

Coarse-scale PDEs from fine-scale observations via machine learning

Seungjoon Lee,¹ Mahdi Kooshkbaghi,² Konstantinos Spiliotis,³ Constantinos I. Siettos,⁴ and Ioannis G. Kevrekidis^{1, a)}

¹⁾*Department of Chemical and Biomolecular Engineering, Johns Hopkins University*

²⁾*Program in Applied and Computational Mathematics, Princeton University*

³⁾*Institute of Mathematics, University of Rostock*

⁴⁾*Dipartimento di Matematica e Applicazioni “Renato Caccioppoli”, Università degli Studi di Napoli Federico II*

(Dated: 13 September 2019)

Complex spatiotemporal dynamics of physicochemical processes are often modeled at a microscopic level (through e.g. atomistic, agent-based or lattice models) based on first principles. Some of these processes can also be successfully modeled at the macroscopic level using e.g. partial differential equations (PDEs) describing the evolution of the right few macroscopic observables (e.g. concentration and momentum fields). Deriving good macroscopic descriptions (the so-called “closure problem”) is often a time-consuming process requiring deep understanding/intuition about the system of interest. Recent developments in data science provide alternative ways to effectively extract/learn accurate macroscopic descriptions approximating the underlying microscopic observations. In this paper, we introduce a data-driven framework for the identification of unavailable coarse-scale PDEs from microscopic observations via machine learning algorithms. Specifically, using Gaussian Processes, Artificial Neural Networks, and/or Diffusion Maps, the proposed framework uncovers the relation between the relevant macroscopic space fields and their time evolution (the right-hand-side of the explicitly unavailable macroscopic PDE). Interestingly, several choices equally representative of the data can be discovered. The framework will be illustrated through the data-driven discovery of macroscopic, concentration-level PDEs resulting from a fine-scale, Lattice Boltzmann level model of a reaction/transport process. Once the coarse evolution law is identified, it can be simulated to produce long-term macroscopic predictions. Different features (pros as well as cons) of alternative machine learning algorithms for performing this task (Gaussian Processes and Artificial Neural Networks), are presented and discussed.

Keywords: Gaussian process regression; Artificial neural networks; Data mining; System identification, Manifold learning; Diffusion maps

The behavior of microscopic complex systems is often described in terms of effective, macroscopic governing equations, leading to simple and efficient prediction. Yet, the discovery/derivation of such macroscopic governing equations generally relies on deep understanding and prior knowledge about the system, as well as extensive and time-consuming mathematical justification. Recent developments in data-driven computational approaches suggest alternative ways towards uncovering useful coarse-scale governing equations directly from fine scale data. Interestingly, even deciding what the “right” coarse-scale variables are, may present a significant challenge. In this paper, we introduce and implement a framework for systematically extracting coarse-scale observables from microscopic/fine scale data and for discovering the underlying governing equations using machine learning techniques (e.g. Gaussian processes and artificial neural networks) enhanced by feature selection methods. Intrinsic representations of the coarse-scale behavior via manifold learning techniques (in particular, Diffusion Maps), generating alternative possible forms of the governing equations is also explored and discussed.

I. INTRODUCTION

The successful description of the spatiotemporal evolution of complex systems typically relies on detailed mathematical models operating at a fine scale (e.g. molecular dynamics, agent-based, stochastic or lattice-based methods). Such microscopic, first principles models, keeping track of the interactions between huge numbers microscopic level degrees of freedom, typically lead to prohibitive computational cost for large-scale spatiotemporal simulations.

To address this issue (and since we are typically interested in macro-scale features -pressure drops, reaction rates- rather than the position and velocity of each individual molecule), reduced, coarse-scale models are developed and used, leading to significant computational savings in large-scale spatiotemporal simulations¹.

Macroscopically, the fine scale processes may often be successfully modeled using partial differential equations (PDEs) in terms of the right macroscopic observables (“coarse variables”: not molecules and their velocities, say, but rather pressure drops and momentum fields). Deriving the macroscopic PDE that effectively models the microscopic physics (the so-called “closure problem”) requires, however, deep understanding/intuition about the complex system of interest and often extensive mathematical operations; the discovery of macroscopic governing equations is typically a difficult and time-consuming process.

To bypass the first principles discovery of a macroscopic PDE directly, several data-driven approaches provide ways

^{a)}Also at Department of Applied Mathematics and Statistics; Department of Medicine, Johns Hopkins University; Electronic mail: yannisk@jhu.edu

to effectively determine good coarse observables and the approximate coarse-scale relations between them from simulation data. In the early '90s researchers (including our group) employed artificial neural networks for system identification (both lumped and distributed)²⁻⁶. Projective time integration in dynamical systems⁷ and fluid dynamics^{8,9} also provides a good data-driven approximation of long-time prediction based not on closed-form equations, but rather on a “black box” simulator. Furthermore, the equation-free framework for designing fine scale computational experiments and systematically processing their results through “coarse-time steps” has proven its usefulness / computational efficiency in analyzing macroscopic bifurcation diagrams^{10,11}. The easy availability of huge simulation data sets, and recent developments in the efficient implementation of machine learning algorithms, has made the revisiting of the identification of nonlinear dynamical systems from simulation time series an attractive -and successful- endeavor. Working with observations at the macroscopic level, hidden macroscopic PDEs can be recovered directly by artificial neural networks⁶, (see also¹²). Sparse identification of nonlinear dynamics (SINDy)¹³ as well as Gaussian processes^{14,15} have also been successfully used, resulting in *explicit* data-driven PDEs. All these approaches rely on macroscopic observations.

In this paper, we discuss the identification of unavailable coarse-scale PDEs *from fine scale observations* through a combination of machine learning and manifold learning algorithms. Specifically, using Gaussian Processes, Artificial Neural Networks, and/or Diffusion Maps, and starting with candidate coarse fields (e.g. densities), our procedure extracts relevant macroscopic features (e.g. coarse derivatives) from the data, and then uncovers the relations between these macroscopic features and their time evolution (the right-hand-side of the explicitly unavailable macroscopic PDE).

To effectively reduce the input data domain, we employ two feature selection methods: (1) a sensitivity analysis via automatic relevance determination (ARD)¹⁶⁻¹⁸ in Gaussian processes and (2) a manifold learning technique, Diffusion Maps¹⁹. Having selected the relevant macro features in terms of which the evolution can be modelled, we employ two machine learning algorithms to approximate a “good” right-hand-side of the underlying PDEs: (1) Gaussian process regression and (2) artificial neural networks.

Our framework is illustrated through the data-driven discovery of the macroscopic, concentration-level PDE resulting from a fine-scale, Lattice Boltzmann (LB) model of a reaction/transport process (the FitzHugh-Nagumo process in one spatial dimension). Long-term macroscopic prediction is enabled by numerical simulation of the coarse-scale PDE *identified from the Lattice-Boltzmann data*. Different possible feature combinations (leading to different realizations of the same evolution) will also be discussed.

The remainder of the paper is organized as follows: In section II, we present an overview of our proposed framework and briefly review theoretical concepts of Gaussian process regression, artificial neural networks, and Diffusion Maps. Two methods for feature selection are also presented. In section III, we describe two simulators at different scales: (1) the

FitzHugh-Nagumo model at the macro-scale and (2) its Lattice Boltzmann realization at the micro-scale. In section IV, we demonstrate the effectiveness of our proposed framework and discuss the advantages and challenges of different feature selection methods and regression models for performing this task. In section V, we summarize our results and discuss open issues for further development of the data-driven discovery of the underlying coarse PDE from microscopic observations.

II. FRAMEWORK FOR RECOVERING A COARSE-SCALE PDE VIA MACHINE LEARNING

A. Overview

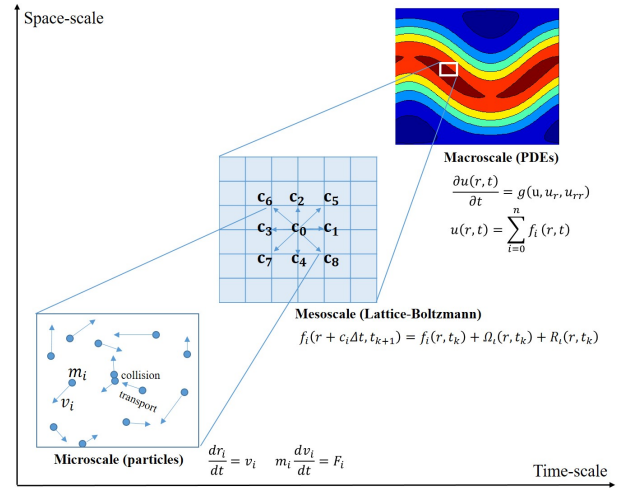


FIG. 1. Schematic illustration of the extraction of coarse-scale observables u from microscopic observations. Through a Lattice Boltzmann model (here, D2Q9), we obtain particle distribution functions (f_i) on a given lattice. Using the zeroth moment field of the particle distribution function at the grid point x_n , we extract the coarse observable u (in this paper, we have two coarse observables, u and v , which represent the density of the activator and the inhibitor, respectively).

The workflow of our framework for recovering hidden coarse-scale PDEs from microscopic observations is schematically shown in figures 1 and 2. Specifically, this framework consists of two subsections: (1) computing coarse-scale observables and (2) identifying coarse-scale PDEs and then numerically simulating them.

To clarify the algorithm, consider a single field (the activator u ; later in this paper we will use two densities, u and v , for the activator and the inhibitor, respectively). As shown in figure 1, we compute the coarse-scale observable (here the u concentration field) through the zeroth moment of the microscopic LB simulation (averaging the particle distribution functions (f_i) on a given lattice point, see section IIIB for more details).

Given the coarse-scale observable we estimate its time-derivative and several of its spatial derivatives (e.g. u_t , u_x ,

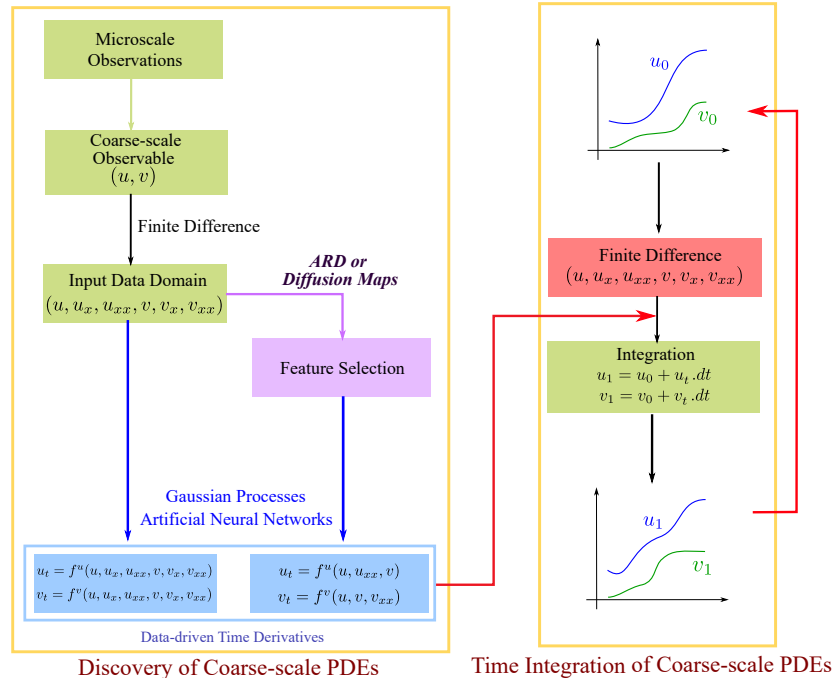


FIG. 2. Workflow for uncovering coarse-scale PDEs. First, we compute macroscopic variables u and v from the Lattice Boltzmann simulation data (see equation (18) and figure 1) and estimate their spatial derivatives (e.g. by finite difference schemes on the lattice). After that, we employ machine learning algorithms (here, Gaussian process regression or artificial neural networks) to identify “proper” time derivatives u_t and v_t from an original input data domain directly (no feature selection among several spatial derivatives) or from a reduced input data domain (feature selection among several spatial derivatives) using ARD in Gaussian processes or Diffusion Maps. We then simulate the identified coarse-scale PDE for given coarse initial conditions $(\mathbf{u}_0, \mathbf{v}_0)$.

and u_{xx}), typically using finite difference schemes in time and space as necessary. A PDE of the form $u_t = L(u) = F(u, u_x, u_{xx}, \dots)$ is a relation between the time-derivative and a number of spatial derivatives; this relation holds at every moment in time and every point in space. For a simple reaction diffusion equation, say $u_t = u_{xx} - ku$, the data triplets (u, u_t, u_{xx}) will in general lie on a two-dimensional manifold in three-dimensional space, since u_t is a function of u and u_{xx} . Knowing that this manifold is two-dimensional suggests (in the spirit of the Whitney and Takens embedding theorems^{20,21}) that any five generic observables suffice to create an embedding - and thus learn u_t , a function on the manifold, as a function of these five observables. One might choose, for example, as observables the values of u at any five spatial points at a given time moment, possibly the five points used in a finite difference stencil for estimating spatial derivatives. In the study of time series through delay embeddings one uses observations on a temporal stencil; it is interesting that here one might use observations on a spatial stencil - encoding information analogous to spatial derivatives (see¹²). Motivated by this perspective, in order to learn the time derivative u_t , we use an original input data domain including several (say, all up to some order) spatial derivatives. We also consider the selection of a reduced input data domain via two feature selection methods: (1) a sensitivity analysis by automatic relevance determination (ARD) in Gaussian processes^{16,18,22} and (2) a manifold learning approach, Diffusion Maps^{23,24}, with a regression loss (see section IV B in more details). Then, we

consider two different machine learning methods (Gaussian process regression and artificial neural networks) to learn u_t based on the selected feature input data domain.

After training, simulation of the learned coarse-scale PDE given a coarse initial condition $u_0(x, t), v_0(x, t)$ can proceed with any acceptable discretization scheme in time and space (from simple finite differences to, say, high order spectral or finite element methods).

B. Gaussian process regression

One of the two approaches we employ to extract dominant features and uncover the RHS of coarse-scale PDEs is Gaussian process regression. In Gaussian processes, to represent a probability distribution over target functions (here, the time derivative), we assume that our observations are a set of random variables whose finite collections have a multivariate Gaussian distribution with an *unknown* mean (usually set to zero) and an *unknown* covariance matrix K . This covariance matrix is commonly formulated by a Euclidean distance-based kernel function κ in the input space, whose hyperparameters are optimized by training data. Here, we employ a radial basis kernel function (RBF), which is the *de facto* default kernel function in Gaussian process regression, with

ARD¹⁷.

$$K_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j; \theta) = \theta_0 \exp \left(-\frac{1}{2} \sum_{l=1}^k \frac{(x_{i,l} - x_{j,l})^2}{\theta_l} \right). \quad (1)$$

where $\theta = [\theta_0, \dots, \theta_k]^T$ is a $k+1$ dimensional vector of hyperparameters and k is the number of dimensions of the input data domain. The optimal hyperparameter set θ^* can be obtained by minimizing a negative log marginal likelihood with the training data set $\{\mathbf{x}, \mathbf{y}\}$:

$$\begin{aligned} \theta^* &= \underset{\theta}{\operatorname{argmin}} \{ -\log p(\mathbf{y}|\mathbf{x}, \theta) \} \\ &= \frac{1}{2} \mathbf{y}^T (K + \sigma^2 I)^{-1} \mathbf{y} + \frac{1}{2} \log |(K + \sigma^2 I)| + \frac{N}{2} \log 2\pi \end{aligned} \quad (2)$$

where N is the number of training data points, σ^2 and I represent the variance of the (Gaussian) observation noise and an $N \times N$ identity matrix, respectively.

To find the Gaussian distribution of the function values (here the time derivative) at test data points, we represent the multivariate Gaussian distribution with the covariance matrix constructed by equation (1) as

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{y}^* \end{bmatrix} = N \left(\mathbf{0}, \begin{bmatrix} K + \sigma^2 I & K_* \\ K_*^T & K_{**} \end{bmatrix} \right), \quad (3)$$

where \mathbf{y}^* is a predictive distribution for test data \mathbf{x}^* , K_* represents a covariance matrix between training and test data while K_{**} represents a covariance matrix between test data.

Finally, we represent a Gaussian distribution for time derivatives at the test point in terms of a predictive mean and its variance, through conditioning a multivariate Gaussian distribution as

$$\bar{\mathbf{y}}^* = K_*(K + \sigma^2 I)^{-1} \mathbf{y}, \quad (4)$$

$$K(\mathbf{y}^*) = K_{**} - K_*^T (K + \sigma^2 I)^{-1} K_*, \quad (5)$$

and we assign the predictive mean ($\bar{\mathbf{y}}^*$) as the estimated time derivative for the corresponding data point.

C. Artificial neural networks

Next, we consider (artificial, possibly deep) neural networks (ANN or NN or DNN) for identifying the RHS of coarse-scale PDEs. Generally, neural networks consist of an input layer, one or more hidden layers, and an output layer, all comprised of several computational neurons, typically fully connected by weights (ω), biases (b), and an activation function ($\psi(\cdot)$). Macroscopic observables and their spatial derivatives are assigned at the input layer, while the corresponding time derivative is obtained at the output layer (here we are considering only first order PDEs in time; higher order equations, like the wave equation, involving second derivatives in time can also be accounted for within the framework). In (feed-forward) neural networks, a universal approximation

theorem²⁵ guarantees that for a single hidden layer with (sufficient) finite number of neurons, an approximate realization \tilde{y} of the target function, y can be found. Here, approximation implies that the target and learned functions are sufficiently close in an appropriately chosen norm ($\forall \epsilon > 0 : |y - \tilde{y}| < \epsilon$). The approximate form of the target function obtained through the feedforward neural net can be written as

$$\tilde{y}(\mathbf{x}) = \sum_{i=1}^N \psi(\omega_i^T \mathbf{x} + b_i). \quad (6)$$

The root-mean-square error cost function

$$E_D = \frac{1}{N} \sum_{j=1}^N (y_j - \tilde{y}(x_j))^2, \quad (7)$$

typically measures the goodness of the approximation.

In order to obtain a *generalizable* network, with good performance on the test data set as well as on the training data set (e.g. preventing overfitting), several regularization approaches have been proposed, mostly relying on modifications of the cost function. Foresee and Hagan²⁶ showed that modifying the cost function by adding the regularization term $E_\omega = \sum_{j=1}^N \omega_j^2$, results in a network that will maximize the posterior probability based on Bayes' rule. We thus trained our network based on a total cost function of the form:

$$E_{total} = \beta_1 E_D + \beta_2 E_\omega, \quad (8)$$

in which β_1 and β_2 are network tuning parameters. Here, we employ Bayesian regularized back-propagation for training, which updates weight and bias values through Levenberg-Marquardt optimization²⁷; we expect that, for our data, comparable results would be obtained through other modern regularization/optimization algorithms.

D. Diffusion Maps

Diffusion Maps (DMAP) have successfully been employed for dimensionality reduction and nonlinear manifold learning^{23,24,28,29}. The Diffusion Maps algorithm can guarantee, for data lying on a smooth manifold -and at the limit of infinite data- that the eigenvectors of the large normalized kernel matrices constructed from the data converge to the eigenfunctions of the Laplace-Beltrami operator on the manifold on which the data lie. These eigenfunctions can also provide nonlinear parametrizations (i.e. sets of coordinates) for such Riemannian manifolds. To approximate the Laplace-Beltrami operator from scattered data points on the manifold, a normalized diffusion kernel matrix between observation (data) points is commonly used:

$$\mathbf{W}_{ij} = \exp \left(-\frac{\|\mathbf{y}_i - \mathbf{y}_j\|^2}{\epsilon} \right), \quad (9)$$

where \mathbf{y}_i are real-valued observations and ϵ is the kernel width. After that, one obtains a normalized matrix $\mathbf{W}^{(\alpha)}$ by

$$\mathbf{W}^{(\alpha)} = \mathbf{D}^{-\alpha} \mathbf{W} \mathbf{D}^{-\alpha}, \quad (10)$$

where \mathbf{D} is a diagonal matrix whose i^{th} entry is the sum of corresponding row of W . Here, $\alpha \in \{0, 1\}$ is a tuning parameter: $\alpha = 0$ corresponds to the classical normalized graph Laplacian^{30,31} while $\alpha = 1$, which takes into account the local data density, yields the Laplace-Beltrami operator²⁴; in this paper, we set $\alpha = 1$. Then, $\tilde{\mathbf{W}}$ is calculated simply as:

$$\tilde{\mathbf{W}} = \tilde{\mathbf{D}}^{-1} \mathbf{W}^{(\alpha)}, \quad (11)$$

where $\tilde{\mathbf{D}}$ is a diagonal matrix whose i^{th} entry is the sum of corresponding row of $\mathbf{W}^{(\alpha)}$.

Finally, an embedding of the manifold is constructed by the first few (say m) nontrivial eigenvectors of $\tilde{\mathbf{W}}$,

$$\mathbf{y}_i \mapsto (\lambda_1^t \phi_{1,i}, \dots, \lambda_m^t \phi_{m,i}), \quad i = 1, \dots, n, \quad (12)$$

where t corresponds to the number of diffusion steps (here $t = 0$) with descending ordered eigenvalues λ_i .

E. Feature selection

Describing the coarse-scale spatiotemporal dynamics in the form of a PDE, involves learning the local field time-derivatives as a function of a few, relevant local field spatial derivatives. Starting with a “full” input data domain consisting of all local field values as well as all their coarse-scale spatial derivatives (up to some order), we must extract the few “relevant” spatial derivatives as dominant features of this input data domain. Such feature selection will typically reduce the dimensionality of the input data domain. Among various feature selection methods, we employ two algorithms based on (1) sensitivity analysis via ARD in Gaussian processes^{16,18,22} and (2) manifold parametrization through output-informed Diffusion Maps¹⁹.

First, we employ sensitivity analysis via automatic relevance determination (ARD) in Gaussian processes, which effectively reduces the number of input data dimensions. In Gaussian processes, we obtain the optimal hyperparameter set θ^* by minimizing a negative log marginal likelihood (see equation (2)). ARD assigns a different hyperparameter θ_i for each input dimension d_i . As can be seen in equation (1), a large value of θ_i nullifies the difference between target function values along the d_i dimension, allowing us to designate this dimension as “insignificant”. Practically, we select the input dimensions with relatively small θ_j to build a reduced input data domain, which can still successfully represent the approximation of right-hand-side on the underlying PDEs.

Alternatively, we employ a manifold learning technique to find the intrinsic representation of the coarse-scale PDE, and then examine the relation between these intrinsic coordinates and given input features (spatial field derivatives). Specifically, Diffusion Maps will provide an intrinsic parametrization off the combined input-output data domain (here, $\{u_t, u, u_x, u_{xx}, v, v_x, v_{xx}\}$ for u and $\{v_t, u, u_x, u_{xx}, v, v_x, v_{xx}\}$ for v). Selecting leading intrinsic coordinates, we can then find the lowest-dimensional embedding space for the PDE manifold (the manifold embodying u_t and v_t as a function of the embedding intrinsic coordinates.) We then test several

combinations of subsets of the input domain coordinates (spatial derivatives) as to their ability to parametrize the discovered intrinsic embedding coordinates. Each such set of such inputs, that successfully parametrize the intrinsic embedding coordinates, provides us a new possibility of learning a PDE formulation that describes the spatiotemporal dynamics of our observation data set.

In principle, any subset of intrinsic coordinates that successfully parametrizes the manifold can be used to learn functions on the manifold, and in particular u_t and v_t . The success of any particular subset of leading intrinsic coordinates in so describing u_t and v_t is confirmed through regression, via a mean-squared-error loss (L).

Next, we investigate which set of features of the input domain (which sets of spatial derivatives) can be best used to parametrize the intrinsic embedding (and thus learn the PDE right-hand-side). One can find the subset of features from a user-defined dictionary (here spatial derivatives) to parametrize the intrinsic embedding coordinates through a linear Group LASSO³². In this paper, we examine several combinations of input domain variables, and find subsets that can minimally parametrize the intrinsic embedding; this is quantified through a total regression loss (L_T) based on a mean-squared-error as

$$L_T = \left(\sum_{j=1}^d L_{\phi_j}^2 \right)^{\frac{1}{2}}, \quad (13)$$

where L_{ϕ_j} represents a regression loss for representing the intrinsic coordinate ϕ_j using selected input features and d represents the number of intrinsic coordinates we chose.

ARD for Gaussian processes suggests the “best” input domain subset in terms of which we will try and predict u_t and v_t . In the manifold learning context, we may find several different input subsets capable of parametrizing the manifold on which the observed behavior lies. Different minimal parametrizing subsets will lead to different (but, in principle, on the data, equivalent) right-hand-sides for the PDE evolution. One expects that some of them will be “better conditioned” (have better Lipschitz constants) than others.

III. DIFFERENT SCALE SIMULATORS FOR ONE-DIMENSIONAL REACTION-DIFFUSION SYSTEMS

A. Macro-scale simulator: FitzHugh-Nagumo model

To describe a one-dimensional reaction-diffusion system that involves an activator u and an inhibitor v , the FitzHugh-Nagumo model consists of two coupled reaction-diffusion partial differential equations:

$$\begin{aligned} \frac{\partial u}{\partial t} &= D^u \frac{\partial^2 u}{\partial x^2} + u - u^3 - v, \\ \frac{\partial v}{\partial t} &= D^v \frac{\partial^2 v}{\partial x^2} + \varepsilon(u - a_1 v - a_0), \end{aligned} \quad (14)$$

where a_1 and a_0 are model parameters, ε represents a kinetic bifurcation parameter, and D^u and D^v represent diffusion co-

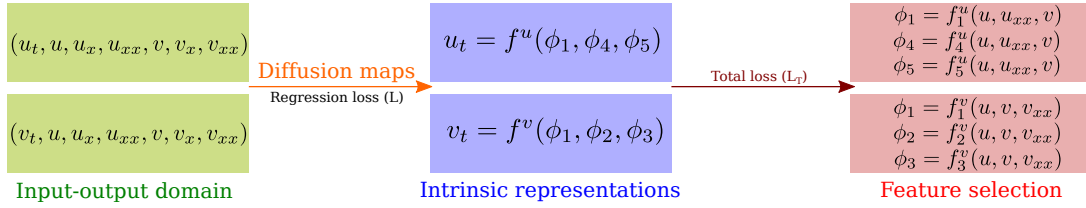


FIG. 3. Input feature selection via output-informed Diffusion Maps. Diffusion Maps provide intrinsic coordinatization of the output (the time derivatives) from combined input-output data. Guided by a regression loss (L), we find a low-dimensional intrinsic embedding space in which u_t (and v_t) can be represented as a function of just a few intrinsic diffusion map coordinates. After that, we search and find minimal subsets of the input data domain that can parametrize the selected intrinsic coordinates (e.g. ϕ_1, ϕ_4, ϕ_5) as quantified by a small total regression loss (see equation (13)).

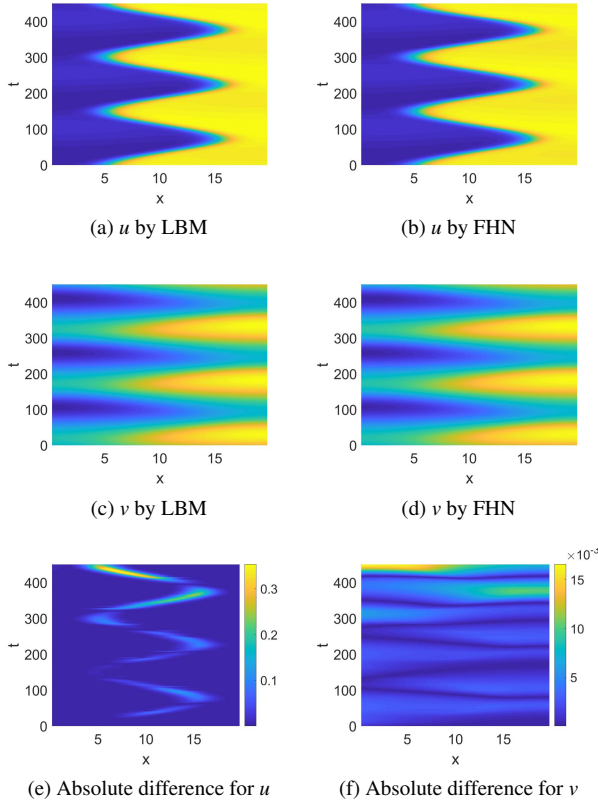


FIG. 4. Spatiotemporal behavior of u and v simulated by the Lattice-Boltzmann model and by the FitzHugh-Nagumo PDE. (a) and (c): u and v from the Lattice Boltzmann model (LBM). (b) and (d): u and v from the FitzHugh-Nagumo PDE. (e) and (f): Normalized absolute difference between the simulations of the two models.

efficients for u and v , respectively. Here, we set these parameters to $a_1 = 2$, $a_0 = -0.03$, $\varepsilon = 0.01$, $D^u = 1$, and $D^v = 4^{11}$. We discretize a spatial domain on $[0, 20]$ with $\Delta x = 0.2$ and a time domain on $[0, 450]$ with $\Delta t = 0.001$, respectively. We impose homogeneous Neumann boundary conditions at both boundaries and solve these equations (for various initial conditions) numerically via the finite element method using the COMSOL software.

B. Micro-scale simulator: the Lattice Boltzmann model

We also introduce a Lattice Boltzmann model (LBM)^{33,34}, which can be thought of as a mesoscopic numerical scheme for describing spatiotemporal dynamics using finite-difference-type discretizations of Boltzmann-BGK equations³⁵, retaining certain advantages of microscopic particle models. In this paper, the Lattice Boltzmann model is our fine scale “microscopic simulator” and its results are considered to be “the truth” from which the coarse-scale PDE will be learned.

The time evolution of the particle distribution function on a given lattice can be described by

$$f_i^l(x_{j+i}, t_{k+1}) = f_i^l(x_j, t_k) + \Omega_i^l(x_j, t_k) + R_i^l(x_j, t_k) \quad l \in \{u, v\}, \quad (15)$$

where a superscript l indicates the activator u and the inhibitor v , and Ω_i^l represents a collision term defined by Bhatnagar-Gross-Krook (BGK)³⁵:

$$\Omega_i^l(x_j, t_k) = -\omega^l(f_i^l(x_j, t_k) - f_i^{l,eq}(x_j, t_k)), \quad (16)$$

where ω^l represents a BGK relaxation coefficient defined as³⁶

$$\omega^l = \frac{2}{1 + 3D^l \frac{\Delta t}{\Delta x^2}}. \quad (17)$$

To compute our coarse-scale observables u and v , we employ the D1Q3 model, which uses three distribution functions on the one-dimensional lattice as (f_{-1}^l, f_0^l, f_1^l) for each density (totalling 6 distribution functions). Through the zeroth moment (in the velocity directions) of the overall distribution function, finally we compute the coarse-scale observable u and v as

$$u(x_j, t_k) = \sum_{i=-1}^1 f_i^u(x_j, t_k), \quad (18)$$

$$v(x_j, t_k) = \sum_{i=-1}^1 f_i^v(x_j, t_k).$$

Based on spatially uniform Local diffusion equilibrium, for which f_i^{eq} is homogeneous in all velocity directions, the

weights are chosen all equal to 1/3:

$$\begin{aligned} f_i^{u,eq}(x_j, t_k) &= \frac{1}{3}u(x_i, t_k), \\ f_i^{v,eq}(x_j, t_k) &= \frac{1}{3}v(x_i, t_k). \end{aligned} \quad (19)$$

Thus, the reaction terms R_i^l in equation (15) are modeled by

$$\begin{aligned} R_i^u(x_j, t_k) &= \frac{1}{3}\Delta t(u(x_j, t_k) - u(x_j, t_k)^3 - v(x_j, t_k)), \\ R_i^v(x_j, t_k) &= \frac{1}{3}\Delta t\varepsilon(u(x_j, t_k) - a_1v(x_j, t_k)^3 - a_0). \end{aligned} \quad (20)$$

All model parameters ($a_0, a_1, \varepsilon, D^u, D^v$) are the same as the FHN PDE. The corresponding spatiotemporal behavior of these coarse observables u and v is shown in figures 4(a) and (c) while the FHN PDE simulation for the same coarse initial conditions is shown in figures 4(b) and (d).

IV. RESULTS

A. Learning without feature selection

We begin by considering our proposed framework without feature selection, so as to later contrast with the results including feature selection. The data come from the fine-scale Lat-

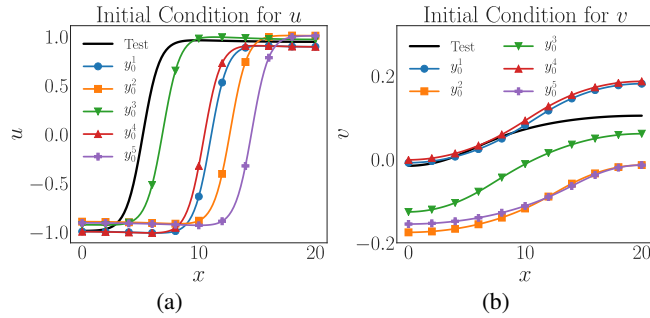


FIG. 5. Five different coarse initial conditions for training and a test coarse initial condition (colored in black). (a) Coarse initial conditions for u . (b) Coarse initial conditions for v . Five initial conditions are randomly chosen near the stable periodic solution.

tice Boltzmann simulation. For the parameter values selected, the long-term dynamics of the LB simulation lie, for all practical purposes, on a stable time-periodic solution. To predict the coarse time derivatives u_t and v_t , we collect training data from five different initial conditions near this stable periodic solution (see figure 5) with the following LB spatiotemporal discretization – in space, 99 discretized points on $[0.2, 19.8]$ with $dx = 0.2$; and in time, 451 discretized points on $[0, 450]$ with $dt = 1$ for each initial condition. Since our data come from the fine scale LB code, we need to initialize at the fine, LB scale of particle distribution functions (and not just of the concentrations u and v). To initialize the particle distribution functions in the Lattice Boltzmann model we apply the equal weights rule, 1/3 for f_{-1} , f_0 , and f_1 , motivated by near-equilibrium

considerations. We do expect that such initialization features will soon be “forgotten” as higher distribution moments become quickly slaved to the lower (here the zeroth) moments (see for example³⁷). To ensure that our results are not affected by the initialization details, we only start collecting training data after relaxation by short time simulation (here, 2000 time steps with $\Delta t = 0.001$ or $t = 2$), see appendix A. We estimate the local coarse fields and their (several) spatial and temporal derivatives through finite differences, and then apply machine learning algorithms (here Gaussian processes as well as neural networks) to learn the time derivatives of the activator u_t and the inhibitor v_t using as input variables the local u , v and all their spatial derivatives up to and including order two ($u, u_x, u_{xx}, v, v_x, v_{xx}$).

$$\begin{aligned} u_t(x, t) &= f^u(u, u_x, u_{xx}, v, v_x, v_{xx}), \\ v_t(x, t) &= f^v(u, u_x, u_{xx}, v, v_x, v_{xx}). \end{aligned} \quad (21)$$

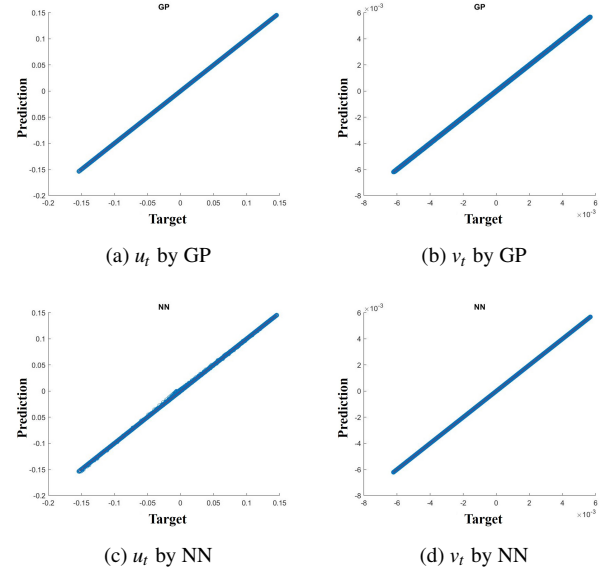


FIG. 6. No feature selection: $u_t = f^u(u, u_x, u_{xx}, v, v_x, v_{xx})$ and $v_t = f^v(u, u_x, u_{xx}, v, v_x, v_{xx})$. Regression results of the two methods for time derivatives: Gaussian processes (GP) and neural networks (NN).

Specifically, for the neural networks approach, we build two different networks, one for the prediction of the activator and one for the inhibitor. For both the activator and the inhibitor, we set use hidden layers consisting of 6 and 6 neurons using a hyperbolic tangent sigmoid activation function; as mentioned above, we use Levenberg-Marquardt optimization with a Bayesian regularization (see section IIC). Both networks use the mean-squared-error as their loss function. For Gaussian processes, we employ a radial basis kernel function with ARD (see equation (1)). Regression results obtained by each the two methods for the time derivatives in the training data set are shown in figure 6. Both methods provide good approximations of the target time derivatives u_t and v_t . Given the test coarse initial condition (black curves in figure 5), simulation results *with the learned PDE* from $t = 0$ to $t = 450$ with

$\Delta t = 0.001$ and their normalized absolute differences from the “ground truth” LB simulations are shown in figures 7. The order of magnitude of these absolute differences for both models is the same as those between the LB FHN and the explicitly known FHN PDE (see figures 4(e) and (f)).

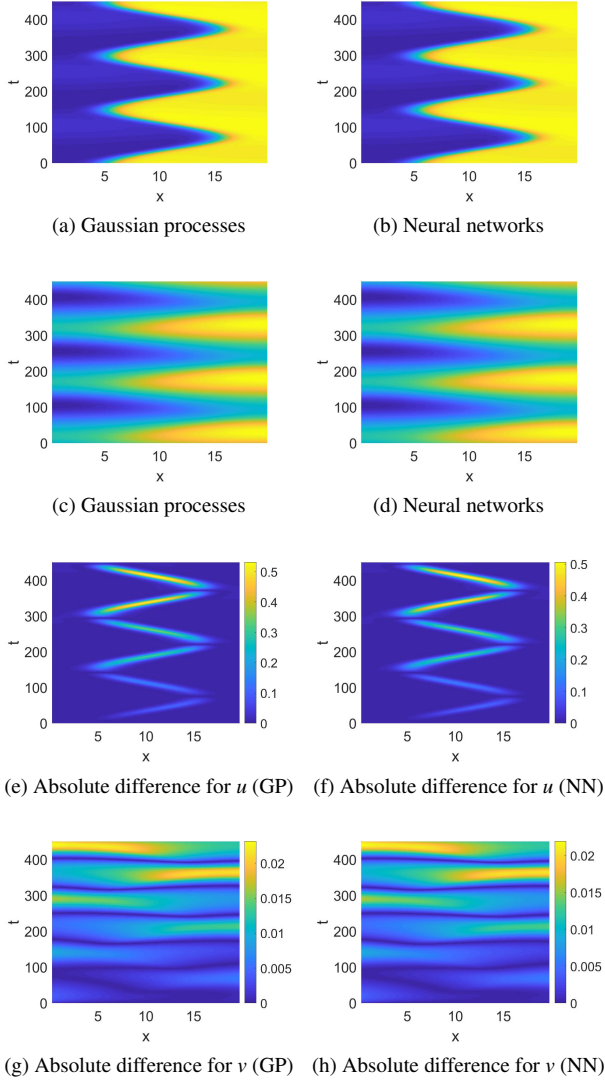


FIG. 7. No feature selection: $u_t = f^u(u, u_x, u_{xx}, v, v_x, v_{xx})$ and $v_t = f^v(u, u_x, u_{xx}, v, v_x, v_{xx})$. (a)-(d): Simulation results of the two methods for u and v . (e)-(h): The normalized absolute differences from the “ground truth” LB simulations for u and v .

B. Learning with feature selection

Now, we consider the possibility of feature selection, in an attempt to learn the RHS of coarse-scale PDEs with a minimal number of input domain variables (spatial derivatives). First, we apply the sensitivity analysis via ARD in the case of Gaussian process approximation. The optimal ARD weights (θ^*) for u_t and v_t are tabulated in table I. u_t has three rela-

tively small weights for (u, u_{xx}, v) and v_t has also three relatively small weights for (u, v, v_{xx}) . It is interesting to observe that the selected features via ARD are the same as those in the explicitly known FHN PDE (see equation (14)). This shows that ARD can effectively guide in selecting the appropriate dimensionality of the input data domain, resulting here in the same spatial derivative choices as in the explicitly known FHN PDE. Now, we use the reduced input data domain (u, u_{xx}, v) for

TABLE I. Optimal ARD weights (θ^* for u_t and v_t in equation (2)). As mentioned in section II E, features which have relatively small ARD weights can be regarded as dominant features for the target functions u_t and v_t .

	u	u_x	u_{xx}	v	v_x	v_{xx}
u_t	5.28E+00	4.23E+06	9.13E+02	2.13E+03	5.32E+08	4.78E+07
v_t	1.33E+02	6.69E+06	1.94E+06	5.09E+02	4.20E+06	1.75E+02

u_t and (u, v, v_{xx}) for v_t to recover the RHS of the coarse-scale PDEs as

$$\begin{aligned} u_t(x, t) &= f_1^u(u, u_{xx}, v), \\ v_t(x, t) &= f_1^v(u, v, v_{xx}). \end{aligned} \quad (22)$$

Regression results of our two methods for the time derivatives are shown in figure 8. Results of long time simulation of the

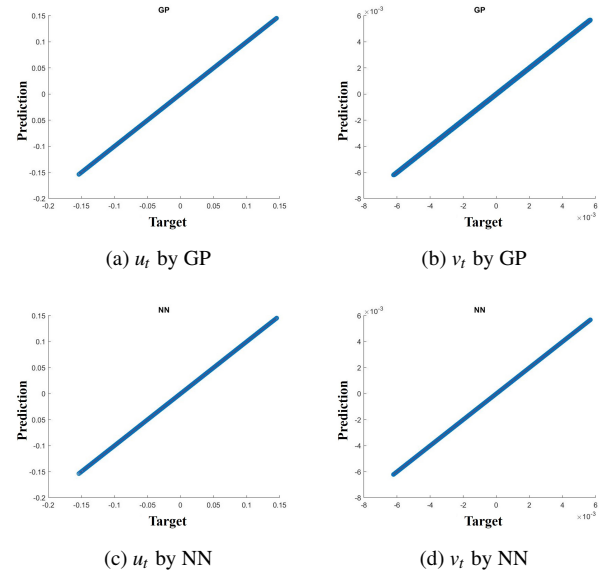


FIG. 8. Feature selection 1: $u_t = f_1^u(u, u_{xx}, v)$ and $v_t = f_1^v(u, v, v_{xx})$. These selected variables are the same as those that appear in the right-hand-side of the explicitly known FHN PDE. Regression results of the two methods for time derivatives: Gaussian processes (GP) and neural networks (NN).

learned PDEs by each method, from $t = 0$ to $t = 450$, as well as normalized absolute differences from the simulation of the “ground truth” LB are shown in figure 9.

The two machine learning methods operating with a reduced input data domain can still provide good approximations of the time derivatives and of the resulting dynamics.

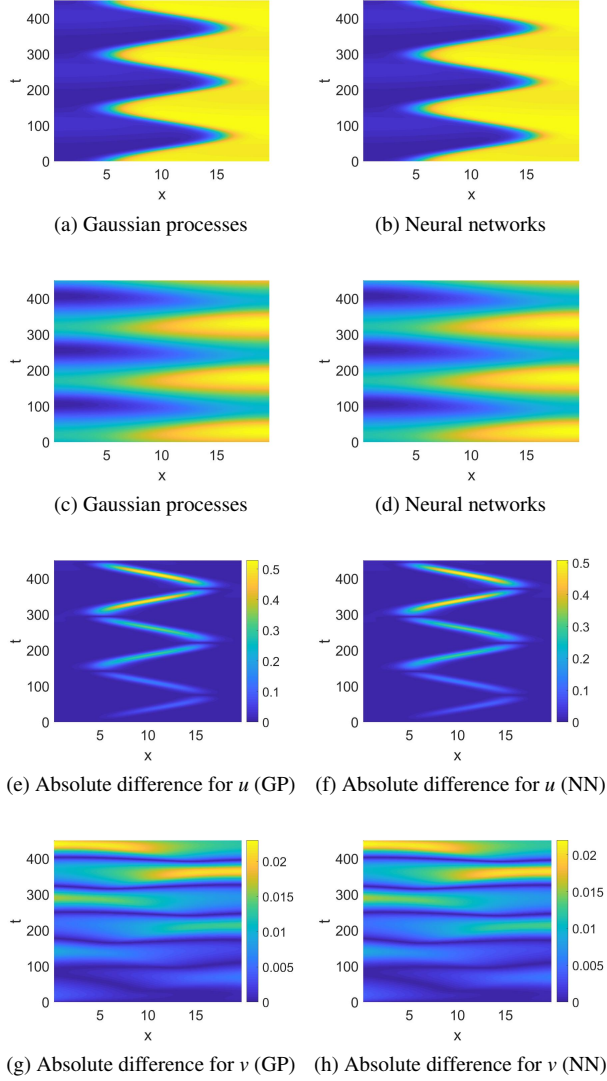


FIG. 9. Feature selection 1: $u_t = f^u(u, u_{xx}, v)$ and $v_t = f^v(u, v, v_{xx})$. (a)-(d): Simulation results of the two methods for u and v . (e)-(h): The normalized absolute differences from the “ground truth” LB simulations for u and v .

The order of magnitude of these absolute differences is effectively the same as the difference of the FHN LB from the explicitly known FHN PDE. It is, therefore, clear that our framework effectively recovers the coarse-scale PDE from fine scale observation data; the difference is that the right hand-side of the PDE is now given in terms of the ANN right-hand-side, or in terms of the observed data and the GP kernel/hyperparameters, rather than the simple algebraic formula of equation (14).

Next, we consider an alternative approach for feature selection, via our manifold learning technique, Diffusion Maps. The best candidate set among different combinations of intrinsic coordinates (varying the number of leading intrinsic dimensions and recording the corresponding Gaussian Process regression loss) are shown in table II. Since this three-

TABLE II. The best candidates and the corresponding regression loss (L) for u_t and v_t with respect to the number of Diffusion map coordinates

Optimal intrinsic coordinates		Regression Loss (L)	
u_t	v_t	u_t	v_t
1d (ϕ_5^u)	(ϕ_2^v)	4.60E-04	7.69E-06
2d (ϕ_1^u, ϕ_5^u)	(ϕ_1^v, ϕ_2^v)	1.40E-06	1.50E-06
3d $(\phi_1^u, \phi_4^u, \phi_5^u)$	$(\phi_1^v, \phi_2^v, \phi_3^v)$	2.18E-08	4.74E-08
4d $(\phi_1^u, \phi_3^u, \phi_4^u, \phi_5^u)$	$(\phi_1^v, \phi_2^v, \phi_3^v, \phi_4^v)$	1.64E-08	5.71E-09

TABLE III. The best candidates and corresponding total loss for $u_t = (\phi_1^u, \phi_4^u, \phi_5^u)$ and $v_t = (\phi_1^v, \phi_2^v, \phi_3^v)$ with respect to the number of features.

$u_t = (\phi_1^u, \phi_4^u, \phi_5^u)$		$v_t = (\phi_1^v, \phi_2^v, \phi_3^v)$	
Features	Total Loss (L_T)	Features	Total Loss (L_T)
1d (u)	6.51E-05	(u)	7.93E-05
2d (u, v)	1.65E-08	(u, v)	1.49E-05
3d (u, u_{xx}, v)	6.52E-09	(u, v, v_{xx})	3.32E-07
(u, u_x, v)	7.39E-09	(u, u_x, v_{xx})	6.21E-07
4d (u, u_x, u_{xx}, v)	2.68E-09	(u, v, v_x, v_{xx})	4.47E-09

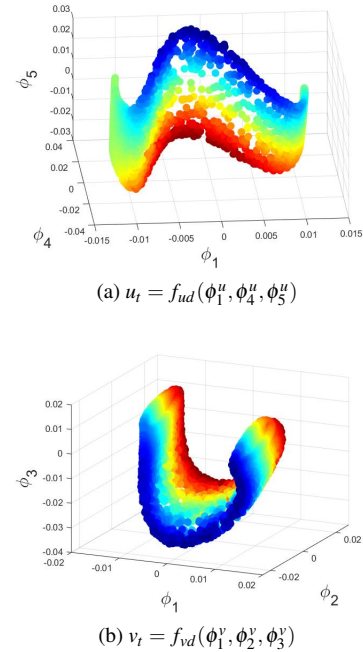


FIG. 10. Three leading Diffusion map coordinates: Colors represent u_t in (a) and v_t in (b).

dimensional intrinsic embedding space exhibits a (tiny) regression loss of order 10^{-8} , we choose an input domain for u_t consisting of $(\phi_1^u, \phi_4^u, \phi_5^u)$ as shown in figure 10(a). For v_t , by the same token, we choose the three-dimensional embedding space consisting of $(\phi_1^v, \phi_2^v, \phi_3^v)$ as shown in figure 10(b). Based on these identified intrinsic embedding spaces, we examined several subsets of input domain features (spatial derivatives) using the total loss of equation (13). “Good” subsets of input features (those that result in small regres-

sion losses with minimal input dimension) are presented in table III. Clearly, different choices of such input feature subsets can give comparable total losses; this suggests that we may construct different right-hand-sides of the unknown coarse-scale PDE that are comparably successful in representing the observed dynamics.

The good candidates for u_t and v_t identified this way, consisting of three input features, are (u, u_{xx}, v) and (u, v, v_{xx}) ; they are the same as those found from GP via ARD, and also the same as the ones in the explicitly known FHN PDE. Interestingly, another possible alternative candidate set is also identified: (u, u_x, v) for u_t and (u, u_x, v_{xx}) for v_t .

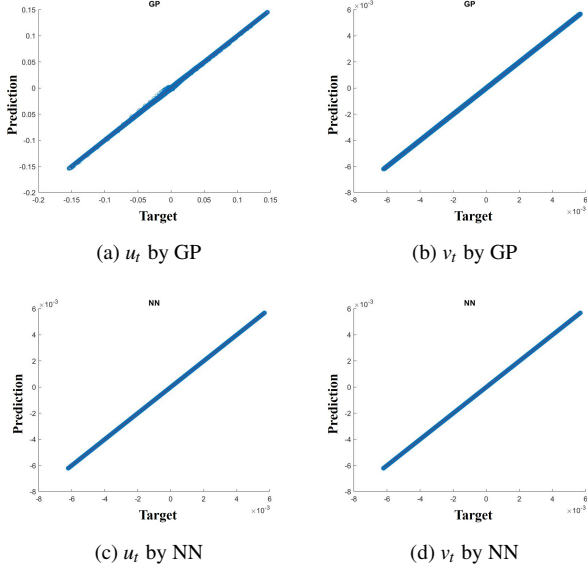


FIG. 11. Feature selection 2: $u_t = f_2^u(u, u_x, v)$ and $v_t = f_2^v(u, v, v_{xx})$. Regression results of the two methods for time derivatives: Gaussian processes (GP) and neural networks (NN).

Using these alternative candidate feature sets, we model different “versions” of what, on the data, is effectively the same coarse-scale PDE. The “alternative” version of the PDE can be symbolically written as

$$\begin{aligned} u_t(x, t) &= f_2^u(u, u_x, v), \\ v_t(x, t) &= f_2^v(u, v, v_{xx}), \end{aligned} \quad (23)$$

and the corresponding regression results of the time derivatives are shown in figure 11. Specifically, we use the first spatial derivative u_x instead of the second spatial derivative u_{xx} for learning u_t .

As shown in figure 12, both models provide good predictions of the “ground truth” LB simulations; we observe, however, that the accuracy of the neural network based predictions is enhanced. These results confirm that, on the data, alternative coarse-scale PDE forms can provide successful macroscopic description.

To further explore this possibility of alternative PDE forms that represent the observed data with *qualitatively comparable accuracy*, we also explored the efficacy of a third coarse-scale

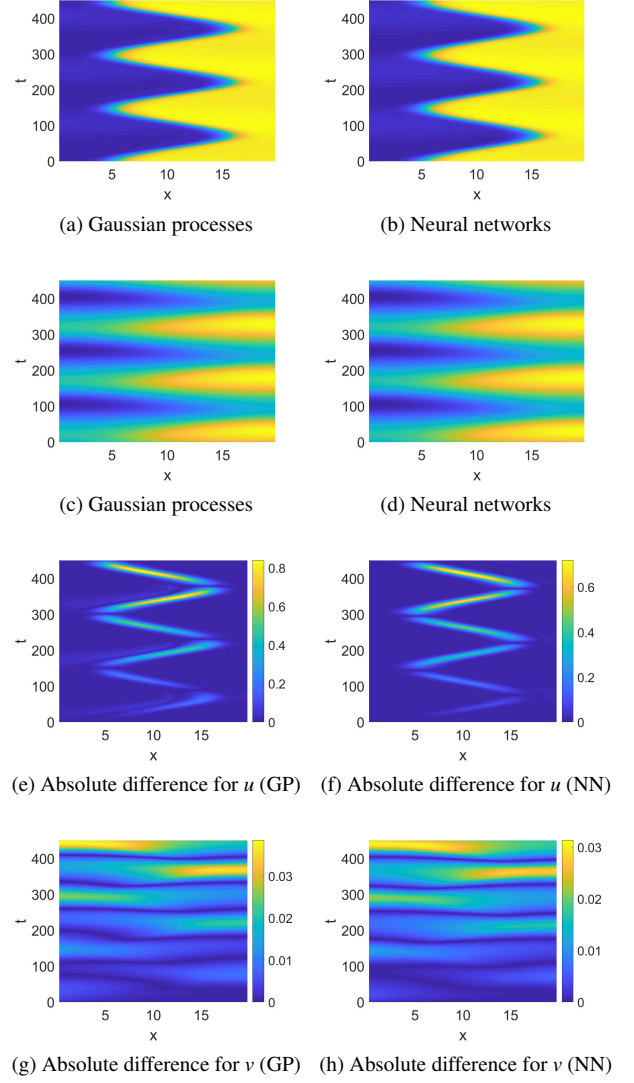


FIG. 12. Feature selection 2: $u_t = f_2^u(u, u_x, v)$ and $v_t = f_2^v(u, v, v_{xx})$. (a)-(d): Simulation results of the two methods for u and v . (e)-(h): The normalized absolute differences from the “ground truth” LB simulations for u and v .

PDE description, in terms of a yet different input feature set: (u, u_{xx}, v) for u_t and (u, u_x, v_{xx}) for v_t , so that the PDE can symbolically be written as

$$\begin{aligned} u_t(x, t) &= f_3^u(u, u_{xx}, v), \\ v_t(x, t) &= f_3^v(u, u_x, v_{xx}). \end{aligned} \quad (24)$$

The corresponding prediction results of the time derivatives are shown in figure 13.

As shown in figure 13, both regression methods provide an inaccurate approximation of v_t near $v_t = 0$; the order of magnitude of this error is 10^{-3} . The long term prediction results are not as accurate representations of the ground truth LB simulation as the previous two coarse-scale PDE realizations; yet they may still be qualitatively informative. Normalized absolute differences of long-time simulation for both machine

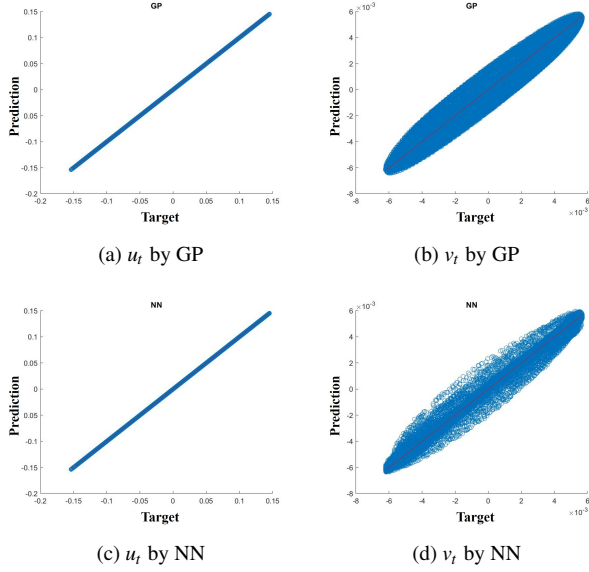


FIG. 13. Feature selection 3: $u_t = f_3^u(u, u_{xx}, v)$ and $v_t = f_3^v(u, u_x, v_{xx})$. Regression results of the two methods for time derivatives: Gaussian processes (GP) and neural networks (NN).

learning methods are shown in figure 14. As was the case in the previous alternative PDE realizations, the NN model appears more accurate than the GP one.

To compare our identified coarse-scale PDEs with the explicitly known FHN PDE (see equations (14)), we also compare the predictions of our coarse-scale PDEs to those of the FHN PDE via mean normalized absolute differences for the test coarse initial condition followed from $t = 0$ to $t = 450$ as

$$\begin{aligned} \text{MNAD}_u &= \frac{1}{N_T} \sum_{i=1}^{99} \sum_{j=0}^{450} \frac{|u(i, j) - u_f(i, j)|}{\max(u_f) - \min(u_f)}, \\ \text{MNAD}_v &= \frac{1}{N_T} \sum_{i=1}^{99} \sum_{j=0}^{450} \frac{|v(i, j) - v_f(i, j)|}{\max(v_f) - \min(v_f)}, \end{aligned} \quad (25)$$

where N_T is a total number of data points and u_f and v_f represent simulation results of the FHN PDE, respectively. The comparison of these representative simulation of our various coarse-scale PDEs is summarized in table IV. The differences across our various coarse-scale identified PDEs are of order 10^{-2} and below, comparable to the difference between each of them and the FHN PDE. Specifically, ‘feature selection 1’ (figure 9), whose variables are the same as those of the explicit FHN PDE, provides the best PDE realization via *both* the GP and the NN models.

V. CONCLUSION

In this paper, we demonstrated the data-driven discovery of macroscopic, concentration-level PDEs for reaction/transport processes resulting from fine-scale observations (here, from simulations of a Lattice Boltzmann mesoscopic model).

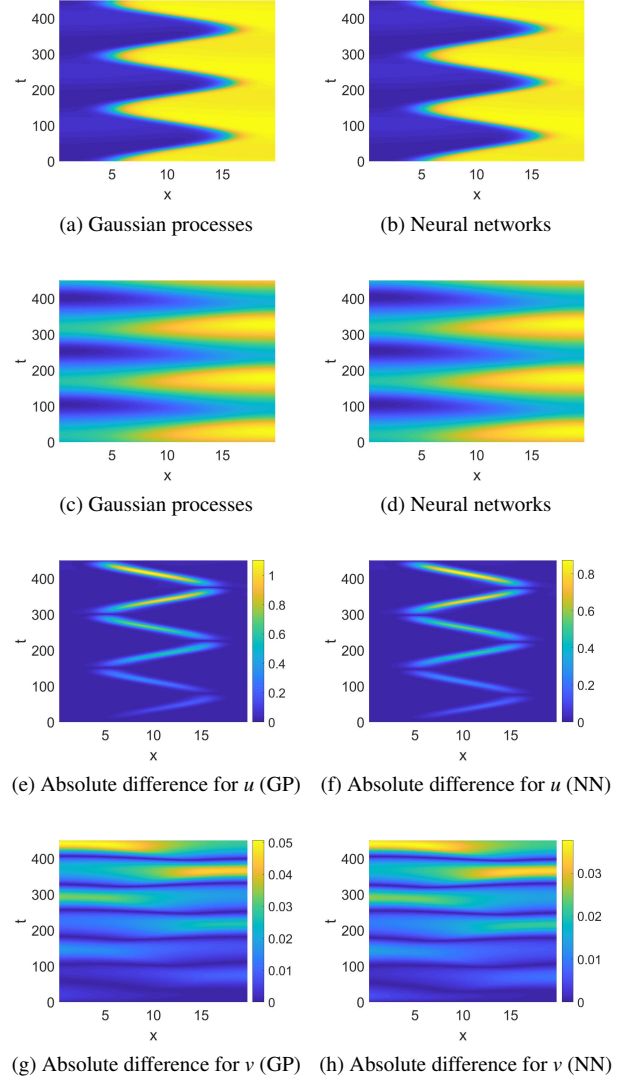


FIG. 14. Feature selection 3: $u_t = f^u(u, u_{xx}, v)$ and $v_t = f^v(u, u_x, v_{xx})$. (a)-(d): Simulation results of the two methods for u and v . (e)-(h): The normalized absolute differences from the ‘ground truth’ LB simulations for u and v .

TABLE IV. Mean normalized absolute difference (MNAD) for different coarse-scale PDEs. ‘GP’ and ‘NN’ represent ‘Gaussian processes’ and ‘Neural networks’, respectively.

	MNAD _u	MNAD _v
No Feature selection with GP	1.59E-02	1.62E-02
No Feature selection with NN	1.53E-02	1.56E-02
Feature selection 1 with GP	1.58E-02	1.62E-02
Feature selection 1 with NN	1.54E-02	1.57E-02
Feature selection 2 with GP	2.39E-02	2.20E-02
Feature selection 2 with NN	2.00E-02	2.11E-02
Feature selection 3 with GP	3.20E-02	3.31E-02
Feature selection 3 with NN	2.08E-02	2.16E-02

Long-term macroscopic prediction is then obtained by simulation of the identified (via machine-learning methods) coarse-scale PDE. We explored the effect of input feature selection capability on the effectiveness of our framework to identify the underlying macroscopic PDEs.

Our framework suggests four different PDEs (one without and three with feature selection), all comparable with the explicit FitzHugh-Nagumo PDE *on the data*: all of them provide good approximations of sample coarse-scale dynamic trajectories. The FHN PDE terms have a well-established mechanistic physical meaning (reaction and diffusion); it would be interesting to explore if any physical meaning can be ascribed to our alternative parametrizations of the right-hand-side of the coarse-scale evolution PDE.

In our framework, we employed finite differences to estimate spatial derivatives in the formulation of the PDE. Instead of numerical spatial derivatives, we may use the values of coarse observables at neighboring points directly to uncover the coarse evolution law. The effect of this alternative embedding for the PDE right-hand-side, explored in¹², on the accuracy of the identified model predictions, is the subject of ongoing research.

We believe that the framework we presented is easily generalizable to multiscale and/or multifidelity data. Here we worked across a single scale gap and a single fine-scale simulator providing the ground truth. We envision that data fusion tools can be combined with the approach to exploit data at several cascaded scales, and taking advantage of simulation data from several heterogeneous simulators^{9,38}.

ACKNOWLEDGMENTS

S. Lee, M. Kooshkbaghi, and I. G. Kevrekidis gratefully acknowledge partial support by NIH and by DARPA. Discussions of the authors with Dr. Felix Dietrich are gratefully acknowledged.

Appendix A: “Healing” for the Lattice Boltzmann model

An important assumption underlying our work is that the fine-scale model can “close” at a coarse-scale level. In our particular case, this means that, even though our fine scale Lattice Boltzmann model (LBM) evolves particle distribution functions, one can be predictive at the coarse level of the zeroth moments of these functions, the concentrations u and v of the activator and the inhibitor. The hypothesis that allows this reduction is that the problem is *singularly perturbed* in time: higher order moments of these distribution functions become quickly slaved to the “slow”, governing, zeroth order moment fields. Yet, while initializing the FHN PDE only requires spatial profiles of u and v at the initial time, initializing the full FHN LBM requires initial conditions for *all* the evolving particle distributions. Constructing such detailed fine scale initializations consistent with the coarse initial conditions is an important conceptual component of equation-free computation; the term used is “lifting”^{7,39}.

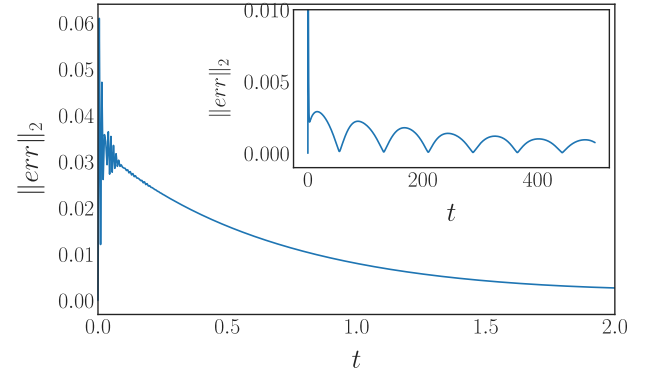


FIG. 15. Evolution of the L^2 norm (see equation (A2)) of the coarse difference between trajectories with the same coarse but different fine initial conditions. After the initial small (but violent) oscillation in error abates (for $t \lesssim 0.1$), the perturbed system relaxes to the vicinity of the base solution over $t \approx 2$.

Here, in lifting the coarse-scale observable (a concentration field ρ) to the microscopic description (particle distribution function f) on each node, we employ an equal weight rule, i.e., the three particle distribution function values at the node x_n are chosen to be

$$f_{-1}(x_n) = f_0(x_n) = f_1(x_n) = \frac{\rho(x_n)}{3}. \quad (\text{A1})$$

This equal weight choice (the local, spatially uniform diffusive equilibrium distribution) is not, in general, consistent with the (spatially nonuniform, and not simply diffusive) macroscopic PDE model (here, the FitzHugh-Nagumo PDEs); yet we expect that the fine scale simulation features will become rapidly slaved to the local concentration field³⁷. To estimate the appropriate slaving/relaxation time, we compare the L^2 norm of the density predicted by two differently initialized LBM simulations: one that lies on the long-term stable limit cycle, and one that results from it by retaining the coarse state, but perturbing the associated fine scale states according to the local diffusive equilibrium $\frac{1}{3}$ rule (for $\varepsilon = 0.01$ in equation (20)).

To explore the slaving time scale, we trace the L^2 norm of the difference between the simulations resulting from these two initializations. This L^2 norm is defined as

$$\|err\|_2 = \|\rho_{eq}(x, t) - \rho(x, t)\|_2, \quad (\text{A2})$$

where ρ_{eq} and ρ represent the density with equal weights and the density without equal weights (reference solution), respectively. As shown in figure 15, after a fast transient oscillation of L^2 (for $t < 0.1$), the norm decays smoothly until $t \approx 2$. There is still a small inherent bias (the trajectory will come back to a *nearby* point along the limit cycle); this does not affect our estimate of the slaving time. We therefore chose a relaxation time $t = 2$ (or 2000 LB time steps) and we started collecting coarse observations as training data from our various initializations only after $t = 2$.

¹W. Noid, “Perspective: Coarse-grained models for biomolecular systems,” The Journal of chemical physics **139**, 09B201.1 (2013).

- ²J. Hudson, M. Kube, R. Adomaitis, I. Kevrekidis, A. Lapedes, and R. Farber, "Nonlinear signal processing and system identification: applications to time series from electrochemical reactions," *Chemical Engineering Science* **45**, 2075–2081 (1990).
- ³K. Krischer, R. Rico-Martinez, I. Kevrekidis, H. Rotermund, G. Ertl, and J. Hudson, "Model identification of a spatiotemporally varying catalytic reaction," *AIChE Journal* **39**, 89–98 (1993).
- ⁴R. Rico-Martinez, J. Anderson, and I. Kevrekidis, "Continuous-time nonlinear signal processing: a neural network based approach for gray box identification," in *Proceedings of IEEE Workshop on Neural Networks for Signal Processing* (IEEE, 1994) pp. 596–605.
- ⁵J. Anderson, I. Kevrekidis, and R. Rico-Martinez, "A comparison of recurrent training algorithms for time series analysis and system identification," *Computers & chemical engineering* **20**, S751–S756 (1996).
- ⁶R. Gonzalez-Garcia, R. Rico-Martinez, and I. Kevrekidis, "Identification of distributed parameter systems: A neural net based approach," *Computers & chemical engineering* **22**, S965–S968 (1998).
- ⁷I. G. Kevrekidis, C. W. Gear, J. M. Hyman, P. G. Kevrekidis, O. Runborg, and C. Theodoropoulos, "Equation-free, coarse-grained multiscale computation: Enabling microscopic simulators to perform system-level analysis," *Commun. Math. Sci.* **1**, 715–762 (2003).
- ⁸S. Sirisup, G. E. Karniadakis, D. Xiu, and I. G. Kevrekidis, "Equation-free/galerkin-free pod-assisted computation of incompressible flows," *Journal of Computational Physics* **207**, 568–587 (2005).
- ⁹S. Lee, I. G. Kevrekidis, and G. E. Karniadakis, "A resilient and efficient cfd framework: Statistical learning tools for multi-fidelity and heterogeneous information fusion," *Journal of Computational Physics* **344**, 516–533 (2017).
- ¹⁰C. Siettos, C. Gear, and I. Kevrekidis, "An equation-free approach to agent-based computation: Bifurcation analysis and control of stationary states," *EPL (Europhysics Letters)* **99**, 48007 (2012).
- ¹¹C. Theodoropoulos, Y.-H. Qian, and I. G. Kevrekidis, "coarse stability and bifurcation analysis using time-steppers: A reaction-diffusion example," *Proceedings of the National Academy of Sciences* **97**, 9840–9843 (2000).
- ¹²Y. Bar-Sinai, S. Hoyer, J. Hickey, and M. P. Brenner, "Learning data-driven discretizations for partial differential equations," *Proceedings of the National Academy of Sciences* **116**, 15344–15349 (2019).
- ¹³S. H. Rudy, S. L. Brunton, J. L. Proctor, and J. N. Kutz, "Data-driven discovery of partial differential equations," *Science Advances* **3**, e1602614 (2017).
- ¹⁴M. Raissi, P. Perdikaris, and G. E. Karniadakis, "Machine learning of linear differential equations using gaussian processes," *Journal of Computational Physics* **348**, 683–693 (2017).
- ¹⁵M. Raissi and G. E. Karniadakis, "Hidden physics models: Machine learning of nonlinear partial differential equations," *Journal of Computational Physics* **357**, 125–141 (2018).
- ¹⁶Y. A. Qi, T. P. Minka, R. W. Picard, and Z. Ghahramani, "Predictive automatic relevance determination by expectation propagation," in *Proceedings of the twenty-first international conference on Machine learning* (ACM, 2004) p. 85.
- ¹⁷C. E. Rasmussen and C. K. I. Williams, *Gaussian processes for machine learning* (MIT Press, 2006).
- ¹⁸D. P. Wipf and S. S. Nagarajan, "A new view of automatic relevance determination," in *Advances in neural information processing systems* (2008) pp. 1625–1632.
- ¹⁹A. Holiday, M. Kooshkbaghi, J. M. Bello-Rivas, C. W. Gear, A. Zagaris, and I. G. Kevrekidis, "Manifold learning for parameter reduction," *Journal of computational physics* **392**, 419–431 (2019).
- ²⁰H. Whitney, "Differentiable manifolds," *Annals of Mathematics* **37**, 645–680 (1936).
- ²¹F. Takens, "Detecting strange attractors in turbulence," in *Dynamical Systems and Turbulence, Warwick 1980*, edited by D. Rand and L.-S. Young (Springer Berlin Heidelberg, Berlin, Heidelberg, 1981) pp. 366–381.
- ²²C. K. Williams and C. E. Rasmussen, "Gaussian processes for regression," in *Advances in neural information processing systems* (1996) pp. 514–520.
- ²³R. R. Coifman, S. Lafon, A. B. Lee, M. Maggioni, B. Nadler, F. Warner, and S. W. Zucker, "Geometric diffusions as a tool for harmonic analysis and structure definition of data: Diffusion maps," *Proceedings of the National Academy of Sciences* **102**, 7426–7431 (2005), <https://www.pnas.org/content/102/21/7426.full.pdf>.
- ²⁴R. R. Coifman and S. Lafon, "Diffusion maps," *Applied and computational harmonic analysis* **21**, 5–30 (2006).
- ²⁵G. Cybenko, "Approximation by superpositions of a sigmoidal function," *Mathematics of control, signals and systems* **2**, 303–314 (1989).
- ²⁶F. D. Foresee and M. T. Hagan, "Gauss-newton approximation to bayesian learning," in *Proceedings of International Conference on Neural Networks (ICNN'97)*, Vol. 3 (IEEE, 1997) pp. 1930–1935.
- ²⁷M. T. Hagan and M. B. Menhaj, "Training feedforward networks with the marquardt algorithm," *IEEE transactions on Neural Networks* **5**, 989–993 (1994).
- ²⁸B. Nadler, S. Lafon, R. R. Coifman, and I. G. Kevrekidis, "Diffusion maps, spectral clustering and reaction coordinates of dynamical systems," *Applied and Computational Harmonic Analysis* **21**, 113–127 (2006).
- ²⁹B. Nadler, S. Lafon, R. Coifman, and I. G. Kevrekidis, "Diffusion maps—a probabilistic interpretation for spectral embedding and clustering algorithms," in *Principal manifolds for data visualization and dimension reduction* (Springer, 2008) pp. 238–260.
- ³⁰M. Belkin and P. Niyogi, "Laplacian eigenmaps and spectral techniques for embedding and clustering," in *Advances in neural information processing systems* (2002) pp. 585–591.
- ³¹M. Belkin and P. Niyogi, "Laplacian eigenmaps for dimensionality reduction and data representation," *Neural computation* **15**, 1373–1396 (2003).
- ³²M. Meila, S. Koelle, and H. Zhang, "A regression approach for explaining manifold embedding coordinates," *arXiv preprint arXiv:1811.11891* (2018).
- ³³S. Chen and G. D. Doolen, "Lattice boltzmann method for fluid flows," *Annual review of fluid mechanics* **30**, 329–364 (1998).
- ³⁴S. Succi, *The lattice Boltzmann equation: for fluid dynamics and beyond* (Oxford university press, 2001).
- ³⁵P. L. Bhatnagar, E. P. Gross, and M. Krook, "A model for collision processes in gases. i. small amplitude processes in charged and neutral one-component systems," *Physical review* **94**, 511 (1954).
- ³⁶Y. Qian and S. Orszag, "Scalings in diffusion-driven reaction $a + b \rightarrow c$: Numerical simulations by lattice bgk models," *Journal of Statistical Physics* **81**, 237–253 (1995).
- ³⁷P. Van Leemput, K. Lust, and I. G. Kevrekidis, "Coarse-grained numerical bifurcation analysis of lattice boltzmann models," *Physica D: Nonlinear Phenomena* **210**, 58–76 (2005).
- ³⁸S. Lee, F. Dietrich, G. E. Karniadakis, and I. G. Kevrekidis, "Linking gaussian process regression with data-driven manifold embeddings for nonlinear data fusion," *Interface Focus* **9**, 20180083 (2019).
- ³⁹I. G. Kevrekidis and G. Samaey, "Equation-free multiscale computation: Algorithms and applications," *Annual review of physical chemistry* **60**, 321–344 (2009).