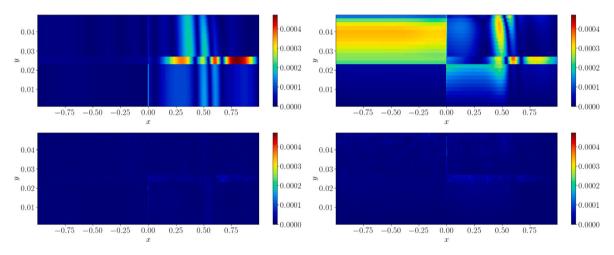


Fig. 9. Top left: *u* relative error for LS-ROM with collocation HR, WFPC, 48 DoF; Top right: *v* relative error for LS-ROM with collocation HR, WFPC, 48 DoF; Bottom left: *u* relative error for NM-ROM with collocation HR, WFPC, 48 DoF; Bottom right: *v* relative error for NM-ROM with collocation HR, WFPC, 48 DoF.

of magnitude higher error than the other configurations. Since SRPC performs as well or worse compared to WFPC for the cases tested, we only consider WFPC in the remainder of this section.

7.7. Pareto fronts

Next we compute Pareto fronts to compare LS-ROM and NM-ROM with WFPC while varying different parameters. The relative error reported is the same as defined in equation (55). Fig. 10 shows the Pareto fronts for varying $(n_i^{\Omega}, n_i^{\Gamma})$ for both the non-HR and

Table 3
Relative error and speedup for LS-ROM and NM-ROM with and without HR applied to the WFPC and SRPC formulations. We use $N_i^B = 100$ HR nodes per subdomain in the HR case.

	Constraints	n_i^{Ω}	\widetilde{n}_{j}^{p}	n_i^{Γ}	Total DoF	Error	Speedup	Error (HR)	Speedup (HR)
LS-ROM	WFPC	4	_	2	24	4.12×10^{-2}	32.1	3.45×10^{-2}	352.7
		6	-	3	36	2.06×10^{-2}	48.7	1.78×10^{-2}	340.0
		8	-	4	48	1.98×10^{-2}	30.0	1.44×10^{-2}	347.6
		10	-	5	60	1.50×10^{-2}	16.3	1.16×10^{-2}	329.6
		16	-	8	96	2.66×10^{-3}	18.3	3.23×10^{-3}	280.4
	SRPC	6	1	5	44	5.17×10^{-2}	24.5	5.12×10^{-2}	315.6
		8	2	7	60	3.75×10^{-2}	22.0	4.22×10^{-2}	313.9
		10	3	9	76	1.87×10^{-2}	19.4	2.94×10^{-2}	262.6
		16	4	11	108	2.00×10^{-2}	17.8	3.37×10^{-2}	253.8
NM-ROM	WFPC	4	_	2	24	6.94×10^{-3}	22.7	7.04×10^{-3}	37.4
		6	-	3	36	2.42×10^{-3}	26.2	2.60×10^{-3}	44.7
		8	-	4	48	1.28×10^{-3}	21.7	1.64×10^{-3}	43.9
		10	-	5	60	1.09×10^{-3}	15.0	1.19×10^{-3}	43.6
		16	-	8	96	7.87×10^{-4}	13.9	9.80×10^{-4}	37.5
	SRPC	6	1	5	44	2.75×10^{-2}	27.6	3.17×10^{-2}	41.1
		8	2	7	60	1.19×10^{-3}	28.8	1.70×10^{-3}	42.6
		10	3	9	76	1.46×10^{-3}	17.0	2.54×10^{-3}	43.1
		16	4	11	108	1.11×10^{-3}	16.8	2.45×10^{-3}	47.1

Table 4Relative error and speedup for LS-ROM and NM-ROM with and without HR and different subdomain configurations for the WFPC and SRPC formulations. We use $n_i^{\Omega} = 8$ for all cases, $n_i^{\Gamma} = 4$ for the WFPC cases, $\tilde{n}_i^{R} = 2$ for the SRPC cases, and $N_i^{B} = 100$ for the HR cases.

	Constraints	n_Q^x	n_{Ω}^{y}	# subdomains	Error	Speedup	Error (HR)	Speedup (HR)
	WFPC $ \begin{array}{ccccccccccccccccccccccccccccccccccc$			2	6.36×10^{-3}	25.1	6.64×10^{-3}	285.7
							1.44×10^{-2}	347.6
		4		4			7.47×10^{-3}	373.1
LS-ROM		35.8	4.21×10^{-2}	259.2				
L3-KOW		2	1	2	6.85×10^{-3}	30.5	30.5 9.49×10^{-3} 22.0 4.22×10^{-2}	293.0
	SRPC	2	2	4	3.75×10^{-2}	22.0	4.22×10^{-2}	313.9
		4	1	4	1.04×10^{-2}	24.6	5.94×10^{-2}	287.5
		4	2	8	4.96×10^{-2}	12.2	5.19×10^{-2}	181.6
NM-ROM	WFPC	2	1	2	1.34×10^{-3}	16.8	1.36×10^{-3}	30.5
		2	2	4	1.28×10^{-3}	21.7	1.64×10^{-3}	43.9
		4	1	4	3.14×10^{-3}	27.8	4.98×10^{-3}	38.6
		4	2	8	4.82×10^{-3}	26.3	5.98×10^{-3}	40.4
	SRPC	2	1	2	1.00×10^{-3}	17.0	1.37×10^{-3}	35.9
		2	2	4	1.19×10^{-3}	28.8	1.70×10^{-3}	42.6
		4	1	4	1.68×10^{-3}	27.4	2.12×10^{-3}	39.1
		4	2	8	1.67×10^{-2}	24.2	2.39×10^{-2}	32.5

HR cases. In the HR case, we use $N_i^B = 100$ for each subdomain. We see that LS-ROM wins in terms of relative wall time, while NM-ROM attains an order of magnitude lower error in comparison to LS-ROM in each case.

Fig. 11 shows the Pareto front for varying number of HR nodes per subdomain, N_i^B . In this case, $(n_i^\Omega, n_i^\Gamma) = (8, 4)$ for each experiment. As in the case for varying (n_i^Ω, n_i^Γ) , NM-ROM attains an order of magnitude lower relative error. Moreover, both the relative error and relative wall time for NM-ROM remains small for each value of N_i^B , whereas the relative error and relative wall time for LS-ROM has more variability with respect to number of HR nodes.

Fig. 12 shows the Pareto fronts for varying number of training snapshots in the non-HR and HR cases. We used $(n_i^{\Omega}, n_i^{\Gamma}) = (8, 4)$ for each experiment, and used $N_i^B = 100$ HR nodes in the HR case. For LS-ROM, the whole training set was used to compute the POD bases, whereas for NM-ROM, the displayed number of snapshots underwent a random 90-10 split for training and validation, respectively. In the case of LS-ROM, the relative error remained constant for the number of training snapshots used, whereas the relative error for NM-ROM decreased as more training snapshots were used.

8. Conclusion

In this work, we detail the first application of NM-ROM with HR to a DD problem. We extend the DD framework of [16] and compute ROMs using the NM-ROM [55] approach on each subdomain. Furthermore, we apply HR to NM-ROM on each subdomain, which informs the use of shallow, sparse autoencoders, as in [56]. We detail how to implement an inexact Lagrange–Newton SQP method to solve the constrained least-squares formulation of the ROM, where the inexactness comes from using a Gauss–Newton approximation of the residual terms, and from neglecting second-order decoder derivative terms coming from the compatibility

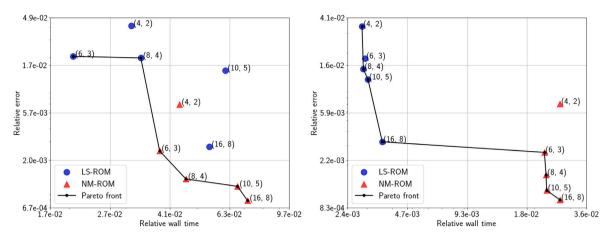


Fig. 10. Left: Pareto front for LS-ROM and NM-ROM without HR with varying $(n_i^{\Omega}, n_i^{\Gamma})$ for WFPC formulation; Right: Pareto front for LS-ROM and NM-ROM with varying $(n_i^{\Omega}, n_i^{\Gamma})$ and $N_i^{B} = 100$ HR nodes per subdomain for WFPC formulation.

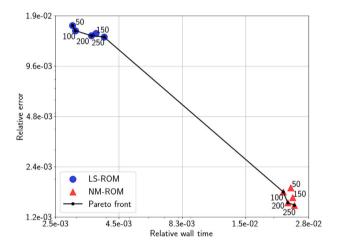


Fig. 11. Pareto front for LS-ROM and NM-ROM with $(n_i^D, n_i^T) = (8, 4)$ and varying number of HR nodes per subdomain N_i^B for WFPC formulation.

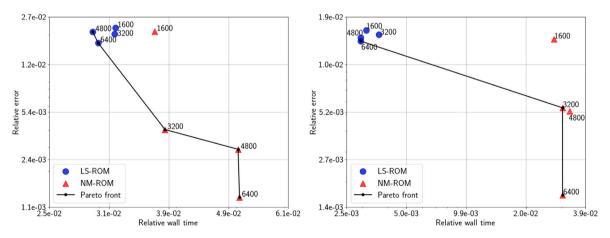


Fig. 12. Left: Pareto front for LS-ROM and NM-ROM with $(n_i^0, n_i^\Gamma) = (8, 4)$ and varying number of training snapshots for WFPC formulation; Right: Pareto front for LS-ROM and NM-ROM with $(n_i^0, n_i^\Gamma) = (8, 4)$, $N_i^B = 100$ HR nodes per subdomain, and varying number of training snapshots for WFPC formulation.

constraints. We then provide a convergence result for the SQP solver used using standard theory of inexact Newton's method. We also provide *a priori* and *a posteriori* error bounds between the FOM and ROM solutions.

From our numerical experiments on the 2D steady-state Burgers' equation, we showed that using the DD NM-ROM approach significantly decreases the number of required NN parameters per subdomain compared to the monolithic single-domain NM-ROM. We also showed that DD NM-ROM achieves an order of magnitude lower relative error than DD LS-ROM in nearly all cases tested. Furthermore, in the non-HR case, NM-ROM is faster than LS-ROM in some instances. We also saw that NM-ROM is more robust than LS-ROM with respect to changes in subdomain configuration. In some cases, the subdomain configuration increased the LS-ROM relative error by an order of magnitude. While LS-ROM with HR achieves much higher speedup than NM-ROM with HR, NM-ROM is still the clear winner in terms of ROM accuracy for a given ROM size. Moreover, HR allows NM-ROM to gain an extra 15-20× speedup compared to the non-HR cases. While the speedup is not as drastic as for LS-ROM, these speedup gains for NM-ROM are the highest that have been achieved for NM-ROM to our knowledge. Our results indicate that NM-ROM should be the preferred approach for problems where ROM accuracy for a given ROM size is more important than speedup.

Although NM-ROM performs better than LS-ROM in our experiments, LS-ROM still attains a relatively low relative error. This indicates that the model problem considered may still be too benign to expose the advantages that NM-ROM has over LS-ROM, particularly when applied to problems with slowly decaying Kolmogorov *n*-width. Therefore, in future work, we plan to apply DD NM-ROM to more complicated problems, including those with slowly decaying Kolmogorov *n*-width, as well as to time-dependent problems. Furthermore, the speedup of the DD NM-ROM is highly dependent on the SQP solver used. Thus, it will be important to investigate the use of other optimization algorithms for the solution of (13) and examine their effects on computational speedup. Other directions for future research include a greedy sampling strategy for the parameter space *D* when choosing which FOM snapshots to compute for NM-ROM training, implementing a "bottom-up" training strategy that uses subdomain snapshots rather than full-domain snapshots for training, applying the DD NM-ROM framework to decomposable or component-based systems, and applying NM-ROM to other DD approaches such as the Schwarz method. Finally error estimates based on the first order necessary optimality conditions, as outlined in the last paragraph of Section 6 is also part of future research.

CRediT authorship contribution statement

Alejandro N. Diaz: Writing – review & editing, Writing – original draft, Visualization, Validation, Software, Methodology, Formal analysis, Data curation, Conceptualization. Youngsoo Choi: Writing – review & editing, Supervision, Resources, Methodology, Investigation, Conceptualization. Matthias Heinkenschloss: Writing – review & editing, Writing – original draft, Supervision, Methodology, Formal analysis.

Declaration of competing interest

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: Alejandro N. Diaz reports financial support was provided by National Defense Science and Engineering Graduate Fellowship. Youngsoo Choi reports financial support was provided by US Department of Energy. Matthias Heinkenschloss reports financial support was provided by Air Force Office of Scientific Research. If there are other authors, they declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

Data will be made available on request.

Acknowledgments

This work was performed at Lawrence Livermore National Laboratory. A. N. Diaz was supported for this work by a Defense Science and Technology Internship (DSTI) at Lawrence Livermore National Laboratory and a 2021 National Defense Science and Engineering Graduate Fellowship, United States. Y. Choi was supported for this work by the US Department of Energy under the Mathematical Multifaceted Integrated Capability Centers – DoE Grant DE – SC0023164; The Center for Hierarchical and Robust Modeling of Non-Equilibrium Transport (CHaRMNET) and partially by LDRD (21-SI-006). M. Heinkenschloss was supported by AFOSR, United States Grant FA9550-22-1-0004 at Rice University. Lawrence Livermore National Laboratory is operated by Lawrence Livermore National Security, LLC, for the U.S. Department of Energy, National Nuclear Security Administration, United States under Contract DE-AC52-07NA27344; IM release number: LLNL-JRNL-849457.

References

- [1] Y. Maday, E.M. Rønquist, A reduced-basis element method, J. Sci. Comput. 17 (1-4) (2002) 447-459, http://dx.doi.org/10.1023/A:1015197908587.
- [2] Y. Maday, E.M. Rønquist, The reduced basis element method: application to a thermal fin problem, SIAM J. Sci. Comput. 26 (1) (2004) 240–258, http://dx.doi.org/10.1137/S1064827502419932.
- [3] L. Iapichino, A. Quarteroni, G. Rozza, A reduced basis hybrid method for the coupling of parametrized domains represented by fluidic networks, Comput. Methods Appl. Mech. Engrg. 221/222 (2012) 63–82, http://dx.doi.org/10.1016/j.cma.2012.02.005.
- [4] P.F. Antonietti, P. Pacciarini, A. Quarteroni, A discontinuous Galerkin reduced basis element method for elliptic problems, ESAIM Math. Model. Numer. Anal. 50 (2) (2016) 337–360, http://dx.doi.org/10.1051/m2an/2015045.

- [5] J.L. Eftang, D.B.P. Huynh, D.J. Knezevic, E.M. Ronquist, A.T. Patera, Adaptive port reduction in static condensation, IFAC Proc. Vol. 45 (2) (2012) 695–699, http://dx.doi.org/10.3182/20120215-3-AT-3016.00123. 7th Vienna International Conference on Mathematical Modelling.
- [6] D.B.P. Huynh, D.J. Knezevic, A.T. Patera, A static condensation reduced basis element method: approximation and *a posteriori* error estimation, ESAIM Math. Model. Numer. Anal. 47 (1) (2013) 213–251, http://dx.doi.org/10.1051/m2an/2012022.
- [7] J.L. Eftang, A.T. Patera, Port reduction in parametrized component static condensation: approximation and *a posteriori* error estimation, Internat. J. Numer. Methods Engrg. 96 (5) (2013) 269–302, http://dx.doi.org/10.1002/nme.4543.
- [8] M. Buffoni, H. Telib, A. Iollo, Iterative methods for model reduction by domain decomposition, Comput. Fluids 38 (6) (2009) 1160–1167, http://dx.doi.org/10.1016/j.compfluid.2008.11.008.
- [9] J. Barnett, I. Tezaur, A. Mota, The Schwarz alternating method for the seamless coupling of nonlinear reduced order models and full order models, 2022, http://dx.doi.org/10.48550/ARXIV.2210.12551, arXiv:2210.12551.
- [10] A. de Castro, P. Bochev, P. Kuberry, I. Tezaur, Explicit synchronous partitioned scheme for coupled reduced order models based on composite reduced bases, Comput. Methods Appl. Mech. Engrg. 417 (2023) 116398, http://dx.doi.org/10.1016/j.cma.2023.116398.
- [11] A. Iollo, G. Sambataro, T. Taddei, A one-shot overlapping Schwarz method for component-based model reduction: application to nonlinear elasticity, Comput. Methods Appl. Mech. Engrg. 404 (2023) 115786, http://dx.doi.org/10.1016/j.cma.2022.115786, 32.
- [12] K. Smetana, T. Taddei, Localized model reduction for nonlinear elliptic partial differential equations: localized training, partition of unity, and adaptive enrichment, SIAM J. Sci. Comput. 45 (3) (2023) A1300–A1331, http://dx.doi.org/10.1137/22M148402X.
- [13] K. Sun, R. Glowinski, M. Heinkenschloss, D.C. Sorensen, Domain decomposition and model reduction of systems with local nonlinearities, in: K. Kunisch, G. Of, O. Steinbach (Eds.), Numerical Mathematics and Advanced Applications, in: ENUMATH 2007, Springer-Verlag, Heidelberg, 2008, pp. 389–396, http://dx.doi.org/10.1007/978-3-540-69777-0 46.
- [14] S. McBane, Y. Choi, Component-wise reduced order model lattice-type structure design, Comput. Methods Appl. Mech. Engrg. 381 (2021) 113813, http://dx.doi.org/10.1016/j.cma.2021.113813, 28.
- [15] S. McBane, Y. Choi, K. Willcox, Stress-constrained topology optimization of lattice-like structures using component-wise reduced order models, Comput. Methods Appl. Mech. Engrg. 400 (2022) 115525, http://dx.doi.org/10.1016/j.cma.2022.115525, 25.
- [16] C. Hoang, Y. Choi, K. Carlberg, Domain-decomposition least-squares Petrov-Galerkin (DD-LSPG) nonlinear model reduction, Comput. Methods Appl. Mech. Engrg. 384 (2021) 113997, http://dx.doi.org/10.1016/j.cma.2021.113997, 41.
- [17] K. Li, K. Tang, T. Wu, Q. Liao, D3M: A deep domain decomposition method for partial differential equations, IEEE Access 8 (2020) 5283–5294, http://dx.doi.org/10.1109/ACCESS.2019.2957200.
- [18] W. Li, X. Xiang, Y. Xu, Deep domain decomposition method: Elliptic problems, in: J. Lu, R. Ward (Eds.), Proceedings of the First Mathematical and Scientific Machine Learning Conference, in: Proceedings of Machine Learning Research, vol. 107, PMLR, 2020, pp. 269–286, URL https://proceedings.mlr.press/v107/li20a.html.
- [19] Q. Sun, X. Xu, H. Yi, Domain decomposition learning methods for solving elliptic problems, 2022, http://dx.doi.org/10.48550/arXiv.2207.10358, arXiv preprint arXiv:2207.10358.
- [20] S. Li, Y. Xia, Y. Liu, Q. Liao, A deep domain decomposition method based on fourier features, J. Comput. Appl. Math. 423 (2023) 114963, http://dx.doi.org/10.1016/j.cam.2022.114963.
- [21] 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, K. Willcox (Eds.), Model Reduction and Approximation: Theory and Algorithms, Computational Science and Engineering, SIAM, Philadelphia, 2017, pp. 65–136, http://dx.doi.org/10.1137/1.9781611974829.ch2.
- [22] A. Quarteroni, A. Manzoni, F. Negri, Reduced Basis Methods for Partial Differential Equations. An Introduction, in: Unitext, vol. 92, Springer, Cham, 2016, http://dx.doi.org/10.1007/978-3-319-15431-2.
- [23] M. Hinze, S. Volkwein, Proper orthogonal decomposition surrogate models for nonlinear dynamical systems: Error estimates and suboptimal control, in: P. Benner, V. Mehrmann, D.C. Sorensen (Eds.), Dimension Reduction of Large-Scale Systems, in: Lecture Notes in Computational Science and Engineering, vol. 45, Springer-Verlag, Heidelberg, 2005, pp. 261–306, http://dx.doi.org/10.1007/3-540-27909-1_10.
- [24] M. Gubisch, S. Volkwein, Chapter 1: Proper orthogonal decomposition for linear-quadratic optimal control, in: P. Benner, A. Cohen, M. Ohlberger, K. Willcox (Eds.), Model Reduction and Approximation: Theory and Algorithms, in: Computational Science and Engineering, SIAM, Philadelphia, 2017, pp. 3–64, http://dx.doi.org/10.1137/1.9781611974829.ch1.
- [25] S.W. Cheung, Y. Choi, D.M. Copeland, K. Huynh, Local lagrangian reduced-order modeling for the Rayleigh-Taylor instability by solution manifold decomposition, J. Comput. Phys. 472 (2023) 111655, http://dx.doi.org/10.1016/j.jcp.2022.111655.
- [26] D.M. Copeland, S.W. Cheung, K. Huynh, Y. Choi, Reduced order models for lagrangian hydrodynamics, Comput. Methods Appl. Mech. Engrg. 388 (2022) 114259, http://dx.doi.org/10.1016/j.cma.2021.114259.
- [27] K. Carlberg, Y. Choi, S. Sargsyan, Conservative model reduction for finite-volume models, J. Comput. Phys. 371 (2018) 280–314, http://dx.doi.org/10. 1016/j.jcp.2018.05.019.
- [28] A.C. Antoulas, Approximation of Large-Scale Dynamical Systems, in: Advances in Design and Control, vol. 6, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2005, http://dx.doi.org/10.1137/1.9780898718713.
- [29] P. Benner, T. Breiten, Chapter 6: Model order reduction based on system balancing, in: P. Benner, A. Cohen, M. Ohlberger, K. Willcox (Eds.), Model Reduction and Approximation: Theory and Algorithms, in: Computational Science and Engineering, SIAM, Philadelphia, 2017, pp. 261–295, http://dx.doi.org/10.1137/1.9781611974829.ch6.
- [30] A.C. Antoulas, C.A. Beattie, S. Gugercin, Interpolatory Model Reduction, in: Computational Science & Engineering, vol. 21, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2020, http://dx.doi.org/10.1137/1.9781611976083.
- [31] C. Gu, QLMOR: A projection-based nonlinear model order reduction approach using quadratic-linear representation of nonlinear systems, IEEE Trans. Comput.-Aided Des. Integr. Circuits Syst. 30 (9) (2011) 1307–1320, http://dx.doi.org/10.1109/TCAD.2011.2142184.
- [32] P. Benner, T. Breiten, Two-sided projection methods for nonlinear model order reduction, SIAM J. Sci. Comput. 37 (2) (2015) B239–B260, http://dx.doi.org/10.1137/14097255X.
- [33] A.J. Mayo, A.C. Antoulas, A framework for the solution of the generalized realization problem, Linear Algebra Appl. 425 (2-3) (2007) 634-662, http://dx.doi.org/10.1016/j.laa.2007.03.008.
- [34] A.C. Antoulas, I.V. Gosea, A.C. Ionita, Model reduction of bilinear systems in the Loewner framework, SIAM J. Sci. Comput. 38 (5) (2016) B889–B916, http://dx.doi.org/10.1137/15M1041432.
- [35] I.V. Gosea, A.C. Antoulas, Data-driven model order reduction of quadratic-bilinear systems, Numer. Linear Algebra Appl. 25 (6) (2018) e2200, http://dx.doi.org/10.1002/nla.2200.
- [36] Y. Choi, P. Brown, W. Arrighi, R. Anderson, K. Huynh, Space-time reduced order model for large-scale linear dynamical systems with application to boltzmann transport problems, J. Comput. Phys. 424 (2021) 109845, http://dx.doi.org/10.1016/j.jcp.2020.109845.
- [37] Y. Kim, K. Wang, Y. Choi, Efficient space–time reduced order model for linear dynamical systems in python using less than 120 lines of code, Mathematics 9 (14) (2021) 1690, http://dx.doi.org/10.3390/math9141690.
- [38] Y. Choi, K. Carlberg, Space-time least-squares Petrov-Galerkin projection for nonlinear model reduction, SIAM J. Sci. Comput. 41 (1) (2019) A26-A58, http://dx.doi.org/10.1137/17M1120531.

- [39] M. Ohlberger, S. Rave, Reduced basis methods: Success, limitations and future challenges, Proc. Conf. Algoritmy (2016) 1–12, URL http://www.iam.fmph.uniba.sk/amuc/ois/index.php/algoritmy/article/view/389.
- [40] N.J. Nair, M. Balajewicz, Transported snapshot model order reduction approach for parametric, steady-state fluid flows containing parameter-dependent shocks, Internat. J. Numer. Methods Engrg. 117 (12) (2019) 1234–1262, http://dx.doi.org/10.1002/nme.5998.
- [41] A. Iollo, D. Lombardi, Advection modes by optimal mass transfer, Phys. Rev. E 89 (2014) 022923, http://dx.doi.org/10.1103/PhysRevE.89.022923.
- [42] N. Cagniart, Y. Maday, B. Stamm, Model order reduction for problems with large convection effects, in: B.N. Chetverushkin, W. Fitzgibbon, Y.A. Kuznetsov, P. Neittaanmäki, J. Periaux, O. Pironneau (Eds.), Contributions to Partial Differential Equations and Applications, in: Comput. Methods Appl. Sci., vol. 47, Springer, Cham, 2019, pp. 131–150, http://dx.doi.org/10.1007/978-3-319-78325-3_10.
- [43] J. Reiss, P. Schulze, J. Sesterhenn, V. Mehrmann, The shifted proper orthogonal decomposition: a mode decomposition for multiple transport phenomena, SIAM J. Sci. Comput. 40 (3) (2018) A1322–A1344, http://dx.doi.org/10.1137/17M1140571.
- [44] G. Welper, Interpolation of functions with parameter dependent jumps by transformed snapshots, SIAM J. Sci. Comput. 39 (4) (2017) A1225–A1250, http://dx.doi.org/10.1137/16M1059904.
- [45] R. Mojgani, M. Balajewicz, Lagrangian basis method for dimensionality reduction of convection dominated nonlinear flows, 2017, http://dx.doi.org/10.48550/arXiv.1701.04343, arXiv:1701.04343v1.
- [46] M. Dihlmann, M. Drohmann, B. Haasdonk, Model reduction of parametrized evolution problems using the reduced basis method with adaptive time-partitioning, in: International Conference on Adaptive Modeling and Simulation, ADMOS 2011, 2011, p. 64.
- [47] M. Drohmann, B. Haasdonk, M. Ohlberger, Adaptive reduced basis methods for nonlinear convection–diffusion equations, in: J. Fort, J. Fürst, J. Halama, R. Herbin, F. Hubert (Eds.), Finite Volumes for Complex Applications VI Problems & Perspectives: FVCA 6, International Symposium, Prague, June 6–10, 2011, Springer, Berlin, Heidelberg, 2011, pp. 369–377, http://dx.doi.org/10.1007/978-3-642-20671-9_39.
- [48] T. Taddei, S. Perotto, A. Quarteroni, Reduced basis techniques for nonlinear conservation laws, ESAIM Math. Model. Numer. Anal. 49 (3) (2015) 787–814, http://dx.doi.org/10.1051/m2an/2014054.
- [49] B. Peherstorfer, D. Butnaru, K. Willcox, Online Adaptive Model Reduction for Nonlinear Systems via Low-Rank Updates, SIAM J. Sci. Comput. 37 (4) (2015) A2123–A2150, http://dx.doi.org/10.1137/140989169.
- [50] D. Amsallem, M.J. Zahr, C. Farhat, Nonlinear model order reduction based on local reduced-order bases, Internat. J. Numer. Methods Engrg. 92 (10) (2012) 891–916, http://dx.doi.org/10.1002/nme.4371.
- [51] J. Barnett, C. Farhat, Quadratic approximation manifold for mitigating the Kolmogorov barrier in nonlinear projection-based model order reduction, J. Comput. Phys. 464 (2022) 111348, http://dx.doi.org/10.1016/j.jcp.2022.111348.
- [52] R. Geelen, S. Wright, K. Willcox, Operator inference for non-intrusive model reduction with nonlinear manifolds, 2022, arXiv:2205.02304v1.
- [53] K. Kashima, Nonlinear model reduction by deep autoencoder of noise response data, in: 2016 IEEE 55th Conference on Decision and Control, CDC, 2016, pp. 5750–5755, http://dx.doi.org/10.1109/CDC.2016.7799153.
- [54] D. Hartman, L.K. Mestha, A deep learning framework for model reduction of dynamical systems, in: 2017 IEEE Conference on Control Technology and Applications, CCTA, 2017, pp. 1917–1922, http://dx.doi.org/10.1109/CCTA.2017.8062736.
- [55] K. Lee, K.T. Carlberg, Model reduction of dynamical systems on nonlinear manifolds using deep convolutional autoencoders, J. Comput. Phys. 404 (2020) 108973, http://dx.doi.org/10.1016/j.jcp.2019.108973, 32.
- [56] Y. Kim, Y. Choi, D. Widemann, T. Zohdi, A fast and accurate physics-informed neural network reduced order model with shallow masked autoencoder, J. Comput. Phys. 451 (2022) 110841, http://dx.doi.org/10.1016/j.jcp.2021.110841, 29.
- [57] Y. Kim, Y. Choi, D. Widemann, T. Zohdi, Efficient nonlinear manifold reduced order model, 2020, http://dx.doi.org/10.48550/arXiv.2011.07727, arXiv preprint arXiv:2011.07727.
- [58] F. Romor, G. Stabile, G. Rozza, Non-linear manifold reduced-order models with convolutional autoencoders and reduced over-collocation method, J. Sci. Comput. 94 (3) (2023) 74, http://dx.doi.org/10.1007/s10915-023-02128-2.
- [59] T. Taddei, X. Xu, L. Zhang, A non-overlapping optimization-based domain decomposition approach to component-based model reduction of incompressible flows, 2023, http://dx.doi.org/10.48550/arXiv.2310.20267, arXiv.
- [60] C. Farhat, S. Grimberg, A. Manzoni, A. Quarteroni, Computational bottlenecks for PROMs: precomputation and hyperreduction, in: P. Benner, S. Grivet-Talocia, A. Quarteroni, G. Rozza, W. Schilders, L.M. Silveira (Eds.), Model Order Reduction, in: Snapshot-Based Methods and Algorithms, vol. 2, Walter de Gruyter & Co., Berlin, 2021, pp. 181–243, http://dx.doi.org/10.1515/9783110671490-005.
- [61] R. Everson, L. Sirovich, The Karhunen-Loéve procedure for gappy data, J. Opt. Soc. Amer. 12 (8) (1995) 1657-1664.
- [62] Y. Choi, D. Coombs, R. Anderson, SNS: A solution-based nonlinear subspace method for time-dependent model order reduction, SIAM J. Sci. Comput. 42 (2) (2020) A1116-A1146.
- [63] J.T. Lauzon, S.W. Cheung, Y. Shin, Y. Choi, D.M. Copeland, K. Huynh, S-OPT: A points selection algorithm for hyper-reduction in reduced order models, 2022, arXiv preprint arXiv:2203.16494.
- [64] P.T. Boggs, J.W. Tolle, Sequential quadratic programming, in: Acta Numerica, 1995, Cambridge Univ. Press, Cambridge, 1995, pp. 1–51, http://dx.doi.org/10.1017/s0962492900002518.
- [65] J. Nocedal, S.J. Wright, Numerical Optimization, second ed., Springer Verlag, Berlin, Heidelberg, New York, 2006, http://dx.doi.org/10.1007/978-0-387-40065-5
- [66] M. Benzi, G.H. Golub, J. Liesen, Numerical solution of saddle point problems, Acta Numer. 14 (2005) 1–137, http://dx.doi.org/10.1017/ S0962492904000212.
- [67] H.G. Bock, E. Kostina, J.P. Schlöder, Numerical methods for parameter estimation in nonlinear differential algebraic equations, GAMM-Mitt. 30 (2007) 376–408, http://dx.doi.org/10.1002/gamm.200790024.
- [68] M. Heinkenschloss, Mesh independence for nonlinear least squares problems with norm constraints, SIAM J. Optim. 3 (1993) 81–117, http://dx.doi.org/ 10.1137/0803005.
- [69] G. Cybenko, Approximation by superpositions of a sigmoidal function, Math. Control Signals Systems 2 (4) (1989) 303–314, http://dx.doi.org/10.1007/BF02551274.
- [70] A. Pinkus, Approximation theory of the MLP model in neural networks, in: Acta Numerica, 1999, in: Acta Numer., vol. 8, Cambridge Univ. Press, Cambridge, 1999, pp. 143–195, http://dx.doi.org/10.1017/S0962492900002919.
- [71] K.T. Carlberg, C. Farhat, J. Cortial, D. Amsallem, The GNAT method for nonlinear model reduction: Effective implementation and application to computational fluid dynamics and turbulent flows, J. Comput. Phys. 242 (2013) 623–647, http://dx.doi.org/10.1016/j.jcp.2013.02.028.
- [72] A. Schmidt, D. Wittwar, B. Haasdonk, Rigorous and effective a-posteriori error bounds for nonlinear problems—application to RB methods, Adv. Comput. Math. 46 (2) (2020) 32, http://dx.doi.org/10.1007/s10444-020-09741-x, 30.