

# Machine Learning for Prediction with Missing Dynamics

John Harlim

Department of Mathematics  
Department of Meteorology and Atmospheric Science  
Institute for CyberScience  
The Pennsylvania State University, PA 16802, USA  
jharlim@psu.edu

Shixiao W. Jiang

Department of Mathematics  
The Pennsylvania State University, PA 16802, USA  
suj235@psu.edu

Senwei Liang

Department of Mathematics  
National University of Singapore, Singapore  
liangsenwei@u.nus.edu

Haizhao Yang

Department of Mathematics  
Purdue University, IN 47907, USA  
haizhao@purdue.edu

October 9, 2019

## Abstract

This article presents a general framework for recovering missing dynamical systems using available data and machine learning techniques. The proposed framework reformulates the prediction problem as a supervised learning problem to approximate a map that takes the memories of the resolved and identifiable unresolved variables to the missing components in the resolved dynamics. We demonstrate the effectiveness of the proposed framework with a theoretical guarantee of a path-wise convergence of the resolved variables up to finite time and numerical tests on prototypical models in various scientific domains. These include the 57-mode barotropic stress models with multiscale interactions that mimic the blocked and unblocked patterns observed in the atmosphere, the nonlinear Schrödinger equation which found many applications in physics such as optics and Bose-Einstein-Condense, the Kuramoto-Sivashinsky equation which spatiotemporal chaotic pattern formation models trapped ion mode in plasma and phase dynamics in reaction-diffusion systems. While many machine learning techniques can be used to validate the proposed framework, we found that recurrent neural networks outperform kernel regression methods in terms of recovering the trajectory of the resolved components and the equilibrium one-point and two-point statistics. This superb performance suggests that recurrent neural networks are an effective tool for recovering the missing dynamics that involves approximation of high-dimensional functions.

**Keywords.** Closure Modeling, Missing dynamics, Machine Learning, Kernel Mean Embedding of Conditional Distribution, Long Short Term Memory

**AMS subject classifications:** 44A55, 65R10 and 65T50.

## 1 Introduction

The problem of missing dynamics is ubiquitous in any scientific domain that concerns with prediction

through computational models. This long-standing problem has been posted under various names, including model error, sub-grid scale parameterization, closure modeling [15, 28, 6, 25, 32, 35, 37, 40, 67]. Another relevant topic of broad interest is the reduced-order modeling whose ultimate goal is to systematically deduce a computationally efficient model to predict the evolution of the resolved variables when the full underlying model is too expensive to solve [20, 44, 45, 65, 12, 13, 22, 24, 29]. In our context, the proposed framework adopted here does not require any knowledge of the full equations that govern the underlying dynamical systems. The proposed approach that we consider assumes that only the dynamical components corresponding to the resolved variables are given. As in [26], the missing components will be learned from a historical time series of the resolved variables and the identifiable unresolved variables, where the latter serves as external forces to the dynamics of the relevant components. In essence, the proposed approach is to learn a dynamical model for the time-dependent external forces (or the identifiable unresolved variables).

The success of deep learning as a supervised learning algorithm has drawn tremendous interest on reduced-order modeling applications. A closely related approach to the modeling framework in this paper is presented in [53]. They proposed a Feedforward Neural Network (FNN) as a representation of the dynamics of the irrelevant variables. These authors also provide a linear control theory perspective to justify the identifiability of their dynamical representation on a class of nonlinear systems with a dual linear closure. In this article, we consider the closure modeling framework with a nonparametric formulation and provide a theoretical guarantee of path-wise convergence of the estimation of the resolved variables for the first time. For discrete dynamics obtained from a temporal discretization of differential equations, we found that when the unresolved variables are fully identifiable, the error rate is  $O(T^2 \Delta^2 \epsilon)$ , under mild assumptions. Here,  $T > 0$  denotes the prediction time index,  $\Delta$  denotes the discrete time step, and  $\epsilon > 0$  denotes the approximation error of the missing dynamics. Recalling the theory of nonparametric regression [62], if the missing dynamics is a function of a Sobolev class,  $H^\beta$ , where  $\beta > 0$  denotes the regularity parameter, the learning rate  $\epsilon$  of any nonparametric regression algorithm with i.i.d data has an optimal global error rate of an order  $\epsilon = O(N^{-\frac{\beta}{2\beta+d}})$ , where  $N$  denotes the length of training data and  $d$  denotes the dimension of the domain of the function. Hence, nonparametric regression algorithms suffer from the curse of dimensionality unless when the regularity parameter  $\beta = O(d)$ . However, even for the case of  $\beta = O(d)$ , there are no efficient tools to carry out the computation for high-dimensional problems.

Fortunately, recent advances in the theoretical analysis of deep neural networks show that they can avoid the curse of dimensionality in terms of approximation error in both the case of sufficiently smooth functions [5, 47, 49] and general continuous functions [48]. Also, there is no curse of dimensionality of deep neural networks in terms of generalization error when the target functions admit sufficient smoothness [17], when the data are sampled on a low-dimensional manifold [52], or in the case of classification functions [9]. While the generalization error for deep neural networks on general functions is an open problem, empirical numerical evidence has indicated that deep neural networks together with their stochastic training algorithms (e.g., batch-based stochastic gradient descent) are automatic tools that can identify the “low complexity” of the underlying systems, e.g., the smoothness or the low-dimensional domain that leads to no curse of dimensionality. In particular, recurrent neural networks as a special case of deep neural networks has the potential to avoid the curse of dimension when learning and predicting discrete dynamical systems with low complexity structures.

While the closure modeling framework can be numerically realized using any approximation/regression methods, we will consider a special type of recurrent neural networks called the Long-Short-Term-Memory (LSTM). We will show that this approach can overcome the curse of dimension suffered by the standard nonparametric regression method such as the kernel mean embedding approximation used in [26]. Our choice of using the LSTM is also encouraged by the success of it in recent closure modeling applications as proposed in [39, 64, 46]. We should stress that these existing approaches [39, 64, 46] share a similarity, that is, they specify the closure model as a function of only the memory of the resolved variables and motivate

their framework using a heuristic connection with the Mori-Zwanzig formalism [50, 68]. In contrast, we will demonstrate that it is critical for the closure model to also depend on the memory of the identifiable unresolved variables in addition to the resolved components. We will demonstrate the effectiveness of our framework on several tough prototype complex systems that arise in geophysical fluid dynamics, optics, quantum fluid such as Bose-Einstein-Condensate, and plasma physics, in addition to theoretical justification.

## 2 Data-Driven Modeling for Missing Dynamics

Throughout this paper, we will describe the closure modeling approach in the context of discrete maps that naturally arise from numerical discretization of partial or ordinary differential equations. We will discuss the stochastic case in the next section. Let the resolved,  $x_t \in \mathcal{X}$ , and unresolved,  $y_t \in \mathcal{Y}$ , variables be the solution of the following deterministic discrete dynamical systems,

$$x_{t+1} = \mathcal{F}(x_t, y_t), \quad y_{t+1} = \mathcal{G}(x_t, y_t), \quad (1)$$

given initial conditions  $x_0, y_0$ . Furthermore, we assume that the full system is ergodic with invariant measure  $\mu$  and the maps  $\mathcal{F}$  and  $\mathcal{G}$  are globally Lipschitz in  $x$  and  $y$ .

Our goal is to predict  $\{x_t : t \in \mathbb{N}\}$  and its statistics, such as, the mean, covariance, and auto-correlation functions, with the following constraints. We assume that  $\mathcal{F}$  is given, but  $\mathcal{G}$  is not available. Basically, the absence of  $\mathcal{G}$  means that we are missing the unresolved dynamics for  $y$ . Our goal is to reconstruct the missing dynamics in (1) from the given historical data  $\{x_t, \theta_t\}_{t=1, \dots, N}$ , where  $\theta_t := \theta(x_t, y_t)$  is the identifiable unresolved variable. To be precise,  $\theta_t$  is the component of the unresolved variables that can be identified from  $\mathcal{F}(x, y) := \mathcal{F}(x, \theta(x, y))$  in (1) and observations  $\{x_t\}$ . This restriction is motivated by practical constraints where only the resolved variables are observed. Given  $\{x_t\}$ , one can extract  $\{\theta_t\}$  using a regression fitting [14]. In the case when the observed  $x_t$  is noisy, one can also use appropriate filtering methods [25, 7]. In the numerical simulations shown below, we assume that a historical timeseries of  $\{x_t, \theta_t\}_{t=1, \dots, N}$  is available to us.

Define  $z_{t,m} := (x_{t-m:t}, \theta_{t-m:t}) \in \mathcal{Z}$  with  $x_{t-m:t} := (x_{t-m}, x_{t-m+1}, \dots, x_t)$  and  $\theta_{t-m:t} := (\theta_{t-m}, \theta_{t-m+1}, \dots, \theta_t)$  for some integer  $m \geq 0$  which characterizes the memory length. We consider a general closure model of the following form,

$$\hat{x}_{t+1} = \mathcal{F}(\hat{x}_t, \hat{\theta}_t), \quad \hat{\theta}_{t+1} = \mathbb{E}^\epsilon[\Theta_{t+1} | \hat{z}_{t,m}] + \xi_{t+1}, \quad (2)$$

where  $\hat{\cdot}$  is used to denote the numerical approximation of the corresponding variable in the closure model. In (2), the noise  $\xi_t$  is added to account for the residual due to misspecification of hypothesis space of  $\mathbb{E}[\Theta_{t+1} | \cdot]$ . In fact, we shall see from our numerical experiments below that this additional noise is not needed for the deterministic problems when LSTM is used. For simplicity, we only consider  $\xi_t \sim \Xi$  to be Gaussian with variance,

$$\begin{aligned} \mathbb{E}[\Xi^2] &:= \mathbb{E}[(\Theta_{t+1} - \mathbb{E}^\epsilon[\Theta_{t+1} | Z_{t,m}])^2] \\ &= \mathbb{E}[(\mathbb{E}[\Theta_{t+1} | Z_{t,m}] - \mathbb{E}^\epsilon[\Theta_{t+1} | Z_{t,m}])^2]. \end{aligned} \quad (3)$$

Here, the notation  $\mathbb{E}^\epsilon[\Theta_{t+1} | \cdot]$  is to denote a numerical approximation to the conditional expectation,  $\mathbb{E}[\Theta_{t+1} | \cdot]$ , defined with respect to the invariant conditional distribution of the full dynamics. This conditional expectation will be estimated from the historical solutions of the ergodic system in (1) since its explicit expression is unknown in general. Essentially, the closure model is reformulated as a supervised learning problem to learn the map  $\hat{z}_{t,m} \mapsto \mathbb{E}[\Theta_{t+1} | \hat{z}_{t,m}]$ .

## 2.1 Fully identifiable unresolved variables

To give an intuition, suppose that the entire unresolved variables can be identified, that is,  $\theta(x, y) = y$ . In this case, it is clear that  $\mathbb{E}[\Theta_{t+1}|z_{t,0}] = \mathbb{E}[Y_{t+1}|x_t, y_t] = \mathcal{G}(x_t, y_t)$  such that one can rewrite (1) as,

$$x_{t+1} = \mathcal{F}(x_t, y_t), \quad y_{t+1} = \mathbb{E}[Y_{t+1}|x_t, y_t]. \quad (4)$$

If  $\mathbb{E}^\epsilon[Y_{t+1}|\cdot]$  is a consistent estimator to  $\mathbb{E}[Y_{t+1}|\cdot]$  with variance error rate  $\epsilon^2$ , it is clear from (3) that  $\mathbb{E}[\Xi^2] = O(\epsilon^2)$ . In this case, we can show that

**Theorem 1.** *Let  $x_{t+1}$  be the solutions of (4) and  $\hat{x}_{t+1}$  be the solutions of,*

$$\hat{x}_{t+1} = \mathcal{F}(\hat{x}_t, \hat{y}_t), \quad \hat{y}_{t+1} = \mathbb{E}^\epsilon[Y_{t+1}|\hat{x}_t, \hat{y}_t] + \xi_{t+1}, \quad (5)$$

*under the same initial conditions  $x_0 = \hat{x}_0, y_0 = \hat{y}_0$ . Suppose that the variance error in (3) are of  $O(\epsilon^2)$  and  $\mathcal{F}$  and  $\mathcal{G}$  are globally Lipschitz in  $x$  and  $y$ . Then,*

$$\mathbb{E}\left[\max_{t \in \{0, \dots, T\}} |\hat{x}_t - x_t|\right] = O(a^T \epsilon). \quad (6)$$

*for some constant  $a > 1$  that is independent of  $T$  and  $\epsilon$ .*

*Proof.* See SI Appendix A. □

When the discrete dynamical system is a result of the Euler-Maruyama discretization on a system of stochastic differential equations,

$$\begin{aligned} dx &= f(x, y) dt + \sigma_x dW_{x,t}, \\ dy &= g(x, y) dt + \sigma_y dW_{y,t}, \end{aligned}$$

where  $dW_{x,t}$  and  $dW_{y,t}$  denote independent Gaussian white noises, we have:

$$\begin{aligned} x_{t+1} &= x_t + f(x_t, y_t)\Delta + \sigma_x \Delta^{1/2} \xi_{x,t+1}, \\ y_{t+1} &= y_t + g(x_t, y_t)\Delta + \sigma_y \Delta^{1/2} \xi_{y,t+1}, \end{aligned} \quad (7)$$

where  $\Delta$  denotes the time step. Here,  $\xi_x, \xi_y \sim \mathcal{N}(0, I)$  are samples of the standard Gaussian white noises. When  $g$  and  $\sigma_y$  are unknown, we can directly estimate these terms and obtain a sharper estimate:

**Theorem 2.** *Let  $x_{t+1}$  be the solutions of (7) and  $\hat{x}_{t+1}$  be the solutions of,*

$$\begin{aligned} \hat{x}_{t+1} &= \hat{x}_t + f(\hat{x}_t, \hat{y}_t)\Delta + \sigma_x \Delta^{1/2} \xi_{x,t+1}, \\ \hat{y}_{t+1} &= \hat{y}_t + \Delta \mathbb{E}^\epsilon[Y_{t+1}^\Delta | \hat{x}_t, \hat{y}_t] + \hat{\sigma}_y \Delta^{1/2} \xi_{y,t+1}, \end{aligned} \quad (8)$$

*under the same initial conditions  $x_0 = \hat{x}_0, y_0 = \hat{y}_0$ . Here, we have defined  $Y_{t+1}^\Delta := \frac{Y_{t+1} - Y_t}{\Delta}$  and the noise variance,*

$$\hat{\sigma}_y^2 \Delta := \mathbb{E}\left[\left(Y_{t+1} - Y_t - \Delta \mathbb{E}^\epsilon[Y_{t+1}^\Delta | X_t, Y_t]\right)^2\right]$$

*is estimated from the training data. Suppose that the learning variance error rate is,*

$$\mathbb{E}[(\mathbb{E}[Y_{t+1}^\Delta | X_t, Y_t] - \mathbb{E}^\epsilon[Y_{t+1}^\Delta | X_t, Y_t])^2] \leq C\epsilon^2.$$

*Let  $f$  and  $g$  be Lipschitz continuous in  $x$  and  $y$ . Then,*

$$\mathbb{E}\left[\max_{t \in \{0, \dots, T\}} |\hat{x}_t - x_t|\right] = O(\epsilon T^2 \Delta^2).$$

*Proof.* See SI Appendix B. □

This result suggests that the solution of the proposed approximate dynamics in (8) converges path-wise to that of (7) up to finite time. The convergence rate suggests that one can achieve a path-wise prediction with an accuracy of order learning rate error,  $\epsilon$ , up to order-one model unit time,  $(T\Delta)^2 = O(1)$ . In other words, the length of accurate path-wise prediction is inversely proportional to the square root of the learning error rate,  $T\Delta \approx \epsilon^{-1/2}$ .

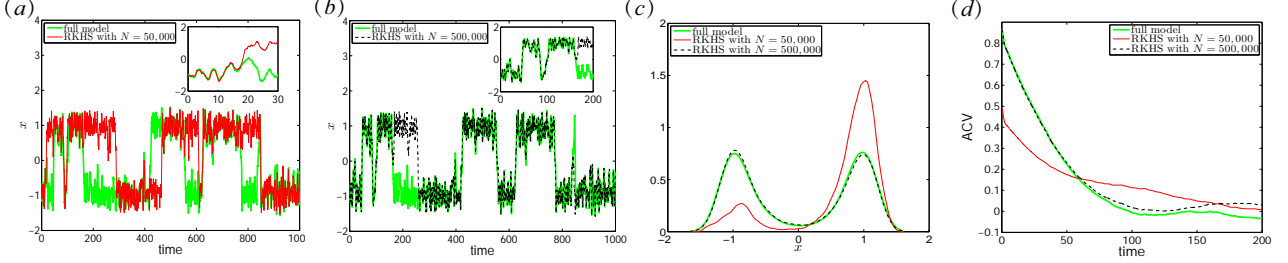


Figure 1: The first two panels from the left display comparison of the trajectories between the full model (green) and RKHS closure model using training dataset  $N = 50,000$  (red) and  $N = 500,000$  (dash black), respectively. The last two panels compare PDFs and auto-covariance functions (ACVs) among the full model (green), RKHS closure models trained with  $N = 50,000$  (red) and  $N = 500,000$  (dash black).

## 2.2 Parametric versus nonparametric closure models

The essence of parametric closure modeling is to simply specify  $\mathbb{E}^\epsilon$  in (2) with a specific choice of function  $\mathcal{P}(\hat{z}_{t,m}; W)$  that depends on a finite-dimensional parameter  $W$ . The choice of ansatz  $\mathcal{P}$  is usually based on physical intuition; see e.g., [40, 25, 6, 35, 37]. Once the model is specified, the hyper-parameter  $W$  can be obtained by regressing the pairs  $\{z_{t,m}, \theta_{t+1}\}$ . Subsequently, the variance of  $\Xi$  is estimated as in (3). In the later section, we will consider the Long-Short-Term-Memory model for  $\mathcal{P}$ .

In [26], a non-parametric closure model is considered. Specifically, the conditional expectation in (2) is estimated using the kernel mean embedding of conditional distribution formula [61, 60]. In the implementation, they assume that  $\mathbb{E}[\Theta_{t+1}|\cdot]$  belongs to the reproducing kernel Hilbert space (RKHS)  $\mathcal{H} \subset L^2(\mathcal{Z}, q)$  induced by orthonormal basis functions  $\{\varphi_k\}$  of this  $L^2$ -space, weighted by an arbitrary positive  $q \in L^1(\mathcal{Z})$  [26]. Then, for any  $z_t \in \mathcal{Z}$ , we can represent:

$$\mathbb{E}[\Theta_{t+1}|z_t] = \sum_{k=0}^{\infty} c_k \varphi_k(z_t), \quad (9)$$

where the coefficients  $c_k$  are precomputed from the available training data  $\{\theta_t, x_t\}$ ; see [26] for details. The key point is that this nonparametric formulation turns the problem of choosing the closure model into a problem of constructing basis functions of a Hilbert space (i.e., choosing an appropriate hypothesis space). Any parametric model  $\mathcal{P}(z; W)$  is a result of a finite truncation of a particular choice of basis functions  $\{\varphi_k\}$ .

Now, let us demonstrate the effectiveness of this approach in the following simple yet nontrivial example.

**Example:** Consider a Langevin dynamics

$$\begin{aligned} dx &= y dt, \\ dy &= (-\nabla V(x) - \gamma y) dt + \sigma_y dW_t, \end{aligned} \quad (10)$$

Table 1: Comparisons of mean exit time  $\bar{\tau}_0$  and reaction rate  $\nu_R$  between the full model and closure models.

	True	RKHS $N = 50,000$	RKHS $N = 500,000$
$\bar{\tau}_0$	99.2	69.1	102.7
$\nu_R$	0.0079	0.0040	0.0075

where  $x \in \mathbb{R}$  is the displacement,  $y \in \mathbb{R}$  is the velocity,  $V(x) = -x^2/2 + x^4/4$  is the double-well potential,  $\gamma = 1$  is the damping coefficient,  $dW$  is a standard Gaussian white noise, and  $\sigma_y = 3\sqrt{2}/10$  is the driving strength. We observe the trajectories of the variables  $(x_t, y_t)$  at every time step  $\Delta = 0.01$ , obtained using the Euler-Maruyama discretization scheme. In our previous notation,  $\theta(y) = y$  and we consider a closure model in (2) with  $\mathbb{E}[\Theta_{t+1}|\hat{z}_{t,0}] = \mathbb{E}[Y_{t+1}|\hat{x}_t, \hat{y}_t]$ . This example is nontrivial due to the transition state induced by the double-well potential  $V(x)$ . Note that only when the driving strength  $\sigma_y$  is in a reasonable region, the transition state phenomenon can be observed for this double-well potential system (see Fig. 1 for example).

In Fig. 1, we compare the result obtained using the RKHS approximation in (9) with the true trajectories and the statistics from the full model in (10). In this comparison, we apply the formula in (9) with a tensor product of  $50 \times 50$  Hermite polynomials. We compare the prediction of the trajectory up to finite-time, marginal density of  $x$ , and auto-covariance function  $\mathbb{E}[x_\tau x_0]$  of the true dynamics in (10) with those from the closure models, trained using  $N = 5 \times 10^4$  and  $5 \times 10^5$  data points. Notice that using large enough training data, we are not only accurately recovering the trajectory path-wise longer in time but also the density and auto-covariance function. Compare to the optimal learning rate  $\epsilon = O(N^{-\frac{\beta}{2\beta+d}})$  of [62] for very smooth function with  $\beta = \infty$ , the empirical prediction length (as shown in Fig. 1) is on the order of the theoretical prediction length  $T\Delta = \epsilon^{-1/2} \approx 14.95$  for  $N = 5 \times 10^4$  and is slightly longer than the conservative estimate  $T\Delta = \epsilon^{-1/2} \approx 26.59$  for  $N = 5 \times 10^5$ . In Table 1, we also see the agreement of several statistics that are commonly used to characterize metastable dynamics. When the training data set is large, we see a relatively accurate estimation of the mean exit time  $\bar{\tau}_0$  of a particle to escape one of the wells [19] and the reaction rate  $\nu_R$  (the number of trajectories to escape a well in the time interval  $T$  as  $T \rightarrow \infty$ ) [63].

### 2.3 Partially identifiable unresolved variables

While the closure model in (5) is theoretically consistent and the example above shows a very promising approach, in real applications, the function  $\theta(x, y) \neq y$  since the unresolved variables,  $y$ , are usually not identifiable from the data  $\{x_t\}$  and the map  $\mathcal{F}(x_t, \theta(x_t, y_t))$ . Even if the full data of  $y$  are available, they are very high-dimensional relative to  $\theta$ .

In this case, let us rewrite the dynamics for the identified unresolved variable,  $\theta_{t+1}$ , as a function of  $(x_{t-m:t}, \theta_{t-m:t})$ . To do this, we consider a Koopman operator  $S : \mathcal{H} \rightarrow \mathcal{H}$  defined as,  $S\theta(x_t, y_t) = \theta(x_{t+1}, y_{t+1})$ , where  $\theta \in \mathcal{H}$  belongs to a Hilbert space of functions that takes value in  $\mathcal{X} \times \mathcal{Y}$ , equipped with an inner product weighted by the invariant measure  $\mu$ . Let  $P : \mathcal{H} \rightarrow \mathcal{V}$  be an orthogonal projection operator to the space of functions that depends only on  $(x, \theta)$  and  $Q := I - P$  be the projection map to the orthogonal space  $\mathcal{V}^\perp \subset \mathcal{H}$ . Here, following the notation in [34], we define the map  $\pi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{X} \times \Theta$  as,  $\pi(x, y) = (x, \theta)$ . Applying the Dyson formula [16, 34],

$$S^{t+1} = \sum_{k=0}^t S^{t-k} P S (Q S)^k + (Q M)^{t+1},$$

to  $\pi$  and evaluate at  $(x_0, y_0)$ , we can rewrite the dynamics of the  $\theta$  component as,

$$\theta_{t+1} = \sum_{k=0}^t P(F_k \circ \mathcal{G})(x_{t-k}, \theta_{t-k}) + F_{n+1}(x_0, y_0),$$

where  $F_n := (QS)^n \pi$ . While the dynamics of  $\theta$  depends on the whole history  $(x_{0:t}, \theta_{0:t})$  and initial condition  $(x_0, y_0)$ , in practice, the length of the memory is often finite,  $0 < m < t$ , and it can be estimated as shown in [22, 54]. In fact, for nonlinear systems with linear dual closure, one can determine the minimum length of  $m$  to guarantee the identifiability of  $\theta_{t+1}$  from  $z_{t,m} := (x_{t-m:t}, \theta_{t-m:t}) \in \mathcal{Z}$  [39]. In the remainder of this paper, we assume that the underlying dynamics have a finite memory length  $0 < m < t$ ,

$$\begin{aligned} x_{t+1} &= \mathcal{F}(x_t, \theta_t), \\ \theta_{t+1} &= \sum_{k=0}^m P(F_k \circ \mathcal{G})(x_{t-k}, \theta_{t-k}). \end{aligned} \tag{11}$$

With this assumption, we have  $\mathbb{E}[\Theta_{t+1}|z_{t,m}] = \sum_{k=0}^m P(F_k \circ \mathcal{G})(x_{t-k}, \theta_{t-k})$ . Since orthogonal projection operators  $P, Q$  and the Koopman operator  $S$  are all bounded linear operators, it is clear that if  $\mathcal{G}$  is Lipschitz, then  $P(F_k \circ \mathcal{G})$  are also Lipschitz continuous and the following result holds.

**Theorem 3.** *Let  $x_{t+1}$  and  $\hat{x}_{t+1}$  be the solutions of (11) and (2), respectively, under the same initial conditions  $x_{-m:0} = \hat{x}_{-m:0}, \theta_{-m:0} = \hat{\theta}_{-m:0}$ . Let the equality in (3) holds and of order- $\epsilon^2$ . For  $\mathcal{F}$  and  $\mathcal{G}$  Lipschitz continuous in  $x$  and  $y$ , then*

$$\mathbb{E}\left[\max_{t \in \{0, \dots, T+1\}} |\hat{x}_t - x_t|\right] = O(a^T \epsilon)$$

where  $a > 1$  is a constant that is independent of  $T$  and  $\epsilon$ .

*Proof.* See SI Appendix C. □

This error rate is rather conservative since it depends on  $T$  exponentially. One might achieve an improved error rate by analyzing the eigenvalue problem corresponding to the autoregressive model of order- $m$ , which bounds the dynamical equation for the errors between (11) and (2).

For dynamical systems driven by stochastic noises, one can rewrite the full dynamics as an autonomous dynamical system by augmenting  $(x_n, y_n)$  with the entire history of the noises. See [30] for the details or the Appendix of [34] for the key idea. Subsequently, one can apply the Dyson formula to the resulting autonomous dynamics, defined on appropriate state space that includes  $\mathcal{X} \times \mathcal{Y}$  and the space of the history of the noises, and derive an analogous representation as in (11). We suspect that the result is not different from that in Theorem 2 and thus we will not present this derivation.

### 3 Deep Learning via Long-Short-Term-Memory

As a nonlinear type parametric regression method, deep learning outperforms kernel methods including the RKHS approach in terms of generalization error when the target functions are sufficiently smooth [17]. Though it is theoretically unclear whether there are any advantages of using deep learning over other non-parametric regression methods for general continuous functions in terms of overcoming the curse of dimension [5, 47, 49, 48], deep learning has practical advantages over the RKHS approaches. A significant challenge with the RKHS approximation in (9) is that there is no a priori guideline for choosing the appropriate hypothesis space. If the orthogonal basis is used, it is practically difficult to even construct these basis functions on very high-dimensional variables  $z \in \mathcal{Z}$ . On the other hand, if arbitrary radial functions are used as a basis, the evaluation of the resulting model on a new point  $z$  requires evaluating the basis functions on  $\|z - z_{t,m}\|$  for all training data  $t = 1, \dots, N$ , making the prediction with (2) becomes too costly since we need to evaluate the conditional expectation in (9) on a new point in each iteration. In contrast, deep learning as a nonlinear parametric regression method is not hampered by these issues, since it is practically just

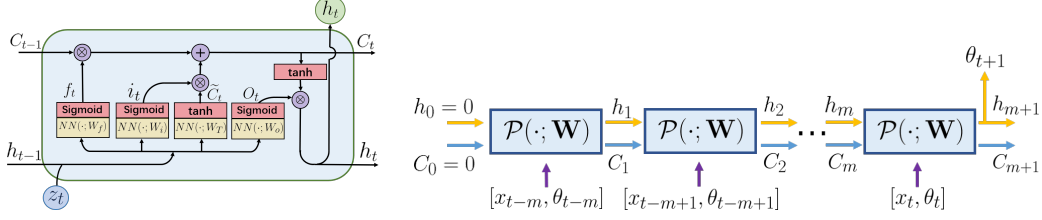


Figure 2: Left: the basic computational flow of a LSTM recurrence.  $+$  and  $\otimes$  are element-wise addition and multiplication respectively. Right: a sequence of LSTM cells applied compositionally.

a nonlinear interpolation technique using a composition of nonlinear activation functions and linear transforms. Of course, the main issue with nonlinear regression is whether one can obtain the minimum on such a non-convex optimization problem in the training phase. Recent advances in optimization theory show that simple gradient descent can identify a local minimizer with an arbitrarily small loss and a generalization error without the curse of dimensionality when the network size is sufficiently large for classification problems [9]. Though there is no existing optimization theory that guarantees good local minimizers in general settings, motivated by many positive numerical results shown in other closure modeling approaches [64, 46], we will consider realizing the closure model in (2) using recurrent neural networks.

As a special case of recurrent neural networks, Long-Short-Term-Memory (LSTM) is capable of learning multi-scale temporal effects and hence is adopted in our method. The computational flow of the LSTM consists of a sequence of computational cells, each of which is

$$\begin{aligned} f_t &= \sigma \circ NN(h_{t-1}, z_t; W_f), & i_t &= \sigma \circ NN(h_{t-1}, z_t; W_i) \\ o_t &= \sigma \circ NN(h_{t-1}, z_t; W_o), & \tilde{C}_t &= \tanh(NN(h_{t-1}, z_t; W_T)), \\ C_t &= f_t \otimes C_{t-1} + i_t \otimes \tilde{C}_t, & h_t &= o_t \otimes \tanh(C_t), \end{aligned}$$

where  $\sigma$  denotes the sigmoid function,  $\otimes$  is the pointwise product, and  $NN$  denotes a fully connected network which stacks layers of linear transformation and nonlinear activation function. See Fig. 2 (left) for an illustration of an LSTM cell. For simplicity, let us denote the above flow as  $(h_t, C_t) = \mathcal{P}(z_t, h_{t-1}, C_{t-1}; W)$  with parameters  $W$ , inputs  $(z_t, h_{t-1}, C_{t-1})$ , and outputs  $(h_t, C_t)$ . LSTM cells can be applied compositionally and we denote the LSTM sequence with  $m + 1$  cells as  $(h_{m+1}, C_{m+1}) = \mathcal{P}_m(\{z_t\}_{t=1}^{m+1}, h_0, C_0; W)$  (see Fig. 2 (right) for an illustration).

Now let us apply the LSTM to approximate the closure model in (2) with the given training data  $\{z_{t,m}, \theta_t\}$ , where  $z_{t,m} := (x_{t-m:t}, \theta_{t-m:t}) \in \mathcal{Z}$ . We train an  $(m + 1)$ -cell LSTM  $(h_{m+1}, C_{m+1}) = \mathcal{P}_m(z_{t,m}, h_0, C_0; W)$  with an input in the  $j$ -th cell as  $(x_{t-m+j-1}, \theta_{t-m+j-1})$  such that  $h_{m+1}$  predicts  $\theta_{t+1}$  well. The parameters  $h_0$  and  $C_0$  are set to be 0 for simplicity and  $W$  is identified via minimizing a mean squared error (MSE) function specified below. In what follows, we adopt the notation  $h_{m+1} = \mathcal{P}_m(z_{t,m}; W)$  for simplicity. Define the MSE loss as

$$\mathcal{L}(W) := \frac{1}{N - m - 1} \sum_{s=m+1}^{N-1} (\mathcal{P}_m(z_{s,m}; W) - \theta_{s+1})^2, \quad (12)$$

i.e., we aim to identify a predictor  $\mathcal{P}_m(z_{s,m}; W)$  such that it can predict the  $(s + 1)$ -th sample in the time series given  $(m + 1)$  preceding samples. Minimizing (12) can be achieved efficiently via a mini-batch stochastic gradient descent (SGD) and backpropagation through time (BPTT) [51, 58, 66]. Though the global minimizer of the above highly non-convex optimization might not be available, empirically numerical results [55] and partial theoretical analysis show that gradient-based algorithms are able to provide a minimizer  $W^*$  with a reasonably good generalization capacity in the case of over-parametrized networks [3, 1, 11]. Once  $W^*$  has been identified,  $\theta_{t+1} = \mathcal{P}_m(z_{t,m}; W^*)$  is applied instead of the conditional expectation in (2).



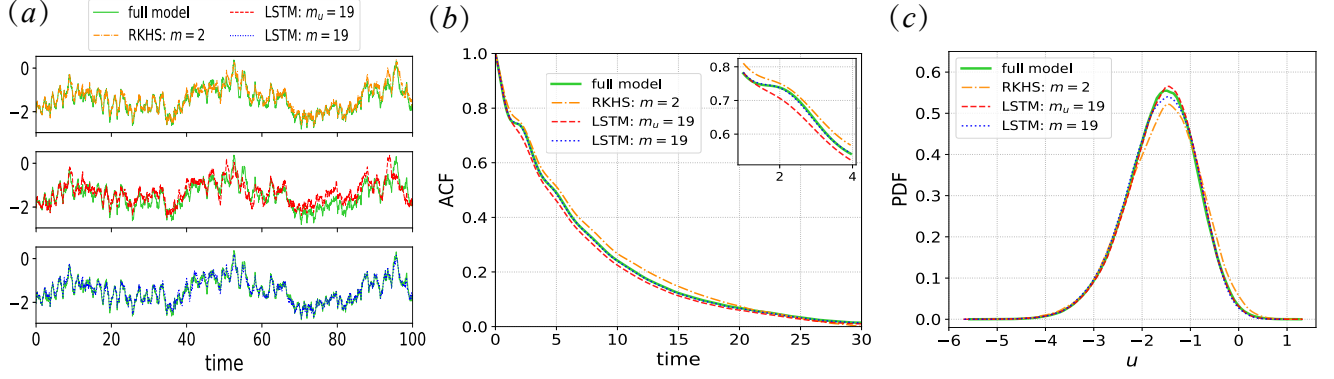


Figure 3: Comparison of (a) trajectories, (b) ACFs, and (c) densities, between the full and closure models for the large-scale mean velocity  $u$  in the regime of  $H = 5\sqrt{2}/4$  and  $\bar{d} = 0.1$ .

The computational cost of the proposed framework mainly consists of two parts: training and evaluation. In the training part, suppose the SGD with a batch size  $K$  is applied to minimize the loss in (12), then the computational cost for evaluating an approximate loss and gradient via BPTT is essentially  $O(Km)$  matrix-vector multiplications with a matrix of size  $d_L \times d_L$ , where  $d_L$  is the dimension of the hidden layer of  $NN$ . According to the approximation theory of DNNs [38, 59],  $d_L$  is required to be larger than  $d$ ; based on the optimization analysis of DNNs [8, 2, 9], a larger  $d_L$  admits a simpler optimization problem in deep learning. In the numerical examples shown below, we empirically set  $d_L = 500$  for problems with dimensions  $d = 40, 80$ , and set  $d_L = 1000$  for a problem that involves  $d = 480$ . Hence, the total computational cost is  $O(Kmd_L^2)$  for one iteration in the SGD, which has been parallelized via GPU computing in standard machine learning packages. Recent theoretical analysis shows that the convergence of SGD is linear under mild conditions [3]. Therefore, the upper bound of the iteration number in the SGD to guarantee a loss less than  $\epsilon$  is  $O(\log(\frac{1}{\epsilon}))$ . In our numerical examples below, the total numbers of iterations required for accurate performance are  $4 \times 10^3$  for the first two problems and  $4 \times 10^4$  for the last problem. In the evaluation part, the cost estimation for predicting one data sample with  $m$  pairs of historical data is  $O(md_L^2)$ , which has been parallelized via GPU.

## 4 Numerical Examples

In this section, we numerically demonstrate the effectiveness of our proposed closure framework on severely truncated dynamical systems in three prototypical applications. First, the topographic mean flow interaction that mimics the blocked and unblocked patterns observed in the atmosphere [10, 23, 42, 18] is considered. In our test, we will use the stochastic version of the 57-mode model studied in [57]. Second, the non-linear Schrödinger equation that finds many applications in optics and Bose-Einstein-Condensate (see the references in [27]) is studied. Finally, the Kuramoto-Shivashinsky equation with a spatiotemporal chaotic pattern formation with applications in trapped ion modes in plasma [33] and phase dynamics in reaction-diffusion systems [31]. We shall see that the closure models in these three examples progressively involve approximations of functions of dimensions 40 to 480.

## 4.1 Topographic mean flow interaction

We consider the topographic mean flow interaction that solves a barotropic quasi-geostrophic equation with a large-scale zonal mean flow  $u(t)$  on a two-dimensional  $2\pi \times 2\pi$  periodic domain, formulated as in [57]:

$$\begin{aligned} \frac{du}{dt} + \oint \frac{\partial h}{\partial x} \psi &= -\bar{d}u + \sigma\mu^{-1/2}\dot{W}_0, \\ \frac{\partial \omega}{\partial t} + \nabla^\perp \psi \cdot \nabla q + u \frac{\partial q}{\partial x} + \beta \frac{\partial \psi}{\partial x} &= -\mathcal{D}\psi + \Sigma \dot{W}. \end{aligned} \quad (13)$$

Here,  $q = \omega + h$  is the small-scale potential vorticity which is advected by the velocity field  $\mathbf{v} = \nabla^\perp \psi \equiv (-\partial_y \psi, \partial_x \psi)$ ;  $\omega = \Delta \psi$  and  $\psi$  are the relative potential vorticity and the stream function, respectively;  $h(\mathbf{x}) = h(x, y)$  is the topography. The parameter  $\beta$  is associated with the  $\beta$ -plane approximation to the Coriolis force. The integral in (13) is a two-dimensional integral over a periodic box of  $[-\pi, \pi] \times [-\pi, \pi]$ . On the right hand side of (13), the dissipation and forcing operators are applied on both the small and the large scales. On the small scale, the dissipation operator is in the form of  $\mathcal{D} = -\bar{d}\Delta$  with  $\bar{d} \geq 0$  and  $\Delta$  the Laplace operator corresponding to the Ekman drag dissipation. On the large scale, operator  $-\bar{d}u$  represents the momentum damping. The forcing terms are represented by random Gaussian white noises (e.g. unresolved baroclinic instability processes on small scales, random wind stress, etc), where  $W(t)$  and  $W_0(t)$  are standard Wiener processes;  $\sigma\mu^{-1/2} > 0$  is a constant amplitude and  $\Sigma$  is spatially dependent.

Following [42, 18, 57], we construct a set of special solutions to (13), which inherit the nonlinear coupling of the small-scale vortical modes with the large-scale mean flow via topographic stress. Consider the truncated spectral expansion of the state variables for  $\psi$  and  $\omega$  with high wavenumber truncation  $1 \leq |\mathbf{k}| \leq K$  using standard Fourier basis  $\exp(i\mathbf{k} \cdot \mathbf{x})$  with  $\mathbf{k} = (k_x, k_y)$ . We can rewrite (13) for the large-scale mean flow  $u(t)$  in a truncated Fourier form, as:

$$\frac{du}{dt} = i \sum_{1 \leq |\mathbf{k}| \leq K} \frac{k_x}{|\mathbf{k}|^2} \hat{h}_{\mathbf{k}}^* \omega_{\mathbf{k}} - \bar{d}(u - u_{eq}) + \sigma\mu^{-1/2}\dot{W}_t. \quad (14)$$

Here,  $\hat{h}_{\mathbf{k}}$  and  $\omega_{\mathbf{k}}$  are the Fourier transform of the topography  $h(\mathbf{x})$  and the relative potential vorticity  $\omega$ , respectively;  $u_{eq} = -\beta/\mu$  is the equilibrium mean of  $u(t)$ . The parameter  $\sigma$  is chosen such that  $\sigma_{eq}^2 = \frac{\sigma^2}{2\bar{d}} = 1$ . The parameters  $\beta = 1$  and  $\mu = 2$  are fixed in our simulation. More details can be found in SI Appendix D.

In our implementation, we consider the ground truth as the solution corresponding to the truncation  $1 \leq |\mathbf{k}| \leq K$  with  $K = 17$  such that there are 57 degrees of freedom for integers  $\mathbf{k} = (k_x, k_y)$ . In this topographic 57-mode model, we use the standard 4th order Runge-Kutta method for the time integration up to  $5 \times 10^7$  time iterations with a time step  $\delta t = 2.5\text{E-}3$ , which is small enough to capture the small-scale dynamics. Here, the initial condition,  $\psi(\mathbf{x}, 0)$ , is a sample of Gaussian distribution with random phases and amplitudes consistent with the ensemble mean and enstrophy as in [41]. The observed data are recorded at every 20 time steps, that is, we observe the data at every  $\Delta = 0.05$  time unit. Taking half of this data set for training,  $N = 1.25 \times 10^6$  samples. For the topography  $h(\mathbf{x})$ , we use a simple layered topography with variation only in the  $x$ -direction,  $h(\mathbf{x}) = H(\cos(x) + \sin(x))$ , where  $H$  denotes the topography amplitude.

We now present the closure model for the large-scale mean flow  $u(t)$  in Eq. (14). The application of the Euler-Maruyama scheme for the large-scale mean flow  $\hat{u}(t)$  gives

$$\begin{aligned} \hat{u}_{t+1} &= \hat{u}_t + \Delta \hat{\theta}_t - \Delta \bar{d}(\hat{u}_t - u_{eq}) + \sqrt{\Delta} \sigma \mu^{-1/2} \eta_{t+1}, \\ \hat{\theta}_{t+1} &= \mathbb{E}^\epsilon [\Theta_{t+1} | \hat{u}_{t-m:t}, \hat{\theta}_{t-m:t}] + \xi_{t+1}, \end{aligned} \quad (15)$$

where the time step  $\Delta = 0.05$  and the identifiable unresolved variable  $\hat{\theta}_t$  is an estimator of  $\theta_t = i \sum_{1 \leq |\mathbf{k}| \leq N} \frac{k_x}{|\mathbf{k}|^2} \hat{h}_{\mathbf{k}}^* \omega_{\mathbf{k}}$ . In this case,  $\theta$  is a function of the unresolved variables alone. The noises  $\xi_t$  are i.i.d. standard Gaussian while

the noises  $\eta_t$  are Gaussian with variance estimated as in (3). We will approximate the conditional expectation  $\mathbb{E}^\epsilon$  in (15) with the LSTM model with  $m = 19$ , which involves an approximation of a forty dimensional function. We will also include an experiment mimicking the existing approach in [39, 64, 46], where the conditional expectation in (15) is replaced with  $\mathbb{E}^\epsilon[\Theta_{t+1}|\hat{u}_{t-m_u:t}]$ , a function that depends only on the memory of the resolved scale with memory length  $m_u = 19$ . In this case, the conditional expectation is a twenty dimensional function. In addition, we also report the RKHS approximation to (15) with  $m = 2$ , which involves only an approximation of a six-dimensional function  $(\hat{u}_{t-2:t}, \hat{\theta}_{t-2:t}) \mapsto \mathbb{E}^\epsilon[\Theta_{t+1}|\hat{u}_{t-2:t}, \hat{\theta}_{t-2:t}]$ . Here, we use a tensor product of Hermite polynomials to represent this six dimensional function. The curse of dimensionality makes the RKHS model with orthogonal polynomials prohibitive in higher dimensions.

In Fig. 3, we present the short-time predictions and the long-time statistics of the large-scale mean flow  $u$  in a regime when  $H = 5\sqrt{2}/4$  and  $\bar{d} = 0.1$  (a regime considered in [57]). We would like to emphasize that we verify the closure model in (15) using initial conditions not in the training data set. More numerical results in other regimes of coupling  $H$  and damping  $\bar{d}$  can be found in SI Appendix D. In Fig. 3(a), the comparison of trajectories between the full model and the three closure models are presented. We see that the short-time prediction by the LSTM with both the resolved and identifiable unresolved variables ( $m = 19$ ) is the most accurate; however, the LSTM without the unresolved variables ( $m_u = 19$ ) proposed in [39, 64, 46] makes the worst prediction, which justifies the importance of considering both the resolved and identifiable unresolved variables in the closure model proposed in this paper. In Fig. 3 (b) and (c), we also show the comparison of the Auto-Correlation Functions (ACFs) and the equilibrium distribution of the full model and three prediction methods, which verify the importance of considering both the resolved and identifiable unresolved variables in the closure model as well.

## 4.2 Nonlinear Schrödinger equation

We consider the following nonlinear Schrödinger (NLS) equation,

$$i\frac{\partial u}{\partial t} = -\frac{\partial^2 u}{\partial x^2} + |u|^2 u, \quad (16)$$

with  $t$  as time and  $x$  as one-dimensional space. The periodic boundary condition is applied for the domain  $x \in [0, 2\pi]$ . Numerically, we generate the truth by integrating (16) using the pseudospectral methods in space [4] with finite wavenumbers  $|k| \leq K$  and Strang's splitting method in time [4]. The initial condition is generated by a Monte Carlo algorithm. Here, the number of modes  $K = 32$  and the observation time interval  $\Delta = 0.02$ . The observation data length is  $10^6$ . Taking half of this data set for training,  $N = 5 \times 10^5$  samples.

In this example, we are interested in constructing a closure model for the dynamics of the zeroth mode  $u_0$  of the NLS equation. In particular, the closure model consists of coupling a discrete approximation to the dynamics of the zeroth mode,

$$\frac{d\hat{u}_0}{dt} = -i(|\hat{u}_0|^2 \hat{u}_0 + \hat{\theta}_1 \hat{u}_0 + \hat{\theta}_2), \quad (17)$$

and

$$\hat{\theta}_{t+1} = \mathbb{E}^\epsilon[\Theta_{t+1}|\hat{u}_{0,t-m:t}, \hat{\theta}_{t-m:t}]. \quad (18)$$

Here  $\hat{\theta}_1$  and  $\hat{\theta}_2$  denote the approximation to the identifiable unresolved variables,  $\theta_1 = 2 \sum_{k \neq 0} |u_k|^2$  and  $\theta_2 = \sum_{k_1, k_2 \in \mathbb{Z}} u_{k_1} u_{k_2} u_{k_1+k_2}^* - 2 \sum_{k \neq 0} |u_k|^2 u_0 - |u_0|^2 u_0$ , respectively. We should point that  $\theta_1, \theta_2$  are independent of  $u_0$ . See SI Appendix E for detailed derivation of (17). The notation  $\hat{\theta}_t := (\hat{\theta}_{1,t}, \hat{\theta}_{2,t})$  in (18) is to denote the discrete estimate at time index  $t$ . Numerically, we use a splitting method to approximate the solution of (17) (see SI Appendix E for detail).

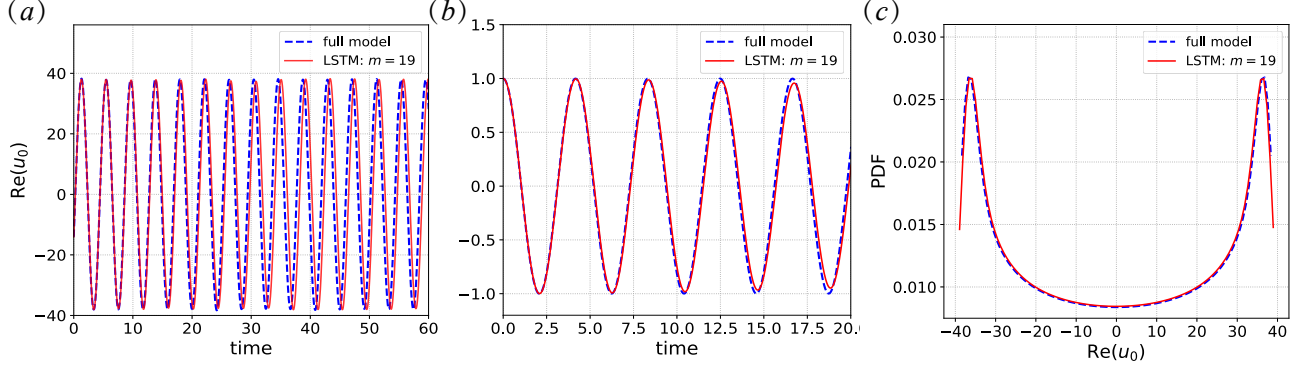


Figure 4: Comparison of (a) trajectories, (b) ACFs, and (c) densities of the full and closure LSTM models in the high temperature regime with  $k_B T = 10$ . The time step for the closure model is  $\Delta = 0.02$  and the LSTM method uses  $m = 19$  memory terms for  $(\theta, u_0)$  in Eq. (17).

In this study, we fix the memory length to be  $m = 19$  in the LSTM method, resulting in an approximation of 80-dimensional function  $\mathbb{E}^\epsilon[\Theta_{t+1} | \cdot]$ . No residual term is added in (18). We now present the numerical results from a training data set that is obtained by simulating initial conditions sampled from the Gibbs distribution  $\pi = \exp(-\frac{E}{k_B T})$ , where  $E$  denotes the Hamiltonian of the ODE system resulting from the Fourier Galerkin truncation on (16);  $k_B$  and  $T$  denote the Boltzmann constant and temperature, respectively. Qualitatively, both the amplitude and frequency of the solutions increase as a function of the temperature. In Fig. 4, we show numerical results in a high-temperature regime with  $k_B T = 10$ . For the short-time forecasting, we observe from Fig. 4(a) that the path-wise solution of  $u_0$  is well captured. For the long-time statistics, we observe from Fig. 4(b) and (c) that the ACF and the density of the true  $u_0$  are also well reproduced. Therefore, for the first mode  $u_0$  of the NLS equation, the proposed closure model using the LSTM method can reasonably replicate the short-time forecasting skill and long-time statistics in the high-temperature regime. Numerical results in the low-temperature regime are reported in SI Appendix E and we can draw the same conclusion.

We should point out that the resulting model is only valid in predicting the evolution of the system on the same energy level, since the underlying Hamiltonian system is not ergodic. This implies that the verification will only be valid to predict the evolution of the system with initial conditions sampled from the same Gibbs distribution where the training data is generated from.

### 4.3 The Kuramoto-Sivashinsky equation

We consider the Kuramoto-Sivashinsky equation (KSE) on an  $L$ -periodic domain, the Fourier representation of which can be written as

$$\frac{d}{dt}v_k = (q_k^2 - q_k^4)v_k - \frac{iq_k}{2} \sum_{l=-\infty}^{\infty} v_l v_{k-l}, \quad (19)$$

where  $q_k = 2\pi k/L$  with  $k \in \mathbb{Z}$ , and  $v_k$  denotes the  $k$ th Fourier mode.

In our numerical implementation, we let the full dynamics to be the Galerkin truncation of (19) for  $|k| \leq K/2$ , where  $K = 96$ . Notice that in the linearized equations of (19), each Fourier mode has an eigenvalue  $q_k^2 - q_k^4$  so that high  $k$  modes with  $|q_k| > 1$  are linearly stable whereas low  $k$  modes with  $|q_k| \leq 1$  are not. We set the spatial length  $L = 2\pi/\sqrt{0.085}$  so that the number of linearly unstable modes is  $\lfloor 1/\sqrt{0.085} \rfloor = 3$ . In this case, the energy is transferred from the linearly unstable low 3 modes to the damped high  $K/2 - 3 = 45$

modes through the nonlinear terms so that the KSE is well-posed and the solutions remain globally bounded in time [21]. This regime is exactly the same as the one considered in [36, 34].

We predict the six leading modes of the KSE with the following partial dynamics,

$$\frac{d}{dt}\hat{v}_k = (q_k^2 - q_k^4)\hat{v}_k - \frac{iq_k}{2} \sum_{1 \leq |l|, |k-l| \leq 6} \hat{v}_l \hat{v}_{k-l} + \hat{\theta}_k, \quad k = 1, \dots, 6. \quad (20)$$

In this case, since the nonlinear terms in (20) only involve summation of terms that are restricted to  $1 \leq |l|, |k-l| \leq 6$ , the identifiable unresolved variables,  $\theta_k$ , depends also on the resolved modes, in addition to the unresolved modes. The proposed closure model is to concatenate the numerical discretization of (20) with the discrete nonparametric closure model,

$$\hat{\theta}_{t+1} = \mathbb{E}^\epsilon[\Theta_{t+1}|\hat{v}_{t-m:t}, \hat{\theta}_{t-m:t}], \quad (21)$$

where  $\hat{\theta}_t = (\hat{\theta}_{1,t}, \dots, \hat{\theta}_{6,t}) \in \mathbb{C}^6$  and  $\hat{v}_t = (\hat{v}_{1,t}, \dots, \hat{v}_{6,t}) \in \mathbb{C}^6$ . In our numerical experiment, we set  $m = 19$  such that  $\mathbb{E}^\epsilon$  in (21) is a function that maps a real-valued vector of size  $(19 + 1) \times 12 \times 2 = 480$  to a 12-dimensional vector consisting of the real and imaginary values of  $\hat{\theta}_{t+1}$ . To evolve the dynamics in (20)-(21), we discretize (20) with the midpoint rule and a time step  $\Delta$ .

In our numerical experiment, the true time series for training are obtained by integrating the full dynamics, that is, (19) truncated on  $1 \leq |k| \leq 48$  with a time step  $\delta t = 0.005$ . We observe only the first 6 modes at a time step  $\Delta = 0.05$ . The size of the training data set is  $N = 2.5 \times 10^5$ . The identifiable unresolved variable,  $\theta_t$ , is estimated by fitting the time series  $v_t$  to the dynamics in (20). Subsequently, we use the pair  $\{\theta_t, v_t\}$  to train the LSTM model for (21).

Fig. 5(a) and (b) display the comparison of the short-time spatiotemporal manifestation between the full model and the closure model. One can see that the spatio-temporal pattern of the proposed closure model is consistent with that of the full KS model up to time  $t = 50$ . In fact, we can see a strong similarity in the spatio-temporal patterns for time  $t > 50$ . A close inspection shows an accurate path-wise prediction of  $\text{Re}(v_1)$  up to time 50 (see Fig. 5(c)) and an accurate recovery of long-time statistics including the ACFs, cross-correlation functions (CCFs) between  $|v_4|^2$  and  $|v_1|^2$ , PDFs, and energy spectra (see Figs. 5(d)-(g)). Complete results of all the other Fourier modes are shown in SI Appendix F.

While such an accurate recovery in path-wise and statistical prediction has also been achieved with the NARMAX parametric closure in [36, 34], careful choice of parametric ansatz is necessary with the NARMAX model. Here, an accurate recovery is obtained with a much simpler nonparametric model in (21).

## 5 Summary

We have presented a general nonparametric framework for prediction with missing dynamics. The proposed framework reformulates the closure model as a supervised learning problem in which the task is to approximate a high-dimensional map that takes the history of resolved and identifiable unresolved variables to the missing components in the resolved dynamics. Mathematically, we validate the approach with a theoretical guarantee for the path-wise convergence of the resolved variables. Numerically, we demonstrate the effectiveness of our framework in replicating severely truncated complex nonlinear problems arise in many applications. While the framework can be realized using any machine learning technique, we found that the LSTM as a special class of RNN is robust for this particular task.

From the positive numerical tests, several open questions deserve further investigation. For example, justifying the existence of the equilibrium distribution of the closure model; demonstrating the convergence to the underlying equilibrium distribution; clarifying the condition under which we can achieve a stable closure model.

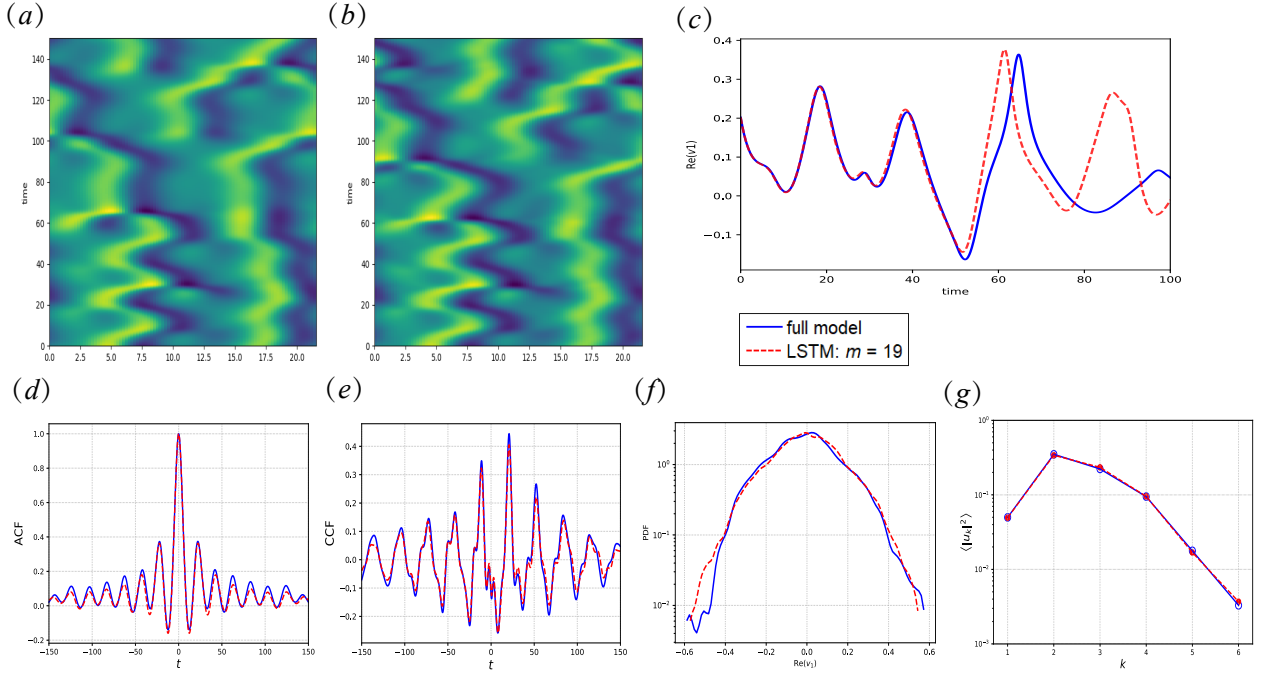


Figure 5: Comparison of spatiotemporal manifestation of KS solutions starting from the same initial conditions between (a) the full model and (b) the closure model using the LSTM method. Also plotted are the comparisons of (c) trajectories and (d) autocorrelation functions (ACFs) for  $\text{Re}(v_1)$ , (e) cross correlation functions (CCFs) between  $|v_1|^2$  and  $|v_4|^2$ , (f) probability density functions (PDFs) for  $\text{Re}(v_1)$ , and (g) the energy spectra  $\langle |v_k|^2 \rangle$  for the KS solutions between the full and the closure LSTM models.

**Acknowledgments.** The research of J.H. was partially supported by the ONR Grant N00014-16-1-2888 and NSF Grant DMS-1854299. J.H. thanks Di Qi for sharing the codes for the 57-mode barotropic stress model. S. L. and H. Y. gratefully acknowledge the support of National Supercomputing Center Singapore (NSCC) and High-Performance Computing (HPC) of the National University of Singapore for providing computational resources, and the support of NVIDIA Corporation with the donation of the Titan Xp GPU used for this research. H. Y. is partially supported by the startup of the Department of Mathematics at Purdue University.

## Appendix A: Proof of Theorem 1

Before we prove Theorem 1 in the main text, let us review the following bound which will be used below as well as in the proof of Theorem 3.

**Lemma 1.** *Let  $\alpha, c > 0$  be real numbers and  $m, T \geq 0$  be integers. Suppose that,*

$$E_{T+1} \leq \alpha \sum_{j=T-m}^T E_j + c,$$

*If  $E_j = 0$  for  $j = -m, \dots, 0$ , then for all integer  $T \geq 0$ .*

$$E_{T+1} \leq c(1 + \alpha)^T.$$

*Proof.* We proceed by induction. one can verify that,  $E_1 \leq c$ ,  $E_2 \leq c(1 + \alpha)$  and so on. In fact, we can verify for  $j = 0, \dots, m$  one by one that

$$E_j \leq c(1 + \alpha)^{j-1}. \quad (22)$$

By induction, for  $T \geq m$ , we have

$$E_{T+1} \leq c\alpha \sum_{j=T-m}^T (1 + \alpha)^{j-1} + c \leq c\alpha \sum_{j=1}^T (1 + \alpha)^{j-1} + c = c(1 + \alpha)^T. \quad (23)$$

□

Now we proceed with the proof of Theorem 1. As we mentioned before, since

$$\mathbb{E}[y_{t+1}|x_t, y_t] = \mathcal{G}(x_t, y_t), \quad (24)$$

we can rewrite the full dynamics as,

$$\begin{aligned} x_{t+1} &= \mathcal{F}(x_t, y_t), \\ y_{t+1} &= \mathbb{E}[Y_{t+1}|x_t, y_t]. \end{aligned}$$

We consider an approximate dynamics given as,

$$\begin{aligned} \hat{x}_{t+1} &= \mathcal{F}(\hat{x}_t, \hat{y}_t) \\ \hat{y}_{t+1} &= \mathbb{E}^\epsilon[Y_{t+1}|\hat{x}_t, \hat{y}_t] + \xi_{t+1}, \end{aligned}$$

where  $\xi_{t+1} \sim \Xi$  are Gaussian white noises with variance,

$$\mathbb{E}[\Xi^2] := \mathbb{E}[(Y_{t+1} - \mathbb{E}^\epsilon[Y_{t+1}|X_t, Y_t])^2] = \mathbb{E}[(\mathbb{E}[Y_{t+1}|X_t, Y_t] - \mathbb{E}^\epsilon[Y_{t+1}|X_t, X_t])^2] = O(\epsilon^2). \quad (25)$$

Define  $E_{x,t} := |x_{t+1} - \hat{x}_{t+1}|$  and  $E_{y,t} := |y_{t+1} - \hat{y}_{t+1}|$ , using the consistency in (24) and the Lipschitz conditions of  $\mathcal{F}$  and  $\mathcal{G}$ , we deduce

$$\begin{aligned} E_{y,t+1} &\leq |\mathbb{E}[Y_{t+1}|x_t, y_t] - \mathbb{E}^\epsilon[Y_{t+1}|\hat{x}_t, \hat{y}_t]| + |\xi_{t+1}| \\ &\leq |\mathbb{E}[Y_{t+1}|x_t, y_t] - \mathbb{E}[Y_{t+1}|\hat{x}_t, \hat{y}_t]| + |\mathbb{E}[Y_{t+1}|\hat{x}_t, \hat{y}_t] - \mathbb{E}^\epsilon[Y_{t+1}|\hat{x}_t, \hat{y}_t]| + |\xi_{t+1}| \\ &< |\mathcal{G}(x_t, y_t) - \mathcal{G}(\hat{x}_t, \hat{y}_t)| + |\mathbb{E}[Y_{t+1}|\hat{x}_t, \hat{y}_t] - \mathbb{E}^\epsilon[Y_{t+1}|\hat{x}_t, \hat{y}_t]| + |\xi_{t+1}| \\ &\leq L_1 E_{y,t} + L_2 E_{x,t} + |\mathbb{E}[Y_{t+1}|\hat{x}_t, \hat{y}_t] - \mathbb{E}^\epsilon[Y_{t+1}|\hat{x}_t, \hat{y}_t]| + |\xi_{t+1}| \end{aligned}$$



Define  $E_{x,T+1}^* := \mathbb{E}[\max_{t=\{0,\dots,T+1\}} E_{x,t}]$  and  $E_{y,T+1}^* := \mathbb{E}[\max_{t=\{0,\dots,T+1\}} E_{y,t}]$ . Then, by the Burkholder-Davis-Gundy inequality [56],

$$E_{y,T+1}^* \leq L_1 E_{y,T}^* + L_2 E_{x,T}^* + C\epsilon, \quad (26)$$

in which we have used (25) to bound the last two terms. This bound can be explicitly written as,

$$\begin{aligned} E_{y,T+1}^* &\leq L_1^{T+1} E_{y,0}^* + \sum_{j=0}^T L_1^j (L_2 E_{x,T-j}^* + C\epsilon) \\ &\leq L_1^{T+1} E_{y,0}^* + (L_2 E_{x,T}^* + C\epsilon) \sum_{j=0}^T L_1^j \\ &= (L_2 E_{x,T}^* + C\epsilon) \frac{L_1^{T+1} - 1}{L_1 - 1}. \end{aligned} \quad (27)$$

where we have used the fact that  $L_2 E_{x,t}^* + C\epsilon$  is non-decreasing to get the second inequality and  $E_{y,0}^* = 0$  to obtain the last equality.

Using similar algebra, we have

$$E_{x,T+1}^* \leq L_3 E_{x,T}^* + L_4 E_{y,T}^* \quad (28)$$

Inserting (27) into (28), we obtain

$$\begin{aligned} E_{x,T+1}^* &\leq L_3 E_{x,T}^* + L_4 (L_2 E_{x,T-1}^* + C\epsilon) \frac{L_1^T - 1}{L_1 - 1} \\ &\leq \alpha (E_{x,T}^* + E_{x,T-1}^*) + CL_4 \epsilon \frac{L_1^T - 1}{L_1 - 1} \end{aligned} \quad (29)$$

where we have define  $\alpha = \max\{L_3, L_2 L_4 \frac{L_1^T - 1}{L_1 - 1}\}$ .

Given  $E_{x,0}^* = 0$ , we apply the bound in Lemma 1 for  $m = 1$ ,

$$E_{x,T+1}^* \leq CL_4 \epsilon \frac{L_1^T - 1}{L_1 - 1} (1 + \alpha)^T = O(a^T \epsilon),$$

for some constant  $a > 1$  and the proof is complete.

## Appendix B: Proof of Theorem 2

Let  $Y_{t+1}^\Delta := \frac{Y_{t+1} - Y_t}{\Delta}$  such that,

$$\mathbb{E}[Y_{t+1}^\Delta | x_t, y_t] = g(x_t, y_t). \quad (30)$$

With this definition, we can rewrite the full dynamics as,

$$\begin{aligned} x_{t+1} &= x_t + f(x_t, y_t) \Delta + \Delta^{1/2} \sigma_x \xi_{x,t+1}, \\ y_{t+1} &= y_t + \mathbb{E}[Y_{t+1}^\Delta | x_t, y_t] \Delta + \Delta^{1/2} \sigma_y \xi_{y,t+1}. \end{aligned}$$

We consider an approximate dynamics given as,

$$\begin{aligned} \hat{x}_{t+1} &= \hat{x}_t + f(\hat{x}_t, \hat{y}_t) \Delta + \Delta^{1/2} \sigma_x \xi_{x,t+1}, \\ \hat{y}_{t+1} &= \hat{y}_t + \mathbb{E}^\epsilon[Y_{t+1}^\Delta | \hat{x}_t, \hat{y}_t] \Delta + \Delta^{1/2} \hat{\sigma}_y \xi_{y,t+1}. \end{aligned}$$

First, notice that

$$\begin{aligned}
\Delta \hat{\sigma}_y^2 &= \mathbb{E}[(Y_{t+1} - Y_t - \Delta \mathbb{E}^\epsilon[Y_{t+1}^\Delta | X_t, Y_t])^2] \\
&\leq \mathbb{E}[(Y_{t+1} - Y_t - \Delta \mathbb{E}[Y_{t+1}^\Delta | X_t, Y_t])^2 + \Delta^2(\mathbb{E}[Y_{t+1}^\Delta | X_t, Y_t] - \mathbb{E}^\epsilon[Y_{t+1}^\Delta | X_t, Y_t])^2 \dots \\
&\quad + 2\Delta \mathbb{E}[(Y_{t+1} - Y_t - \Delta \mathbb{E}[Y_{t+1}^\Delta | X_t, Y_t])(\mathbb{E}[Y_{t+1}^\Delta | X_t, Y_t] - \mathbb{E}^\epsilon[Y_{t+1}^\Delta | X_t, Y_t])] \\
&= \Delta \sigma_y^2 + O(\Delta^2 \epsilon^2),
\end{aligned} \tag{31}$$

where the last term vanishes since the mean of  $y_{t+1} - y_t - \Delta \mathbb{E}[Y_{t+1}^\Delta | x_t, y_t] = \Delta^{1/2} \sigma_y \xi_{x,t+1}$  is zero.

Define  $E_{x,t+1} := |x_{t+1} - \hat{x}_{t+1}|$  and  $E_{y,t+1} := |y_{t+1} - \hat{y}_{t+1}|$ , using the consistency in (30) and the Lipschitz conditions of  $f$  and  $g$ , we deduce

$$\begin{aligned}
E_{y,t+1} &\leq E_{y,t} + \Delta |\mathbb{E}[Y_{t+1}^\Delta | x_t, y_t] - \mathbb{E}^\epsilon[Y_{t+1}^\Delta | \hat{x}_t, \hat{y}_t]| + \Delta^{1/2} |\sigma_y - \hat{\sigma}_y| |\xi_{y,t+1}| \\
&\leq E_{y,t} + \Delta |\mathbb{E}[Y_{t+1}^\Delta | x_t, y_t] - \mathbb{E}[Y_{t+1}^\Delta | \hat{x}_t, \hat{y}_t]| + \Delta |\mathbb{E}[Y_{t+1}^\Delta | \hat{x}_t, \hat{y}_t] - \mathbb{E}^\epsilon[Y_{t+1}^\Delta | \hat{x}_t, \hat{y}_t]| + \Delta^{1/2} |\sigma_y - \hat{\sigma}_y| |\xi_{y,t+1}| \\
&< E_{y,t} + \Delta |g(x_t, y_t) - g(\hat{x}_t, \hat{y}_t)| + \Delta |\mathbb{E}[Y_{t+1}^\Delta | \hat{x}_t, \hat{y}_t] - \mathbb{E}^\epsilon[Y_{t+1}^\Delta | \hat{x}_t, \hat{y}_t]| + \Delta^{1/2} |\sigma_y - \hat{\sigma}_y| |\xi_{y,t+1}| \\
&\leq (1 + \Delta \ell) E_{y,t} + \Delta \ell E_{x,t} + \Delta |\mathbb{E}[Y_{t+1}^\Delta | \hat{x}_t, \hat{y}_t] - \mathbb{E}^\epsilon[Y_{t+1}^\Delta | \hat{x}_t, \hat{y}_t]| + \Delta^{1/2} |\sigma_y - \hat{\sigma}_y| |\xi_{y,t+1}|.
\end{aligned}$$

where  $\ell = O(1)$  denotes the largest Lipschitz constant in all directions. Define  $E_{x,T+1}^* := \mathbb{E}[\max_{t=\{0,\dots,T+1\}} E_{x,t}]$  and  $E_{y,T+1}^* := \mathbb{E}[\max_{t=\{0,\dots,T+1\}} E_{y,t}]$ . Then, by the Burkholder-Davis-Gundy inequality [56], we have

$$E_{y,T+1}^* \leq (1 + \Delta \ell) E_{y,T}^* + \Delta \ell E_{x,T}^* + C \Delta \epsilon, \tag{32}$$

which have used (31). Concatenate this with

$$E_{x,T+1}^* \leq (1 + \Delta \ell) E_{x,T}^* + \Delta \ell E_{y,T}^*,$$

we have

$$E_{T+1}^* \leq (I + A + A^2 + \dots + A^T) b,$$

where

$$E_{T+1}^* = \begin{pmatrix} E_{x,T+1}^* \\ E_{y,T+1}^* \end{pmatrix}, \quad A = \begin{pmatrix} 1 + \Delta \ell & \Delta \ell \\ \Delta \ell & 1 + \Delta \ell \end{pmatrix}, \quad b = \begin{pmatrix} 0 \\ C \Delta \epsilon \end{pmatrix}.$$

Using the fact that,

$$A = \frac{1}{2} \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 1 + 2\ell \Delta \end{pmatrix} \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix},$$

one can deduce that,

$$E_{x,T+1}^* \leq C \Delta \epsilon \left( -(T+1) + \frac{1 - (1 + 2\ell \Delta)^{T+1}}{-2\ell \Delta} \right) = O(\epsilon \Delta^2 T^2).$$

and the proof is complete.

### Appendix C: Proof of Theorem 3

In this case, we have

$$\mathbb{E}[\Theta_{t+1}|z_{t,m}] = P\mathcal{G}(x_t, \theta_t) + \sum_{k=1}^m P(F_k \circ \mathcal{G})(x_{t-k}, \theta_{t-k}). \quad (33)$$

where  $z_{t,m} = (x_{t-m:t}, \theta_{t-m:t})$ .

Define  $E_{\theta,t} := |\theta_{t+1} - \hat{\theta}_{t+1}|$  and  $E_{x,t} := |x_{t+1} - \hat{x}_{t+1}|$ . Since the orthogonal projection operators  $P, Q$  and the Koopman operator  $S$  are all bounded linear operator, for Lipschitz continuous  $\mathcal{F}$  and  $\mathcal{G}$ , it is clear that  $P\mathcal{G}$  and  $P(F_k \circ \mathcal{G})$  are all Lipschitz in  $x$  and  $\theta$ . Thus, we have

$$\begin{aligned} E_{\theta,t+1} &\leq |\mathbb{E}[\Theta_{t+1}|z_{t,m}] - \mathbb{E}^\epsilon[\Theta_{t+1}|\hat{z}_{t,m}] + |\xi_{t+1}| \\ &\leq |\mathbb{E}[\Theta_{t+1}|z_{t,m}] - \mathbb{E}[\Theta_{t+1}|\hat{z}_{t,m}]| + |\mathbb{E}[\Theta_{t+1}|\hat{z}_{t,m}] - \mathbb{E}^\epsilon[\Theta_{t+1}|\hat{z}_{t,m}]| + |\xi_{t+1}| \\ &\leq \sum_{k=0}^m |P(F_k \circ \mathcal{G})(x_{t-k}, \theta_{t-k}) - P(F_k \circ \mathcal{G})(\hat{x}_{t-k}, \hat{\theta}_{t-k})| + |\mathbb{E}[\Theta_{t+1}|\hat{z}_{t,m}] - \mathbb{E}^\epsilon[\Theta_{t+1}|\hat{z}_{t,m}]| + |\xi_{t+1}| \\ &\leq \sum_{s=t-m}^t K_{s-(t-m)} E_{\theta,s} + \sum_{s=t-m}^t L_{s-(t-m)} E_{x,s} + |\mathbb{E}[\Theta_{t+1}|\hat{z}_{t,m}] - \mathbb{E}^\epsilon[\Theta_{t+1}|\hat{z}_{t,m}]| + |\xi_{t+1}|, \end{aligned} \quad (34)$$

where  $K_s, L_s$  are Lipschitz constants.

Define  $E_{\theta,T+1}^* := \mathbb{E}[\max_{t=\{0,\dots,T+1\}} E_{\theta,t}]$  and  $E_{x,T+1}^* := \mathbb{E}[\max_{t=\{0,\dots,T+1\}} E_{x,t}]$ . Then we have,

$$E_{\theta,T+1}^* \leq \sum_{s=T-m}^T K_{s-(T-m)} E_{\theta,s}^* + \sum_{s=T-m}^T L_{s-(T-m)} E_{x,s}^* + C\epsilon, \quad (35)$$

where the expectation of the last term in (34) is bounded using the Burkholder-Davis-Gundy inequality [56]. We should point out that since the expectation in

$$\mathbb{E}[(\mathbb{E}[\Theta_{t+1}|Z_{t,m}] - \mathbb{E}^\epsilon[\Theta_{t+1}|Z_{t,m}])^2] = O(\epsilon^2),$$

is defined with respect to the invariant measure  $\pi$  supported on  $\mathcal{Z}$ , where  $0 < \pi(\mathcal{Z}) < \infty$ , it is clear that expectation of the third term in (34),

$$\mathbb{E}[|\mathbb{E}[\Theta_{t+1}|Z_{t,m}] - \mathbb{E}^\epsilon[\Theta_{t+1}|Z_{t,m}]|] \leq \mathbb{E}[|\mathbb{E}[\Theta_{t+1}|Z_{t,m}] - \mathbb{E}^\epsilon[\Theta_{t+1}|Z_{t,m}]|^2]^{1/2} \pi(\mathcal{Z})^{1/2} = C\epsilon. \quad (36)$$

is also bounded by order- $\epsilon$ .

Let  $0 < K := \max\{K_0, \dots, K_m\}$ , applying the bound in Lemma 1, we can obtain from (35)

$$E_{\theta,T+1}^* \leq \left( \sum_{s=T-m}^T L_{s-(T-m)} E_{x,s}^* + C\epsilon \right) (1 + K)^T. \quad (37)$$

Using similar algebra, we have

$$E_{x,T+1}^* \leq L_{m+1} E_{x,T}^* + K_{m+1} E_{\theta,T}^*, \quad (38)$$

for some constants  $K_{m+1}, L_{m+1} > 0$ . Inserting (37) into (38), let  $0 < L := \max_{j=0,\dots,m} \{L_{m+1}, K_{m+1} L_j (1 +$

$K)^{T-1}\}$ , applying the bound (23), we obtain

$$E_{x,T+1}^* \leq L_{m+1}E_{x,T}^* + K_{m+1}\left(\sum_{s=T-m-1}^{T-1} L_{s-(T-m-1)}E_{x,s}^* + C\epsilon\right)(1+K)^{T-1} \quad (39)$$

$$\leq L \sum_{s=T-m-1}^T E_{x,s}^* + K_{m+1}C\epsilon(1+K)^{T-1} \quad (40)$$

$$\leq K_{m+1}C\epsilon(1+K)^{T-1}(1+L)^T \quad (41)$$

$$= O(a^T\epsilon),$$

for some  $a > 1$  and the proof is completed.

## Appendix D: 57-mode barotropic stress equation

The 57-mode barotropic stress equation describes the anisotropic interaction between large and small scales through the geophysical effects from  $u(t)$ ,  $\beta$ , and the topography  $h(\mathbf{x})$  [57]:

$$\begin{aligned} \frac{du}{dt} + \oint \frac{\partial h}{\partial x} \psi &= -\bar{d}u + \sigma\mu^{-1/2}\dot{W}_0, \\ \frac{\partial \omega}{\partial t} + \nabla^\perp \psi \cdot \nabla q + u \frac{\partial q}{\partial x} + \beta \frac{\partial \psi}{\partial x} &= -\mathcal{D}\psi + \Sigma \dot{W}. \end{aligned} \quad (42)$$

On the right hand side of (42), the dissipation and forcing operators are applied on both the small and the large scales. On the small scale, the dissipation operator is in the form of  $\mathcal{D} = -\bar{d}\Delta$  with  $\bar{d} \geq 0$  and  $\Delta$  the Laplace operator corresponding to the Ekman drag dissipation. On the large scale, the constant damping  $-\bar{d}u$  in the large-scale mean flow represents the momentum damping. The forcing terms are represented by random Gaussian white noises (e.g. unresolved baroclinic instability processes on small scales, random wind stress, etc), where  $W(t)$  and  $W_0(t)$  are standard Wiener processes;  $\sigma\mu^{-1/2} > 0$  is a constant amplitude and  $\Sigma(\mathbf{x}) := \sum_{1 \leq |\mathbf{k}| \leq K} \sigma_{\mathbf{k}} e^{i\mathbf{k} \cdot \mathbf{x}}$ .

In the absence of dissipation and forcing on the right hand sides of (42), two quantities conserve in the statistical theory, one being the total kinetic energy  $E$ , and the other being the large-scale enstrophy  $\mathcal{E}$ , defined as

$$E = \frac{1}{2}u^2 + \frac{1}{2} \oint |\nabla \psi|^2, \quad \mathcal{E} = \beta u + \frac{1}{2} \oint (\omega + h)^2. \quad (43)$$

Details about the two conserved quantities can be found in [42, 43].

To construct a set of special solutions of (42), we consider the truncated spectral expansion of the state variables of interest with high wavenumber truncation  $K$  under standard Fourier basis  $\exp(i\mathbf{k} \cdot \mathbf{x})$  with  $\mathbf{k} = (k_x, k_y)$  based on the periodic boundary condition

$$\begin{aligned} \omega &= \sum_{1 \leq |\mathbf{k}| \leq K} \hat{\omega}_{\mathbf{k}} e^{i\mathbf{k} \cdot \mathbf{x}}, \\ \psi &= \sum_{1 \leq |\mathbf{k}| \leq K} \hat{\psi}_{\mathbf{k}} e^{i\mathbf{k} \cdot \mathbf{x}} = - \sum_{1 \leq |\mathbf{k}| \leq K} \frac{\hat{\omega}_{\mathbf{k}}}{|\mathbf{k}|^2} e^{i\mathbf{k} \cdot \mathbf{x}}, \end{aligned}$$

where  $\psi$  is the Galerkin truncated stream function, and  $\omega = \Delta\psi$  is the small-scale truncated relative vorticity. Consider the Galerkin projection operator  $\mathcal{P}_K$  on a subspace with high wavenumbers truncation  $1 \leq |\mathbf{k}| \leq K$ , we can obtain the barotropic equations (42) in Fourier form:

$$\begin{aligned}
\frac{du}{dt} &= i \sum_{1 \leq |\mathbf{k}| \leq K} \frac{k_x}{|\mathbf{k}|^2} \hat{h}_{\mathbf{k}}^* \hat{\omega}_{\mathbf{k}} - \bar{d} (u - u_{eq}) + \sigma \mu^{-1/2} \dot{W}_t, \\
\frac{d\hat{\omega}_{\mathbf{k}}}{dt} &= \mathcal{P}_{K,\mathbf{k}} (\nabla^\perp \psi \cdot \nabla q) + ik_x \left( \frac{\beta}{|\mathbf{k}|^2} - u \right) \hat{\omega}_{\mathbf{k}} - ik_x \hat{h}_{\mathbf{k}} u - \bar{d} (\hat{\omega}_{\mathbf{k}} - \hat{\omega}_{eq,\mathbf{k}}) + \sigma_{\mathbf{k}} \dot{W}_{\mathbf{k},t}, \quad 1 \leq |\mathbf{k}| \leq K,
\end{aligned} \tag{44}$$

Here,  $\hat{\omega}_{eq,\mathbf{k}} = -|\mathbf{k}|^2 \hat{h}_{\mathbf{k}} / (\mu + |\mathbf{k}|^2)$  is the mean relative vorticity,  $\sigma_{\mathbf{k}} = \sigma (1 + \mu |\mathbf{k}|^{-2})^{-1/2}$  is the forcing strength for each mode  $\mathbf{k}$ , and  $u_{eq} = -\beta/\mu$  is the equilibrium mean flow. The topography  $h(\mathbf{x})$  with  $\mathbf{x} = (x, y)$  is given by  $h(\mathbf{x}) = \sum_{1 \leq |\mathbf{k}| \leq K} \hat{h}_{\mathbf{k}} e^{i\mathbf{k} \cdot \mathbf{x}}$ . The parameter  $\sigma$  is chosen such that  $\sigma_{eq}^2 = \frac{\sigma^2}{2\bar{d}} = 1$  and the parameters  $\beta = 1$  and  $\mu = 2$  are fixed in our simulation.

In our implementation, we consider the truncation  $|\mathbf{k}| \leq K$  with  $K = 17$  which corresponds to 57 degrees of freedom for integers  $\mathbf{k} = (k_x, k_y)$ . For the space integration, we apply the pseudo-spectral method to the topographic barotropic model (42). For the nonlinear advection term,  $\mathcal{P}_{K,\mathbf{k}} (\nabla^\perp \psi \cdot \nabla q)$ , the 2/3 rule is applied for de-aliasing [57]. The noise is added using the standard Euler-Maruyama scheme. The true time series are generated from initial condition  $\hat{\psi}_{\mathbf{k}}(0)$ , drawn randomly from a Gaussian distribution with mean and variance,

$$\text{Mean}(\hat{\psi}_{\mathbf{k}}) = -\frac{h_{\mathbf{k}}}{\mu + |\mathbf{k}|^2}, \quad \text{Var}(\hat{\psi}_{\mathbf{k}}) = \frac{1}{|\mathbf{k}|^2(\mu + |\mathbf{k}|^2)}, \tag{45}$$

which are the ensemble mean of the small scale and the quadratic invariants of the enstrophy [41].

For the topography  $h(\mathbf{x})$ , we use the simple layered topography with variation only in  $x$ -direction

$$h(\mathbf{x}) = H (\cos(x) + \sin(x)) = \hat{h}_{(1,0)} e^{i(1,0) \cdot \mathbf{x}} + \hat{h}_{(-1,0)} e^{-i(1,0) \cdot \mathbf{x}},$$

where  $\hat{h}_{(1,0)} = H/2 - H/2i$  and  $\hat{h}_{(-1,0)} = \hat{h}_{(1,0)}^*$ , and  $H$  denotes the topography amplitude. Three regimes with different topographic strength  $H$  and damping  $\bar{d}$  are studied:

- Regime 1: weak coupling  $H = 3\sqrt{2}/4$  and strong damping  $\bar{d} = 0.5$  are applied to both small and large scales.
- Regime 2: intermediate coupling  $H = 5\sqrt{2}/4$  and intermediate damping  $\bar{d} = 0.1$  are applied to both small and large scales (shown in the main text).
- Regime 3: strong coupling  $H = 7\sqrt{2}/4$  and weak damping  $\bar{d} = 0.05$  are applied to both small and large scales. In the following, we show the numerical results in regime 1 and regime 3 for the large-scale mean flow  $u(t)$ .

In this example, we focus on the short-time forecasting and long-time statistics of the large-scale mean flow  $u(t)$  in the topographic barotropic equations (42). To compare the path-wise trajectories for short-time forecasting, we need the same noise realization,  $\eta_{t+1}$ , of the random variable  $\dot{W}_{t+1}$  in the full and closure models. However, we notice that the full model and the closure models are integrated with different time steps. The full model is integrated with a relatively small time step  $\delta t = 2.5\text{E-}3$  in order to resolve all the small-scale vorticity modes. Nevertheless, the closure models are integrated with a relatively large time step  $\Delta = 0.05$  for resolving only the large-scale mean flow  $u(t)$ . To compare the path-wise trajectories, we first generate a realization of the white noise  $\dot{W}_{t+1}$  from the identifiable variables of the full model using a finite difference method

$$\eta_{t+1} := \frac{u_{t+1} - [u_t + \Delta\theta_t - \Delta\bar{d}(u_t - u_{eq})]}{\sqrt{\Delta}\sigma\mu^{-1/2}},$$

where  $u_t$  and  $\theta_t$  are the identifiable variables from the dataset for verification of the full model observed with time interval  $\Delta = 0.05$ . Using  $\eta_{t+1}$  in (46) as a realization of  $\dot{W}_{t+1}$ , we can run the following closure model

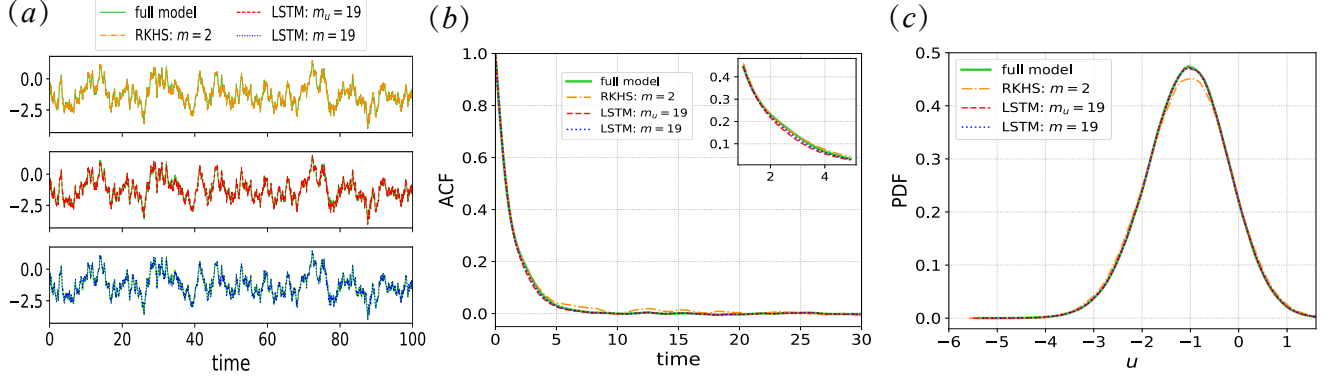


Figure 6: (Color online) Comparison of (a) trajectories, (b) ACFs, and (c) densities, between the full model and closure models for the large-scale mean velocity  $u$  in the regime of weak coupling  $H = 3\sqrt{2}/4$  and strong damping  $\bar{d} = 0.5$ .

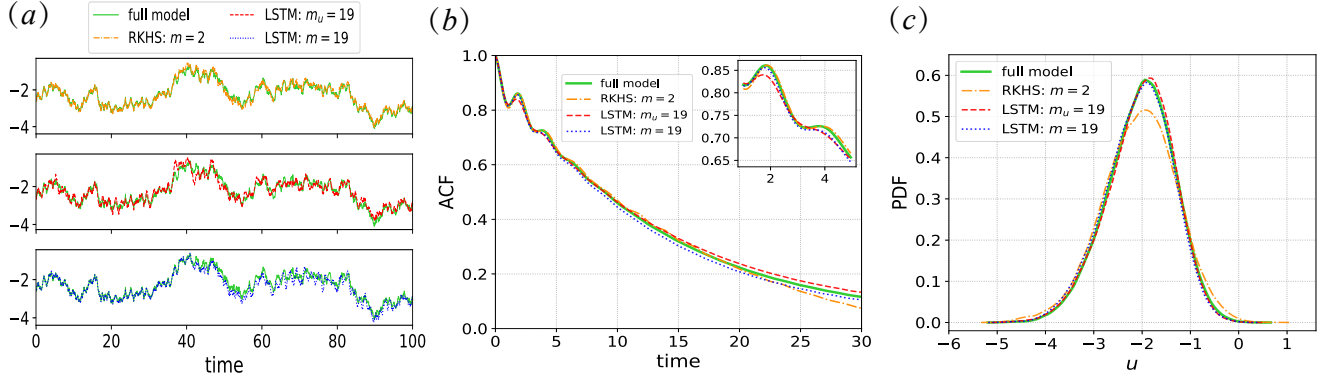


Figure 7: (Color online) Comparison of (a) trajectories, (b) ACFs, and (c) densities, between the full model and closure models for the large-scale mean velocity  $u$  in the regime of strong coupling  $H = 7\sqrt{2}/4$  and weak damping  $\bar{d} = 0.05$ .

with  $\Delta = 0.05$  starting at the same initial condition as the full model,

$$\begin{aligned}\hat{u}_{t+1} &= \hat{u}_t + \Delta \hat{\theta}_t - \Delta d (\hat{u}_t - u_{eq}) + \sqrt{\Delta} \sigma \mu^{-1/2} \eta_{t+1}, \\ \hat{\theta}_{t+1} &= \mathbb{E}^\epsilon [\Theta_{t+1} | \hat{u}_{t-m:t}, \hat{\theta}_{t-m:t}] + \xi_{t+1}.\end{aligned}\tag{46}$$

Beside comparing trajectories, we also verify the closure model in the estimation of two standard long-time statistical quantities, ACF's and PDF's. The auto-correlation function (ACF) for the large-scale mean flow  $u$  is calculated as  $\langle U_t U_0 \rangle / \langle U_0 U_0 \rangle$ , where  $U_t = u_t - \langle u_t \rangle$  with  $\langle \cdot \rangle$  being the temporal average over  $1.25 \times 10^6$  data for verification, which is different from the  $N = 1.25 \times 10^6$  training dataset. The probability density function (PDF) for  $u$  is obtained from the same verification dataset using the kernel density estimation (KDE) method.

Figures 6 and 7 display the comparison of short-time forecasting and long-time statistics in the weak coupling regime and strong coupling regime, respectively. In the weak coupling regime, one can see from Fig. 6 that the short-time predictions are all excellent among 3 closure models. For the long-time statistics, both the LSTM methods provide a better approximation than the RKHS model. This is because not enough memory terms are used in the RKHS model. In the strong coupling regime, one can see from Fig. 7 that

the LSTM method with  $m = 19$  is the best approximation. This is because that the variance of the training residual is small about  $10^{-5}$  for the LSTM with  $m = 19$  and the variance is relatively large about  $10^{-1}$  for the LSTM with  $m_u = 19$ . This result confirms the robustness of our framework with a closure model that depends on, both, the memories of the resolved and identifiable unresolved variables.

## Appendix E: Nonlinear Schrödinger equation

In this section, we discuss some details for the NLS equation and then show the numerical results in a low-temperature regime. The solution of the NLS equation can be expanded by the Fourier series,

$$u(x, t) = \sum_{k \in \mathbb{Z}} u_k(t) e^{ikx}. \quad (47)$$

Then, the NLS equation for the Fourier modes can be written as

$$\frac{du_k}{dt} = -i\omega_k u_k - i \sum_{k_1 \in \mathbb{Z}} \sum_{k_2 \in \mathbb{Z}} u_{k_1} u_{k_2} u_{k_1+k_2-k}^*, \quad (48)$$

where the dispersion relation is given by

$$\omega_k = k^2. \quad (49)$$

Notice that for the zeroth mode  $k = 0$ , the dispersion relation vanishes,  $\omega_k = 0$ , at the right hand side of (48). Writing the zeroth mode, we have,

$$\frac{du_0}{dt} = -i \sum_{k_1 \in \mathbb{Z}} \sum_{k_2 \in \mathbb{Z}} u_{k_1} u_{k_2} u_{k_1+k_2}^* = -i(|u_0|^2 u_0 + 2 \sum_{k \neq 0} |u_k|^2 u_0 + (\sum_{k_1, k_2 \in \mathbb{Z}} u_{k_1} u_{k_2} u_{k_1+k_2}^* - 2 \sum_{k \neq 0} |u_k|^2 u_0 - |u_0|^2 u_0)). \quad (50)$$

For our closure model, we assume that the identifiable unresolved variables are defined as,

$$\begin{aligned} \theta_1 &= 2 \sum_{k \neq 0} |u_k|^2, \\ \theta_2 &= \sum_{k_1, k_2 \in \mathbb{Z}} u_{k_1} u_{k_2} u_{k_1+k_2}^* - 2 \sum_{k \neq 0} |u_k|^2 u_0 - |u_0|^2 u_0, \end{aligned}$$

which are functions of only the unresolved variables  $\{u_k\}_{k \neq 0}$ . Thus, the dynamics of the resolved variable  $u_0$  can be rewritten as,

$$\frac{du_0}{dt} = -i \sum_{k_1 \in \mathbb{Z}} \sum_{k_2 \in \mathbb{Z}} u_{k_1} u_{k_2} u_{k_1+k_2}^* = -i(|u_0|^2 u_0 + \theta_1 u_0 + \theta_2). \quad (51)$$

The closure model is obtained by concatenating a discretization on (51) with time step  $\Delta$  with a map,

$$\mathbb{E}^\epsilon[\Theta_{t+1} | \cdot] : \mathbb{R}^{(m+1) \times 4} \rightarrow \mathbb{R}^2,$$

defined as,

$$\theta_{t+1} = (\theta_{1,t+1}, \theta_{2,t+1}) = \mathbb{E}^\epsilon[\Theta_{t+1} | u_{0,t-m:t}, \theta_{t-m:t}]. \quad (52)$$

Here  $\theta_{t+1} := (\theta_{1,t+1}, \theta_{2,t+1})$  denotes the values of  $\theta_1$  and  $\theta_2$  at discrete time  $t+1$ . In our numerical experiment, we will set  $m = 19$  and train the map in (52) with the LSTM model using a time series of length  $N = 5 \times 10^5$ . We discretize (51) with the following time-splitting scheme. We use Euler scheme to solve,

$$\frac{du_0}{dt} = -i(\theta_1 u_0 + \theta_2),$$

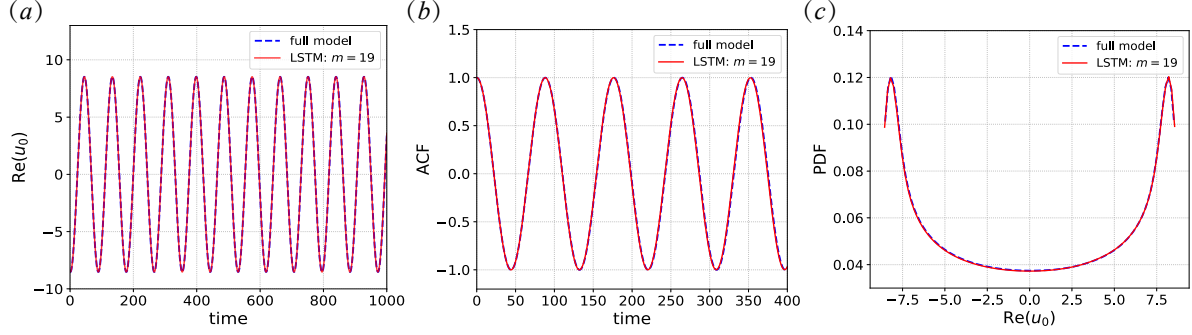


Figure 8: (Color online) Comparison of (a) trajectories, (b) ACFs, and (c) densities of the full model and closure models in the low-temperature regime with  $k_B T = 10^{-1}$ . The time step for the closure model is  $\Delta t = 0.02$  and the LSTM method uses  $m = 19$  for  $(\theta, u_0)$  in the closure model.

since we only have discrete estimates of  $\theta_1$  and  $\theta_2$ , and we use the explicit solution  $u_0(t) = |u_0(t_0)| \exp(-i|u_0(t_0)|t)$  for the nonlinear ODE,  $du_0/dt = -i|u_0|^2 u_0$ .

As mentioned in the main text, we simulate the initial conditions by sampling from the Gibbs distribution  $\pi = \exp(-\frac{E}{k_B T})$ , where  $E$  denotes the Hamiltonian of the ODE system resulting from the Fourier representation (48);  $k_B$  and  $T$  denote the Boltzmann constant and temperature, respectively. In this case, the Hamiltonian is given by  $E = E_2 + E_4$ ,

$$E_2 = \sum_{k \in \mathbb{Z}} \omega_k |u_k|^2, \quad E_4 = \frac{1}{2} \sum_{k_1 \in \mathbb{Z}} \sum_{k_2 \in \mathbb{Z}} \sum_{k_3 \in \mathbb{Z}} u_{k_1} u_{k_2} u_{k_3}^* u_{k_1+k_2-k_3}^*.$$

We now show the numerical results for the zeroth mode  $u_0$  in a low-temperature regime with  $k_B T = 10^{-1}$ . Notice that the amplitude and the frequency of the solution as shown in Fig. 8(a) are much smaller and slower, respectively, compared to those in the case of higher temperature (see Fig. 4 in the main text). Since smaller time steps are required for accurate solutions with higher amplitude as well as the faster frequency, the problem is numerically stiff as the temperature increases.

We observe from Fig. 8(a) that the path-wise trajectory of  $\text{Re}(u_0)$  is well captured for sufficiently long time. For the statistics, we observe from Figs. 8(b) and (c) that the ACF and PDF of the true  $u_0$  are also accurately reproduced. Here, the ACF for the  $\text{Re}(u_0)$  is calculated as  $\langle V_t V_0 \rangle / \langle V_0 V_0 \rangle$ , where  $V_t = \text{Re}(u_{0,t}) - \langle \text{Re}(u_{0,t}) \rangle$  with  $\langle \cdot \rangle$  being the temporal average over  $5 \times 10^5$  verification data, which is different from the  $N = 5 \times 10^5$  training dataset. The PDF for  $\text{Re}(u_0)$  is obtained from the same verification dataset using the kernel density estimation (KDE) method.

## Appendix F: Kuramoto-Sivashinsky equation

In this section, we report the full results on the numerical experiment discussed in the main text.

In particular, Figure 9 displays the results for the comparison of trajectories, ACFs, CCFs, and PDFs for all the Fourier modes  $v_1, \dots, v_6$ . Here, the CCF's are defined as the cross correlation functions between  $|v_k|^2$  and  $|v_4|^2$ . The ACF's and CCF's are estimated by Monte-Carlo estimation over  $2.5 \times 10^5$  number of dataset. The PDF's are computed using the KDE method. We can see that the path-wise trajectories can be well captured up to time 50 for all modes. Moreover, we can see that ACFs, CCFs, and PDFs can be well reproduced by the LSTM for all modes. Therefore, the closure model using the LSTM can provide an accurate recovery for both the short-time forecasting and the long-time statistics of the KSE.



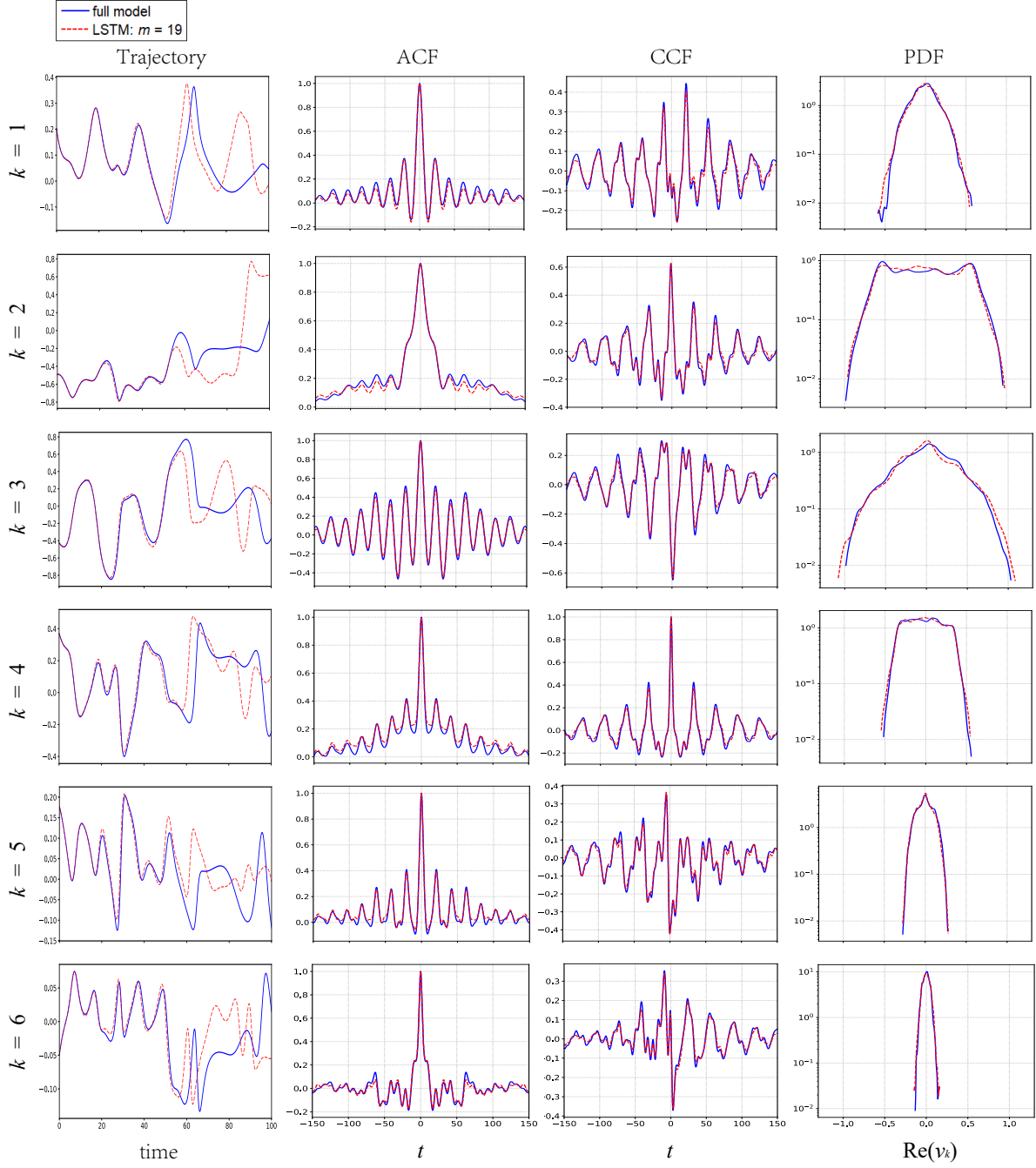


Figure 9: (Color online) Comparison of trajectories for  $\text{Re}(v_k)$ , ACFs for  $\text{Re}(v_k)$ , CCFs between  $|v_k|^2$  and  $|v_4|^2$ , and PDFs for  $\text{Re}(v_k)$  between the full model and the closure model. Solid blue line corresponds to the full model and dashed red line corresponds to the closure model.

## References

- [1] Z. Allen-Zhu and Y. Li. Can SGD learn recurrent neural networks with provable generalization? *CoRR*, abs/1902.01028, 2019.
- [2] Z. Allen-Zhu, Y. Li, and Y. Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. *CoRR*, abs/1811.04918, 2018.
- [3] Z. Allen-Zhu, Y. Li, and Z. Song. On the convergence rate of training recurrent neural networks. *CoRR*, abs/1810.12065, 2018.
- [4] W. Bao, S. Jin, and P. A. Markowich. Numerical study of time-splitting spectral discretizations of nonlinear schrödinger equations in the semiclassical regimes. *SIAM Journal on Scientific Computing*, 25(1):27–64, 2003.
- [5] A. R. Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information Theory*, 39(3):930–945, May 1993.
- [6] T. Berry and J. Harlim. Linear Theory for Filtering Nonlinear Multiscale Systems with Model Error. *Proc. Roy. Soc. A* 20140168, 2014.
- [7] T. Berry and J. Harlim. Semiparametric modeling: Correcting low-dimensional model error in parametric models. *J. Comput. Phys.*, 308:305–321, 2016.
- [8] A. Brutzkus, A. Globerson, E. Malach, and S. Shalev-Shwartz. SGD learns over-parameterized networks that provably generalize on linearly separable data. *CoRR*, abs/1710.10174, 2017.
- [9] Y. Cao and Q. Gu. A generalization theory of gradient descent for learning over-parameterized deep relu networks. *CoRR*, abs/1902.01384, 2019.
- [10] G. F. Carnevale and J. S. Frederiksen. Nonlinear stability and statistical mechanics of flow over topography. *Journal of Fluid Mechanics*, 175:157–181, 1987.
- [11] M. Chen, X. Li, and T. Zhao. On generalization bounds of a family of recurrent neural networks, 2019.
- [12] A. Chorin, O. Hald, and R. Kupferman. Optimal prediction with memory. *Physica D: Nonlinear Phenomena*, 166(3):239–257, 2002.
- [13] A. Chorin and P. Stinis. Problem reduction, renormalization, and memory. *Communications in Applied Mathematics and Computational Science*, 1(1):1–27, 2007.
- [14] A. J. Chorin and F. Lu. Discrete approach to stochastic parametrization and dimension reduction in nonlinear dynamics. *Proc. Nat. Acad. Sci.*, 112(32):9804–9809, 2015.
- [15] D. Crommelin and E. Vanden-Eijnden. Subgrid-scale parameterization with conditional markov chains. *Journal of the Atmospheric Sciences*, 65(8):2661–2675, 2008.
- [16] E. Darve, J. Solomon, and A. Kia. Computing generalized langevin equations and generalized fokker-planck equations. *Proceedings of the National Academy of Sciences*, 106(27):10884–10889, 2009.
- [17] W. E, C. Ma, and L. Wu. A priori estimates of the generalization error for two-layer neural networks. *ArXiv*, abs/1810.06397, 2018.

- [18] C. Franzke, I. Horenko, A. J. Majda, and R. Klein. Systematic metastable atmospheric regime identification in an agcm. *Journal of the Atmospheric Sciences*, 66(7):1997–2012, 2009.
- [19] T. Gao, J. Duan, X. Li, and R. Song. Mean exit time and escape probability for dynamical systems driven by lvy noises. *SIAM Journal on Scientific Computing*, 36(3):A887–A906, 2014.
- [20] D. Givon, R. Kupferman, and A. Stuart. Extracting macroscopic dynamics: model problems and algorithms. *Nonlinearity*, 17(6):R55, 2004.
- [21] J. Goodman. Stability of the kuramoto-sivashinsky and related systems. *Communications on Pure and Applied Mathematics*, 47(3):293–306, 1994.
- [22] A. Gouasmi, E. J. Parish, and K. Duraisamy. A priori estimation of memory effects in reduced-order models of nonlinear systems using the mori–zwanzig formalism. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 473(2205):20170385, 2017.
- [23] M. J. Grote, A. J. Majda, and C. G. Ragazzo. Dynamic mean flow and small-scale interaction through topographic stress. *Journal of Nonlinear Science*, 9(1):89–130, 1999.
- [24] J. Harlim and X. Li. Parametric reduced models for the nonlinear Schrödinger equation. *Phys. Rev. E.*, 91:053306, 2015.
- [25] J. Harlim, A. Mahdi, and A. Majda. An ensemble Kalman filter for statistical estimation of physics constrained nonlinear regression models. *J. Comput. Phys.*, 257, Part A:782–812, 2014.
- [26] S. Jiang and J. Harlim. Modeling of missing dynamical systems: Deriving parametric models using a nonparametric framework. *submitted to J. Nonlinear Sci.*, 2019.
- [27] P. Kevrekidis and D. Frantzeskakis. Solitons in coupled nonlinear schrödinger models: a survey of recent developments. *Reviews in Physics*, 1:140–153, 2016.
- [28] B. Khouider, J. A. Biello, and A. J. Majda. A stochastic multicloud model for tropical convection. *Comm. Math. Sci.*, 8:187–216, 2010.
- [29] D. Kondrashov, M. D. Chekroun, and M. Ghil. Data-driven non-markovian closure models. *Physica D: Nonlinear Phenomena*, 297:33–55, 2015.
- [30] H. Kunita. *Stochastic flows and stochastic differential equations*, volume 24. Cambridge university press, 1997.
- [31] Y. Kuramoto and T. Tsuzuki. Persistent Propagation of Concentration Waves in Dissipative Media Far from Thermal Equilibrium. *Progress of Theoretical Physics*, 55(2):356–369, 02 1976.
- [32] F. Kwasniok. Data-based stochastic subgrid-scale parametrization: an approach using cluster-weighted modelling. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 370(1962):1061–1086, 2012.
- [33] R. E. LaQuey, S. Mahajan, P. Rutherford, and W. Tang. Nonlinear saturation of the trapped-ion mode. *Physical Review Letters*, 34(7):391, 1975.
- [34] K. K. Lin and F. Lu. Data-driven model reduction, wiener projections, and the mori-zwanzig formalism. *arXiv preprint arXiv:1908.07725*, 2019.

- [35] F. Lu, K. Lin, and A. Chorin. Comparison of continuous and discrete-time data-based modeling for hypoelliptic systems. *Communications in Applied Mathematics and Computational Science*, 11(2):187–216, 2016.
- [36] F. Lu, K. K. Lin, and A. J. Chorin. Data-based stochastic model reduction for the kuramoto–sivashinsky equation. *Physica D: Nonlinear Phenomena*, 340:46–57, 2017.
- [37] F. Lu, X. Tu, and A. J. Chorin. Accounting for model error from unresolved scales in ensemble kalman filters by stochastic parameterization. *Monthly Weather Review*, 145(9):3709–3723, 2017.
- [38] Z. Lu, H. Pu, F. Wang, Z. Hu, and L. Wang. The expressive power of neural networks: A view from the width. *CoRR*, abs/1709.02540, 2017.
- [39] C. Ma and J. Wang. Model reduction with memory and the machine learning of dynamical systems. *arXiv preprint arXiv:1808.04258*, 2018.
- [40] A. Majda and J. Harlim. Physics constrained nonlinear regression models for time series. *Nonlinearity*, 26:201–217, 2013.
- [41] A. Majda, I. Timofeyev, and E. Vanden-Eijnden. Systematic strategies for stochastic mode reduction in climate. *Journal of the Atmospheric Sciences*, 60:1705–1722, 2003.
- [42] A. Majda and X. Wang. *Nonlinear dynamics and statistical theories for basic geophysical flows*. Cambridge University Press, UK, 2006.
- [43] A. J. Majda. Statistical energy conservation principle for inhomogeneous turbulent dynamical systems. *Proceedings of the National Academy of Sciences*, 112(29):8937–8941, 2015.
- [44] A. J. Majda, I. Timofeyev, and E. V. Eijnden. Models for stochastic climate prediction. *Proceedings of the National Academy of Sciences*, 96(26):14687–14691, 1999.
- [45] A. J. Majda, I. Timofeyev, and E. Vanden Eijnden. A mathematical framework for stochastic climate models. *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, 54(8):891–974, 2001.
- [46] R. Maulik, A. Mohan, B. Lusch, S. Madireddy, and P. Balaprakash. Time-series learning of latent-space dynamics for reduced-order model closure. *arXiv preprint arXiv:1906.07815*, 2019.
- [47] H. Montanelli and Q. Du. New error bounds for deep networks using sparse grids. 2017.
- [48] H. Montanelli and H. Yang. Error bounds for deep relu networks using the kolmogorov–arnold superposition theorem. *arXiv:1906.11945v1 [math.NA]*, 2019.
- [49] H. Montanelli, H. Yang, and Q. Du. Deep relu networks overcome the curse of dimensionality for bandlimited functions. *arXiv:1903.00735 [math.NA]*, 2019.
- [50] H. Mori. Transport, collective motion, and Brownian motion. *Prog. Theor. Phys.*, 33:423 – 450, 1965.
- [51] M. C. Mozer. Backpropagation. chapter A Focused Backpropagation Algorithm for Temporal Pattern Recognition, pages 137–169. L. Erlbaum Associates Inc., Hillsdale, NJ, USA, 1995.
- [52] R. Nakada and M. Imaizumi. Adaptive approximation and estimation of deep neural network to intrinsic dimensionality. *arXiv:1907.02177 [stat.ML]*, 2019.

- [53] S. Pan and K. Duraisamy. Data-driven discovery of closure models. *SIAM Journal on Applied Dynamical Systems*, 17(4):2381–2413, 2018.
- [54] E. J. Parish and K. Duraisamy. A dynamic subgrid scale model for large eddy simulations based on the mori–zwanzig formalism. *Journal of Computational Physics*, 349:154–175, 2017.
- [55] R. Pascanu, T. Mikolov, and Y. Bengio. On the difficulty of training recurrent neural networks. In *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28*, ICML’13, pages III–1310–III–1318. JMLR.org, 2013.
- [56] G. Pavliotis and A. Stuart. *Multiscale Methods: Averaging and Homogenization*, volume 53 of *Texts in Applied Mathematics*. Springer, 2000.
- [57] D. Qi and A. J. Majda. Low-dimensional reduced-order models for statistical response and uncertainty quantification: Barotropic turbulence with topography. *Physica D: Nonlinear Phenomena*, 343:7–27, 2017.
- [58] A. J. Robinson and F. Fallside. The utility driven dynamic error propagation network. Technical Report CUED/F-INFENG/TR.1, Engineering Department, Cambridge University, Cambridge, UK, 1987.
- [59] Z. Shen, H. Yang, and S. Zhang. Deep network approximation characterized by number of neurons. *arXiv:1906.05497 [math.NA]*, 2019.
- [60] L. Song, K. Fukumizu, and A. Gretton. Kernel embeddings of conditional distributions: A unified kernel framework for nonparametric inference in graphical models. *IEEE Signal Process. Mag.*, 30(4):98–111, 2013.
- [61] L. Song, J. Huang, A. Smola, and K. Fukumizu. Hilbert space embeddings of conditional distributions with applications to dynamical systems. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 961–968. ACM, 2009.
- [62] C. J. Stone. Optimal global rates of convergence for nonparametric regression. *The annals of statistics*, pages 1040–1053, 1982.
- [63] E. Vanden-Eijnden. Transition-path theory and path-finding algorithms for the study of rare events. *Annual review of physical chemistry*, 61:391–420, 2010.
- [64] P. R. Vlachas, W. Byeon, Z. Y. Wan, T. P. Sapsis, and P. Koumoutsakos. Data-driven forecasting of high-dimensional chaotic systems with long short-term memory networks. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 474(2213):20170844, 2018.
- [65] E. Weinan, B. Engquist, X. Li, W. Ren, and E. Vanden-Eijnden. Heterogeneous multiscale methods: a review. *Commun. Comput. Phys*, 2(3):367–450, 2007.
- [66] P. J. Werbos. Generalization of backpropagation with application to a recurrent gas market model. *Neural Networks*, 1(4):339 – 356, 1988.
- [67] D. S. Wilks. Effects of stochastic parametrizations in the lorenz’96 system. *Quarterly Journal of the Royal Meteorological Society*, 131(606):389–407, 2005.
- [68] R. Zwanzig. Statistical mechanics of irreversibility. *Lectures in Theoretical Physics*, 3:106–141, 1961.