

# Autoencoding Reduced Order Models for Control through the Lens of Dynamic Mode Decomposition

Anonymous authors

Paper under double-blind review

## Abstract

Modeling and control of high-dimensional dynamical systems often involve some dimensionality reduction techniques to construct a lower-order model that makes the associated task computationally feasible or less demanding. In recent years, two techniques have become widely popular for analysis and [reduced order modeling](#) of high-dimensional dynamical systems: (1) dynamic mode decomposition and (2) deep autoencoding learning. This paper establishes a connection between dynamic mode decomposition and autoencoding learning for controlled dynamical systems. Specifically, we first show that an optimization objective for learning a linear autoencoding reduced order model can be formulated such that its solution closely resembles the solution obtained by the *dynamic mode decomposition with control* algorithm. The linear autoencoding architecture is then extended to a deep autoencoding architecture to learn a nonlinear reduced order model. Finally, the learned reduced order model is used to design a controller utilizing stability-constrained deep neural networks. The studied framework does not require knowledge of the governing equations of the underlying system and learns the model and controller solely from time series data of observations and actuations. We provide empirical analysis on modeling and control of spatiotemporal high-dimensional systems, including fluid flow control.

## 1 Introduction

Designing controllers for high-dimensional dynamical systems remains a challenge as many control algorithms become computationally prohibitive in high dimensions. Typically, a *reduce-then-design* approach (Atwell et al. (2001)) is used in practice, which involves two steps: (1) develop a reduced order model (ROM) using dimensionality reduction techniques and (2) design a controller for that reduced order model (Figure 1a). For controlled dynamical systems, the [reduced order modeling](#) approaches either combine analytical techniques with empirical approximation (Willcox & Peraire (2002)) or are purely data-driven (Juang & Pappa (1985); Juang et al. (1993); Proctor et al. (2016)). Among these, the dynamic mode decomposition (DMD) based methods have become widely popular in recent years due to a strong connection between DMD and Koopman operator theory (Rowley et al. (2009)). Another recent research trend involves utilizing deep neural networks (DNNs), particularly autoencoders, for modeling and control of high-dimensional dynamical systems (Lusch et al. (2018); Erichson et al. (2019); Eivazi et al. (2020); Morton et al. (2018); Bounou et al. (2021); Chen et al. (2021); Ping et al. (2021)).

In this paper, we provide a perspective connecting DMD and autoencoding reduced order models for controlled dynamical systems and present a framework to learn control policies for such systems by means of the DNN-based reduced order models. We first formulate an objective function for data-driven learning of controlled dynamical systems in a linear autoencoding configuration. We analytically show that the associated objective function encourages a linear ROM that closely resembles the lower-order model obtained using the *dynamic mode decomposition with control* (DMDc) algorithm (Proctor et al. (2016)). The linear autoencoding architecture is designed in such a way that its components can be replaced with DNNs and the corresponding objective function can be optimized by gradient descent to obtain a nonlinear ROM. The architecture with DNN components, DeepROM, closely resembles the aforementioned deep autoencoding architectures used in recent literature for the prediction and control of dynamical systems. The learned

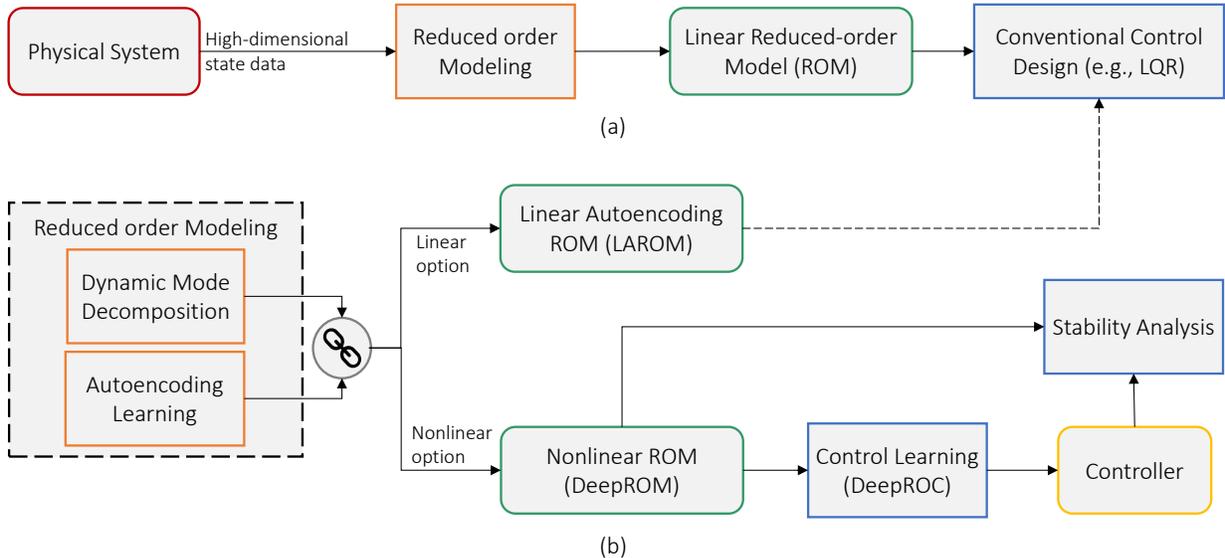


Figure 1: (a) Reduce-then-design paradigm for designing control for high-dimensional systems, (b): our work in the context. The rounded-corner rectangles denote the actual physical system (red outlined) or its models (green outlined) or a controller (yellow outlined). The sharp-corner rectangles indicate the techniques/procedure to obtain the models or controllers. Among those sharp-corner rectangles, the ones outlined in orange are associated with modeling, while the ones outlined in blue are associated with control. This work proposes an autoencoding learning framework that establishes a *link* with dynamic mode decomposition. The similarity with dynamic mode decomposition for control is shown through a linear autoencoding reduced order model while the prediction and control performance is evaluated using a deep autoencoding reduced order model. The dashed arrow in the figure represents the potential control methods that can be applied to the linear ROM.

DNN-based reduced order model is then used in a control learning framework, deep reduced order control (DeepROC), to design a controller for the underlying system. The control policy is learned by jointly training two DNNs: one stability-constrained DNN predicts a target closed-loop dynamics for the learned ROM while the other DNN serves as a controller to achieve that target dynamics. We analytically show that keeping the joint learning objective within a sufficiently small value implies stability for the closed-loop ROM in terms of *ultimate boundedness*, i.e., trajectories starting close to the desired state stay close to the desired state. The overall workflow of this paper is shown in Figure 1(b). We provide empirical analysis using examples of spatiotemporal physical processes such as reaction–diffusion and fluid flow systems. In summary, our contributions are as follows:

- For controlled dynamical systems, we show that an objective function can be formulated in a linear autoencoding configuration and optimized by gradient descent such that the corresponding linear ROM closely resembles the ROM obtained using the DMDc algorithm.
- We extend the linear autoencoding configuration to a deep autoencoding configuration to learn a DNN-based nonlinear ROM.
- We analytically show that a DNN controller can be trained such that the closed-loop trajectories of the learned ROM remain ultimately bounded.
- We empirically show the similarity of the linear autoencoding reduced order model (LAROM) with DMDc and evaluate the prediction performance of DeepROM and control performance of DeepROC in experiments with high-dimensional systems, including suppressing vortex shedding in fluid flow.

## 2 Related Work

In recent years, deep learning has seen widespread application in scientific and engineering problems, including understanding complex dynamics of large-scale or high-dimensional systems and solving associated computational tasks. The majority of the research in this area focuses on the modeling and prediction of such complex dynamics using deep neural networks (DNNs) (Xingjian et al. (2015); Long et al. (2018); Raissi (2018); Seo et al. (2019); Ayed et al. (2019); Donà et al. (2020)) and has found application in several fields including fluid flow (Erichson et al. (2019); Eivazi et al. (2020); Srinivasan et al. (2019)), biochemical and electric power systems (Yeung et al. (2019)), climate and ocean systems (Scher (2018); Ren et al. (2021); Yang et al. (2017); De Bézenac et al. (2019)), and structural analysis Zhang et al. (2020), just to name a few.

A second line of research, though relatively less prevalent than modeling and prediction, is utilizing [deep learning](#) for controlling high-dimensional systems [and aligns closely with our work](#). [Most of these works focus on fluid flow control tasks](#). Rabault et al. (2019); Tang et al. (2020) applied deep reinforcement learning in active flow control for vortex shedding suppression and used a system-specific reward function involving lift and drag. Ma et al. (2018) used an autoencoder for encoding high-dimensional fluid state to low-dimensional features and trained RL agents with those features to control rigid bodies in a coupled 2D system involving both fluid and rigid bodies. Garnier et al. (2021) also used an autoencoder for feature extraction to train an RL agent for controlling the position of a small cylinder in a two-cylinder fluid system to reduce the total drag. Beintema et al. (2020) used deep reinforcement learning with system-specific rewards to control a Rayleigh–Bénard system with the aim of reducing convective effects. Model-free deep reinforcement learning methods have high sample complexity necessitating a large number of interactions with the environment and often require system-specific reward construction. [In contrast, we consider learning the control policies offline with limited pre-collected data](#). [Moreover, we utilize standard distance-based metrics with respect to a hypothesized target dynamics instead of relying on any system-specific rewards or loss functions](#). [Model-free RL methods require running numerical solvers in every iteration to provide feedback to the agents, which is computationally expensive](#). The same concern arises for the methods involving differentiable simulators. For example, Holl et al. (2020) used a differentiable partial differential equation (PDE) solver to generate gradient-based feedback for a predictor-corrector scheme to control PDE-driven high-dimensional systems. Takahashi et al. (2021) too integrated a differentiable simulator with DNNs to learn control in coupled solid-fluid systems. [In comparison, our method avoids the need for computationally intensive simulators during the learning process as it learns from pre-collected data in an offline manner](#).

The alternative to model-free methods takes the traditional approach: develop a model first and then use that to design controllers. Deep learning is now being used in this process by developing frameworks like DeepMPC (Lenz et al. (2015)) which incorporates DNN features in model predictive control (MPC). Bieker et al. (2020) and Morton et al. (2018) utilized the DeepMPC framework for fluid flow control. Bieker et al. (2020) used a recurrent neural network to model the dynamics of only control-relevant quantities (i.e. lift and drag) of the system, which is then employed in an MPC framework for the flow control tasks. Morton et al. (2018) followed the method proposed by Takeishi et al. (2017) and used DNN-based embedding to first learn a linear reduced order model in Koopman invariant subspace and then incorporate it in the MPC framework. Similar approaches have been applied to other high-dimensional control tasks like control from video input (Bounou et al. (2021)), automatic generation control in wind farms in the presence of dynamic wake effect (Chen et al. (2021)), and transient stabilization in power grids (Ping et al. (2021)). [These model-based methods constrain the latent dynamic models to be linear that work well within a short time window](#). Khodkar et al. (2019) showed that the linear combination of a finite number of dynamic modes may not accurately represent the long-term nonlinear characteristics of complex dynamics and adding nonlinear forcing terms yields better prediction accuracy (Eivazi et al. (2020)). [The linear ROMs are needed to be updated with online observations during operation for better prediction accuracy](#). Accordingly, the aforementioned model-based approaches utilize the MPC framework to optimize the control policy online using the updated dynamic model. Running online optimization at every step may not be computationally feasible in many scenarios. Conversely, we investigate if a nonlinear ROM provides a more accurate prediction over a longer time window so that an offline control learning method can be used.

### 3 Problem and Preliminaries

#### 3.1 Problem statement

Consider a time-invariant controlled dynamical system

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

where  $\mathbf{x}(t) \in \mathbb{X} \subset \mathbb{R}^{d_x}$ ,  $d_x \gg 1$  and  $\mathbf{u}(t) \in \mathbb{U} \subset \mathbb{R}^{d_u}$  are the system state and the actuation (or control input), respectively, at time  $t$ . Our objective is to learn a feedback controller  $\mathbf{u} = \pi(\mathbf{x})$  for this high-dimensional ( $d_x \gg 1$ ) system of (1) to stabilize it at a desired state in a data-driven reduce-then-design approach when the nonlinear function  $f : \mathbb{X} \times \mathbb{U} \rightarrow \mathbb{R}^{d_x}$  is unknown. We assume that we have observations from the system consisting of time series data  $\mathbf{x}(t_i)$ ,  $i = 0, 1, \dots, n$  subjected to random values of actuations  $\mathbf{u}(t_i)$ ,  $i = 0, 1, \dots, (n-1)$ .

Note, we use  $v$  (in place of  $v(t)$  for brevity) as notation for any continuous-time variable (e.g., system state, control input), whereas  $v(t_i)$  is used to denote their discrete sample at time instance  $t_i$ .

#### 3.2 Stabilization of controlled systems

We assume that the system we are aiming to stabilize at an equilibrium point is *locally stabilizable*. Suppose the function  $f$  in (1) is locally Lipschitz and  $(\mathbf{x} = \mathbf{0}, \mathbf{u} = \mathbf{0})$  is an equilibrium pair of the system, i.e.,  $f(\mathbf{0}, \mathbf{0}) = \mathbf{0}$ . The system (1) is said to be *locally stabilizable* with respect to the equilibrium pair if there exists a locally Lipschitz function

$$\pi : \mathbb{X}_0 \rightarrow \mathbb{U}, \quad \pi(\mathbf{0}) = \mathbf{0},$$

defined on some neighborhood  $\mathbb{X}_0 \subset \mathbb{X}$  of the origin  $\mathbf{x} = \mathbf{0}$  for which the closed-loop system

$$\frac{d\mathbf{x}}{dt} = f(\mathbf{x}, \pi(\mathbf{x})) \quad (2)$$

is locally *asymptotically stable*, i.e.  $\|\mathbf{x}(t_0)\| < \delta$  implies  $\lim_{t \rightarrow \infty} \mathbf{x}(t) = \mathbf{0}$  (Sontag (2013)). [We discuss the criteria for stability and asymptotic stability in the following paragraph.](#)

Stability of the closed-loop system  $\frac{d\mathbf{x}}{dt} = f(\mathbf{x}, \pi(\mathbf{x})) = h(\mathbf{x})$  at equilibrium points can be analyzed using the method of Lyapunov. Let  $\mathcal{V} : \mathbb{X} \rightarrow \mathbb{R}$  be a continuously differentiable function such that

$$\mathcal{V}(\mathbf{0}) = 0, \quad \text{and} \quad \mathcal{V}(\mathbf{x}) > 0 \quad \forall \mathbf{x} \in \mathbb{X} \setminus \{\mathbf{0}\}, \quad (3)$$

and the time derivative of  $\mathcal{V}$  along the trajectories

$$\frac{d\mathcal{V}}{dt} = \nabla \mathcal{V}(\mathbf{x})^\top \frac{d\mathbf{x}}{dt} = \nabla \mathcal{V}(\mathbf{x})^\top h(\mathbf{x}) \leq 0 \quad \forall \mathbf{x} \in \mathbb{X}. \quad (4)$$

Then, the equilibrium point  $\mathbf{x} = \mathbf{0}$  is *stable*, i.e., for each  $\epsilon > 0$ , there exists a  $\delta = \delta(\epsilon) > 0$  such that  $\|\mathbf{x}(t_0)\| < \delta$  implies  $\|\mathbf{x}(t)\| < \epsilon$ ,  $\forall t > t_0$ . The function  $\mathcal{V}$  with the above properties is called a Lyapunov function. If  $\frac{d\mathcal{V}}{dt} < 0$  in [some subset](#)  $\mathbb{X}_s \subset \mathbb{X} \setminus \{\mathbf{0}\}$ , then  $\mathbf{x} = \mathbf{0}$  is *locally asymptotically stable*. Moreover, if there exist positive constants  $c_1, c_2, c_3$  and  $c_4$  such that

$$c_1 \|\mathbf{x}\|^2 \leq \mathcal{V}(\mathbf{x}) \leq c_2 \|\mathbf{x}\|^2, \quad (5)$$

and

$$\nabla \mathcal{V}(\mathbf{x})^\top h(\mathbf{x}) \leq -c_3 \|\mathbf{x}\|^2, \quad \forall \mathbf{x} \in \mathbb{X}_s, \quad (6)$$

then  $\mathbf{x} = \mathbf{0}$  is *exponentially stable*, i.e., there exist positive constants  $\delta, \lambda$  and  $\gamma$  such that  $\|\mathbf{x}(t)\| \leq \lambda \|\mathbf{x}(t_0)\| e^{-\gamma(t-t_0)}$ ,  $\forall \|\mathbf{x}(t_0)\| < \delta$  (Khalil (2002)).

In this paper, we assume that the system we are aiming to stabilize at an equilibrium point is stabilizable in the sense of the aforementioned definition and criteria, i.e., there exists a continuously differentiable function

$\mathcal{V}$  and a Lipschitz continuous control law  $\pi$  such that criteria (3) and (4) are conformed. To ensure the stability of the closed-loop reduced order model at equilibrium, we utilize a target dynamics hypothesis that is exponentially stable at the origin, i.e., (5) and (6) are satisfied as well for this target dynamics. Detail on the target dynamics hypothesis is discussed in subsection 4.2.

Though the above formulation is for stabilization at an equilibrium point  $\mathbf{x} = \mathbf{0}$ , the same can be used to stabilize the system at any arbitrary point  $\mathbf{x}_{ss}$ . In that case, a steady-state control input  $\mathbf{u}_{ss}$  is required that can maintain the equilibrium at  $\mathbf{x}_{ss}$ , i.e.,  $f(\mathbf{x}_{ss}, \mathbf{u}_{ss}) = \mathbf{0}$ . The change of variables  $\mathbf{x}_e = \mathbf{x} - \mathbf{x}_{ss}$ ,  $\mathbf{u}_e = \mathbf{u} - \mathbf{u}_{ss}$  leads to a transformed system where we can apply the aforementioned formulation of stabilization. The overall control, in this case,  $\mathbf{u} = \mathbf{u}_e + \mathbf{u}_{ss}$  comprises a feedback component  $\mathbf{u}_e$  and a feedforward component  $\mathbf{u}_{ss}$  (Khalil (2002)).

### 3.3 Dynamic mode decomposition with control

DMD (Schmid (2010)) is a data-driven method that reconstructs the underlying dynamics using only a time series of snapshots from the system. DMD computes a modal decomposition where each mode is associated with an oscillation frequency and decay/growth rate. DMD has become a widely used technique for spectral analysis of dynamical systems. DMDc (Proctor et al. (2016)) is an extension of DMD for dynamical systems with control. DMDc seeks best-fit linear operators  $\mathbf{A}$  and  $\mathbf{B}$  between successive observed states and the actuations:

$$\hat{\mathbf{x}}(t_{i+1}) = \mathbf{A}\mathbf{x}(t_i) + \mathbf{B}\mathbf{u}(t_i), \quad i = 0, 1, \dots, n-1, \quad (7)$$

where  $\hat{\mathbf{x}}(t)$  denotes an approximation of  $\mathbf{x}(t)$ ,  $\mathbf{A} \in \mathbb{R}^{d_x \times d_x}$ , and  $\mathbf{B} \in \mathbb{R}^{d_x \times d_u}$ . Direct analysis of (7) could be computationally prohibitive for  $d_x \gg 1$ . DMDc leverages dimensionality reduction to compute a ROM

$$\mathbf{x}_{\text{R,DMDc}}(t_i) = \mathbf{E}_{\text{DMDc}}\mathbf{x}(t_i), \quad (8a)$$

$$\mathbf{x}_{\text{R,DMDc}}(t_{i+1}) = \mathbf{A}_{\text{R,DMDc}}\mathbf{x}_{\text{R,DMDc}}(t_i) + \mathbf{B}_{\text{R,DMDc}}\mathbf{u}(t_i), \quad i = 0, 1, \dots, n-1, \quad (8b)$$

which retains the dominant dynamic modes of (7). Here,  $\mathbf{x}_{\text{R,DMDc}}(t_i) \in \mathbb{R}^{r_x}$  is the reduced state, where  $r_x \ll d_x$ , and  $\mathbf{E}_{\text{DMDc}} \in \mathbb{R}^{r_x \times d_x}$ ,  $\mathbf{A}_{\text{R,DMDc}} \in \mathbb{R}^{r_x \times r_x}$ ,  $\mathbf{B}_{\text{R,DMDc}} \in \mathbb{R}^{r_x \times d_u}$ . The full state is reconstructed from the reduced state using the transformation  $\hat{\mathbf{x}}(t_i) = \mathbf{D}_{\text{DMDc}}\mathbf{x}_{\text{R,DMDc}}(t_i)$ , where  $\mathbf{D}_{\text{DMDc}} \in \mathbb{R}^{d_x \times r_x}$ . DMDc computes truncated singular value decomposition (SVD) of the data matrices  $\mathbf{Y} = [\mathbf{x}(t_1), \mathbf{x}(t_2), \dots, \mathbf{x}(t_n)] \in \mathbb{R}^{d_x \times n}$  and  $\mathbf{\Omega} = [\boldsymbol{\omega}(t_0), \boldsymbol{\omega}(t_1), \dots, \boldsymbol{\omega}(t_{n-1})] \in \mathbb{R}^{(d_x+d_u) \times n}$ ,  $\boldsymbol{\omega}(t_i) = [\mathbf{x}(t_i)^\top, \mathbf{u}(t_i)^\top]^\top \in \mathbb{R}^{d_x+d_u}$  as follows:

$$\mathbf{Y} = \hat{\mathbf{U}}_{\mathbf{Y}} \hat{\boldsymbol{\Sigma}}_{\mathbf{Y}} \hat{\mathbf{V}}_{\mathbf{Y}}^\top, \quad \mathbf{\Omega} = \hat{\mathbf{U}}_{\boldsymbol{\Omega}} \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\Omega}} \hat{\mathbf{V}}_{\boldsymbol{\Omega}}^\top, \quad (9)$$

where  $\hat{\mathbf{U}}_{\mathbf{Y}} \in \mathbb{R}^{d_x \times r_x}$ ,  $\hat{\boldsymbol{\Sigma}}_{\mathbf{Y}} \in \mathbb{R}^{r_x \times r_x}$ ,  $\hat{\mathbf{V}}_{\mathbf{Y}} \in \mathbb{R}^{n \times r_x}$ ,  $\hat{\mathbf{U}}_{\boldsymbol{\Omega}} \in \mathbb{R}^{(d_x+d_u) \times r_{xu}}$ ,  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\Omega}} \in \mathbb{R}^{r_{xu} \times r_{xu}}$ , and  $\hat{\mathbf{V}}_{\boldsymbol{\Omega}} \in \mathbb{R}^{n \times r_{xu}}$ .  $r_x < \min(d_x, n)$  and  $r_{xu} < \min(d_x + d_u, n)$  denote the truncation dimensions of SVDs. Utilizing the SVDs of (9) the parameters of the ROM (8) is obtained as

$$\mathbf{E}_{\text{DMDc}} = \hat{\mathbf{U}}_{\mathbf{Y}}^\top, \quad \mathbf{D}_{\text{DMDc}} = \hat{\mathbf{U}}_{\mathbf{Y}}, \quad (10a)$$

$$\mathbf{A}_{\text{R,DMDc}} = \hat{\mathbf{U}}_{\mathbf{Y}}^\top \mathbf{Y} \hat{\mathbf{V}}_{\boldsymbol{\Omega}} \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\Omega}}^{-1} \hat{\mathbf{U}}_{\boldsymbol{\Omega},1}^\top \hat{\mathbf{U}}_{\mathbf{Y}}, \quad \mathbf{B}_{\text{R,DMDc}} = \hat{\mathbf{U}}_{\mathbf{Y}}^\top \mathbf{Y} \hat{\mathbf{V}}_{\boldsymbol{\Omega}} \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\Omega}}^{-1} \hat{\mathbf{U}}_{\boldsymbol{\Omega},2}^\top, \quad (10b)$$

where  $\hat{\mathbf{U}}_{\boldsymbol{\Omega},1} \in \mathbb{R}^{d_x \times r_{xu}}$ ,  $\hat{\mathbf{U}}_{\boldsymbol{\Omega},2} \in \mathbb{R}^{d_u \times r_{xu}}$ , and  $\hat{\mathbf{U}}_{\boldsymbol{\Omega}}^\top = [\hat{\mathbf{U}}_{\boldsymbol{\Omega},1}^\top \quad \hat{\mathbf{U}}_{\boldsymbol{\Omega},2}^\top]$ .

## 4 Method

As mentioned earlier, in the reduce-then-design approach, we first need to develop a ROM and then design a controller using that ROM. A controller designed for the ROM is expected to perform well in the full system only if the ROM effectively captures the dynamic characteristics of the underlying system. In this section, we first describe how to design a ROM that effectively captures the relation between successive observations and actuation. Next, we delineate the process for learning controllers utilizing the learned ROM.

## 4.1 Learning a reduced order model

DMDc can extract the dominant modes of underlying dynamics in a reduced order model (Proctor et al. (2016)). In order to develop a nonlinear ROM utilizing DNNs that effectively capture the underlying dynamics, we first investigate if we can obtain a linear ROM similar to DMDc, in a gradient descent arrangement. Specifically, we analyze optimization objectives that encourage a DMDc-like solution for a [reduced order modeling](#) problem using linear networks (single layer without nonlinear activation). Consider the following [reduced order modeling](#) problem

$$\mathbf{x}_R(t_i) = \mathbf{E}_x \mathbf{x}(t_i), \quad \mathbf{x}_R(t_{i+1}) = \mathbf{A}_R \mathbf{x}_R(t_i) + \mathbf{B}_R \mathbf{u}(t_i), \quad \hat{\mathbf{x}}(t_i) = \mathbf{D}_x \mathbf{x}_R(t_i), \quad i = 0, 1, \dots, n-1, \quad (11)$$

where the linear operators  $\mathbf{E}_x \in \mathbb{R}^{r_x \times d_x}$  and  $\mathbf{D}_x \in \mathbb{R}^{d_x \times r_x}$  projects and reconstructs back, respectively, the high-dimensional system state to and from a low-dimensional feature  $\mathbf{x}_R \in \mathbb{R}^{r_x}$ . The linear operators  $\mathbf{A}_R \in \mathbb{R}^{r_x \times r_x}$  and  $\mathbf{B}_R \in \mathbb{R}^{r_x \times d_u}$  describe the relations between successive reduced states and actuations. We refer to this reduced order model with linear networks as linear autoencoding ROM or LAROM. In the following, we first analyze the solution of the optimization objective of LAROM for a fixed *encoder*  $\mathbf{E}_x$ . Then we establish a connection between the solution of LAROM and the solution of DMDc, and further discuss the choice of the encoder to promote similarity between the two. Finally, we extend the linear model to a DNN-based model, which we refer to as DeepROM.

### 4.1.1 Analysis of the linear reduced order model for a fixed encoder

The DMDc algorithm essentially solves for  $\tilde{\mathbf{G}} \in \mathbb{R}^{r_x \times (d_x + d_u)}$  to minimize  $\frac{1}{n} \sum_{i=0}^{n-1} \|\mathbf{E}_x \mathbf{x}(t_{i+1}) - \tilde{\mathbf{G}} \boldsymbol{\omega}(t_i)\|^2$  for a fixed projection matrix  $\mathbf{E}_x = \mathbf{E}_{\text{DMDc}} = \hat{\mathbf{U}}_Y^\top$ . Here,  $\boldsymbol{\omega}(t_i)$  is the concatenated vector of state and actuation as defined in section 3.3. The optimal solution  $\tilde{\mathbf{G}}_{\text{opt}}$  is then partitioned as  $[\tilde{\mathbf{A}} \quad \tilde{\mathbf{B}}]$  such that  $\tilde{\mathbf{A}} \in \mathbb{R}^{r_x \times d_x}$ ,  $\tilde{\mathbf{B}} \in \mathbb{R}^{r_x \times d_u}$ . Finally,  $\tilde{\mathbf{A}}$  is post-multiplied with the reconstruction operator  $\mathbf{D}_{\text{DMDc}} = \hat{\mathbf{U}}_Y$  to get the ROM components  $\mathbf{A}_{R, \text{DMDc}}$  and  $\mathbf{B}_{R, \text{DMDc}}$ . Details of this process along with the proofs are given in appendix A.5. Note, the final step of this process offers dimensionality reduction only for the linear case, not in the case when the projection and reconstruction operators are nonlinear (e.g. DNNs). Therefore, we use an alternative formulation with the following results to design a loss function that encourages a DMDc-like solution for (11) and also offers dimensionality reduction when nonlinear components are used.

**Theorem 4.1.1.** *Consider the following objective function*

$$L_{\text{pred}}(\mathbf{E}_x, \mathbf{G}) = \frac{1}{n} \sum_{i=0}^{n-1} \|\mathbf{E}_x \mathbf{x}(t_{i+1}) - \mathbf{G} \mathbf{E}_{x\mathbf{u}} \boldsymbol{\omega}(t_i)\|^2, \quad (12)$$

where  $\mathbf{G} = [\mathbf{A}_R \quad \mathbf{B}_R] \in \mathbb{R}^{r_x \times (r_x + d_u)}$ ,  $\mathbf{E}_{x\mathbf{u}} = \begin{bmatrix} \mathbf{E}_x & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d_u} \end{bmatrix} \in \mathbb{R}^{(r_x + d_u) \times (d_x + d_u)}$ ,  $\mathbf{I}_{d_u}$  being the identity matrix of order  $d_u$ . For any fixed matrix  $\mathbf{E}_x$ , the objective function  $L_{\text{pred}}$  is convex in the coefficients of  $\mathbf{G}$  and attains its minimum for any  $\mathbf{G}$  satisfying

$$\mathbf{G} \mathbf{E}_{x\mathbf{u}} \boldsymbol{\Omega} \boldsymbol{\Omega}^\top \mathbf{E}_{x\mathbf{u}}^\top = \mathbf{E}_x \mathbf{Y} \boldsymbol{\Omega}^\top \mathbf{E}_{x\mathbf{u}}^\top, \quad (13)$$

where  $\mathbf{Y}$  and  $\boldsymbol{\Omega}$  are the data matrices as defined in section (3.3). If  $\mathbf{E}_x$  has full rank  $r_x$ , and  $\boldsymbol{\Omega} \boldsymbol{\Omega}^\top$  is non-singular, then  $L_{\text{pred}}$  is strictly convex and has a unique minimum for

$$\mathbf{G} = [\mathbf{A}_R \quad \mathbf{B}_R] = \mathbf{E}_x \mathbf{Y} \boldsymbol{\Omega}^\top \mathbf{E}_{x\mathbf{u}}^\top (\mathbf{E}_{x\mathbf{u}} \boldsymbol{\Omega} \boldsymbol{\Omega}^\top \mathbf{E}_{x\mathbf{u}}^\top)^{-1}. \quad (14)$$

*Proof sketch.* This can be proved by a method similar to the one used for deriving the solution of linear autoencoder in (Baldi & Hornik (1989)). For any fixed  $\mathbf{E}_x$ , the objective function of (12) can be written as  $L_{\text{pred}}(\mathbf{E}_x, \mathbf{G}) = \|\text{vec}(\mathbf{E}_x \mathbf{Y}) - (\boldsymbol{\Omega}^\top \mathbf{E}_{x\mathbf{u}}^\top \otimes \mathbf{I}_{r_x}) \text{vec}(\mathbf{G})\|^2$ , where  $\otimes$  denotes the Kronecker product and  $\text{vec}(\cdot)$  denotes vectorization of a matrix. Optimizing this linear least-square problem, we get (13) and (14), given the stated conditions are satisfied. The complete proof is given in appendix A.1.

**Remark.** For a unique solution, we assume that  $\mathbf{E}_x$  has full rank. The other scenario, i.e.,  $\mathbf{E}_x$  is rank-deficient suggests poor utilization of the hidden units of the model. In that case, the number of hidden units (which represents the dimension of the reduced state) should be decreased. The assumption that the covariance matrix  $\mathbf{\Omega}\mathbf{\Omega}^\top$  is invertible can be ensured when  $n \geq d_x + d_u$ , by removing any linearly dependent features in system state and actuation. When  $n < d_x + d_u$ , the covariance matrix  $\mathbf{\Omega}\mathbf{\Omega}^\top$  is not invertible. However, similar results can be obtained by adding  $\ell_2$  regularization (for the coefficients/entries of  $\mathbf{G}$ ) to the objective function. Proof of this is given in appendix A.4.

#### 4.1.2 The connection between the solutions of the linear autoencoding model and DMDC

The connection between the ROM obtained by minimizing  $L_{\text{pred}}$  (for a fixed  $\mathbf{E}_x$ ), i.e., (14) and the DMDC ROM of (10b) is not readily apparent. To interpret the connection, we formulate an alternative representation of (14) utilizing the SVD and the Moore-Penrose inverse of matrices. This alternative representation leads to the following result.

**Corollary 4.1.1.1.** *Consider the (full) SVD of the data matrix  $\mathbf{\Omega}$  given by  $\mathbf{\Omega} = \mathbf{U}_\Omega \mathbf{\Sigma}_\Omega \mathbf{V}_\Omega^\top$ , where  $\mathbf{U}_\Omega \in \mathbb{R}^{(d_x+d_u) \times (d_x+d_u)}$ ,  $\mathbf{\Sigma}_\Omega \in \mathbb{R}^{(d_x+d_u) \times n}$ , and  $\mathbf{V}_\Omega \in \mathbb{R}^{n \times n}$ . If  $\mathbf{E}_x = \widehat{\mathbf{U}}_Y^\top$  and  $\mathbf{\Omega}\mathbf{\Omega}^\top$  is non-singular, then the solution for  $\mathbf{G} = [\mathbf{A}_R \ \mathbf{B}_R]$  corresponding to the unique minimum of  $L_{\text{pred}}$  can be expressed as*

$$\mathbf{A}_R = \widehat{\mathbf{U}}_Y^\top \mathbf{Y} \mathbf{V}_\Omega \mathbf{\Sigma}^* \mathbf{U}_{\Omega,1}^\top \widehat{\mathbf{U}}_Y, \quad \text{and} \quad \mathbf{B}_R = \widehat{\mathbf{U}}_Y^\top \mathbf{Y} \mathbf{V}_\Omega \mathbf{\Sigma}^* \mathbf{U}_{\Omega,2}^\top, \quad (15)$$

where  $[\mathbf{U}_{\Omega,1}^\top \ \mathbf{U}_{\Omega,2}^\top] = \mathbf{U}_\Omega^\top$  with  $\mathbf{U}_{\Omega,1} \in \mathbb{R}^{d_x \times (d_x+d_u)}$ ,  $\mathbf{U}_{\Omega,2} \in \mathbb{R}^{d_u \times (d_x+d_u)}$ , and  $\mathbf{\Sigma}^* = \lim_{\varepsilon \rightarrow 0} (\mathbf{\Sigma}_\Omega^\top \mathbf{U}_{\Omega,1}^\top \widehat{\mathbf{U}}_Y \widehat{\mathbf{U}}_Y^\top \mathbf{U}_{\Omega,1} \mathbf{\Sigma}_\Omega + \mathbf{\Sigma}_\Omega^\top \mathbf{U}_{\Omega,2}^\top \mathbf{U}_{\Omega,2} \mathbf{\Sigma}_\Omega + \varepsilon^2 \mathbf{I}_n)^{-1} \mathbf{\Sigma}_\Omega^\top$ .

*Proof sketch.* This can be derived by plugging  $\mathbf{E}_x = \widehat{\mathbf{U}}_Y^\top$  into (14), and using the SVD definition and the limit definition (Albert (1972)) of the Moore-Penrose inverse. The complete proof is given in appendix A.3 that uses some preliminary results presented in appendix A.2.

**Remark.** It can be verified easily that if we use the truncated SVD (as defined by 9), instead of the full SVD, for  $\mathbf{\Omega}$  in corollary 4.1.1.1, we get an approximation of (15):

$$\widehat{\mathbf{A}}_R = \widehat{\mathbf{U}}_Y^\top \mathbf{Y} \widehat{\mathbf{V}}_\Omega \widehat{\mathbf{\Sigma}}^* \widehat{\mathbf{U}}_{\Omega,1}^\top \widehat{\mathbf{U}}_Y, \quad \text{and} \quad \widehat{\mathbf{B}}_R = \widehat{\mathbf{U}}_Y^\top \mathbf{Y} \widehat{\mathbf{V}}_\Omega \widehat{\mathbf{\Sigma}}^* \widehat{\mathbf{U}}_{\Omega,2}^\top, \quad (16)$$

where  $\widehat{\mathbf{\Sigma}}^* = \lim_{\varepsilon \rightarrow 0} (\widehat{\mathbf{\Sigma}}_\Omega^\top \widehat{\mathbf{U}}_{\Omega,1}^\top \widehat{\mathbf{U}}_Y \widehat{\mathbf{U}}_Y^\top \widehat{\mathbf{U}}_{\Omega,1} \widehat{\mathbf{\Sigma}}_\Omega + \widehat{\mathbf{\Sigma}}_\Omega^\top \widehat{\mathbf{U}}_{\Omega,2}^\top \widehat{\mathbf{U}}_{\Omega,2} \widehat{\mathbf{\Sigma}}_\Omega + \varepsilon^2 \mathbf{I}_{r_{xx}})^{-1} \widehat{\mathbf{\Sigma}}_\Omega^\top$ . We can see that (16) has the same form as (10b), except  $\widehat{\mathbf{\Sigma}}_\Omega^{-1}$  is replaced with  $\widehat{\mathbf{\Sigma}}^*$ .

All the aforementioned results are derived for a fixed  $\mathbf{E}_x$  and the relation to the DMDC is specific to the case  $\mathbf{E}_x = \widehat{\mathbf{U}}_Y^\top$ . Note that the columns of the  $\widehat{\mathbf{U}}_Y$  are the left singular vectors, corresponding to the leading singular values, of  $\mathbf{Y}$ . Equivalently, those are also the eigenvectors, corresponding to the leading eigenvalues, of the covariance matrix  $\mathbf{Y}\mathbf{Y}^\top$ .  $L_{\text{pred}}$  alone does not constrain  $\mathbf{E}_x$  to take a similar form and we need another loss term to encourage such form for the encoder. To this end, we follow the work of Baldi & Hornik (1989) on the similarity between principle component analysis and linear autoencoders, optimized with the following objective function:

$$L_{\text{recon}}(\mathbf{E}_x, \mathbf{D}_x) = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}(t_i) - \mathbf{D}_x \mathbf{E}_x \mathbf{x}(t_i)\|^2. \quad (17)$$

They showed that all the critical points of  $L_{\text{recon}}$  correspond to projections onto subspaces associated with subsets of eigenvectors of the covariance matrix  $\mathbf{Y}\mathbf{Y}^\top$ . Moreover,  $L_{\text{recon}}$  has a unique global minimum corresponding to the first  $r_x$  (i.e., the desired dimension of the reduced state) number of eigenvectors of  $\mathbf{Y}\mathbf{Y}^\top$ , associated with the leading  $r_x$  eigenvalues. In other words, for any invertible matrix  $\mathbf{C} \in \mathbb{R}^{r_x \times r_x}$ ,  $\mathbf{D}_x = \mathbf{U}_{r_x} \mathbf{C}$  and  $\mathbf{E}_x = \mathbf{C}^{-1} \mathbf{U}_{r_x}^\top$  globally minimizes  $L_{\text{recon}}$ , where  $\mathbf{U}_{r_x}$  denotes the matrix containing leading  $r_x$  eigenvectors of  $\mathbf{Y}\mathbf{Y}^\top$ . Since the left singular vectors of  $\mathbf{Y}$  are the eigenvectors of  $\mathbf{Y}\mathbf{Y}^\top$ , we have  $\mathbf{U}_{r_x} = \widehat{\mathbf{U}}_Y$ . Hence, we consider to utilize  $L_{\text{recon}}$  to promote learning an encoder  $\mathbf{E}_x$  in the form of  $\mathbf{C}^{-1} \widehat{\mathbf{U}}_Y^\top$ .

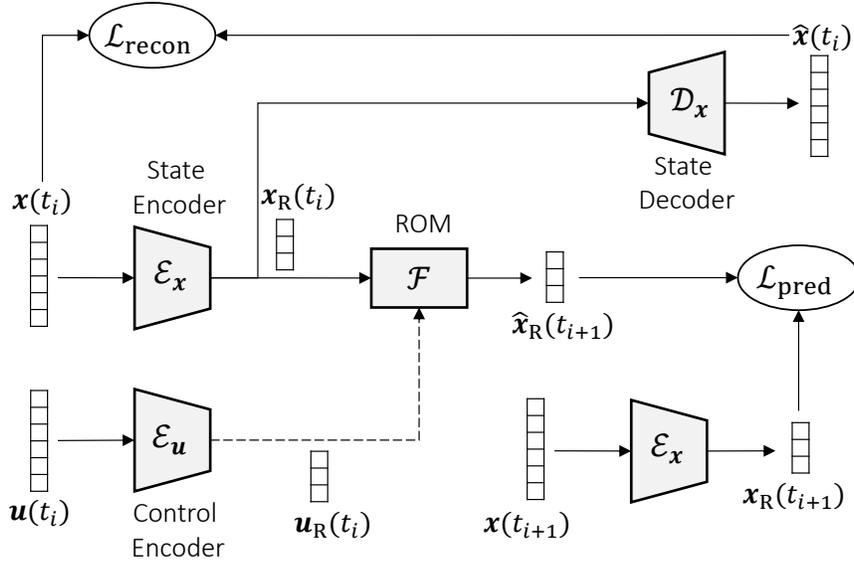


Figure 2: Autoencoding architecture for [reduced order modeling](#). The state encoder  $\mathcal{E}_x$  and control encoder  $\mathcal{E}_u$  reduce the dimension of the state and actuation, respectively. The ROM  $\mathcal{F}$  takes the current reduced state and actuation to predict the next reduced state, which is then uplifted to the full state by the state decoder  $\mathcal{D}_x$ . All modules are trained together using a combined loss involving  $\mathcal{L}_{\text{pred}}$  and  $\mathcal{L}_{\text{recon}}$ . The dashed arrow indicates that the  $\mathcal{E}_u$  is used only when  $d_u \gg 1$ ; otherwise, the actuation is used as a direct input to ROM.

Accordingly, we propose to minimize the following objective function to encourage a DMDc-like solution for LAROM:

$$L(\mathbf{E}_x, \mathbf{D}_x, \mathbf{G}) = L_{\text{pred}}(\mathbf{E}_x, \mathbf{G}) + \beta_1 L_{\text{recon}}(\mathbf{E}_x, \mathbf{D}_x), \quad (18)$$

where  $\beta_1 > 0$  is a tunable hyperparameter.

It is important to note that  $L_{\text{recon}}$  is minimized for any invertible matrix  $\mathbf{C}$ ,  $\mathbf{D}_x = \hat{\mathbf{U}}_Y \mathbf{C}$ , and  $\mathbf{E}_x = \mathbf{C}^{-1} \hat{\mathbf{U}}_Y^\top$ . When optimized using gradient descent, it is highly unlikely to get  $\mathbf{C}$  as the identity matrix like DMDc. Rather, we expect a random  $\mathbf{C}$ . Therefore, we need additional constraints to promote similarity with DMDc. For this purpose, we tie the matrices  $\mathbf{E}_x$  and  $\mathbf{D}_x$  to be the transpose of each other and add a semi-orthogonality constraint  $\beta_4 \|\mathbf{E}_x \mathbf{E}_x^\top - \mathbf{I}_{r_x}\|$ ,  $\beta_4 > 0$  to the optimization objective of (18).

#### 4.1.3 Extending the linear model to a deep model

Here, we discuss the process of extending LAROM to a nonlinear [reduced order modeling](#) framework. We replace all the trainable components of LAROM, i.e.,  $\mathbf{E}_x$ ,  $\mathbf{D}_x$ , and  $\mathbf{G}$ , with DNNs. Specifically, we use an encoding function or *encoder*  $\mathcal{E}_x: \mathbb{X} \rightarrow \mathbb{R}^{r_x}$  and a decoding function or *decoder*  $\mathcal{D}_x: \mathbb{R}^{r_x} \rightarrow \mathbb{X}$  to transform the high-dimensional system state to low-dimensional features and reconstruct it back, respectively, i.e.,

$$\mathbf{x}_R = \mathcal{E}_x(\mathbf{x}), \quad \hat{\mathbf{x}} = \mathcal{D}_x(\mathbf{x}_R), \quad (19)$$

where  $\mathbf{x}_R \in \mathbb{R}^{r_x}$  denotes the reduced state, and  $\hat{\mathbf{x}}$  is the reconstruction of  $\mathbf{x}$ . Unlike the linear case, we use an encoder  $\mathcal{E}_u: \mathbb{U} \rightarrow \mathbb{R}^{r_u}$ ,  $r_u \ll d_u$  for the actuation as well, in cases where the control space is also high-dimensional (for example, distributed control of spatiotemporal PDEs). The control encoder  $\mathcal{E}_u$  maps the high-dimensional actuation to a low-dimensional representation:  $\mathbf{u}_R = \mathcal{E}_u(\mathbf{u})$ , where  $\mathbf{u}_R \in \mathbb{R}^{r_u}$  denotes the encoded actuation. The encoded state and control are then fed to another DNN that represents the reduced order dynamics

$$\frac{d\mathbf{x}_R}{dt} = \mathcal{F}(\mathbf{x}_R, \mathbf{u}_R), \quad (20)$$

where  $\mathcal{F} : \mathbb{R}^{r_x} \times \mathbb{R}^{r_u} \rightarrow \mathbb{R}^{r_x}$ . Given the current reduced state  $\mathbf{x}_R(t_i)$  and control input  $\mathbf{u}_R(t_i)$ , the next reduced state  $\mathbf{x}_R(t_{i+1})$  can be computed by integrating  $\mathcal{F}$  using standard numerical integrator or neural ODE (Chen et al. (2018)):

$$\mathbf{x}_R(t_{i+1}) = \mathbf{x}_R(t_i) + \int_{t_i}^{t_{i+1}} \mathcal{F}(\mathbf{x}_R(t_i), \mathbf{u}_R(t_i)) dt \triangleq \mathcal{G}(\mathbf{x}_R(t_i), \mathbf{u}_R(t_i)). \quad (21)$$

We can say that  $\mathcal{G}$  is the nonlinear counterpart of  $\mathbf{G}$ .

Note, here the ROM is represented as a continuous-time dynamics, unlike the linear case where we used a discrete-time model. We use a discrete-time formulation for LAROM to establish its similarity with DMDC, which is formulated in discrete time. DeepROM can be formulated in a similar fashion as well. However, the specific control learning algorithm we used, which will be discussed in the next subsection, requires vector fields of the learned ROM for training. Therefore, we formulate the ROM in continuous time so that it provides the vector field  $\mathcal{F}(\mathbf{x}_R, \mathbf{u}_R)$  of the dynamics. In cases where only the prediction model is of interest and control learning is not required, a discrete-time formulation should be used for faster training of the ROM.

We train  $\mathcal{E}_x, \mathcal{E}_u, \mathcal{D}_x$ , and  $\mathcal{F}$  by minimizing the following loss function, analogous to (18),

$$\mathcal{L}(\mathcal{E}_x, \mathcal{E}_u, \mathcal{D}_x, \mathcal{F}) = \mathcal{L}_{\text{pred}}(\mathcal{E}_x, \mathcal{E}_u, \mathcal{F}) + \beta_2 \mathcal{L}_{\text{recon}}(\mathcal{E}_x, \mathcal{D}_x), \quad (22)$$

where  $\beta_2 > 0$  is a tunable hyperparameter and  $\mathcal{L}_{\text{pred}}, \mathcal{L}_{\text{recon}}$  are defined as follows,

$$\begin{aligned} \mathcal{L}_{\text{pred}}(\mathcal{E}_x, \mathcal{E}_u, \mathcal{F}) &= \frac{1}{n} \sum_{i=0}^{n-1} \left\| \mathcal{E}_x(\mathbf{x}(t_{i+1})) - \mathcal{G}(\mathcal{E}_x(\mathbf{x}(t_i)), \mathcal{E}_u(\mathbf{u}(t_i))) \right\|^2, \\ \mathcal{L}_{\text{recon}}(\mathcal{E}_x, \mathcal{D}_x) &= \frac{1}{n} \sum_{i=1}^n \left\| \mathbf{x}(t_i) - \mathcal{D}_x \circ \mathcal{E}_x(\mathbf{x}(t_i)) \right\|^2. \end{aligned} \quad (23)$$

Here, the operator  $\circ$  denotes the composition of two functions. In experiments,  $\mathcal{L}_{\text{recon}}$  also includes the reconstruction loss of the desired state where we want to stabilize the system. Figure 2 shows the overall framework for training DeepROM.

## 4.2 Learning control

Once we get a trained ROM of the form (20) using the method proposed in section 4.1, the next goal is to design a controller for the system utilizing that ROM. Since our ROM is represented by DNNs, we need a data-driven method to develop the controller. We adopt the approach presented by Saha et al. (2021) for learning control law for nonlinear systems, represented by DNNs. The core idea of the method is to hypothesize a target dynamics that is exponentially stable at the desired state and simultaneously learn a control policy to realize that target dynamics in closed loop. A DNN is used to represent the vector field  $\mathcal{F}_s : \mathbb{R}^{r_x} \rightarrow \mathbb{R}^{r_x}$  of the target dynamics  $\frac{d\mathbf{x}_R}{dt} = \mathcal{F}_s(\mathbf{x}_R)$ . We use another DNN to represent a controller  $\Pi : \mathbb{R}^{r_x} \rightarrow \mathbb{R}^{d_u}$  that provides the necessary actuation for a given reduced state  $\mathbf{x}_R$ :

$$\mathbf{u} = \Pi(\mathbf{x}_R). \quad (24)$$

This control  $\mathbf{u}$  is then encoded by (trained)  $\mathcal{E}_u$  to its low-dimensional representation  $\mathbf{u}_R$ . Finally, the reduced state  $\mathbf{x}_R$  and actuation  $\mathbf{u}_R$  are fed to the (trained) ROM of (20) to get  $\mathcal{F}(\mathbf{x}_R, \mathbf{u}_R)$ . The overall framework for learning control is shown in Figure 3.

Our training objective is to minimize the difference between  $\mathcal{F}(\mathbf{x}_R, \mathbf{u}_R)$  and  $\mathcal{F}_s(\mathbf{x}_R)$ , i.e.,

$$\mathcal{L}_{\text{ctrl}}(\mathcal{F}_s, \Pi) = \frac{1}{n} \sum_{i=1}^n \left\| \mathcal{F}(\mathcal{E}_x(\mathbf{x}(t_i)), \mathcal{E}_u \circ \Pi \circ \mathcal{E}_x(\mathbf{x}(t_i))) - \mathcal{F}_s \circ \mathcal{E}_x(\mathbf{x}(t_i)) \right\|^2. \quad (25)$$

To minimize the control effort, we add a regularization loss with (25), and the overall training objective for learning control is given by

$$\mathcal{L}_{\text{ctrl,reg}}(\mathcal{F}_s, \Pi) = \mathcal{L}_{\text{ctrl}}(\mathcal{F}_s, \Pi) + \beta_3 \frac{1}{n} \sum_{i=1}^n \left\| \Pi(\mathbf{x}_R(t_i)) \right\|^2, \quad (26)$$

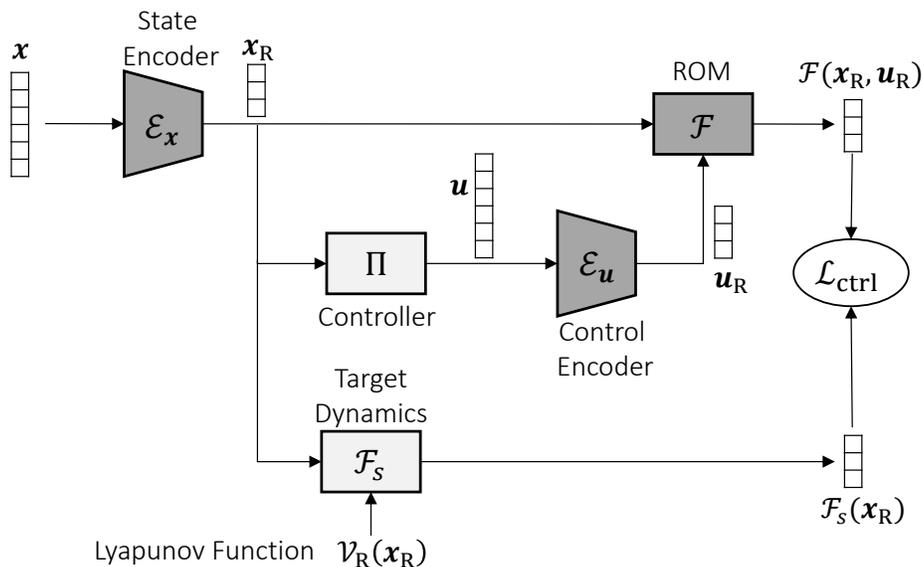


Figure 3: The control learning process. Given a reduced state,  $\mathcal{F}_s$  predicts a target dynamics for the closed-loop system, and the controller  $\Pi$  predicts an actuation to achieve that target. Both the modules are trained jointly using the loss function  $\mathcal{L}_{\text{ctrl}}$ . Parameters of the dark-shaded modules are kept fixed during this process.

where  $\beta_3 > 0$  is a tunable hyperparameter. Here we jointly train the DNNs representing  $\Pi$  and  $\mathcal{F}_s$  only, whereas the previously-trained DNNs for  $\mathcal{E}_x$ ,  $\mathcal{E}_u$ , and  $\mathcal{F}$  are kept frozen. Once all the DNNs are trained, we only need  $\mathcal{E}_x$  and  $\Pi$  during evaluation to generate actuation for the actual system, given a full-state observation:

$$\mathbf{u} = \Pi \circ \mathcal{E}_x(\mathbf{x}) = \pi(\mathbf{x}). \quad (27)$$

As we mentioned earlier, we require the target dynamics, hypothesized by a DNN, to be exponentially stable at the desired state. Without loss of generality, we consider stability at  $\mathbf{x}_R = \mathbf{0}$ . As we mentioned earlier, the system can be stabilized at any desired state by adding a feedforward component to the control. Dynamics represented by a standard neural network is not stable at any equilibrium point, in general. Kolter & Manek (2019) showed that it is possible to design a DNN, by means of Lyapunov functions, to represent a dynamics that is exponentially stable at an equilibrium point. Accordingly, we represent our target dynamics as follows:

$$\frac{d\mathbf{x}_R}{dt} = \mathcal{F}_s(\mathbf{x}_R) = \mathcal{P}(\mathbf{x}_R) - \frac{\text{ReLU}(\nabla \mathcal{V}_R(\mathbf{x}_R)^\top \mathcal{P}(\mathbf{x}_R) + \alpha \mathcal{V}_R(\mathbf{x}_R))}{\|\nabla \mathcal{V}_R(\mathbf{x}_R)\|^2} \nabla \mathcal{V}_R(\mathbf{x}_R), \quad (28)$$

where  $\alpha$  is a positive constant,  $\text{ReLU}(z) = \max(0, z)$ ,  $z \in \mathbb{R}$ , and  $\mathcal{V}_R : \mathbb{R}^{r_\omega} \rightarrow \mathbb{R}$  is a candidate Lyapunov function, i.e., satisfies the criteria similar to (3) and (5). We use

$$\mathcal{V}_R(\mathbf{x}_R) = \mathbf{x}_R^\top \mathbf{K} \mathbf{x}_R, \quad (29)$$

where  $\mathbf{K} \in \mathbb{R}^{r_\omega \times r_\omega}$  is a positive definite matrix.

Though the efficacy of learning control by minimizing the difference with respect to a target dynamics is experimentally demonstrated by Saha et al. (2021), the stability of the closed-loop system subjected to the learned control law has not been studied analytically. Here, we present a result that shows that if we can minimize  $\mathcal{L}_{\text{ctrl}}$  such that the difference between the target dynamics and the closed-loop dynamics is sufficiently small for all  $\mathbf{x}_R \in \mathbb{X}_R \subset \mathbb{R}^{r_\omega}$ , then the trajectories of the closed-loop ROM starting sufficiently close to the origin remains close to the origin, i.e., *ultimately bounded* (Khalil (2002)). [Boundedness of the closed-loop ROM trajectories under the proposed control policy is a necessary but not sufficient requirement for the stability of the original system.](#)

**Theorem 4.2.1.** Consider the target dynamics defined by (28) and the candidate Lyapunov function defined by (29). Suppose the difference between the target dynamics and the closed-loop dynamics satisfies

$$\|\mathcal{F}(\mathbf{x}_R, \mathcal{E}_u \circ \Pi(\mathbf{x}_R)) - \mathcal{F}_s(\mathbf{x}_R)\| \leq \delta < \frac{\alpha\theta\lambda_{\min}(\mathbf{K})}{2\lambda_{\max}(\mathbf{K})} \sqrt{\frac{\lambda_{\min}(\mathbf{K})}{\lambda_{\max}(\mathbf{K})}} \eta, \quad (30)$$

for all  $\mathbf{x}_R \in \mathbb{X}_R = \{\mathbf{x}_R \in \mathbb{R}^{r^*} \mid \|\mathbf{x}_R\| < \eta\}$  and  $0 < \theta < 1$ . Then, for all initial points satisfying  $\|\mathbf{x}_R(t_0)\| < \sqrt{\frac{\lambda_{\min}(\mathbf{K})}{\lambda_{\max}(\mathbf{K})}} \eta$ , the solution of the closed-loop ROM  $\frac{d\mathbf{x}_R}{dt} = \mathcal{F}(\mathbf{x}_R, \mathcal{E}_u \circ \Pi(\mathbf{x}_R))$  satisfies

$$\|\mathbf{x}_R(t)\| \leq \lambda e^{-\gamma(t-t_0)} \|\mathbf{x}_R(t_0)\|, \quad \forall t_0 \leq t < t_c + t_0 \quad (31)$$

and

$$\|\mathbf{x}_R(t)\| \leq \frac{2\delta}{\alpha\theta} \lambda^3, \quad \forall t \geq t_c + t_0 \quad (32)$$

for some finite  $t_c > 0$ , where

$$\gamma = \frac{\alpha(1-\theta)\lambda_{\min}(\mathbf{K})}{2\lambda_{\max}(\mathbf{K})} \quad \text{and} \quad \lambda = \sqrt{\frac{\lambda_{\max}(\mathbf{K})}{\lambda_{\min}(\mathbf{K})}} \quad (33)$$

*Proof Sketch.* This can be proved by first deriving the Lyapunov conditions for the target dynamics (28) (Theorem 1, Kolter & Manek (2019)) and then applying the stability analysis of perturbed systems (Lemma 9.2, Khalil (2002)) and ultimate boundedness (Theorem 4.18, Khalil (2002)) on the closed-loop ROM. A unified proof is provided in appendix A.6.

## 5 Empirical Results

For empirical analysis, we consider modeling and controlling spatiotemporal PDE-driven systems with high-dimensional measurements over discretized space. One of the primary applications of reduced order modeling lies in comprehending the behavior of complex physical processes which are typically characterized by systems of PDEs. Since spatiotemporal PDE-driven systems are infinite-dimensional in their continuous form and high-dimensional when discretized, they are a fitting choice for evaluating our method. The first example investigates a single variable actuation, whereas distributed actuation is considered for the second example.

### 5.1 Baselines

The similarity between DMDc and LAROM is demonstrated using the dynamic modes estimated in respective methods. The prediction performance of DeepROM is compared against DMDc and the Deep Koopman model (Morton et al. (2018)). The Deep Koopman model shares a similar DNN-based autoencoding structure as ours, with the distinction that its (reduced order) dynamic model is linear. The method proposed by Morton et al. (2018) considers a model predictive scenario, where the state/system matrix of the linear reduced order model is updated with online observations during operation while the input/control matrix is kept fixed. However, in contrast to the original method, we keep both matrices fixed during operation (once those are trained) as we consider offline control design in this paper. For the same reason, we apply *linear quadratic regulator* (LQR) on the ROM obtained from the Deep Koopman method, instead of model predictive control, to compare the control performance with our method: DeepROC. The control performance is also compared against the reduced order controller obtained by applying LQR on the ROM derived from DMDc.

Details on the neural network architectures and training settings for the Deep Koopman model are given in appendix D.

## 5.2 Reaction–diffusion system stabilization

For the first experiment, we consider the Newell–Whitehead–Segel reaction-diffusion equation with the Neumann boundary condition

$$\begin{aligned} \frac{\partial q}{\partial t} &= \sigma \nabla^2 q + q(1 - q^2) + \mathbf{1}_{\mathbb{W}} w \quad \text{in } \mathbb{I} \times \mathbb{R}^+, \\ \nabla q(\zeta_l, t) &= \nabla q(\zeta_r, t) = 0, \quad t \in \mathbb{R}^+, \end{aligned} \quad (34)$$

which is used to describe various nonlinear physical systems including Rayleigh–Bénard convection. This example is used by Kalise & Kunisch (2018) to evaluate nonlinear controllers designed from reduced order state space representation. Similar systems are used for modeling problems as well in 1D (Raissi et al. (2019)) and 2D (Li et al. (2020)). In (34),  $q(\zeta, t) \in \mathbb{R}$  denotes the measurement variable such as concentration or temperature at location  $\zeta \in \mathbb{I} \subset \mathbb{R}$  and time  $t$ ;  $\sigma$  denotes the diffusion coefficient;  $w(t) \in \mathbb{R}$  is the actuation at time  $t$  and  $\mathbf{1}_{\mathbb{W}}(\zeta)$  is the indicator function with  $\mathbb{W} \subset \mathbb{I}$ ;  $\zeta_l$  and  $\zeta_r$  denote the boundary points of  $\mathbb{I}$ . (34) is a bistable system with  $\pm 1$  as stable and 0 as unstable equilibria. For the control task, we consider feedback stabilization of (34) at the unstable equilibrium 0, as studied by Kalise & Kunisch (2018). We use  $\mathbb{I} = (-1, 1)$ ,  $\mathbb{W} = (-0.2, 0.2)$ , and  $\sigma = 0.2$ . Details on dataset generation, neural network architectures, and training settings are given in appendix B.

### 5.2.1 Similarity with DMDC

To investigate the similarity DMDC, we first train the LAROM using gradient descent to minimize the objective (18) with the semi-orthogonality regularization and enforcing  $\mathbf{D}_{\mathbf{x}} = \mathbf{E}_{\mathbf{x}}^\top$ , as discussed in 4.1.2.

The dynamic modes for LAROM are computed as  $\varphi_i = \mathbf{D}_{\mathbf{x}} \mathbf{z}_i$ , where  $\mathbf{z}_i$  is the  $i^{\text{th}}$  eigenvector of  $\mathbf{A}_{\mathbf{R}}$ . Similarly, the dynamic modes for DMDC are computed as  $\varphi_{i,\text{DMDC}} = \mathbf{D}_{\text{DMDC}} \mathbf{z}_{i,\text{DMDC}}$ , where  $\mathbf{z}_{i,\text{DMDC}}$  is the  $i^{\text{th}}$  eigenvector of  $\mathbf{A}_{\mathbf{R},\text{DMDC}}$ . Note, these dynamic modes are similar to the ones used in the original DMD algorithm Schmid (2010), not the exact modes obtained in Proctor et al. (2016). Exact modes cannot be computed for LAROM since it does not involve SVD. Modes defined by  $\varphi_{i,\text{DMDC}} = \mathbf{D}_{\text{DMDC}} \mathbf{z}_{i,\text{DMDC}} = \widehat{\mathbf{U}}_{\mathbf{Y}} \mathbf{z}_{i,\text{DMDC}}$  are the orthogonal projection of the exact modes onto the range of  $\mathbf{Y}$  (Theorem 3, Tu et al. (2014)). Figure 4 compares the dynamic modes obtained using DMDC and LAROM for the case when the dimension of the ROMs is 3. It is important to note that the numbering of the modes is arbitrary as the optimal ranking of DMDC modes is not trivial. The correspondence between the DMDC modes and LAROM modes are determined by comparing the eigenvalues of  $\mathbf{A}_{\mathbf{R},\text{DMDC}}$  and  $\mathbf{A}_{\mathbf{R}}$ . Dynamic modes of both methods are similar except for the different signs of the first two modes.

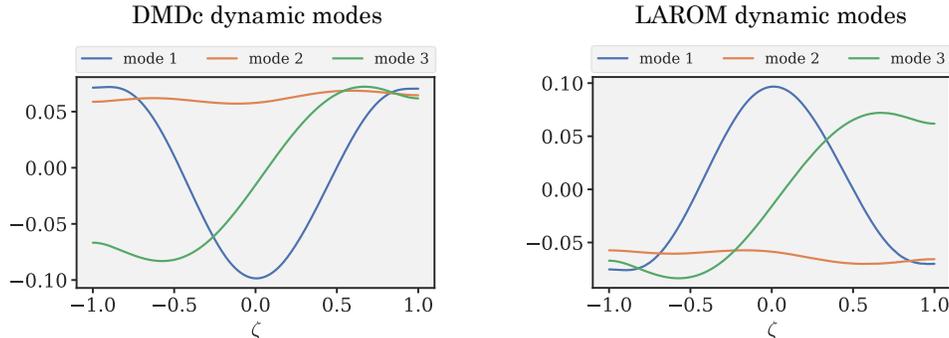


Figure 4: The first three dynamic modes of the reaction–diffusion system, obtained using DMDC and LAROM.

### 5.2.2 Prediction performance of DeepROM

We now compare the performance of DeepROM, Deep Koopman model, and DMDC in the prediction task. Note, this example uses low-dimensional actuation (just a single variable). Accordingly, the control encoder

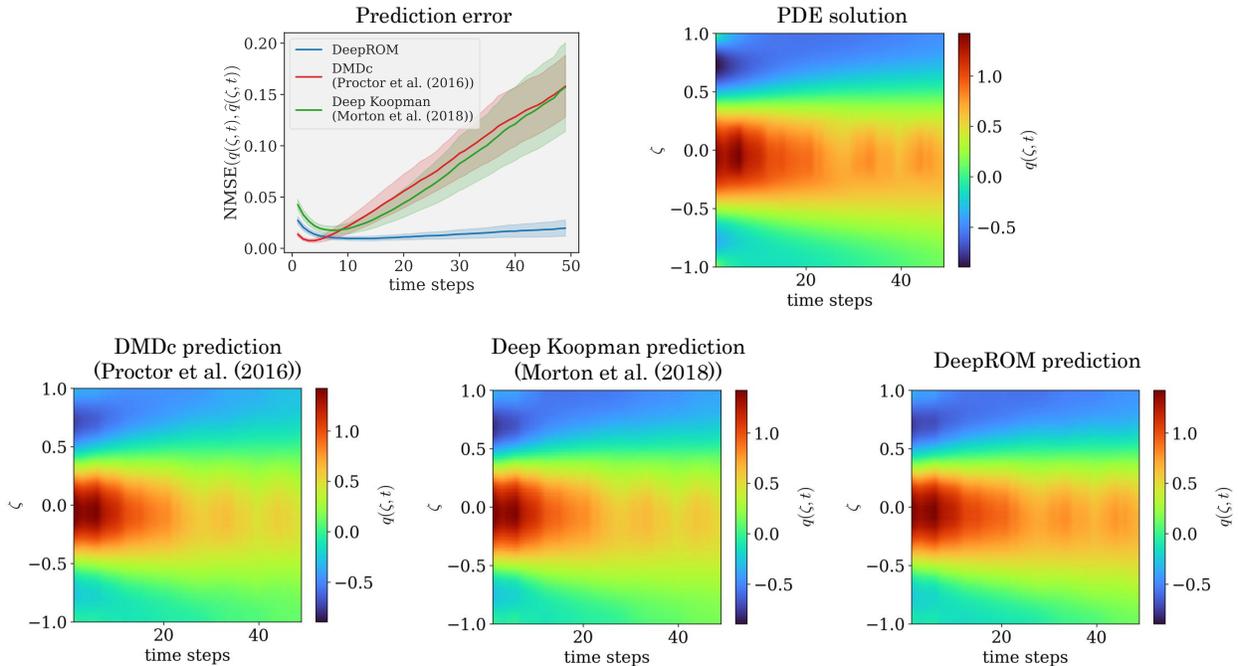


Figure 5: Prediction performance of DMDc, [Deep Koopman](#), and DeepROM in the reaction–diffusion example. The prediction error plot shows the mean error and 95% confidence interval from 100 test sequences and for [Deep Koopman](#) and DeepROM, 3 different training instances. One example sequence is used to visually compare the predictions with the solution from a PDE solver.

$\mathcal{E}_u$  is not used here. Figure 5 shows the quantitative and qualitative comparison of the recursive multi-step predictions obtained using DMDc, [Deep Koopman model](#), and DeepROM. The prediction error is computed as *normalized mean squared error* (NMSE) with respect to the solution obtained using the PDE solver. The prediction error plot shows the mean error and 95% confidence interval from 100 test sequences and for [Deep Koopman](#) and DeepROM, 3 different training instances. The color maps are shown for one example sequence with one training instance. Prediction error increases more quickly for DMDc and [Deep Koopman](#) than DeepROM as the linear ROMs become less accurate in the long term.

### 5.2.3 Control performance of DeepROC

Figure 6 shows the control performance of DeepROC, [Deep Koopman + LQR](#), and DMDc + LQR in the task of stabilizing the system at the unstable equilibrium 0 from an initial state  $2 + \cos(2\pi\zeta)\cos(\pi\zeta)$ . We use the following metrics for comparison:

- (i) mean squared error over time between the controlled solutions and the desired profile
- (ii) differential magnitude that measures the differential changes between the profiles at consecutive time steps. In the steady state, the differential magnitude should be close to zero.
- (iii) the amount of actuation applied

For [Deep Koopman](#) and DeepROC, the plots show the mean values with 1-standard deviation interval from 3 training instances. All methods show similar closed-loop error profiles. However, DeepROC requires significantly less amount of actuation in comparison with DMDc + LQR and [Deep Koopman + LQR](#) to reach a similar steady-state error. DeepROC can account for the decaying nonlinear term  $-q^3$  present in the system (34) and therefore learns to apply less actuation. Figure 7 visually compares the uncontrolled solution and the controlled solutions obtained using [the three](#) methods. When uncontrolled, the system reaches the stable equilibrium at 1, whereas the feedback-controlled system is stabilized at the desired state 0 in both cases.

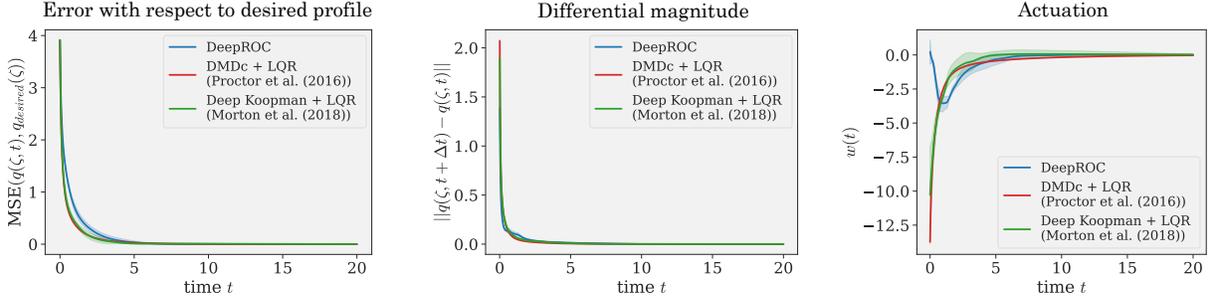


Figure 6: Control performance of DMDc + LQR, **Deep Koopman + LQR**, and DeepROC in the reaction–diffusion example. For **Deep Koopman** and DeepROC, the plots show the mean values with 1-standard deviation interval from 3 training instances.

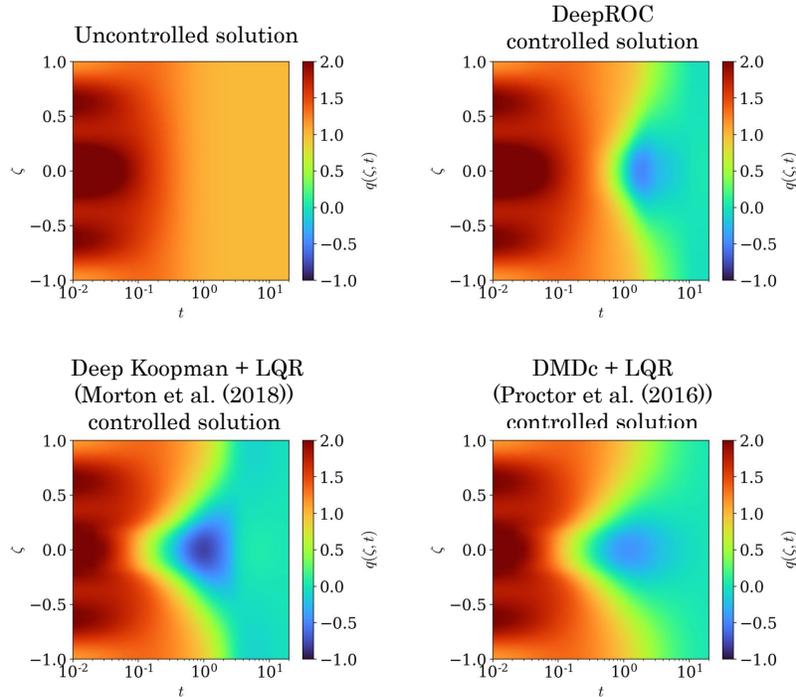


Figure 7: Visual comparison of the uncontrolled solution and the controlled solutions using DeepROC, **Deep Koopman + LQR**, and DMDc + LQR.

### 5.3 Vortex shedding suppression in fluid

In this experiment, we consider modeling and suppressing vortex shedding in two-dimensional incompressible flow past a circular cylinder. This is a well-known problem (Schäfer et al. (1996)) and is of great importance for many engineering applications (Williamson (1996)). Several previous studies on deep learning-based modeling and control have used this system for evaluation (Eivazi et al. (2020); Erichson et al. (2019); Rabault et al. (2019); Tang et al. (2020); Bieker et al. (2020); Morton et al. (2018)). The dynamics is governed by the incompressible Navier-Stokes equations given by

$$\frac{\partial \mathbf{v}}{\partial t} - \nu \nabla^2 \mathbf{v} + (\mathbf{v} \cdot \nabla) \mathbf{v} = -\frac{1}{\rho} \nabla p + \mathbf{1}_W \mathbf{w}, \quad \nabla \cdot \mathbf{v} = \mathbf{0} \quad \text{in } \mathbb{I} \times \mathbb{R}^+, \quad (35)$$

where  $\mathbf{v}(\boldsymbol{\zeta}, t) \in \mathbb{R}^2$  denotes the flow velocity at location  $\boldsymbol{\zeta} \in \mathbb{I} \subset \mathbb{R}^2$  and time  $t$ ,  $p(\boldsymbol{\zeta}, t) \in \mathbb{R}$  denotes the pressure,  $\nu$  denotes the kinematic viscosity and  $\rho$  denotes the density of the fluid.  $\mathbf{w}(\boldsymbol{\zeta}, t)$  is the actuation/force applied to the system and  $\mathbf{1}_{\mathbb{W}}(\boldsymbol{\zeta})$  is the indicator function with  $\mathbb{W} \subset \mathbb{I}$ . We use  $\mathbb{I} = (0, 2.2) \times (0, 0.41)$  and  $\mathbb{W} = (0.11, 0.77) \times (0, 0.41)$ . Density and kinematic viscosity are chosen such that the Reynolds number is  $Re = 50$ , which is just above the cutoff for the onset of the vortex shedding (Williamson (1996)). In this case, vortices are created at the back of the cylinder and are shed periodically from the upper and lower surfaces of the cylinder forming a von Kármán vortex street (Morton et al. (2018)). We use the domain  $\mathbb{W}$  for observation and distributed actuation. The Stokes flow is used as the desired state for the control task. More details on the problem setup, dataset generation, neural network architectures, and training settings are given in appendix C.

### 5.3.1 Similarity with DMDC

To analyze the dynamic modes, we train the LAROM by enforcing  $\mathbf{D}_{\mathbf{x}} = \mathbf{E}_{\mathbf{x}}^{\top}$  and adding the semi-orthogonality constraint to the learning objective, as mentioned in 4.1.2. Figure 8 compares the first two oscillatory dynamic modes obtained using DMDC and LAROM. Only the streamwise components are shown for brevity. Also, complex modes occur in conjugate pairs and only one from each pair is shown. The correspondence between the DMDC modes and LAROM modes are determined by comparing the eigenvalues of  $\mathbf{A}_{R,DMDC}$  and  $\mathbf{A}_R$ . Dynamic modes identified by LAROM are similar to the ones obtained from DMDC, except the real and imaginary components of the first mode are swapped.

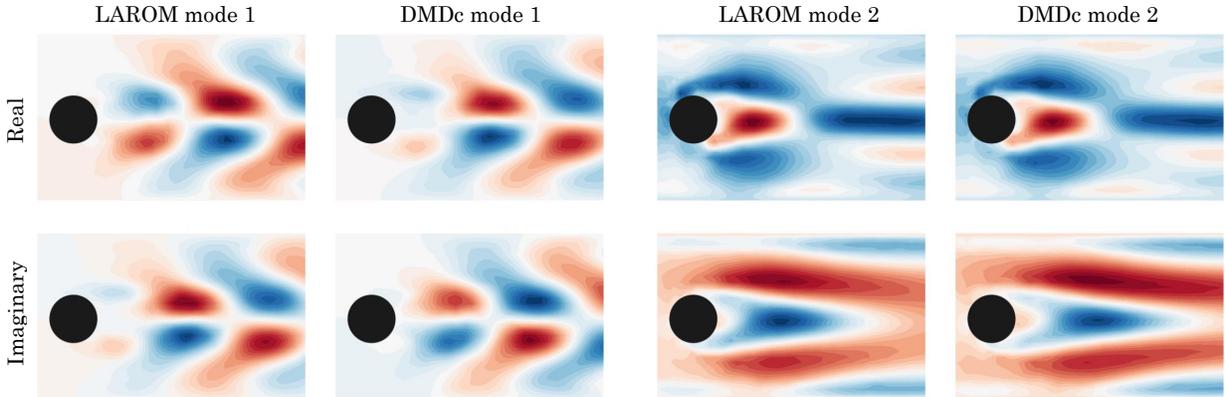


Figure 8: The first two dynamic modes obtained using DMDC and LAROM for the flow past a cylinder system.

### 5.3.2 Prediction performance of DeepROM

Figure 9 shows the quantitative and qualitative comparison of the recursive multi-step predictions, starting from  $t = 0.1$ , obtained using DMDC, Deep Koopman model, and DeepROM. The initial state is chosen at  $t = 0.1$  because the fluid does not reach the observation region  $\mathbb{W}$  before that time. The prediction error is computed as the *mean squared error* (MSE) with respect to the solution obtained using a PDE solver. For Deep Koopman and DeepROM, the prediction error plot shows the mean error and 1-standard deviation interval from 3 training instances. DeepROM shows lower prediction error in comparison with DMDC. The Deep Koopman model shows better prediction performance than DeepROM and DMDC during the initial few steps. However, its accuracy deteriorates rapidly and eventually becomes comparable to that of DMDC. Moreover, unlike DeepROM, DMDC and Deep Koopman model are unable to capture the shedding pattern in multi-step prediction as shown in the contour plots of the velocity magnitude.

### 5.3.3 Control performance of DeepROC

Figure 10 shows the control performance of DeepROC, Deep Koopman + LQR, and DMDC+LQR in the task of suppressing vortex shedding. The controllers of DeepROC and DMDC + LQR directly estimate the high-

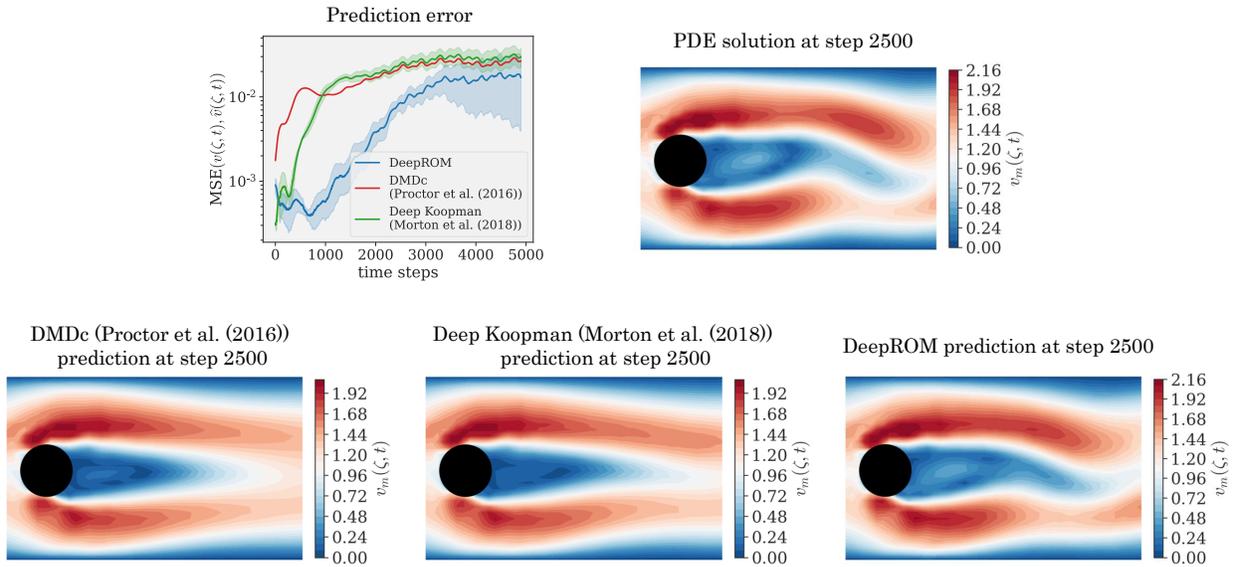


Figure 9: Prediction performance of DMDc, [Deep Koopman](#), and DeepROM in the fluid flow example. For [Deep Koopman](#) and DeepROM, the prediction error plot shows the mean error and 1-standard deviation interval from 3 training instances. Predictions at time step 2500 for the test sequence are visually compared with the solution from a PDE solver.  $v_m$  denotes the velocity magnitude.

dimensional actuation distributed over space. However, the same technique proved ineffective in suppressing the shedding for Deep Koopman + LQR. Therefore, instead of directly estimating the distributed actuation, we utilize a low-dimensional representation of the actuation for Deep Koopman + LQR. We represent the distributed actuation as a linear combination of some space-dependent sinusoidal basis functions. The controller is designed to estimate the coefficients of those basis functions in the linear combination. Details are provided in appendix D.

We use the same metrics as the previous example for comparison except for actuation. Since distributed control is applied in this case, we use the magnitude of the actuation here. For DeepROC and [Deep Koopman + LQR](#), the plots show the mean values with 1-standard deviation interval from 3 training instances. To reach a similar steady-state error, DeepROC takes a longer time compared to DMDc and [Deep Koopman + LQR](#). DeepROM uses the least amount of actuation during the initial few steps, whereas [Deep Koopman + LQR](#) has the least steady-state actuation magnitude. Figure 11 shows the velocity magnitude of the controlled flow for DeepROC, [Deep Koopman + LQR](#), and DMDc+LQR at different times, starting from a von Kármán vortex street pattern. All methods accomplish a similar steady-state flow pattern where vortex shedding has been suppressed.

## 6 Conclusion

We presented a framework for autoencoder-based modeling and control learning for high-dimensional dynamical systems. We showed that autoencoding ROMs are capable of capturing the dominant modes that are essential in analyzing and designing control for the underlying systems. As we showed in experiments, DeepROM offers better prediction accuracy than a linear ROM over a relatively longer prediction horizon when applied to nonlinear systems. However, this advantage does not always translate to significant improvement in control performance. Though the used control learning method theoretically ensures ultimate boundedness for the closed-loop ROM solution, data-driven optimization of the learning objective often makes the models susceptible to distribution shift which can impact the control performance. The control learning process in the DeepROC framework can easily be replaced with other methods like model-based RL or model predictive control. It would be interesting for future work to investigate whether updating both the reduced model and the controller in the MPC framework ensures robustness under distribution shift and

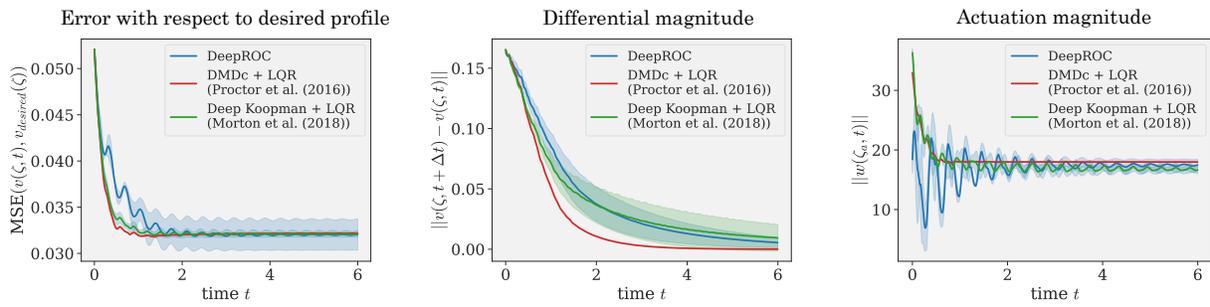


Figure 10: Control performance of DMDc + LQR, Deep Koopman + LQR, and DeepROC in the vortex shedding suppression task. For Deep Koopman + LQR and DeepROC, the plots show the mean values with 1-standard deviation interval from 3 training instances.

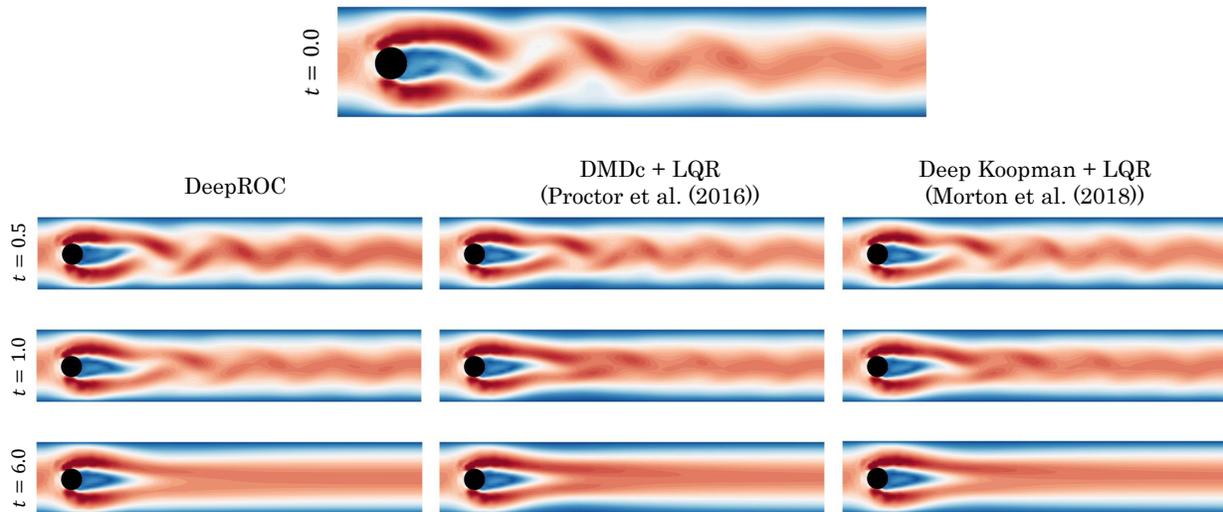


Figure 11: Visual comparison of the velocity magnitude of the flow over time subjected to the controllers obtained using DeepROC, Deep Koopman + LQR, and DMDc + LQR.

offers better control performance. Designing controllers for DNN-based models is a challenging task due to the standard difficulties associated with non-convex optimization. Nevertheless, we envision great prospects in solving many problems of control design for high-dimensional systems utilizing autoencoder-based models as they continue to demonstrate their effectiveness in the analysis and prediction of such systems.

## References

- Arthur Albert. *Regression and the Moore-Penrose Pseudoinverse*. Academic Press, 1972.
- Jeanne A Atwell, Jeffrey T Borggaard, and Belinda B King. Reduced order controllers for burgers' equation with a nonlinear observer. *International Journal of Applied Mathematics and Computer Science*, 11(6): 1311–1330, 2001.
- Ibrahim Ayed, Emmanuel de Bézenac, Arthur Pajot, Julien Brajard, and Patrick Gallinari. Learning dynamical systems from partial observations. *arXiv preprint arXiv:1902.11136*, 2019.
- Pierre Baldi and Kurt Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural networks*, 2(1):53–58, 1989.
- Gerben Beintema, Alessandro Corbetta, Luca Biferale, and Federico Toschi. Controlling rayleigh-bénard convection via reinforcement learning. *Journal of Turbulence*, 21(9-10):585–605, 2020.
- Katharina Bieker, Sebastian Peitz, Steven L Brunton, J Nathan Kutz, and Michael Dellnitz. Deep model predictive flow control with limited sensor data and online learning. *Theoretical and computational fluid dynamics*, 34:577–591, 2020.
- Oumayma Bounou, Jean Ponce, and Justin Carpentier. Online learning and control of dynamical systems from sensory input. In *NeurIPS 2021-Thirty-fifth Conference on Neural Information Processing Systems Year*, 2021.
- Kaixuan Chen, Jin Lin, Yiwei Qiu, Feng Liu, and Yonghua Song. Deep learning-aided model predictive control of wind farms for agc considering the dynamic wake effect. *Control Engineering Practice*, 116: 104925, 2021.
- Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018.
- Emmanuel De Bézenac, Arthur Pajot, and Patrick Gallinari. Deep learning for physical processes: Incorporating prior scientific knowledge. *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12): 124009, 2019.
- Jérémie Donà, Jean-Yves Franceschi, Sylvain Lamprier, and Patrick Gallinari. Pde-driven spatiotemporal disentanglement. *arXiv preprint arXiv:2008.01352*, 2020.
- Hamidreza Eivazi, Hadi Veisi, Mohammad Hossein Naderi, and Vahid Esfahanian. Deep neural networks for nonlinear model order reduction of unsteady flows. *Physics of Fluids*, 32(10):105104, 2020.
- N Benjamin Erichson, Michael Muehlebach, and Michael W Mahoney. Physics-informed autoencoders for lyapunov-stable fluid flow prediction. *arXiv preprint arXiv:1905.10866*, 2019.
- Paul Garnier, Jonathan Viquerat, Jean Rabault, Aurélien Larcher, Alexander Kuhnle, and Elie Hachem. A review on deep reinforcement learning for fluid mechanics. *Computers & Fluids*, 225:104973, 2021.
- Philipp Holl, Nils Thuerey, and Vladlen Koltun. Learning to control pdes with differentiable physics. In *International Conference on Learning Representations*, 2020.
- Jer-Nan Juang and Richard S Pappa. An eigensystem realization algorithm for modal parameter identification and model reduction. *Journal of guidance, control, and dynamics*, 8(5):620–627, 1985.
- Jer-Nan Juang, Minh Phan, Lucas G Horta, and Richard W Longman. Identification of observer/kalman filter markov parameters-theory and experiments. *Journal of Guidance, Control, and Dynamics*, 16(2): 320–329, 1993.
- Dante Kalise and Karl Kunisch. Polynomial approximation of high-dimensional hamilton-jacobi-bellman equations and applications to feedback control of semilinear parabolic pdes. *SIAM Journal on Scientific Computing*, 40(2):A629–A652, 2018.

- Hassan K. Khalil. *Nonlinear systems*. Prentice Hall, third edition, 2002.
- Mohammad Amin Khodkar, Pedram Hassanzadeh, and Athanasios Antoulas. A koopman-based framework for forecasting the spatiotemporal evolution of chaotic dynamics with nonlinearities modeled as exogenous forcings. *arXiv preprint arXiv:1909.00076*, 2019.
- J Zico Kolter and Gaurav Manek. Learning stable deep dynamics models. *Advances in neural information processing systems*, 32, 2019.
- Ian Lenz, Ross A Knepper, and Ashutosh Saxena. Deepmpc: Learning deep latent features for model predictive control. In *Robotics: Science and Systems*, volume 10. Rome, Italy, 2015.
- Angran Li, Ruijia Chen, Amir Barati Farimani, and Yongjie Jessica Zhang. Reaction diffusion system prediction based on convolutional neural network. *Scientific reports*, 10(1):3894, 2020.
- Anders Logg, Kent-Andre Mardal, and Garth Wells. *Automated solution of differential equations by the finite element method: The FEniCS book*, volume 84. Springer Science & Business Media, 2012.
- Zichao Long, Yiping Lu, Xianzhong Ma, and Bin Dong. Pde-net: Learning pdes from data. In *International Conference on Machine Learning*, pp. 3208–3216. PMLR, 2018.
- Bethany Lusch, J Nathan Kutz, and Steven L Brunton. Deep learning for universal linear embeddings of nonlinear dynamics. *Nature communications*, 9(1):4950, 2018.
- Pingchuan Ma, Yunsheng Tian, Zherong Pan, Bo Ren, and Dinesh Manocha. Fluid directed rigid body control using deep reinforcement learning. *ACM Transactions on Graphics (TOG)*, 37(4):1–11, 2018.
- Jan R Magnus and Heinz Neudecker. Symmetry, 0-1 matrices and jacobians: A review. *Econometric Theory*, 2(2):157–190, 1986.
- George Matsaglia and George PH Styán. Equalities and inequalities for ranks of matrices. *Linear and multilinear Algebra*, 2(3):269–292, 1974.
- Jeremy Morton, Antony Jameson, Mykel J Kochenderfer, and Freddie Witherden. Deep dynamical modeling and control of unsteady fluid flows. *Advances in Neural Information Processing Systems*, 31, 2018.
- Zuwei Ping, Zhun Yin, Xiuting Li, Yefeng Liu, and Tao Yang. Deep koopman model predictive control for enhancing transient stability in power grids. *International Journal of Robust and Nonlinear Control*, 31(6):1964–1978, 2021.
- Joshua L Proctor, Steven L Brunton, and J Nathan Kutz. Dynamic mode decomposition with control. *SIAM Journal on Applied Dynamical Systems*, 15(1):142–161, 2016.
- Jean Rabault, Miroslav Kuchta, Atle Jensen, Ulysse Réglade, and Nicolas Cerardi. Artificial neural networks trained through deep reinforcement learning discover control strategies for active flow control. *Journal of fluid mechanics*, 865:281–302, 2019.
- Maziar Raissi. Deep hidden physics models: Deep learning of nonlinear partial differential equations. *The Journal of Machine Learning Research*, 19(1):932–955, 2018.
- Maziar Raissi, Paris Perdikaris, and George E Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational physics*, 378:686–707, 2019.
- Xiaoli Ren, Xiaoyong Li, Kaijun Ren, Junqiang Song, Zichen Xu, Kefeng Deng, and Xiang Wang. Deep learning-based weather prediction: a survey. *Big Data Research*, 23:100178, 2021.
- Clarence W Rowley, Igor Mezić, Shervin Bagheri, Philipp Schlatter, and Dan S Henningson. Spectral analysis of nonlinear flows. *Journal of fluid mechanics*, 641:115–127, 2009.

- Priyabrata Saha, Magnus Egerstedt, and Saibal Mukhopadhyay. Neural identification for control. *IEEE Robotics and Automation Letters*, 6(3):4648–4655, 2021.
- Michael Schäfer, Stefan Turek, Franz Durst, Egon Krause, and Rolf Rannacher. *Benchmark computations of laminar flow around a cylinder*. Springer, 1996.
- Sebastian Scher. Toward data-driven weather and climate forecasting: Approximating a simple general circulation model with deep learning. *Geophysical Research Letters*, 45(22):12–616, 2018.
- Peter J Schmid. Dynamic mode decomposition of numerical and experimental data. *Journal of fluid mechanics*, 656:5–28, 2010.
- Sungyong Seo, Chuizheng Meng, and Yan Liu. Physics-aware difference graph networks for sparsely-observed dynamics. In *International Conference on Learning Representations*, 2019.
- Eduardo D Sontag. *Mathematical control theory: deterministic finite dimensional systems*, volume 6. Springer Science & Business Media, 2013.
- Prem A Srinivasan, L Guastoni, Hossein Azizpour, PHILIPP Schlatter, and Ricardo Vinuesa. Predictions of turbulent shear flows using deep neural networks. *Physical Review Fluids*, 4(5):054603, 2019.
- Tetsuya Takahashi, Junbang Liang, Yi-Ling Qiao, and Ming C Lin. Differentiable fluids with solid coupling for learning and control. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35(7), pp. 6138–6146, 2021.
- Naoya Takeishi, Yoshinobu Kawahara, and Takehisa Yairi. Learning koopman invariant subspaces for dynamic mode decomposition. *Advances in neural information processing systems*, 30, 2017.
- Hongwei Tang, Jean Rabault, Alexander Kuhnle, Yan Wang, and Tongguang Wang. Robust active flow control over a range of reynolds numbers using an artificial neural network trained through deep reinforcement learning. *Physics of Fluids*, 32(5):053605, 2020.
- Jonathan H. Tu, , Clarence W. Rowley, Dirk M. Luchtenburg, Steven L. Brunton, and J. Nathan Kutz and. On dynamic mode decomposition: Theory and applications. *Journal of Computational Dynamics*, 1(2):391–421, 2014. doi: 10.3934/jcd.2014.1.391.
- Karen Willcox and Jaime Peraire. Balanced model reduction via the proper orthogonal decomposition. *AIAA journal*, 40(11):2323–2330, 2002.
- Charles HK Williamson. Vortex dynamics in the cylinder wake. *Annual review of fluid mechanics*, 28(1):477–539, 1996.
- SHI Xingjian, Zhouong Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *Advances in neural information processing systems*, pp. 802–810, 2015.
- Yuting Yang, Junyu Dong, Xin Sun, Estanislau Lima, Quanquan Mu, and Xinhua Wang. A cfcc-lstm model for sea surface temperature prediction. *IEEE Geoscience and Remote Sensing Letters*, 15(2):207–211, 2017.
- Enoch Yeung, Soumya Kundu, and Nathan Hodas. Learning deep neural network representations for koopman operators of nonlinear dynamical systems. In *2019 American Control Conference (ACC)*, pp. 4832–4839. IEEE, 2019.
- Ruiyang Zhang, Yang Liu, and Hao Sun. Physics-informed multi-lstm networks for metamodeling of nonlinear structures. *Computer Methods in Applied Mechanics and Engineering*, 369:113226, 2020.

## Appendices

### A Proofs

This section details the proofs for the results presented in section 4. The proof of theorem 4.1.1 uses the following properties of the rank (denoted by  $\text{rank}(\cdot)$ ), the Kronecker product (denoted by  $\otimes$ ) and vectorization of matrices (denoted by  $\text{vec}(\cdot)$ ). All the definitions and properties are presented in the context of matrices over real numbers.

For any conformable matrices  $\mathbf{D}$  and  $\mathbf{E}$  such that  $\mathbf{E}$  has full row-rank,

$$\text{rank}(\mathbf{D}\mathbf{E}) = \text{rank}(\mathbf{D}). \quad (36a)$$

For any real matrix  $\mathbf{D}$ ,

$$\text{rank}(\mathbf{D}^\top \mathbf{D}) = \text{rank}(\mathbf{D}\mathbf{D}^\top) = \text{rank}(\mathbf{D}^\top) = \text{rank}(\mathbf{D}). \quad (36b)$$

For any matrices (of compatible dimensions)  $\mathbf{D}, \mathbf{E}, \mathbf{F}$ , and  $\mathbf{H}$ ,

$$\text{vec}(\mathbf{D}\mathbf{E}\mathbf{F}^\top) = (\mathbf{F} \otimes \mathbf{D})\text{vec}(\mathbf{E}), \quad (37a)$$

$$(\mathbf{D} \otimes \mathbf{E})^\top = \mathbf{D}^\top \otimes \mathbf{E}^\top, \quad (37b)$$

$$(\mathbf{D} \otimes \mathbf{E})(\mathbf{F} \otimes \mathbf{H}) = (\mathbf{D}\mathbf{F} \otimes \mathbf{E}\mathbf{H}), \quad (37c)$$

whenever these quantities are defined. Furthermore, if  $\mathbf{D}$  and  $\mathbf{E}$  are symmetric and positive semidefinite (resp. positive definite), then  $\mathbf{D} \otimes \mathbf{E}$  is symmetric and positive semidefinite (resp. positive definite), i.e.,

$$\mathbf{D} \succeq 0, \mathbf{E} \succeq 0 \implies (\mathbf{D} \otimes \mathbf{E}) \succeq 0; \quad \mathbf{D} \succ 0, \mathbf{E} \succ 0 \implies (\mathbf{D} \otimes \mathbf{E}) \succ 0. \quad (37d)$$

Proofs of (36) and (37) can be found in (Matsaglia & PH Styan (1974)) and (Magnus & Neudecker (1986)), respectively.

To derive the results presented in corollary (4.1.1.1), we use the following definitions of the Moore-Penrose inverse of a matrix (denoted by  $(\cdot)^+$ ). For any matrix  $\mathbf{D}$  and its (full) SVD, i.e.,  $\mathbf{D} = \mathbf{U}_D \boldsymbol{\Sigma}_D \mathbf{V}_D^\top$ ,

$$\mathbf{D}^+ = (\mathbf{D}^\top \mathbf{D})^{-1} \mathbf{D}^\top, \quad \text{when } (\mathbf{D}^\top \mathbf{D})^{-1} \text{ exists,} \quad (38a)$$

$$\mathbf{D}^+ = \mathbf{D}^\top (\mathbf{D}\mathbf{D}^\top)^{-1}, \quad \text{when } (\mathbf{D}\mathbf{D}^\top)^{-1} \text{ exists,} \quad (38b)$$

$$\mathbf{D}^+ = \mathbf{V}_D \boldsymbol{\Sigma}_D^+ \mathbf{U}_D^\top, \quad (38c)$$

$$\mathbf{D}^+ = \lim_{\varepsilon \rightarrow 0} (\mathbf{D}^\top \mathbf{D} + \varepsilon^2 \mathbf{I})^{-1} \mathbf{D}^\top = \lim_{\varepsilon \rightarrow 0} \mathbf{D}^\top (\mathbf{D}\mathbf{D}^\top + \varepsilon^2 \mathbf{I})^{-1}, \quad (38d)$$

where  $\mathbf{I}$  is the identity matrix of compatible dimension. The proof of (38d) can be found in (Albert (1972)).

To prove Theorem 4.1.1, we use some well-known results, summarized as the following lemma in (Baldi & Hornik (1989)), for linear least-squares optimization.

**Lemma A.0.1.** *The quadratic function  $L(\mathbf{z}) = \|\mathbf{y} - \mathbf{M}\mathbf{z}\|^2 = \mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{M}\mathbf{z} + \mathbf{z}^\top \mathbf{M}^\top \mathbf{M}\mathbf{z}$  is convex, and a point  $\mathbf{z}$  globally minimizes  $L$  if and only if  $\nabla L(\mathbf{z}) = 0$ , or equivalently,  $\mathbf{M}^\top \mathbf{M}\mathbf{z} = \mathbf{M}^\top \mathbf{y}$ . Furthermore, if  $\mathbf{M}^\top \mathbf{M} \succ 0$ , i.e., positive definite, then  $L$  is strictly convex and reaches its unique minimum for  $\mathbf{z} = (\mathbf{M}^\top \mathbf{M})^{-1} \mathbf{M}^\top \mathbf{y}$ .*

#### A.1 Proof of theorem 4.1.1

**Theorem 4.1.1.** *Consider the following objective function*

$$L_{\text{pred}}(\mathbf{E}_x, \mathbf{G}) = \frac{1}{n} \sum_{i=0}^{n-1} \|\mathbf{E}_x \mathbf{x}(t_{i+1}) - \mathbf{G}\mathbf{E}_{xu} \boldsymbol{\omega}(t_i)\|^2, \quad (12)$$

where  $\mathbf{G} = [\mathbf{A}_R \ \mathbf{B}_R] \in \mathbb{R}^{r_x \times (r_x + d_u)}$ ,  $\mathbf{E}_{xu} = \begin{bmatrix} \mathbf{E}_x & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d_u} \end{bmatrix} \in \mathbb{R}^{(r_x + d_u) \times (d_x + d_u)}$ ,  $\mathbf{I}_{d_u}$  being the identity matrix of order  $d_u$ . For any fixed matrix  $\mathbf{E}_x$ , the objective function  $L_{\text{pred}}$  is convex in the coefficients of  $\mathbf{G}$  and attains its minimum for any  $\mathbf{G}$  satisfying

$$\mathbf{G}\mathbf{E}_{xu}\boldsymbol{\Omega}\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top = \mathbf{E}_x\mathbf{Y}\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top, \quad (13)$$

where  $\mathbf{Y}$  and  $\boldsymbol{\Omega}$  are the data matrices as defined in section (3.3). If  $\mathbf{E}_x$  has full rank  $r_x$ , and  $\boldsymbol{\Omega}\boldsymbol{\Omega}^\top$  is non-singular, then  $L_{\text{pred}}$  is strictly convex and has a unique minimum for

$$\mathbf{G} = [\mathbf{A}_R \ \mathbf{B}_R] = \mathbf{E}_x\mathbf{Y}\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top(\mathbf{E}_{xu}\boldsymbol{\Omega}\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top)^{-1}. \quad (14)$$

*Proof.* We can write  $L_{\text{pred}}(\mathbf{E}_x, \mathbf{G})$  as follows,

$$\begin{aligned} L_{\text{pred}}(\mathbf{E}_x, \mathbf{G}) &= \frac{1}{n} \sum_{i=0}^{n-1} \|\mathbf{E}_x\mathbf{x}(t_{i+1}) - \mathbf{G}\mathbf{E}_{xu}\boldsymbol{\omega}(t_i)\|^2 \\ &= \|\text{vec}(\mathbf{E}_x\mathbf{Y}) - \text{vec}(\mathbf{G}\mathbf{E}_{xu}\boldsymbol{\Omega})\|^2 \\ &= \|\text{vec}(\mathbf{E}_x\mathbf{Y}) - (\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top \otimes \mathbf{I}_{r_x})\text{vec}(\mathbf{G})\|^2. \end{aligned} \quad (39)$$

The third equality is obtained using (37a). For fixed  $\mathbf{E}_x$ , we can apply Lemma A.0.1 to (39): (39) is convex in coefficient of  $\mathbf{G}$ , and  $\mathbf{G}$  corresponds to a global minimum of  $L_{\text{pred}}$  if and only if

$$(\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top \otimes \mathbf{I}_{r_x})^\top(\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top \otimes \mathbf{I}_{r_x})\text{vec}(\mathbf{G}) = (\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top \otimes \mathbf{I}_{r_x})^\top\text{vec}(\mathbf{E}_x\mathbf{Y}). \quad (40)$$

Using (37b) and (37c), we can write (40) as

$$(\mathbf{E}_{xu}\boldsymbol{\Omega}\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top \otimes \mathbf{I}_{r_x})\text{vec}(\mathbf{G}) = (\mathbf{E}_{xu}\boldsymbol{\Omega} \otimes \mathbf{I}_{r_x})\text{vec}(\mathbf{E}_x\mathbf{Y}). \quad (41)$$

Applying (37a) on (41), we get  $\mathbf{G}\mathbf{E}_{xu}\boldsymbol{\Omega}\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top = \mathbf{E}_x\mathbf{Y}\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top$ , i.e., (13).

If  $\mathbf{E}_x$  has full rank  $r_x$ , then  $\mathbf{E}_{xu} = \begin{bmatrix} \mathbf{E}_x & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d_u} \end{bmatrix} \in \mathbb{R}^{(r_x + d_u) \times (d_x + d_u)}$  has full rank  $(r_x + d_u)$ . If  $\boldsymbol{\Omega}\boldsymbol{\Omega}^\top \in \mathbb{R}^{(d_x + d_u) \times (d_x + d_u)}$  is non-singular, then  $\boldsymbol{\Omega}$  has full row-rank  $(d_x + d_u)$ . Consequently, using (36a) and (36b), we have

$$\text{rank}(\mathbf{E}_{xu}\boldsymbol{\Omega}\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top) = \text{rank}(\mathbf{E}_{xu}\boldsymbol{\Omega}) = \text{rank}(\mathbf{E}_{xu}) = r_x + d_u. \quad (42)$$

Hence the symmetric positive semidefinite matrix  $\mathbf{E}_{xu}\boldsymbol{\Omega}\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top$  has full rank and therefore positive definite. Using (37b), (37c), and (37d), we can see that  $(\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top \otimes \mathbf{I}_{r_x})^\top(\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top \otimes \mathbf{I}_{r_x}) = (\mathbf{E}_{xu}\boldsymbol{\Omega}\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top \otimes \mathbf{I}_{r_x})$  is positive definite as well. Therefore, by Lemma A.0.1, (39) is strictly convex in the coefficients of  $\mathbf{G}$  and has a unique minimum. Since  $\mathbf{E}_{xu}\boldsymbol{\Omega}\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top \succ 0$ , it is invertible. Hence, from (13), we can say that the unique minimum of (39) is reached at  $\mathbf{G} = \mathbf{E}_x\mathbf{Y}\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top(\mathbf{E}_{xu}\boldsymbol{\Omega}\boldsymbol{\Omega}^\top\mathbf{E}_{xu}^\top)^{-1}$ , i.e., (14).  $\blacksquare$

## A.2 An alternative representation of (14)

Here we provide a possible alternative representation of (14) required to prove corollary 4.1.1.1.

**Lemma A.2.1.** Consider the (full) SVD of the data matrix  $\boldsymbol{\Omega}$  given by  $\boldsymbol{\Omega} = \mathbf{U}_\Omega\boldsymbol{\Sigma}_\Omega\mathbf{V}_\Omega^\top$ , where  $\mathbf{U}_\Omega \in \mathbb{R}^{(d_x + d_u) \times (d_x + d_u)}$ ,  $\boldsymbol{\Sigma}_\Omega \in \mathbb{R}^{(d_x + d_u) \times n}$ , and  $\mathbf{V}_\Omega \in \mathbb{R}^{n \times n}$ . (14) can be expressed as

$$\mathbf{G} = \lim_{\varepsilon \rightarrow 0} \mathbf{E}_x\mathbf{Y}\mathbf{V}_\Omega(\boldsymbol{\Sigma}_\Omega^\top\mathbf{U}_\Omega^\top\mathbf{E}_{xu}^\top\mathbf{E}_{xu}\mathbf{U}_\Omega\boldsymbol{\Sigma}_\Omega + \varepsilon^2\mathbf{I}_n)^{-1}\boldsymbol{\Sigma}_\Omega^\top\mathbf{U}_\Omega^\top\mathbf{E}_{xu}^\top. \quad (43)$$

*Proof.* Replacing  $\boldsymbol{\Omega}$  with its SVD in (14) we get,

$$\begin{aligned} \mathbf{G} &= \mathbf{E}_x\mathbf{Y}\mathbf{V}_\Omega\boldsymbol{\Sigma}_\Omega^\top\mathbf{U}_\Omega^\top\mathbf{E}_{xu}^\top(\mathbf{E}_{xu}\mathbf{U}_\Omega\boldsymbol{\Sigma}_\Omega\mathbf{V}_\Omega^\top\mathbf{V}_\Omega\boldsymbol{\Sigma}_\Omega^\top\mathbf{U}_\Omega^\top\mathbf{E}_{xu}^\top)^{-1} \\ &= \mathbf{E}_x\mathbf{Y}\mathbf{V}_\Omega\boldsymbol{\Sigma}_\Omega^\top\mathbf{U}_\Omega^\top\mathbf{E}_{xu}^\top(\mathbf{E}_{xu}\mathbf{U}_\Omega\boldsymbol{\Sigma}_\Omega\boldsymbol{\Sigma}_\Omega^\top\mathbf{U}_\Omega^\top\mathbf{E}_{xu}^\top)^{-1} \\ &= \mathbf{E}_x\mathbf{Y}\mathbf{V}_\Omega(\mathbf{E}_{xu}\mathbf{U}_\Omega\boldsymbol{\Sigma}_\Omega)^\dagger \end{aligned} \quad (44)$$

The second equality is due to the orthogonality of  $\mathbf{V}_\Omega$ . The third equality is obtained using (38b). Substituting  $(\mathbf{E}_{\mathbf{x}u}\mathbf{U}_\Omega\Sigma_\Omega)^+$  with the *limit definition* (38d) of the Moore-Penrose inverse, we get

$$\mathbf{G} = \lim_{\varepsilon \rightarrow 0} \mathbf{E}_{\mathbf{x}} \mathbf{Y} \mathbf{V}_\Omega (\Sigma_\Omega^\top \mathbf{U}_\Omega^\top \mathbf{E}_{\mathbf{x}u}^\top \mathbf{E}_{\mathbf{x}u} \mathbf{U}_\Omega \Sigma_\Omega + \varepsilon^2 \mathbf{I}_n)^{-1} \Sigma_\Omega^\top \mathbf{U}_\Omega^\top \mathbf{E}_{\mathbf{x}u}^\top. \quad (45)$$

### A.3 Proof of Corollary 4.1.1.1

**Corollary 4.1.1.1.** *Consider the (full) SVD of the data matrix  $\Omega$  given by  $\Omega = \mathbf{U}_\Omega \Sigma_\Omega \mathbf{V}_\Omega^\top$ , where  $\mathbf{U}_\Omega \in \mathbb{R}^{(d_{\mathbf{x}}+d_{\mathbf{u}}) \times (d_{\mathbf{x}}+d_{\mathbf{u}})}$ ,  $\Sigma_\Omega \in \mathbb{R}^{(d_{\mathbf{x}}+d_{\mathbf{u}}) \times n}$ , and  $\mathbf{V}_\Omega \in \mathbb{R}^{n \times n}$ . If  $\mathbf{E}_{\mathbf{x}} = \widehat{\mathbf{U}}_{\mathbf{Y}}^\top$  and  $\Omega\Omega^\top$  is non-singular, then the solution for  $\mathbf{G} = [\mathbf{A}_R \ \mathbf{B}_R]$  corresponding to the unique minimum of  $L_{\text{pred}}$  can be expressed as*

$$\mathbf{A}_R = \widehat{\mathbf{U}}_{\mathbf{Y}}^\top \mathbf{Y} \mathbf{V}_\Omega \Sigma^* \mathbf{U}_{\Omega,1}^\top \widehat{\mathbf{U}}_{\mathbf{Y}}, \quad \text{and} \quad \mathbf{B}_R = \widehat{\mathbf{U}}_{\mathbf{Y}}^\top \mathbf{Y} \mathbf{V}_\Omega \Sigma^* \mathbf{U}_{\Omega,2}^\top, \quad (15)$$

where  $[\mathbf{U}_{\Omega,1}^\top \ \mathbf{U}_{\Omega,2}^\top] = \mathbf{U}_\Omega^\top$  with  $\mathbf{U}_{\Omega,1} \in \mathbb{R}^{d_{\mathbf{x}} \times (d_{\mathbf{x}}+d_{\mathbf{u}})}$ ,  $\mathbf{U}_{\Omega,2} \in \mathbb{R}^{d_{\mathbf{u}} \times (d_{\mathbf{x}}+d_{\mathbf{u}})}$ , and  $\Sigma^* = \lim_{\varepsilon \rightarrow 0} (\Sigma_\Omega^\top \mathbf{U}_{\Omega,1}^\top \widehat{\mathbf{U}}_{\mathbf{Y}} \widehat{\mathbf{U}}_{\mathbf{Y}}^\top \mathbf{U}_{\Omega,1} \Sigma_\Omega + \Sigma_\Omega^\top \mathbf{U}_{\Omega,2}^\top \mathbf{U}_{\Omega,2} \Sigma_\Omega + \varepsilon^2 \mathbf{I}_n)^{-1} \Sigma_\Omega^\top$ .

*Proof.* By the definition of truncated SVD, the columns of  $\widehat{\mathbf{U}}_{\mathbf{Y}}$  are orthonormal. Therefore,  $\widehat{\mathbf{U}}_{\mathbf{Y}}^\top$  has full row-rank  $r_{\mathbf{x}}$ . Hence, by theorem 4.1.1 and lemma A.2.1, if  $\mathbf{E}_{\mathbf{x}} = \widehat{\mathbf{U}}_{\mathbf{Y}}^\top$ , and  $\Omega\Omega^\top$  is non-singular, then the unique minimum of  $L_{\text{pred}}$ , is reached when

$$\mathbf{G} = \widehat{\mathbf{U}}_{\mathbf{Y}}^\top \mathbf{Y} \mathbf{V}_\Omega (\mathbf{E}_{\mathbf{x}u} \mathbf{U}_\Omega \Sigma_\Omega)^+ = \lim_{\varepsilon \rightarrow 0} \widehat{\mathbf{U}}_{\mathbf{Y}}^\top \mathbf{Y} \mathbf{V}_\Omega (\Sigma_\Omega^\top \mathbf{U}_\Omega^\top \mathbf{E}_{\mathbf{x}u}^\top \mathbf{E}_{\mathbf{x}u} \mathbf{U}_\Omega \Sigma_\Omega + \varepsilon^2 \mathbf{I}_n)^{-1} \Sigma_\Omega^\top \mathbf{U}_\Omega^\top \mathbf{E}_{\mathbf{x}u}^\top. \quad (46)$$

Now, substituting  $\mathbf{E}_{\mathbf{x}} = \widehat{\mathbf{U}}_{\mathbf{Y}}^\top$  in  $\mathbf{E}_{\mathbf{x}u}$ , and using the partition  $\mathbf{U}_\Omega^\top = [\mathbf{U}_{\Omega,1}^\top \ \mathbf{U}_{\Omega,2}^\top]$ , where  $\mathbf{U}_{\Omega,1} \in \mathbb{R}^{d_{\mathbf{x}} \times (d_{\mathbf{x}}+d_{\mathbf{u}})}$ ,  $\mathbf{U}_{\Omega,2} \in \mathbb{R}^{d_{\mathbf{u}} \times (d_{\mathbf{x}}+d_{\mathbf{u}})}$ , we get

$$\mathbf{E}_{\mathbf{x}u} \mathbf{U}_\Omega = \begin{bmatrix} \widehat{\mathbf{U}}_{\mathbf{Y}}^\top & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d_{\mathbf{u}}} \end{bmatrix} \begin{bmatrix} \mathbf{U}_{\Omega,1} \\ \mathbf{U}_{\Omega,2} \end{bmatrix} = \begin{bmatrix} \widehat{\mathbf{U}}_{\mathbf{Y}}^\top \mathbf{U}_{\Omega,1} \\ \mathbf{U}_{\Omega,2} \end{bmatrix}, \quad (47)$$

and

$$\mathbf{U}_\Omega^\top \mathbf{E}_{\mathbf{x}u}^\top \mathbf{E}_{\mathbf{x}u} \mathbf{U}_\Omega = \begin{bmatrix} \mathbf{U}_{\Omega,1}^\top \widehat{\mathbf{U}}_{\mathbf{Y}} & \mathbf{U}_{\Omega,2}^\top \end{bmatrix} \begin{bmatrix} \widehat{\mathbf{U}}_{\mathbf{Y}}^\top \mathbf{U}_{\Omega,1} \\ \mathbf{U}_{\Omega,2} \end{bmatrix} = \mathbf{U}_{\Omega,1}^\top \widehat{\mathbf{U}}_{\mathbf{Y}} \widehat{\mathbf{U}}_{\mathbf{Y}}^\top \mathbf{U}_{\Omega,1} + \mathbf{U}_{\Omega,2}^\top \mathbf{U}_{\Omega,2}. \quad (48)$$

Plugging (47) and (48) into (46) leads to

$$\mathbf{G} = \lim_{\varepsilon \rightarrow 0} \widehat{\mathbf{U}}_{\mathbf{Y}}^\top \mathbf{Y} \mathbf{V}_\Omega (\Sigma_\Omega^\top \mathbf{U}_{\Omega,1}^\top \widehat{\mathbf{U}}_{\mathbf{Y}} \widehat{\mathbf{U}}_{\mathbf{Y}}^\top \mathbf{U}_{\Omega,1} \Sigma_\Omega + \Sigma_\Omega^\top \mathbf{U}_{\Omega,2}^\top \mathbf{U}_{\Omega,2} \Sigma_\Omega + \varepsilon^2 \mathbf{I}_n)^{-1} \Sigma_\Omega^\top \begin{bmatrix} \mathbf{U}_{\Omega,1}^\top \widehat{\mathbf{U}}_{\mathbf{Y}} & \mathbf{U}_{\Omega,2}^\top \end{bmatrix}. \quad (49)$$

Defining  $\Sigma^* \triangleq \lim_{\varepsilon \rightarrow 0} (\Sigma_\Omega^\top \mathbf{U}_{\Omega,1}^\top \widehat{\mathbf{U}}_{\mathbf{Y}} \widehat{\mathbf{U}}_{\mathbf{Y}}^\top \mathbf{U}_{\Omega,1} \Sigma_\Omega + \Sigma_\Omega^\top \mathbf{U}_{\Omega,2}^\top \mathbf{U}_{\Omega,2} \Sigma_\Omega + \varepsilon^2 \mathbf{I}_n)^{-1} \Sigma_\Omega^\top$ , we can split (49) into

$$\mathbf{A}_R = \widehat{\mathbf{U}}_{\mathbf{Y}}^\top \mathbf{Y} \mathbf{V}_\Omega \Sigma^* \mathbf{U}_{\Omega,1}^\top \widehat{\mathbf{U}}_{\mathbf{Y}}, \quad \text{and} \quad \mathbf{B}_R = \widehat{\mathbf{U}}_{\mathbf{Y}}^\top \mathbf{Y} \mathbf{V}_\Omega \Sigma^* \mathbf{U}_{\Omega,2}^\top,$$

which is (15). ■

### A.4 The case when $\Omega\Omega^\top$ not invertible

When the covariance matrix  $\Omega\Omega^\top$  is not invertible, which is always true if  $n < d_{\mathbf{x}} + d_{\mathbf{u}}$ , the matrix  $\mathbf{E}_{\mathbf{x}u} \Omega \Omega^\top \mathbf{E}_{\mathbf{x}u}^\top$  is not guaranteed to be invertible. In that case, the minimum of  $L_{\text{pred}}$  corresponds to infinitely many solutions for  $\mathbf{G}$ . However, minimizing  $L_{\text{pred}}$  with added  $\ell_2$  regularization, i.e.,  $L_{\text{pred,reg}}(\mathbf{E}_{\mathbf{x}}, \mathbf{G}) = L_{\text{pred}}(\mathbf{E}_{\mathbf{x}}, \mathbf{G}) + \beta \|\text{vec}(\mathbf{G})\|^2$  provides a unique solution for  $\mathbf{G}$ , for a fixed  $\mathbf{E}_{\mathbf{x}}$ . We have the following result.

**Theorem A.4.1.** *For any fixed matrix  $\mathbf{E}_{\mathbf{x}}$  and  $\beta > 0$ , the objective function  $L_{\text{pred,reg}}(\mathbf{E}_{\mathbf{x}}, \mathbf{G}) = L_{\text{pred}}(\mathbf{E}_{\mathbf{x}}, \mathbf{G}) + \beta \|\text{vec}(\mathbf{G})\|^2$  is strictly convex in the coefficients of  $\mathbf{G}$ , and the global minimum of  $L_{\text{pred,reg}}$  corresponds to the unique solution for  $\mathbf{G}$ , given by*

$$\mathbf{G} = \mathbf{E}_{\mathbf{x}} \mathbf{Y} \Omega^\top \mathbf{E}_{\mathbf{x}u}^\top (\mathbf{E}_{\mathbf{x}u} \Omega \Omega^\top \mathbf{E}_{\mathbf{x}u}^\top + \beta \mathbf{I}_{r_{\mathbf{x}}+d_{\mathbf{u}}})^{-1}. \quad (50)$$

*Proof.*  $L_{\text{pred,reg}}(\mathbf{E}_x, \mathbf{G})$  can be written as, using (37a-c),

$$\begin{aligned} L_{\text{pred,reg}}(\mathbf{E}_x, \mathbf{G}) &= \|\text{vec}(\mathbf{E}_x \mathbf{Y}) - (\boldsymbol{\Omega}^\top \mathbf{E}_{x\mathbf{u}}^\top \otimes \mathbf{I}_{r_x}) \text{vec}(\mathbf{G})\|^2 + \beta \|\text{vec}(\mathbf{G})\|^2 \\ &= \text{vec}(\mathbf{E}_x \mathbf{Y})^\top \text{vec}(\mathbf{E}_x \mathbf{Y}) - 2 \text{vec}(\mathbf{E}_x \mathbf{Y})^\top (\boldsymbol{\Omega}^\top \mathbf{E}_{x\mathbf{u}}^\top \otimes \mathbf{I}_{r_x}) \text{vec}(\mathbf{G}) \\ &\quad + \text{vec}(\mathbf{G})^\top (\mathbf{E}_{x\mathbf{u}} \boldsymbol{\Omega} \boldsymbol{\Omega}^\top \mathbf{E}_{x\mathbf{u}}^\top \otimes \mathbf{I}_{r_x} + \beta \mathbf{I}_{r_x(r_x+d_u)}) \text{vec}(\mathbf{G}) \end{aligned}$$

$\mathbf{E}_{x\mathbf{u}} \boldsymbol{\Omega} \boldsymbol{\Omega}^\top \mathbf{E}_{x\mathbf{u}}^\top$  is a symmetric positive semidefinite matrix, irrespective of whether it has full rank or not. Hence, by (37d),  $\mathbf{E}_{x\mathbf{u}} \boldsymbol{\Omega} \boldsymbol{\Omega}^\top \mathbf{E}_{x\mathbf{u}}^\top \otimes \mathbf{I}_{r_x}$  is symmetric positive semidefinite. Consequently, for any  $\beta > 0$ ,  $\mathbf{E}_{x\mathbf{u}} \boldsymbol{\Omega} \boldsymbol{\Omega}^\top \mathbf{E}_{x\mathbf{u}}^\top \otimes \mathbf{I}_{r_x} + \beta \mathbf{I}_{r_x(r_x+d_u)}$  is positive definite. According to lemma A.0.1,  $L_{\text{pred,reg}}$  is therefore strictly convex in the coefficients of  $\mathbf{G}$  and globally minimized when  $\nabla L_{\text{pred,reg}} = 0$ . The unique solution of (50) can be derived in the same manner as theorem 4.1.1. ■

**Remark.** Replacing  $\boldsymbol{\Omega}$  with its SVD in (50) we get,

$$\mathbf{G} = \mathbf{E}_x \mathbf{Y} \mathbf{V}_\Omega \boldsymbol{\Sigma}_\Omega^\top \mathbf{U}_\Omega^\top \mathbf{E}_{x\mathbf{u}}^\top (\mathbf{E}_{x\mathbf{u}} \mathbf{U}_\Omega \boldsymbol{\Sigma}_\Omega \boldsymbol{\Sigma}_\Omega^\top \mathbf{U}_\Omega^\top \mathbf{E}_{x\mathbf{u}}^\top + \beta \mathbf{I}_{r_x+d_u})^{-1}. \quad (51)$$

In the limit  $\beta \rightarrow 0^+$ , (51) converges to (44).

## A.5 DMDc through a linear autoencoding structure

Here we present a linear autoencoding structure that leads to a linear ROM exactly resembling the DMDc solution when  $\mathbf{E}_x = \widehat{\mathbf{U}}_Y^\top$ . However, its DNN-based nonlinear counterpart does not actually offer dimensionality reduction.

**Theorem A.5.1.** *Consider the following objective function*

$$L_{\text{pred,alt}}(\mathbf{E}_x, \tilde{\mathbf{G}}) = \frac{1}{n} \sum_{i=0}^{n-1} \|\mathbf{E}_x \mathbf{x}(t_{i+1}) - \tilde{\mathbf{G}} \boldsymbol{\omega}(t_i)\|^2, \quad (52)$$

where  $\tilde{\mathbf{G}} \in \mathbb{R}^{r_x \times (d_x+d_u)}$ . For any fixed matrix  $\mathbf{E}_x$ , the objective function  $L_{\text{pred,alt}}$  is convex in the coefficients of  $\tilde{\mathbf{G}}$  and attains its minimum for any  $\tilde{\mathbf{G}}$  satisfying

$$\tilde{\mathbf{G}} \boldsymbol{\Omega} \boldsymbol{\Omega}^\top = \mathbf{E}_x \mathbf{Y} \boldsymbol{\Omega}^\top, \quad (53)$$

where  $\mathbf{Y}$  and  $\boldsymbol{\Omega}$  are the data matrices as defined in section (3.3). If  $\boldsymbol{\Omega} \boldsymbol{\Omega}^\top$  is non-singular, then  $L_{\text{pred,alt}}$  is strictly convex and has a unique minimum for

$$\tilde{\mathbf{G}} = \mathbf{E}_x \mathbf{Y} \boldsymbol{\Omega}^\top (\boldsymbol{\Omega} \boldsymbol{\Omega}^\top)^{-1}. \quad (54)$$

*Proof.* The proof is very similar to the proof of theorem 4.1.1. Using (37a), we can write  $L_{\text{pred,alt}}(\mathbf{E}_x, \tilde{\mathbf{G}})$  as follows,

$$\begin{aligned} L_{\text{pred,alt}}(\mathbf{E}_x, \tilde{\mathbf{G}}) &= \frac{1}{n} \sum_{i=0}^{n-1} \|\mathbf{E}_x \mathbf{x}(t_{i+1}) - \tilde{\mathbf{G}} \boldsymbol{\omega}(t_i)\|^2 \\ &= \|\text{vec}(\mathbf{E}_x \mathbf{Y}) - \text{vec}(\tilde{\mathbf{G}} \boldsymbol{\Omega})\|^2 \\ &= \|\text{vec}(\mathbf{E}_x \mathbf{Y}) - (\boldsymbol{\Omega}^\top \otimes \mathbf{I}_{r_x}) \text{vec}(\tilde{\mathbf{G}})\|^2. \end{aligned} \quad (55)$$

For fixed  $\mathbf{E}_x$ , applying Lemma A.0.1 to (55), we can say  $L_{\text{pred,alt}}$  is convex in the coefficients of  $\tilde{\mathbf{G}}$ , and  $\tilde{\mathbf{G}}$  corresponds to a global minimum of  $L_{\text{pred,alt}}$  if and only if

$$(\boldsymbol{\Omega}^\top \otimes \mathbf{I}_{r_x})^\top (\boldsymbol{\Omega}^\top \otimes \mathbf{I}_{r_x}) \text{vec}(\tilde{\mathbf{G}}) = (\boldsymbol{\Omega}^\top \otimes \mathbf{I}_{r_x})^\top \text{vec}(\mathbf{E}_x \mathbf{Y}). \quad (56)$$

Using (37a-c), we can write (56) as  $\tilde{\mathbf{G}} \boldsymbol{\Omega} \boldsymbol{\Omega}^\top = \mathbf{E}_x \mathbf{Y} \boldsymbol{\Omega}^\top$ , which is (53).

If  $\boldsymbol{\Omega}\boldsymbol{\Omega}^\top$  is non-singular, then it is symmetric positive definite. Using (37b-d), we can see that  $(\boldsymbol{\Omega}^\top \otimes \mathbf{I}_{r_x})^\top (\boldsymbol{\Omega}^\top \otimes \mathbf{I}_{r_x}) = (\boldsymbol{\Omega}\boldsymbol{\Omega}^\top \otimes \mathbf{I}_{r_x})$  is positive definite as well. Therefore, by Lemma A.0.1, (55) is strictly convex in coefficient in  $\tilde{\mathbf{G}}$  and has a unique minimum. In that case, from (53), we can say that the unique minimum of (55) is reached at  $\tilde{\mathbf{G}} = \mathbf{E}_x \mathbf{Y} \boldsymbol{\Omega}^\top (\boldsymbol{\Omega}\boldsymbol{\Omega}^\top)^{-1}$ , i.e., (54). ■

**Corollary A.5.1.1.** *Consider the (full) SVD of the data matrix  $\boldsymbol{\Omega}$  given by  $\boldsymbol{\Omega} = \mathbf{U}_\Omega \boldsymbol{\Sigma}_\Omega \mathbf{V}_\Omega^\top$ , where  $\mathbf{U}_\Omega \in \mathbb{R}^{(d_x+d_u) \times (d_x+d_u)}$ ,  $\boldsymbol{\Sigma}_\Omega \in \mathbb{R}^{(d_x+d_u) \times n}$ , and  $\mathbf{V}_\Omega \in \mathbb{R}^{n \times n}$ . If  $\mathbf{E}_x = \hat{\mathbf{U}}_Y^\top$  and  $\boldsymbol{\Omega}\boldsymbol{\Omega}^\top$  is non-singular, then the solution for  $\tilde{\mathbf{G}}$  corresponding to the unique minimum of  $L_{\text{pred,alt}}$  can be expressed as*

$$\tilde{\mathbf{G}} = \hat{\mathbf{U}}_Y^\top \mathbf{Y} \mathbf{V}_\Omega \boldsymbol{\Sigma}_\Omega^+ \mathbf{U}_\Omega^\top. \quad (57)$$

*Proof.* By theorem A.5.1, if  $\mathbf{E}_x = \hat{\mathbf{U}}_Y^\top$ , and  $\boldsymbol{\Omega}\boldsymbol{\Omega}^\top$  is non-singular, then the unique minimum of  $L_{\text{pred,alt}}$  is reached when

$$\tilde{\mathbf{G}} = \hat{\mathbf{U}}_Y^\top \mathbf{Y} \boldsymbol{\Omega}^\top (\boldsymbol{\Omega}\boldsymbol{\Omega}^\top)^{-1} = \hat{\mathbf{U}}_Y^\top \mathbf{Y} \boldsymbol{\Omega}^+ \quad (58)$$

The second equality is due to (38b). Substituting  $\boldsymbol{\Omega}^+$  with its SVD definition (38c) into (58), we get  $\hat{\mathbf{U}}_Y^\top \mathbf{Y} \mathbf{V}_\Omega \boldsymbol{\Sigma}_\Omega^+ \mathbf{U}_\Omega^\top$ , which is (57). ■

**Remark.** From (52), it can be seen that  $\tilde{\mathbf{G}}$  maps the concatenated vector,  $\boldsymbol{\omega}(t_i)$ , of full state and actuation to the next reduce state  $\mathbf{x}_R(t_{i+1})$ . We can partition (57) as  $\tilde{\mathbf{G}} = \hat{\mathbf{U}}_Y^\top \mathbf{Y} \mathbf{V}_\Omega \boldsymbol{\Sigma}_\Omega^+ [\mathbf{U}_{\Omega,1}^\top \quad \mathbf{U}_{\Omega,2}^\top] = [\tilde{\mathbf{A}} \quad \tilde{\mathbf{B}}]$  to separate out the blocks corresponding to state and actuation. Here,  $\mathbf{U}_{\Omega,1}, \mathbf{U}_{\Omega,2}$  are the same as defined in corollary 4.1.1.1, and  $\tilde{\mathbf{A}} \in \mathbb{R}^{r_x \times d_x}$ ,  $\tilde{\mathbf{B}} \in \mathbb{R}^{r_x \times d_u}$ . Now, if we post-multiply  $\tilde{\mathbf{A}}$  with  $\mathbf{E}_x^\top = \hat{\mathbf{U}}_Y \in \mathbb{R}^{d_x \times r_x}$ , we get a ROM

$$\tilde{\mathbf{A}}_R = \tilde{\mathbf{A}} \hat{\mathbf{U}}_Y = \hat{\mathbf{U}}_Y^\top \mathbf{Y} \mathbf{V}_\Omega \boldsymbol{\Sigma}_\Omega^+ \mathbf{U}_{\Omega,1}^\top \hat{\mathbf{U}}_Y, \quad \tilde{\mathbf{B}}_R = \tilde{\mathbf{B}} = \hat{\mathbf{U}}_Y^\top \mathbf{Y} \mathbf{V}_\Omega \boldsymbol{\Sigma}_\Omega^+ \mathbf{U}_{\Omega,2}^\top, \quad (59)$$

which maps the current reduced state  $\mathbf{x}_R(t_i)$  and actuation  $\mathbf{u}(t_i)$  to the next reduced state  $\mathbf{x}_R(t_{i+1})$ . It can be verified easily that if we use the truncated SVD (as defined by 9), instead of the full SVD, for  $\boldsymbol{\Omega}$  in (58) and follow the similar steps afterward, we get an approximation of (59):

$$\hat{\mathbf{A}}_R = \hat{\mathbf{U}}_Y^\top \mathbf{Y} \hat{\mathbf{V}}_\Omega \hat{\boldsymbol{\Sigma}}_\Omega^{-1} \hat{\mathbf{U}}_{\Omega,1}^\top \hat{\mathbf{U}}_Y = \mathbf{A}_{R,\text{DMDc}}; \quad \hat{\mathbf{B}}_R = \hat{\mathbf{U}}_Y^\top \mathbf{Y} \hat{\mathbf{V}}_\Omega \hat{\boldsymbol{\Sigma}}_\Omega^{-1} \hat{\mathbf{U}}_{\Omega,2}^\top = \mathbf{B}_{R,\text{DMDc}}.$$

In summary, the aforementioned method can be carried out using gradient descent-based optimization and leads to the same ROM as DMDc, when  $\mathbf{E}_x = \hat{\mathbf{U}}_Y^\top$ . However, in this method, the benefit of dimensionality reduction is realized only when linear networks are used. A nonlinear counterpart (a DNN in the context of this paper) of  $\tilde{\mathbf{A}}_R$ , i.e., a nonlinear mapping from  $\mathbb{R}^{r_x}$  to  $\mathbb{R}^{r_x}$ , cannot be pre-computed from a nonlinear counterpart of  $\tilde{\mathbf{G}}$ , unlike the linear case (59). Consequently, we lose the benefit of dimensionality reduction when nonlinear networks are used.

## A.6 Proof of theorem 4.2.1

**Theorem 4.2.1.** *Consider the target dynamics defined by (28) and the candidate Lyapunov function defined by (29). Suppose the difference between the target dynamics and the closed-loop dynamics satisfies*

$$\|\mathcal{F}(\mathbf{x}_R, \mathcal{E}_u \circ \Pi(\mathbf{x}_R)) - \mathcal{F}_s(\mathbf{x}_R)\| \leq \delta < \frac{\alpha\theta\lambda_{\min}(\mathbf{K})}{2\lambda_{\max}(\mathbf{K})} \sqrt{\frac{\lambda_{\min}(\mathbf{K})}{\lambda_{\max}(\mathbf{K})}} \eta, \quad (30)$$

for all  $\mathbf{x}_R \in \mathbb{X}_R = \{\mathbf{x}_R \in \mathbb{R}^{r_x} \mid \|\mathbf{x}_R\| < \eta\}$  and  $0 < \theta < 1$ . Then, for all initial points satisfying  $\|\mathbf{x}_R(t_0)\| < \sqrt{\frac{\lambda_{\min}(\mathbf{K})}{\lambda_{\max}(\mathbf{K})}} \eta$ , the solution of the closed-loop ROM  $\frac{d\mathbf{x}_R}{dt} = \mathcal{F}(\mathbf{x}_R, \mathcal{E}_u \circ \Pi(\mathbf{x}_R))$  satisfies

$$\|\mathbf{x}_R(t)\| \leq \lambda e^{-\gamma(t-t_0)} \|\mathbf{x}_R(t_0)\|, \quad \forall t_0 \leq t < t_c + t_0 \quad (31)$$

and

$$\|\mathbf{x}_R(t)\| \leq \frac{2\delta}{\alpha\theta} \lambda^3, \quad \forall t \geq t_c + t_0 \quad (32)$$

for some finite  $t_c > 0$ , where

$$\gamma = \frac{\alpha(1-\theta)\lambda_{\min}(\mathbf{K})}{2\lambda_{\max}(\mathbf{K})} \quad \text{and} \quad \lambda = \sqrt{\frac{\lambda_{\max}(\mathbf{K})}{\lambda_{\min}(\mathbf{K})}} \quad (33)$$

*Proof.* From the definition of  $\mathcal{V}_R$ , we have

$$\lambda_{\min}(\mathbf{K})\|\mathbf{x}_R\|^2 \leq \mathcal{V}_R(\mathbf{x}_R) \leq \lambda_{\max}(\mathbf{K})\|\mathbf{x}_R\|^2, \quad \forall \mathbf{x}_R \in \mathbb{R}^{r_x}, \quad (60)$$

where  $\lambda_{\min}(\mathbf{K})$  and  $\lambda_{\max}(\mathbf{K})$  denote the smallest and largest eigenvalues, respectively, of  $\mathbf{K}$  and have positive values since the matrix  $\mathbf{K}$  is positive definite. Moreover, the definition of the target dynamics (28) implies

$$\begin{aligned} \nabla \mathcal{V}_R(\mathbf{x}_R)^\top \mathcal{F}_s(\mathbf{x}_R) &= \begin{cases} \nabla \mathcal{V}_R(\mathbf{x}_R)^\top \mathcal{P}(\mathbf{x}_R), & \text{if } \nabla \mathcal{V}_R(\mathbf{x}_R)^\top \mathcal{P}(\mathbf{x}_R) \leq -\alpha \mathcal{V}_R(\mathbf{x}_R) \\ \nabla \mathcal{V}_R(\mathbf{x}_R)^\top \mathcal{P}(\mathbf{x}_R) - \nabla \mathcal{V}_R(\mathbf{x}_R)^\top \frac{\nabla \mathcal{V}_R(\mathbf{x}_R)^\top \mathcal{P}(\mathbf{x}_R) + \alpha \mathcal{V}_R(\mathbf{x}_R)}{\|\nabla \mathcal{V}_R(\mathbf{x}_R)\|^2} \nabla \mathcal{V}_R(\mathbf{x}_R), & \text{otherwise} \end{cases} \\ &= \begin{cases} \nabla \mathcal{V}_R(\mathbf{x}_R)^\top \mathcal{P}(\mathbf{x}_R), & \text{if } \nabla \mathcal{V}_R(\mathbf{x}_R)^\top \mathcal{P}(\mathbf{x}_R) \leq -\alpha \mathcal{V}_R(\mathbf{x}_R) \\ -\alpha \mathcal{V}_R(\mathbf{x}_R), & \text{otherwise} \end{cases} \\ &\leq -\alpha \mathcal{V}_R(\mathbf{x}_R) \\ &\leq -\alpha \lambda_{\min}(\mathbf{K})\|\mathbf{x}_R\|^2, \quad \forall \mathbf{x}_R \in \mathbb{R}^{r_x}. \end{aligned} \quad (61)$$

The last inequality is due to (60).

Now, assume  $\mathcal{F}(\mathbf{x}_R, \mathcal{E}_u \circ \Pi(\mathbf{x}_R)) = \mathcal{H}(\mathbf{x}_R) = \mathcal{F}_s(\mathbf{x}_R) + \mathcal{J}(\mathbf{x}_R)$  for some function  $\mathcal{J} : \mathbb{R}^{r_x} \rightarrow \mathbb{R}^{r_x}$  and consider  $\mathcal{V}_R(\mathbf{x}_R) = \mathbf{x}_R^\top \mathbf{K} \mathbf{x}_R$  as a candidate Lyapunov function for

$$\frac{d\mathbf{x}_R}{dt} = \mathcal{H}(\mathbf{x}_R) = \mathcal{F}_s(\mathbf{x}_R) + \mathcal{J}(\mathbf{x}_R). \quad (62)$$

We have  $\|\nabla \mathcal{V}_R(\mathbf{x}_R)\| = \|2\mathbf{K}\mathbf{x}_R\| \leq 2\lambda_{\max}(\mathbf{K})\|\mathbf{x}_R\|$ . The time-derivative of  $\mathcal{V}_R$  along the trajectories of (62) satisfies

$$\begin{aligned} \frac{d\mathcal{V}_R}{dt} &= \nabla \mathcal{V}_R(\mathbf{x}_R)^\top \mathcal{F}_s(\mathbf{x}_R) + \nabla \mathcal{V}_R(\mathbf{x}_R)^\top \mathcal{J}(\mathbf{x}_R) \\ &\leq -\alpha \lambda_{\min}(\mathbf{K})\|\mathbf{x}_R\|^2 + \|\nabla \mathcal{V}_R(\mathbf{x}_R)\| \|\mathcal{J}(\mathbf{x}_R)\| \\ &\leq -\alpha \lambda_{\min}(\mathbf{K})\|\mathbf{x}_R\|^2 + 2\lambda_{\max}(\mathbf{K})\|\mathbf{x}_R\| \delta, \quad \forall \|\mathbf{x}_R\| < \eta \\ &= -\alpha(1-\theta)\lambda_{\min}(\mathbf{K})\|\mathbf{x}_R\|^2 - \alpha\theta\lambda_{\min}(\mathbf{K})\|\mathbf{x}_R\|^2 + 2\lambda_{\max}(\mathbf{K})\|\mathbf{x}_R\| \delta, \quad 0 < \theta < 1, \forall \|\mathbf{x}_R\| < \eta \\ &\leq -\alpha(1-\theta)\lambda_{\min}(\mathbf{K})\|\mathbf{x}_R\|^2 < 0, \quad \text{when } \eta > \|\mathbf{x}_R\| \geq \frac{2\delta\lambda_{\max}(\mathbf{K})}{\alpha\theta\lambda_{\min}(\mathbf{K})} \triangleq \mu. \end{aligned} \quad (63)$$

The second inequality is obtained using (61) and the third inequality is obtained using (30). Clearly, we have a non-empty region where  $\frac{d\mathcal{V}_R}{dt} < 0$  only when

$$\delta < \frac{\alpha\theta\lambda_{\min}(\mathbf{K})}{2\lambda_{\max}(\mathbf{K})}\eta. \quad (64)$$

Let  $b = \lambda_{\min}(\mathbf{K})\eta^2$  and  $c = \lambda_{\max}(\mathbf{K})\mu^2$ . Consider the sublevel sets  $\chi_b = \{\mathbf{x}_R \in \mathbb{R}^{r_x} \mid \mathcal{V}_R(\mathbf{x}_R) < b\}$  and  $\chi_c = \{\mathbf{x}_R \in \mathbb{R}^{r_x} \mid \mathcal{V}_R(\mathbf{x}_R) \leq c\}$ . It can be easily verified that if  $\delta < \frac{\alpha\theta\lambda_{\min}(\mathbf{K})}{2\lambda_{\max}(\mathbf{K})}\sqrt{\frac{\lambda_{\min}(\mathbf{K})}{\lambda_{\max}(\mathbf{K})}}\eta$ , then  $c < b$ , which implies  $\chi_c \subset \chi_b$ . Note, this condition satisfies the necessary condition (64) for the non-empty region since  $\lambda_{\min}(\mathbf{K}) \leq \lambda_{\max}(\mathbf{K})$ .

For any  $\mathbf{x}_R$  inside  $\chi_b$ , using (60), we have

$$\lambda_{\min}(\mathbf{K})\|\mathbf{x}_R\|^2 \leq \mathcal{V}_R(\mathbf{x}_R) < b = \lambda_{\min}(\mathbf{K})\eta^2, \quad (65)$$

implying  $\|\mathbf{x}_R\| < \eta$ . Similarly, for any  $\mathbf{x}_R$  on the boundary or outside of  $\chi_c$ , we have

$$\lambda_{\max}(\mathbf{K})\mu^2 = c \leq \mathcal{V}_R(\mathbf{x}_R) \leq \lambda_{\max}(\mathbf{K})\|\mathbf{x}_R\|^2, \quad (66)$$

which implies  $\|\mathbf{x}_R\| \geq \mu$ .

Combining (65) and (66) we can say for any  $\mathbf{x}_R$  outside (including the boundary) of  $\chi_c$ , but inside  $\chi_b$ , (63) holds true. For such  $\mathbf{x}_R$  (i.e.  $\mathbf{x}_R \in \chi_b \setminus \chi_c$ ) we have

$$\frac{d\mathcal{V}_R}{dt} \leq -\frac{\alpha(1-\theta)\lambda_{\min}(\mathbf{K})}{\lambda_{\max}(\mathbf{K})}\mathcal{V}_R(\mathbf{x}_R) \triangleq -2\gamma\mathcal{V}_R(\mathbf{x}_R), \quad (67)$$

using (60) and (63).

If the initial point (at time  $t_0$ ) satisfies  $\|\mathbf{x}_R(t_0)\| < \sqrt{\frac{\lambda_{\min}(\mathbf{K})}{\lambda_{\max}(\mathbf{K})}}\eta$ , then by (60),

$$b = \lambda_{\min}(\mathbf{K})\eta^2 > \lambda_{\max}(\mathbf{K})\|\mathbf{x}_R(t_0)\|^2 \geq \mathcal{V}_R(\mathbf{x}_R(t_0)),$$

which implies the initial point  $\mathbf{x}_R(t_0)$  is inside  $\chi_b$ . Assuming an initial point in  $\chi_b \setminus \chi_c$ , and integrating (67) in time interval  $[t_0, t]$ , we get

$$\mathcal{V}_R(\mathbf{x}_R(t)) \leq \mathcal{V}_R(\mathbf{x}_R(t_0))e^{-2\gamma(t-t_0)}. \quad (68)$$

Hence,  $\lambda_{\min}(\mathbf{K})\|\mathbf{x}_R(t)\|^2 \leq \mathcal{V}_R(\mathbf{x}_R(t)) \leq \mathcal{V}_R(\mathbf{x}_R(t_0))e^{-2\gamma(t-t_0)} \leq \lambda_{\max}(\mathbf{K})\|\mathbf{x}_R(t_0)\|^2e^{-2\gamma(t-t_0)}$  as long as  $\mathbf{x}_R(t)$  remains outside of  $\chi_c$ . Since  $\frac{d\mathcal{V}_R}{dt}$  is always negative outside of  $\chi_c$ , any trajectory starting outside of it, must enter  $\chi_c$  in finite time. Let the trajectory starting at  $\mathbf{x}_R(t_0)$  enters  $\chi_c$  for the first time at time  $t_c + t_0$ . Then, we have

$$\|\mathbf{x}_R(t)\| \leq \sqrt{\frac{\lambda_{\max}(\mathbf{K})}{\lambda_{\min}(\mathbf{K})}}e^{-\gamma(t-t_0)}\|\mathbf{x}_R(t_0)\| = \lambda e^{-\gamma(t-t_0)}\|\mathbf{x}_R(t_0)\|, \quad \forall t_0 \leq t < t_c + t_0. \quad (69)$$

Once a trajectory enters  $\chi_c$ , it cannot escape  $\chi_c$  because  $\frac{d\mathcal{V}_R}{dt}$  is negative on the boundary. Therefore, all points of a trajectory after  $t \geq t_c + t_0$  satisfies  $\lambda_{\min}(\mathbf{K})\|\mathbf{x}_R(t)\|^2 \leq \mathcal{V}_R(\mathbf{x}_R(t)) \leq c$ , equivalently,

$$\|\mathbf{x}_R(t)\| \leq \sqrt{\frac{\lambda_{\max}(\mathbf{K})}{\lambda_{\min}(\mathbf{K})}}\mu = \frac{2\delta}{\alpha\theta} \left( \frac{\lambda_{\max}(\mathbf{K})}{\lambda_{\min}(\mathbf{K})} \right)^{3/2} = \frac{2\delta}{\alpha\theta}\lambda^3, \quad \forall t \geq t_c + t_0. \quad (70)$$

From (67), (69) and (70), we have  $\gamma = \frac{\alpha(1-\theta)\lambda_{\min}(\mathbf{K})}{2\lambda_{\max}(\mathbf{K})}$  and  $\lambda = \sqrt{\frac{\lambda_{\max}(\mathbf{K})}{\lambda_{\min}(\mathbf{K})}}$ . ■

## B Details on reaction–diffusion system experiment

### B.1 Dataset

We use FEniCS (Logg et al. (2012)), an open-source computing platform for solving PDEs using the finite element method, with Python interface to generate the dataset. For the reaction-diffusion system of (34), we generate 100 training sequences of length 50 with time step size 0.01 and 256 nodes in  $\mathbb{I}$ . The initial conditions and actuations of these sequences are given by

$$q(\zeta, 0) = |a| \sum_{k=0}^4 b_k T_k(\zeta), \quad \zeta \in \mathbb{I}, \quad (71)$$

and

$$w(t_i) = 10g_i \max_{\zeta} |q(\zeta, t_{i-1})|, \quad i = 1, 2, \dots, 49, \quad (72)$$

where  $T_k$  denotes the  $k^{\text{th}}$  Chebyshev polynomial of the first kind, and  $a \sim \mathcal{N}(0, 1)$ ,  $b_k, g_i \sim \mathcal{U}(-1, 1)$  are chosen randomly. Similarly, 100 sequences are generated for the test set to evaluate the prediction performance.

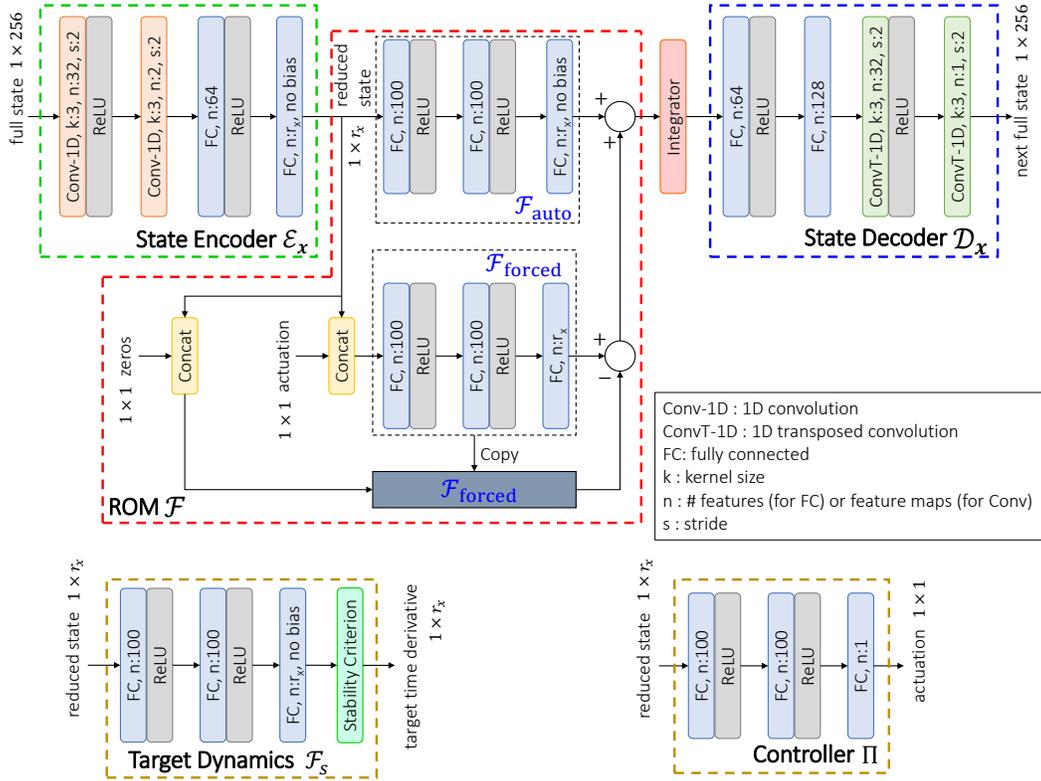


Figure 12: Architectures for all the DNN modules used in the reaction–diffusion experiment. The ‘Copy’ operation denotes the reuse of the same DNN block for zero and nonzero actuation. The ‘Concat’ operator concatenates the input features along the last dimension. **Zeros are concatenated to the reduced state to evaluate the component  $\mathcal{F}_{\text{forced}}(\mathbf{x}_R, \mathbf{0})$ .** The ‘Integrator’ performs the numerical integration for (21). The ‘Stability Criterion’ block implements (28).

## B.2 DNN architectures

Figure 12 shows the DNN architectures used for different modules in the reaction–diffusion experiment. The state encoder comprises 1D convolutional layers, followed by fully connected layers. The state decoder has the reversed order with convolutional layers replaced by transposed convolutional layers. The ROM is designed by breaking the function  $\mathcal{F}$  into two components:  $\mathcal{F}(\mathbf{x}_R, \mathbf{u}_R) = \mathcal{F}_{\text{auto}}(\mathbf{x}_R) + \mathcal{F}_{\text{forced}}(\mathbf{x}_R, \mathbf{u}_R) - \mathcal{F}_{\text{forced}}(\mathbf{x}_R, \mathbf{0})$ .  $\mathcal{F}_{\text{auto}}$  represents the autonomous dynamics that does not depend on the actuation, whereas  $\mathcal{F}_{\text{forced}}$  is responsible for the impact of actuation on dynamics. **The composition  $\mathcal{F}_{\text{forced}}(\mathbf{x}_R, \mathbf{u}_R) - \mathcal{F}_{\text{forced}}(\mathbf{x}_R, \mathbf{0})$  ensures that the component responsible for learning the impact of actuation on the dynamics provides nonzero output only when the actuation is nonzero.** Two multilayer perceptrons (MLPs) are used to implement  $\mathcal{F}_{\text{auto}}$  and  $\mathcal{F}_{\text{forced}}$ . **This specific structure of the ROM is not crucial and a single neural network representing  $\mathcal{F}(\mathbf{x}_R, \mathbf{u}_R)$  works as well. However, we observe better performance in experiments when the aforementioned structure is used.** The output of the ROM is integrated using a numerical integrator to get the next state. The controller is implemented using an MLP. The target dynamics is implemented using another MLP, followed by a stability criterion in the form of (28).

## B.3 Training settings

We use  $r_x = 5$  in the prediction task and  $r_x = 2$  in the control task for **all the methods**. All modules are implemented in PyTorch. In both of the learning phases, learning ROM and learning controller, we use the Adam optimizer with an initial learning rate of 0.001 and apply an exponential scheduler with a decay of

0.99. Modules are trained for 100 epochs in mini-batches of size 32. 10% of the training data is used for validation to choose the best set of models. For DeepROM training, we use  $\beta_2 = 1$  in (22). For learning control, we use  $\beta_3 = 0.2$  in (26),  $\alpha = 0.2$  in (28), and  $\mathbf{K} = 0.5\mathbf{I}_{r_x}$  in (29). Since the learned ROMs from one training instance to another can vary, the hyperparameter pair  $(\alpha, \beta_3)$  may require re-tuning accordingly.

## C Details on vortex shedding suppression experiment

### C.1 Dataset

For the flow past a circular cylinder problem, the geometry and physical parameters of the system are taken from the DFG 2D-2 benchmark (Schäfer et al. (1996)). The geometry is shown in Figure 13. We use the blue-shaded region for observation and actuation. Following the DFG 2D-2 benchmark, we use the no-slip boundary condition of zero velocity for the walls and the cylinder boundary, zero outlet pressure, and the inflow velocity profile (at the inlet) as

$$\mathbf{v}(\boldsymbol{\zeta}, t) = \left( 1.5 \frac{4\zeta_2(0.41 - \zeta_2)}{0.41^2}, 0 \right), \quad (73)$$

where  $\zeta_1$  and  $\zeta_2$  denote the horizontal and vertical coordinates, respectively, of  $\boldsymbol{\zeta}$ . We use kinematic viscosity  $\nu = 0.002$  and density  $\rho = 1$  leading to the Reynolds number  $Re = 50$ . The training sequence of length 5000 is generated in FEniCS with a time step size 0.001 and applying actuations

$$\mathbf{w}(\boldsymbol{\zeta}, t) = a \sum_{k=0}^4 \left[ \sin(k\pi(\zeta_1 - 0.11)/0.66) \quad \sin(k\pi\zeta_2/0.41) \right] \begin{bmatrix} b_{k,1,1} & b_{k,2,1} \\ b_{k,1,2} & b_{k,2,2} \end{bmatrix}, \quad \boldsymbol{\zeta} \in \mathbb{W}, \quad (74)$$

where  $a \sim \mathcal{U}(0, 1)$  and  $b_{k,i,j} \sim \mathcal{U}(-1, 1)$ ,  $i, j = 1, 2$  are chosen randomly. Similarly, a test sequence is generated to evaluate the prediction performance. For learning control, we use the Stokes flow or creeping flow as the desired state, which can be obtained by solving the Stokes equations

$$\nu \nabla^2 \mathbf{v} - \frac{1}{\rho} \nabla p = \mathbf{0}, \quad \nabla \cdot \mathbf{v} = \mathbf{0} \quad \text{in } \mathbb{I} \times \mathbb{R}^+. \quad (75)$$

For training, the flow velocity data from the observation region (blue shaded in Figure 13) are interpolated onto a rectangular uniform grid of size  $32 \times 48$  so that it can be used in standard CNNs.

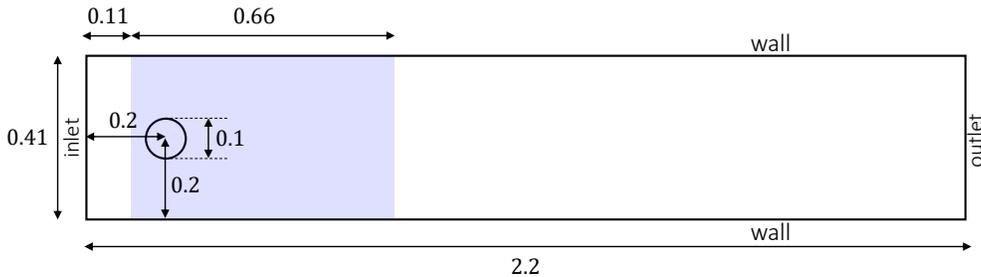


Figure 13: Geometry of the flow past a circular cylinder set-up.

### C.2 DNN architectures

Figure 14 shows the DNN architectures used for different modules in the vortex shedding control experiment. The architectures for the ROM and target dynamics are the same as in the previous example. Moreover, the state encoder and decoder have similar architectures as the previous example except for the 1D convolutions and transposed convolutions are replaced by their 2D counterparts. Here, an additional module is used: the control encoder for encoding the distributed control/actuation. It has the same architecture as the

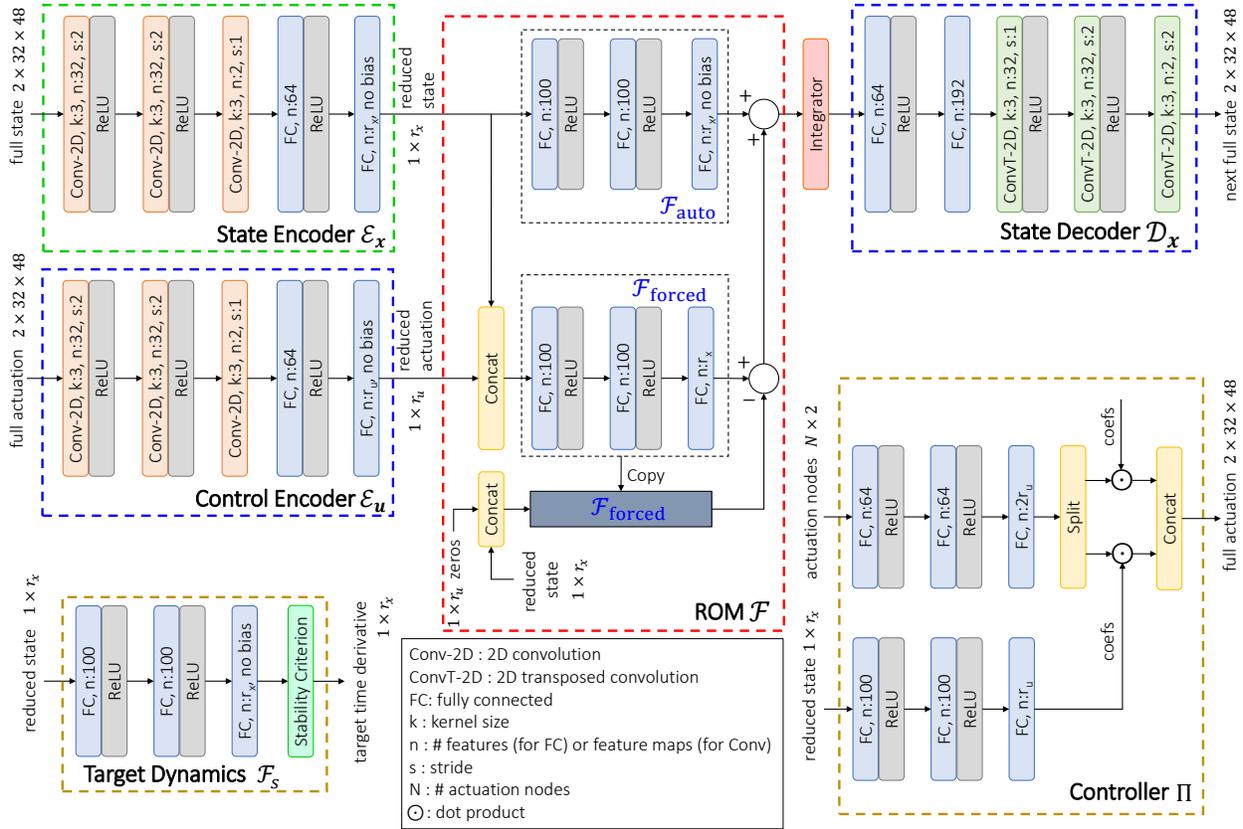


Figure 14: Architectures for all the DNN modules used in the fluid flow experiment. The ‘Split’ operator splits the input features into two vectors, along the last dimension. These split vectors represent the space-dependent polynomial basis associated with the horizontal and vertical components of the actuation.

state encoder. To learn the distributed actuation, we design the controller as a linear combination of space-dependent polynomial basis functions. One MLP is used to learn these space-dependent polynomial basis functions given the locations of the actuation nodes and another MLP is used to learn the corresponding coefficients. The actuation is computed as the dot product of the polynomial basis terms and the coefficient vector. We use this architecture instead of a standard convolutional one because the PDE solver takes the actuation input in a triangular mesh, not in a uniform rectangular grid. The polynomial basis architecture can be used to compute actuation in both uniform rectangular grid during training and triangular mesh during evaluation.

### C.3 Training settings

We use  $r_x = 5$  in both the prediction task and control task for all the methods. All modules are implemented in PyTorch. In both of the learning phases, learning ROM and learning controller, we use the Adam optimizer with an initial learning rate of 0.001 and apply an exponential scheduler with a decay of 0.99. Modules are trained for 100 epochs in mini-batches of size 32. 10% of the training data is used for validation to choose the best set of models. For DeepROM training, we use  $\beta_2 = 1$  in (22). For learning control, we use  $\beta_3 = 2$  in (26),  $\alpha = 0.1$  in (28), and  $\mathbf{K} = 0.5 \mathbf{I}_{r_x}$  in (29). Since the learned ROMs from one training instance to another can vary, the hyperparameter pair  $(\alpha, \beta_3)$  may require re-tuning accordingly.

## D Architecture and training details for the Deep Koopman model

For the encoder and decoder of the Deep Koopman model, we use the same architectures as our state encoder and state decoder. As mentioned in section 5.1, we consider both the system and input matrices of the ROM to be fixed during operation, in contrast to the original method proposed by Morton et al. (2018). Therefore, during training, these matrices are treated as trainable global parameters. Similar to Morton et al. (2018), the input matrix is optimized by gradient descent during training along with the encoder-decoder parameters, whereas the system matrix is obtained using linear least-squares regression. The datasets are divided into staggered 32-step sequences for training, and the model is trained by generating recursive predictions over 32 steps following Morton et al. (2018). We train the model using the Adam optimizer with an initial learning rate of 0.001 and an exponential decay of 0.99 for 200 epochs in mini-batches of size 8. 10% of the training data is used for validation to choose the best set of models.

As mentioned in 5.3.3, we utilize a low-dimensional representation of the distributed actuation for Deep Koopman + LQR, instead of directly estimating the high-dimensional actuation. The distributed actuation is represented as a linear combination of the same space-dependent sinusoidal basis functions used for dataset generation, which are given by (74). The controller is designed to estimate the coefficients  $b_{k,i,j}; i, j = 1, 2; 0 \leq k \leq 4$ .