

Physics-informed learning of governing equations from scarce data

Zhao Chen¹, Yang Liu^{2,*}, and Hao Sun^{1,3,‡}

¹Department of Civil and Environmental Engineering, Northeastern University, Boston, MA 02115, USA

²Department of Mechanical and Industrial Engineering, Northeastern University, Boston, MA 02115, USA

³Department of Civil and Environmental Engineering, MIT, Cambridge, MA 02139, USA

*Corresponding author. E-mail: yang1.liu@northeastern.edu

‡Corresponding author. E-mail: h.sun@northeastern.edu

✉ Author contributions: Y.L. and H.S. contributed to the ideation and design of the research; Z.C. and H.S. performed the research; Z.C., Y.L. and H.S. wrote the paper.

Abstract

Harnessing data to discover the underlying governing laws or equations that describe the behavior of complex physical systems can significantly advance our modeling, simulation and understanding of such systems in various science and engineering disciplines. This work introduces a novel physics-informed deep learning framework to discover governing partial differential equations (PDEs) from scarce and noisy data for nonlinear spatiotemporal systems. In particular, this approach seamlessly integrates the strengths of deep neural networks for rich representation learning, physics embedding, automatic differentiation and sparse regression to (1) approximate the solution of system variables, (2) compute essential derivatives, as well as (3) identify the key derivative terms and parameters that form the structure and explicit expression of the PDEs. The efficacy and robustness of this method are demonstrated, both numerically and experimentally, on discovering a variety of PDE systems with different levels of data scarcity and noise accounting for different initial/boundary conditions. The resulting computational framework shows the potential for closed-form model discovery in practical applications where large and accurate datasets are intractable to capture.

Keywords: governing equation discovery, physics-informed deep learning, PDEs

Current practices on modeling of complex dynamical systems have been mostly rooted in the use of ordinary and/or partial differential equations (ODEs, PDEs) that govern the system behaviors. These governing equations are conventionally obtained from rigorous first principles such as the conservation laws or knowledge-based phenomenological derivations. However, there remain many real-world complex systems underexplored, whose analytical descriptions are undiscovered and parsimonious closed forms of governing equations are unclear or partially unknown. Luckily, observational datasets become increasingly rich and offer an alternative of distilling the underlying equations from data. Harnessing data to uncover the governing laws or equations can significantly advance and transform our modeling, simulation and understanding of complex physical systems in various science and engineering disciplines. For example, obtaining mathematical equations that govern the evolution of sea ice from observational data (e.g., satellite remote sensing images) brings distinct benefits for better understanding and predicting the growth, melt and movement of the Arctic ice pack. Distilling an explicit formulation from field sensing data (e.g., Doppler radar recordings)

will accelerate more accurate prediction of weather and climate patterns. Recently, advances in machine learning theories, computational capacity and data availability kindle significant enthusiasm and efforts towards data-driven discovery of physical laws and governing equations [1–12].

Pioneering contributions by Bongard and Lipson [1] and Schmidt and Lipson [2] leveraged stratified symbolic regression and genetic programming to successfully distil the underlying differential equations that govern nonlinear system dynamics from data. However, this elegant approach doesn't scale up well with the dimensionality of the system, is computationally expensive, and might suffer from overfitting issues. Recently, an impressive breakthrough made by Brunton *et al.* [5] leads to an innovative sparsity-promoting approach called sparse identification of nonlinear dynamics (SINDy), which selects dominant candidate functions from a high-dimensional nonlinear function space based on sparse regression to uncover parsimonious governing equations, ODEs in particular. The sparsity was achieved by a sequential threshold ridge regression (STRidge) algorithm which recursively determines the sparse solution subjected to hard thresholds [5, 6]. Such an approach is capable of balancing the complexity and accuracy of identified models and thus results in parsimony. SINDy has drawn tremendous attention in the past few years, leading to variant algorithms with applications to identify projected low-dimensional surrogate models in the form of first-order ODEs, alternatively with linear embedding [8, 10], for a wide range of nonlinear dynamical systems, such as fluid flows [13, 14], structural systems [15, 16], biological and chemical systems [17–19], active matter [20], predictive control of nonlinear dynamics [21], multi-time-scale systems [22], a predator-prey system [23], and stochastic processes [24], just naming a few among many others. There are also a number of other extensions of SINDy that discover implicit dynamics [17, 25], incorporate physics constraints [13], and embed random sampling to improve the robustness to noise for sparse discovery of high-dimensional dynamics [26]. The convergence and error estimate analyses [27] theoretically sustain the family of SINDy approaches.

The sparsity-promoting paradigm has been later extended for data-driven discovery of spatiotemporal systems governed by PDEs, e.g., the PDE-FIND algorithm [6, 7], where the library of candidate functions is augmented by incorporating spatial partial derivative terms. This method has been further investigated or improved to, for example, obtain parametric PDEs from data [28], discover PDEs enhanced by Bayesian inference [29] and gene expression programming [30], identify diffusion and Navier-Stokes equations based on molecular simulation [31], and learn PDEs for biological transport models [32]. Nevertheless, a critical bottleneck of the SINDy framework, especially for data-driven discovery of PDEs, lies in its strong dependence on both quality and quantity of the measurement data, since numerical differentiation is required to compute the derivatives in order to construct governing equation(s). Especially, the use of finite difference or filtering to calculate derivatives leads to a pivotal challenge that reduces the algorithm robustness. This specially limits the applicability of SINDy in its present form to scenarios given highly incomplete, scarce and noisy data. It is notable that variational system identification [9] shows satisfactory robustness of calculating derivatives based on isogeometric analysis for discovering the weak form of PDEs. However, such an approach doesn't scale down well with respect to the fidelity of available data. Another work [33] shows that weak formulation can significantly improve the discovery robustness against noise, but requires careful design of test functions, which is intractable for high-dimensional spatiotemporal systems.

Graph-based automatic differentiation [34] is well posed to address the above issue, which has been proven successful in deep learning for solving nonlinear PDEs [35–41]. In particular, the deep neural network (DNN) is used to approximate the solution constrained by both the PDE(s) and available data. Latest studies [42, 43] show the potential of using DNNs and automatic differentiation to obtain PDEs from noisy data; yet, false positive identification occurs due to the use of less rigorous sparse regression along with DNN training. Simultaneously optimizing the DNN param-

eters and sparse PDE coefficients poses a significant challenge in finding the global optimum. In this work, we present a novel Physics-informed Deep Learning (PiDL) framework, possessing salient features of interpretability and generalizability, to discover governing PDEs of nonlinear spatiotemporal systems from scarce and noisy data accounting for different initial/boundary conditions. Our methodology integrates the strengths of DNNs for rich representation learning, automatic differentiation for accurate derivative calculation as well as ℓ_0 sparse regression to tackle the fundamental limitation of existing methods that scale poorly with data noise and scarcity. The efficacy and robustness of our method are demonstrated on a variety of PDE systems, both numerically and experimentally.

RESULTS

PiDL with Sparse Regression for PDE Discovery

We consider a multi-dimensional spatiotemporal system whose governing equations can be described by a set of nonlinear, coupled, parameterized PDEs in the general form given by

$$\mathbf{u}_t + \mathcal{F}[\mathbf{u}, \mathbf{u}^2, \dots, \nabla_{\mathbf{x}}\mathbf{u}, \nabla_{\mathbf{x}}^2\mathbf{u}, \nabla_{\mathbf{x}}\mathbf{u} \cdot \mathbf{u}, \dots; \boldsymbol{\lambda}] = \mathbf{p} \quad (1)$$

where $\mathbf{u} = \mathbf{u}(\mathbf{x}, t) \in \mathbb{R}^{1 \times n}$ is the multi-dimensional latent solution (dimension = n) while \mathbf{u}_t is the first-order time derivative term; $t \in [0, T]$ denotes time and $\mathbf{x} \in \Omega$ specifies the space; $\mathcal{F}[\cdot]$ is a complex nonlinear functional of \mathbf{u} and its spatial derivatives, parameterized by $\boldsymbol{\lambda}$; ∇ is the gradient operator with respect to \mathbf{x} ; $\mathbf{p} = \mathbf{p}(\mathbf{x}, t)$ is the source term (note that, in many common cases, $\mathbf{p} = \mathbf{0}$ represents no source input to the system). The PDEs are also subjected to initial and boundary conditions (IBCs), if known, denoted by $\mathcal{I}[\mathbf{x} \in \Omega, t = 0; \mathbf{u}, \mathbf{u}_t] = 0$ and $\mathcal{B}[\mathbf{x} \in \partial\Omega; \mathbf{u}, \nabla_{\mathbf{x}}\mathbf{u}] = 0$. For systems that obey Newton’s second law of motion (e.g., \mathbf{u}_{tt} in wave equations), the governing PDEs can be written in a state-space form of Eq. (1) by defining $\mathbf{v} = \{\mathbf{u}, \mathbf{u}_t\}$ as the solution variable. Our objective is to find the closed form of $\mathcal{F}[\cdot]$ from available spatiotemporal measurements which are assumed to be incomplete, scarce and noisy commonly seen in real-world applications (e.g., when data capture is very costly or the data itself is sparse in nature). We assume that the physical law is governed by only a few important terms which can be selected from a large-space library of candidate functions, where sparse regression can be applied [5–7]. Inherent in this assumption leads to reformulation of Eq. (1) in the following (assuming zero or unknown source for simplicity):

$$\mathbf{u}_t = \boldsymbol{\phi}\boldsymbol{\Lambda} \quad (2)$$

Here, $\boldsymbol{\phi} = \boldsymbol{\phi}(\mathbf{u}) \in \mathbb{R}^{1 \times s}$ is an extensive library of symbolic functions consisting of many candidate terms, e.g., constant, polynomial, and trigonometric terms with respect to each spatial dimension [6, 7], assembled in a row vector given by $\boldsymbol{\phi} = \{1, \mathbf{u}, \mathbf{u}^2, \dots, \mathbf{u}_x, \mathbf{u}_y, \dots, \mathbf{u}^3 \odot \mathbf{u}_{xy}, \dots, \sin(\mathbf{u}), \dots\}$, where \odot represents the element-wise Hadamard product; s denotes the total number of candidate terms in the library; the subscripts in the context of $\{x, y, z\}$ depict the derivatives; $\boldsymbol{\Lambda} \in \mathbb{R}^{s \times n}$ is the sparse coefficient matrix (only the active candidate terms in $\boldsymbol{\phi}$ have non-zero values), e.g., $\boldsymbol{\Lambda} = [\boldsymbol{\lambda}^u \ \boldsymbol{\lambda}^v \ \boldsymbol{\lambda}^w] \in \mathbb{R}^{s \times 3}$ for $\mathbf{u} = \{u, v, w\}$. If there is an unknown source input, the candidate functions for \mathbf{p} can also be incorporated into $\boldsymbol{\phi}$ for discovery (see [Supplementary Note C3](#)). Thus, the discovery problem can then be stated as: given the spatiotemporal measurement data \mathcal{D}_u , find sparse $\boldsymbol{\Lambda}$ such that Eq. (2) holds.

We present an interpretable PiDL paradigm with sparse regression to simultaneously model the system response and identify the parsimonious closed form of the governing PDE(s). The innovative algorithm architecture of this method is shown in Fig. 1, where datasets sampled from two different

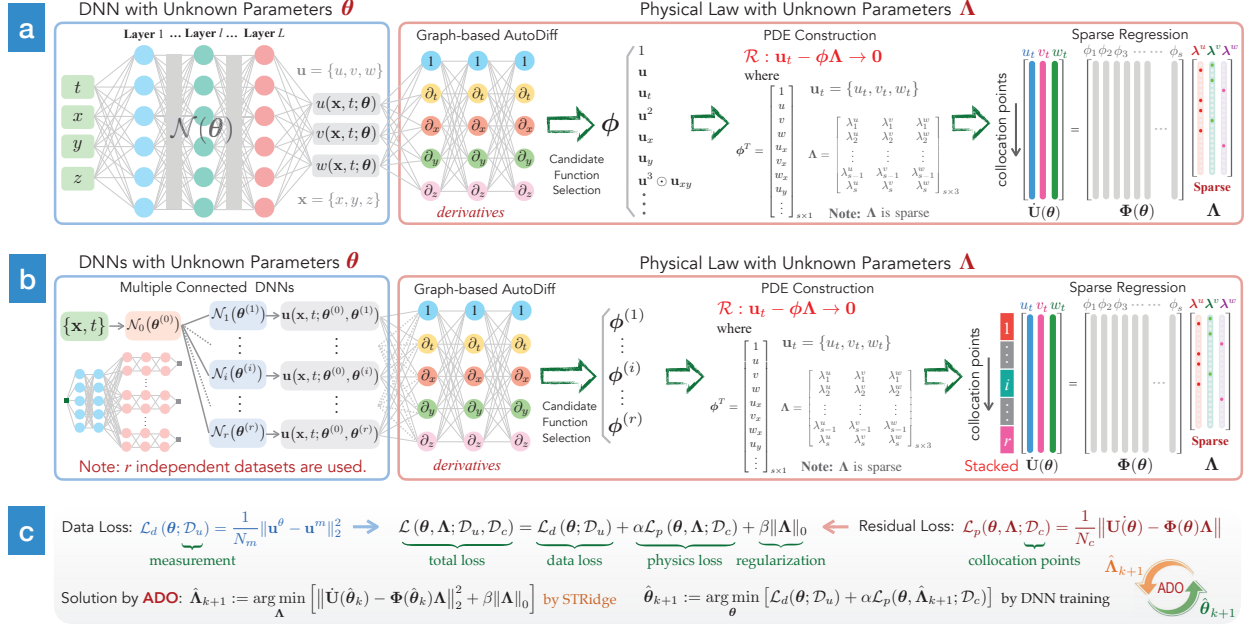


Fig. 1: Schematic architecture of the framework of PiDL with sparse regression for data-driven discovery of PDE(s). (a) the network for one dataset from a single IBC, (b) the “root-branch” network for $r \geq 2$ independent datasets from multiple IBCs, and (c) schematic for training the networks based on alternating direction optimization. The network consists of two components: a DNN governed by the trainable parameters θ , which maps the spatiotemporal coordinates $\{\mathbf{x}, t\}$ to the latent solution $\mathbf{u} = \{u, v, w\}$, and the physical law described by a set of nonlinear PDEs, which are formed by the derivative candidate functions ϕ parameterized by the unknown sparse coefficients Λ . Note that, for the case of multiple independent datasets, the libraries $\phi^{(i)}$ are concatenated to build ϕ for constructing the unified governing PDE(s). The total loss function $\mathcal{L}(\theta, \Lambda; \mathcal{D}_u, \mathcal{D}_c)$ is composed of the data loss $\mathcal{L}_d(\theta, \mathcal{D}_u)$, the physics loss $\alpha \mathcal{L}_p(\theta, \Lambda; \mathcal{D}_c)$, and the ℓ_0 regularization term $\beta \|\Lambda\|_0$ that promotes the sparsity. Here, α and β denote the relative weighting of the loss functions, while \mathcal{D}_u and \mathcal{D}_c represent the measurement data and collocation samples respectively. Note that the physics loss, in a residual form, is only evaluated on the spatiotemporal collocation samples. The colored dots in the sparse coefficients matrix (or vector) on the right denote non-zero values. Simultaneous optimization of the unknown parameters $\{\theta, \Lambda\}$ leads to both the trained DNN for inference of the data-driven full-field solution and the discovered parsimonious closed-form PDEs.

IBC scenarios are considered: (1) one dataset from a single IBC and (2) $r \geq 2$ independent datasets from multiple IBCs. For the case of single dataset, we interpret the latent solution \mathbf{u} by a DNN (denoted by \mathcal{N}), namely, $\mathbf{u}^\theta = \mathbf{u}(\mathbf{x}, t; \theta)$, where θ represents the DNN trainable parameters including weights and biases, as shown in Fig. 1a. When multiple independent datasets are available, a “root-branch” DNN depicted in Fig. 1b is designed to approximate the latent solutions \mathbf{u}_i ($i = 1, \dots, r$) corresponding to different IBCs, viz., $\mathbf{u}_i^\theta = \mathbf{u}(\mathbf{x}, t; \theta^{(0)}, \theta^{(i)})$, where $\theta^{(0)}$ and $\theta^{(i)}$ denote the trainable parameters of the root layers $\mathcal{N}^{(0)}$ and the branch layers $\mathcal{N}^{(i)}$, respectively. Noteworthy, the IBCs are unnecessarily either known *a priori* or measured since the measurement data already reflects the specific IBC (e.g., there exists a one-to-one mapping between the IBC and the PDE solution). The DNN essentially plays a role as a nonlinear functional to approximate the latent solution with the data loss function $\mathcal{L}_d(\theta; \mathcal{D}_u)$. With graph-based automatic differentiation where derivatives on \mathbf{u} are evaluated at machine precision, the library of candidate functions ϕ^θ can be computed from the DNN. For the case of multiple independent datasets, the libraries $\phi^{(i)}$ resulted from the branch nets are concatenated to build ϕ^θ for constructing the unified governing PDE(s). Thus, the sparse representation of the reconstructed PDE(s) can be written in a residual form, namely, $\mathcal{R}^\theta := \mathbf{u}_t^\theta - \phi^\theta \Lambda \rightarrow \mathbf{0}$, where $\mathcal{R}^\theta \in \mathbb{R}^{1 \times n}$ denotes the PDE residuals. The basic concept is to adapt

Table 1: Summary of the PiDL discovery results in the context of accuracy for a range of canonical models.

PDE name	Err. (N-0%)	Err. (N-1%)	Err. (N-10%)	Description of data discretization
Burgers'	0.01±0.01%	0.19±0.11%	1.15±1.20%	$x \in [-8, 8]_{d=256}, t \in [0, 10]_{d=101}$, sub. 1.95%
KS	0.07±0.01%	0.61±0.04%	0.71±0.06%	$x \in [0, 100]_{d=1024}, t \in [0, 100]_{d=251}$, sub. 12.3%
Schrödinger	0.09±0.04%	0.65±0.29%	2.31±0.28%	$x \in [-4.5, 4.5]_{d=512}, t \in [0, \pi]_{d=501}$, sub. 37.5%
NS	0.66±0.72%	0.86±0.63%	1.40±1.83%	$x \in [0, 9]_{d=449}, y \in [-2, 2]_{d=199}, t \in [0, 30]_{d=151}$, sub. 0.22%
λ - ω RD	0.07±0.08%	0.25±0.30%	4.78±3.66%	$x, y \in [-10, 10]_{d=256}, t \in [0, 10]_{d=201}$, sub. 0.29%

Note: The error is defined as the average relative error of the identified non-zero coefficients w.r.t. the ground truth. The percentage values in the parentheses denote the noise levels (e.g., noise free 0%, 1% and 10%) and the subscript d represents the number of discretization. Our method is also compared with SINDy (the PDE-FIND approach presented in [6]) as illustrated in [Supplementary Table S1](#). It is noted that much less measurement data polluted with a higher level of noise are used in our discovery. Gaussian white noise is added to the synthetic response with the noise level defined as the root-mean-square ratio between the noise and the exact solution.

both the DNN trainable parameters θ and the PDE coefficients Λ such that the neural network can fit the measurement data while satisfying the constraints defined by the underlying PDE(s). The PDE residuals will be evaluated on a large number of collocation points $\mathcal{D}_c = \{\mathbf{x}_i, t_i\}_{i=1}^{N_c}$, randomly sampled in the spatiotemporal space, leading to the residual physics loss function $\mathcal{L}_p(\theta, \Lambda; \mathcal{D}_c)$. When multiple IBCs are considered, the measurement data and the collocation points will be stacked when calculating the data loss and the physics loss (based on a unified physics residual formulation $\mathcal{R}^\theta \rightarrow \mathbf{0}$).

The total loss function for training the overall PiDL network is thus composed of the data loss \mathcal{L}_d , the residual physics loss \mathcal{L}_p and a regularization term, expressed as:

$$\mathcal{L}(\theta, \Lambda; \mathcal{D}_u, \mathcal{D}_c) = \mathcal{L}_d(\theta; \mathcal{D}_u) + \alpha \mathcal{L}_p(\theta, \Lambda; \mathcal{D}_c) + \beta \|\Lambda\|_0 \quad (3)$$

where α is the relative weighting of the residual physics loss function; β is the regularization parameter; $\|\cdot\|_0$ represents the ℓ_0 norm. Optimizing the total loss function can produce a DNN that can not only predict the data-driven full-field system response, but also uncover the parsimonious closed-form PDE(s), i.e., $\{\theta^*, \Lambda^*\} := \arg \min_{\{\theta, \Lambda\}} [\mathcal{L}(\theta, \Lambda; \mathcal{D}_u, \mathcal{D}_c)]$, where $\{\theta^*, \Lambda^*\}$ denote the optimal set of parameters. Noteworthy, the total loss function has an implicit complex form, and thus, directly solving the optimization problem is highly intractable since the ℓ_0 regularization makes this problem np -hard. To address this challenge, we present an alternating direction optimization (ADO) algorithm that divides the overall optimization problem into a set of tractable subproblems to sequentially optimize the trainable parameters, as shown in Fig. 1c. Pre-training of PiDL is conducted before running the ADO algorithm for discovery, by simply replacing $\|\Lambda\|_0$ in Eq. (3) with $\|\Lambda\|_1$ where brute-force gradient-based optimization for both θ and Λ becomes applicable. The ℓ_1 -regularized pre-training can accelerate the convergence of ADO by providing an admissible “initial guess”. More detailed formulation and algorithm description are found in [Method](#) and [Supplementary Note A](#).

The synergy of DNN and sparse regression results in the following outcome: the DNN provides accurate modeling of the latent solution, its derivatives and possible candidate function terms as a basis for constructing the governing PDE(s), while the sparsely represented PDE(s) in turn constraints the DNN modeling and projects correct candidate functions, eventually turning the measured system into closed-form PDE(s).

Discovery of Benchmark PDEs with Single Dataset

We observe the efficacy and robustness of our methodology on a group of canonical PDEs used to represent a wide range of physical systems with nonlinear, periodic and/or chaotic behaviors. In particular, we discover the closed forms of Burgers’, Kuramoto-Sivashinsky (KS), nonlinear

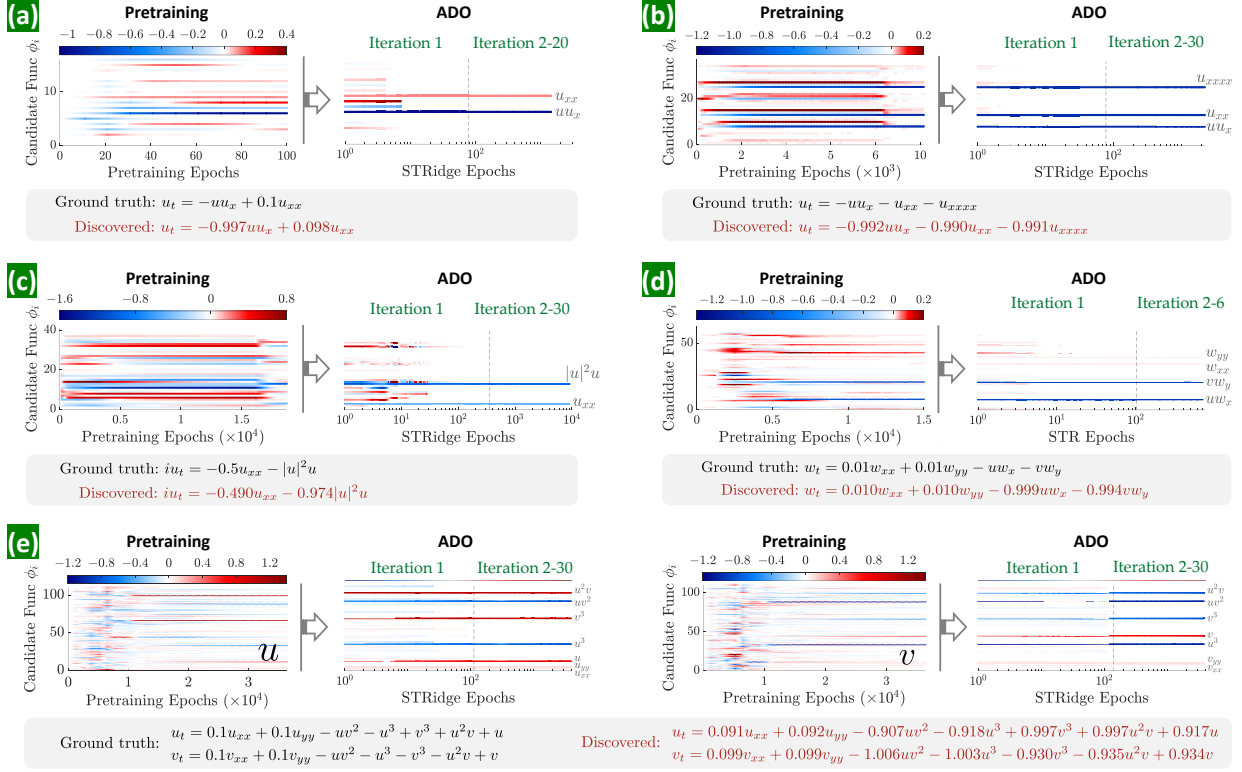


Fig. 2: Discovery of selected benchmark PDEs for sparsely sampled measurement data with 10% noise. (a) Discovered Burgers’ equation: evolution of the sparse coefficients $\Lambda \in \mathbb{R}^{16 \times 1}$ for 16 candidate functions $\phi \in \mathbb{R}^{1 \times 16}$ used to form the PDE, where the color represents the coefficient value. (b) Discovered KS equation: Evolution of the sparse coefficients $\Lambda \in \mathbb{R}^{36 \times 1}$ for 36 candidate functions $\phi \in \mathbb{R}^{1 \times 36}$. (c) Discovered nonlinear Schrödinger equation: evolution of the sparse coefficients $\Lambda \in \mathbb{R}^{40 \times 1}$ for the candidate functions $\phi \in \mathbb{R}^{1 \times 40}$. (d) Discovered NS equation: evolution of the sparse coefficients $\Lambda \in \mathbb{R}^{60 \times 1}$ for 60 candidate functions $\phi \in \mathbb{R}^{1 \times 60}$. (e) Discovered RD equations: evolution of the sparse coefficients $\lambda^u \in \mathbb{R}^{110 \times 1}$ and $\lambda^v \in \mathbb{R}^{110 \times 1}$ ($\Lambda = [\lambda^u \ \lambda^v]$) for 110 candidate functions $\phi \in \mathbb{R}^{1 \times 110}$ used to reconstruct the u -equation and the v -equation, respectively.

Schrödinger, Navier-Stokes (NS), and λ - ω Reaction-Diffusion (RD) equations from scarce and noisy time-series measurements recorded by a number of sensors at fixed locations (data are polluted with Gaussian white noise) from a single IBC. Results are presented in Table 1, Fig. 2 and Fig. 3, which show quite accurate discovery and demonstrate satisfactory performance of the proposed method and its robustness to measurement data scarcity and noise. We also compare our method with SINDy considering different levels of data scarcity and noise (summarized in Supplementary Note B6 and Table S1).

Burgers’ Equation: We first consider a dissipative system with the dynamics governed by a 1D viscous Burgers’ equation expressed as $u_t = -uu_x + \nu u_{xx}$, where ν (equal to 0.1) denotes the diffusion coefficient. The equation describes the decaying stationary viscous shock of a system after a finite period of time, commonly found in simplified fluid mechanics, nonlinear acoustics and gas dynamics. We test the PiDL approach on the recorded traveling shock waves from the solution to Burgers’ equation subjected to a Gaussian initial condition. In particular, 5 sensors are randomly placed at fixed locations among the 256 spatial grids and record the wave for 101 time steps, leading to 1.95% of the dataset used in [6]. A full description of the dataset, design of the library of candidate functions (16 terms) and model training is given in Supplementary Note B.1.1. Fig. 2a shows the discovered Burgers’ equation for a dataset with 10% noise. The evolution of the coefficients $\Lambda \in$

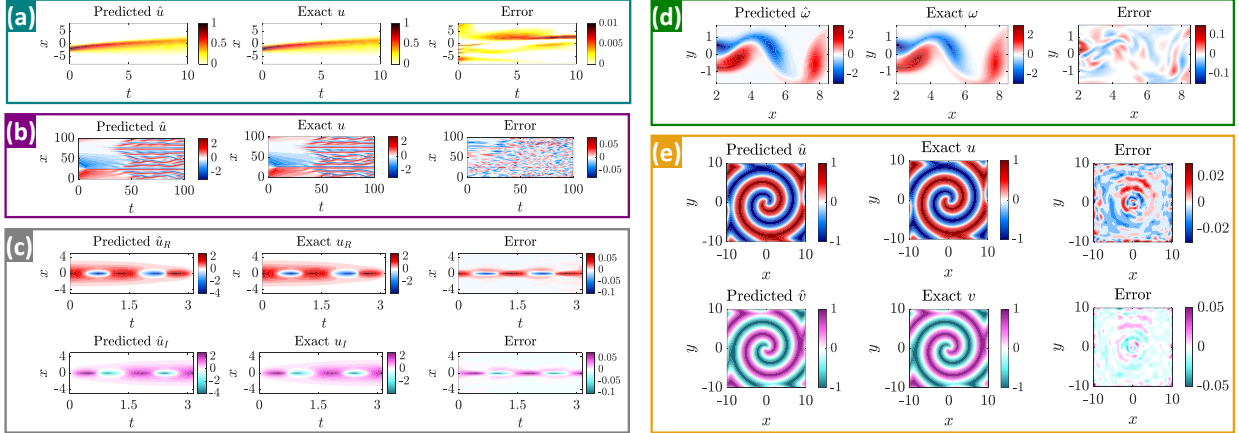


Fig. 3: Predicted responses compared with the exact solutions for selected canonical PDEs. (a) Burgers' equation, (b) KS equation, (c) nonlinear Schrödinger equation, (d) NS equation, and (e) λ - ω RD equations. Note that the sparsely sampled measurement data has 10% noise.

$\mathbb{R}^{16 \times 1}$ illustrates robust convergence to the ground truth (error about 1.2%), resulting in accurate discovery. The trained PiDL properly reproduces the dynamical response from noisy measurements (e.g., the full-field ℓ_2 prediction error is 2.02%) as shown in Fig. 3a. The ADO algorithm converges only after the first alternating iteration and shows capacity to recover the correct sparsity pattern of the PDE. We also discover the Burgers' equation with an unknown/unmeasured source $\sin(x) \sin(t)$, given scarce u -measurement with 10% noise. When discovering the underlying governing equation, the source should be considered and reconstructed concurrently. In this case, we incorporate 14 source candidate functions, composed of $\{\sin(t), \sin(x), \cos(t), \cos(x)\}$ and their combination, into the aforementioned library, resulting in a total of 30 candidate terms for simultaneous discovery of the PDE and reconstruction of the unknown source. The corresponding discovery result is summarized in Extended Data Fig. 1, which includes the discovered equation and source function, the evolution of sparse coefficients $\Lambda \in \mathbb{R}^{30 \times 1}$, and the predicted full-field response. It turns out that both PDE and source terms along with their coefficients are well identified. Nevertheless, if the source is very complex with its general expression or form completely unknown, distinct challenges arise when designing the source candidate functions. This may require an extraordinarily large-space library to retain diversifying representations, and thus pose additional computational complexity for accurate discovery of the PDEs. More discussions are presented in [Supplementary Note C3](#).

Kuramoto-Sivashinsky (KS) Equation: Another dissipative system with intrinsic instabilities is considered, governed by the 1D Kuramoto-Sivashinsky (KS) equation $u_t = -uu_x - u_{xx} - u_{xxxx}$, where the reverse diffusion term $-u_{xx}$ leads to the disruptive behavior while the fourth-order derivative u_{xxxx} introduces chaotic patterns as shown in Fig. 3b, making an ideal test problem for equation discovery. The KS equation is widely used to model the instabilities in laminar flame fronts and dissipative trapped-ion modes among others. We randomly choose 320 points as fixed sensors and record the wave response for 101 time steps, resulting in 12.3% of the dataset used in [6]. A total of 36 candidate functions are employed to construct the underlying PDE. Detail description of this example is found in [Supplementary Note B.1.2](#). It is notable that the chaotic behavior poses significant challenges in approximating the full-field spatiotemporal derivatives, especially the high-order u_{xxxx} , from poorly measured data for discovery of such a PDE. Existing methods (e.g., the family of SINDy methods [6, 7]) eventually fail in this case given very coarse and noisy measurements. Nevertheless, PiDL successfully distills the closed form of the KS equation from subsampled sparse data with 10% noise, shown in Fig. 2b. The evolution of the coefficients $\Lambda \in \mathbb{R}^{36 \times 1}$ in Fig. 2b

illustrates that both the candidate terms and the corresponding coefficients are correctly identified (close to the original parameters; error around 0.7%) within a few ADO iterations. The predicted full-field wave by the trained PiDL also coincides with the exact solution at a relative ℓ_2 error of 1.87% (Fig. 3b).

Nonlinear Schrödinger Equation: In the third example, we discover the nonlinear Schrödinger equation, $iu_t = -0.5u_{xx} - |u|^2u$, where u is a complex field variable. This well-known equation is widely used in modeling the propagation of light in nonlinear optical fibers, Bose-Einstein condensates, Langmuir waves in hot plasmas, and so on. We take 37.5% subsamples (e.g., randomly selected from the spatial grids) of the dataset as shown in Table 1 to construct the PDE using 40 candidate functions $\phi \in \mathbb{R}^{1 \times 40}$. Since the function is complex-valued, we model separately the real part (u_R) and the imaginary part (u_I) of the solution in the output of the DNN, assemble them to obtain the complex solution $u = u_R + iu_I$, and construct the complex-valued candidate functions for PDE discovery. To avoid complex gradients in optimization, we use the modulus $|u|$, instead of the ℓ_2 norm shown in Eq. (5), for the residual physics loss \mathcal{L}_p (see [Supplementary Note B.1.3](#) for more details). Fig. 2c shows the discovered Schrödinger equation for the case of 10% noise. The evolution history of the sparse coefficients $\mathbf{\Lambda} \in \mathbb{R}^{40 \times 1}$ clearly shows the convergence to the actual values (Fig. 2c; error about 4.14%) resulting in accurate closed-form identification of the PDE, while the reconstructed full-field response, for both real and imaginary parts, matches well the exact solution with a slight relative ℓ_2 error of 1% (Fig. 3c).

Navier-Stokes (NS) Equation: We consider a 2D fluid flow passing a circular cylinder with the local rotation dynamics governed by the well-known Navier-Stokes vorticity equation $w_t = -(\mathbf{u} \cdot \nabla)w + \nu \nabla^2 w$, where w is the spatiotemporally variant vorticity, $\mathbf{u} = \{u, v\}$ denotes the fluid velocities, and ν is the kinematic viscosity ($\nu = 0.01$ at Reynolds number 100). We leverage the open simulation data [6] and subsample a dataset of the flow response $\{u, v, w\}$ at 500 spatial locations randomly picked within the indicated region in [Supplementary Fig. S4](#), which record time series for 60 time steps. The resulting dataset is only 10% of that used in [6]. A comprehensive discussion of this example is found in [Supplementary Note B.1.4](#). Fig. 2d summarizes the result of the discovered NS equation for a dataset with 10% noise. It is encouraging that the uncovered PDE expression is almost identical to the ground truth, for both the derivative terms and their coefficients, even under 10% noise corruption. The coefficients $\mathbf{\Lambda} \in \mathbb{R}^{60 \times 1}$, corresponding to 60 candidate functions $\phi \in \mathbb{R}^{1 \times 60}$, converge very quickly to the correct values with precise sparsity right after the first ADO iteration (Fig. 2d). The vorticity patterns and magnitudes are also well predicted as indicated by the snapshot (at $t = 23.8$) shown in Fig. 3d (the full-field ℓ_2 error for all snapshots is about 2.57%). This example provides a compelling test case for the proposed PiDL approach which is capable of discovering the closed-form NS equation with scarce and noisy data.

Reaction-Diffusion (RD) Equations: The examples above are mostly low-dimensional models with limited complexity. We herein consider a λ - ω reaction-diffusion (RD) system in a 2D domain with the pattern forming behavior governed by two coupled PDEs: $u_t = 0.1\nabla^2 u + \lambda(g)u - \omega(g)v$ and $v_t = 0.1\nabla^2 v + \omega(g)u + \lambda(g)v$, where u and v are the two field variables, $g = u^2 + v^2$, $\omega = -g^2$, and $\lambda = 1 - g^2$. The RD equations exhibit a wide range of behaviors including wave-like phenomena and self-organized patterns found in chemical and biological systems. The particular RD equations considered here display spiral waves subjected to periodic boundary conditions. Full details on the dataset, selection of candidate functions and hyperparameter setup of the PiDL model are given in [Supplementary Note B.1.5](#). Fig. 2e shows the evolution of the sparse coefficients $\boldsymbol{\lambda}^u, \boldsymbol{\lambda}^v \in \mathbb{R}^{110 \times 1}$ for 110 candidate functions $\phi \in \mathbb{R}^{1 \times 110}$, given a dataset with 10% noise. Both the sparse terms and the associated coefficients are precisely identified to form the the closed-form equations (as depicted in Fig. 2e). Due to the complexity of the PDEs and the high dimension, slightly more epochs are required in ADO to retain reliable convergence. The predicted response snapshots (e.g., at $t = 2.95$)

by the trained PiDL in Fig. 3e are close to the ground truth. This example shows especially the great ability and robustness of our method for discovering governing PDEs for high-dimensional systems from highly noisy data.

Discovery of PDEs with Multiple Independent Datasets

To demonstrate the “root-branch” network presented in Fig. 1b for discovery of PDE(s) based on multiple independent datasets sampled under different IBCs, we consider (1) the 1D Burgers’ equation with light viscosity that exhibits a shock behavior, and (2) a 2D Fitzhugh-Nagumo (FN) type reaction-diffusion system that describes activator-inhibitor neuron activities excited by external stimulus. The measurement data are sparsely sampled (e.g., time series or snapshots) with 10% noise under three different IBCs. Note that the IBCs are unnecessarily either measured or known *a priori* since the measurements already reflect the specific IBC which holds uniquely one-to-one mapping to the system response. The discovery results are discussed as follows.

Burgers’ Equation with Shock Behavior: In this example, we test the previously discussed Burgers’ equation with a small diffusion/viscosity parameter ($\nu = 0.01/\pi \approx 0.0032$) based on datasets generated by imposing three different IBCs. Such a small coefficient creates shock formation in a compact area with sharp gradient (see Fig. 4c) that could challenge the DNN’s approximation ability and thus affect the discovery. The three initial and Dirichlet boundary conditions include:

$$\begin{aligned} \text{IBC 1: } & u(x, 0) = -\sin(\pi x), u(-1, t) = u(1, t) = 0 \\ \text{IBC 2: } & u(x, 0) = \mathcal{G}(x), u(-1, t) = u(1, t) = 0 \\ \text{IBC 3: } & u(x, 0) = -x^3, u(-1, t) = 1, u(1, t) = -1 \end{aligned}$$

where \mathcal{G} denotes a Gaussian function. Although the measurement datasets for different IBCs exhibit completely distinct system responses, they obey the same underlying PDE, namely, $u_t = -uu_x + 0.0032u_{xx}$. For all IBCs, we assume that there are 30 sensors randomly deployed in space ($x \in [-1, 1]$) measuring the wave traveling (e.g., u) for 500 time instants ($t \in [0, 1]$). A denser sensor grid is needed herein, compared with the previous Burgers’ example, in order to capture the shock behaviors. Fig. 4a shows some of the measurements recorded by 6 typical sensors under 10% noise. A three-branch network ($r = 3$) shown in Fig. 1b is used for discovery. The full description of the dataset, the library of candidate functions (16 terms) and model training is given in [Supplementary Note B.3.1](#). Fig. 4b depicts the evolution of the coefficients ($\mathbf{\Lambda} \in \mathbb{R}^{16 \times 1}$) of candidate functions, where the correct terms in the library (uu_x and u_{xx}) are successfully distilled while other redundant terms are eliminated (e.g., hardly thresholded to zero) by ADO. The coefficients of the active terms are accurately identified as well (in particular the small viscosity parameter that leads to shock formation, e.g., 0.0039). The discovered PDE reads $u_t = -1.006uu_x + 0.0039u_{xx}$. Fig. 4c-d show the predicted responses and errors for three IBC cases, with a stacked full-field ℓ_2 error of 2.24%.

Fitzhugh-Nagumo (FN) Reaction-Diffusion System: We consider the Fitzhugh-Nagumo (FN) type reaction-diffusion system, in a 2D domain $\Omega = [0, 150] \times [0, 150]$ with periodic boundary conditions, whose governing equations are expressed by two coupled PDEs: $u_t = \gamma_u \Delta u + u - u^3 - v + \alpha$ and $v_t = \gamma_v \Delta v + \beta(u - v)$. Here, u and v represent two interactive components/matters (e.g., biological), $\gamma_u = 1$ and $\gamma_v = 100$ are diffusion coefficients, $\alpha = 0.01$ and $\beta = 0.25$ are the coefficients for reaction terms, and Δ is the Laplacian operator. The FN equations are commonly used to describe biological neuron activities excited by external stimulus (α), which exhibit an activator-inhibitor system because one equation boosts the production of both components while the other equation dissipates their new growth. Three random fields are taken as initial conditions to generate three independent datasets for discovery, each of which consists of 31 low-resolution snapshots

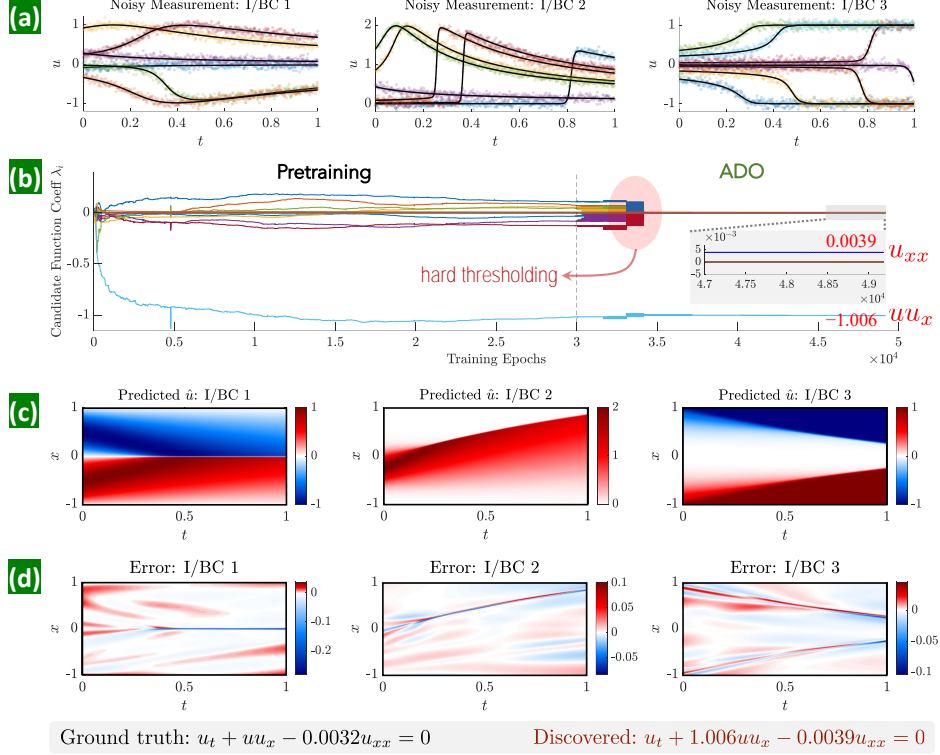


Fig. 4: Discovered Burgers' equation with small viscosity based on datasets sampled under three IBCs with 10% noise. (a) Visualization of noisy measurements for the three IBCs. Note that there are 30 sensors and only a few are illustrated in this figure. (b) Evolution of the sparse coefficients $\Lambda \in \mathbb{R}^{16 \times 1}$ for 16 candidate functions $\phi \in \mathbb{R}^{1 \times 16}$ used to construct the PDE, where the color represents the coefficient value. The correct terms (uu_x and u_{xx}) and their coefficients are successfully identified while other redundant terms are eliminated by ADO. (c-d) The predicted responses and errors for three IBC cases. The ground truth is not listed herein since the visualization is almost indistinguishable from the prediction (see [Supplementary Fig. S.7](#)). The relative full-field ℓ_2 error of the stacked prediction is 2.24%.

(projected into a 31×31 grid) down-sampled from the high-fidelity simulation under a 10% noise condition (see [Extended Data Fig. 2](#)). We assume the diffusion terms (Δu and Δv) are known in the PDEs, whose coefficients (γ_u and γ_v) yet need to be identified. A library with 72 candidate functions ($\phi \in \mathbb{R}^{1 \times 72}$) is designed for discovery of the coupled PDEs (in particular, the nonlinear reaction terms). Similar to the previous example, a root-branch network shown in [Fig. 1b](#) is employed for discovery. More description of the data generation, the specific candidate functions and model training can be found in [Supplementary Note B.3.2](#). [Fig. 5a-b](#) depict the evolution of the sparse coefficients $\lambda^u, \lambda^v \in \mathbb{R}^{72 \times 1}$ for 72 candidate functions. The pretraining step provides a redundant projection of the system onto 72 candidates; however, minor candidates are pruned out right after the first ADO iteration. The rest ADO iterations continue to refine all the trainable parameters including θ , λ^u and λ^v . The finally discovered PDEs are listed in [Fig. 5](#) in comparison with the ground truth. It is seen that the form of the PDEs is precisely uncovered with all correct active terms (including the unknown external stimulus in the first equation). The corresponding identified coefficients are generally close to the ground truth except the diffusion coefficient for v (i.e., γ_v) which seems to be a less sensitive parameter according to our test. It should be noted that, given very scarce and noisy measurement datasets in this example, the ‘‘root-branch’’ DNN is faced with challenges to accurately model the solutions with sharp propagating fronts (see [Fig. 5c](#)). The less accurate solution approximation by DNN then affects the discovery precision. This issue can be

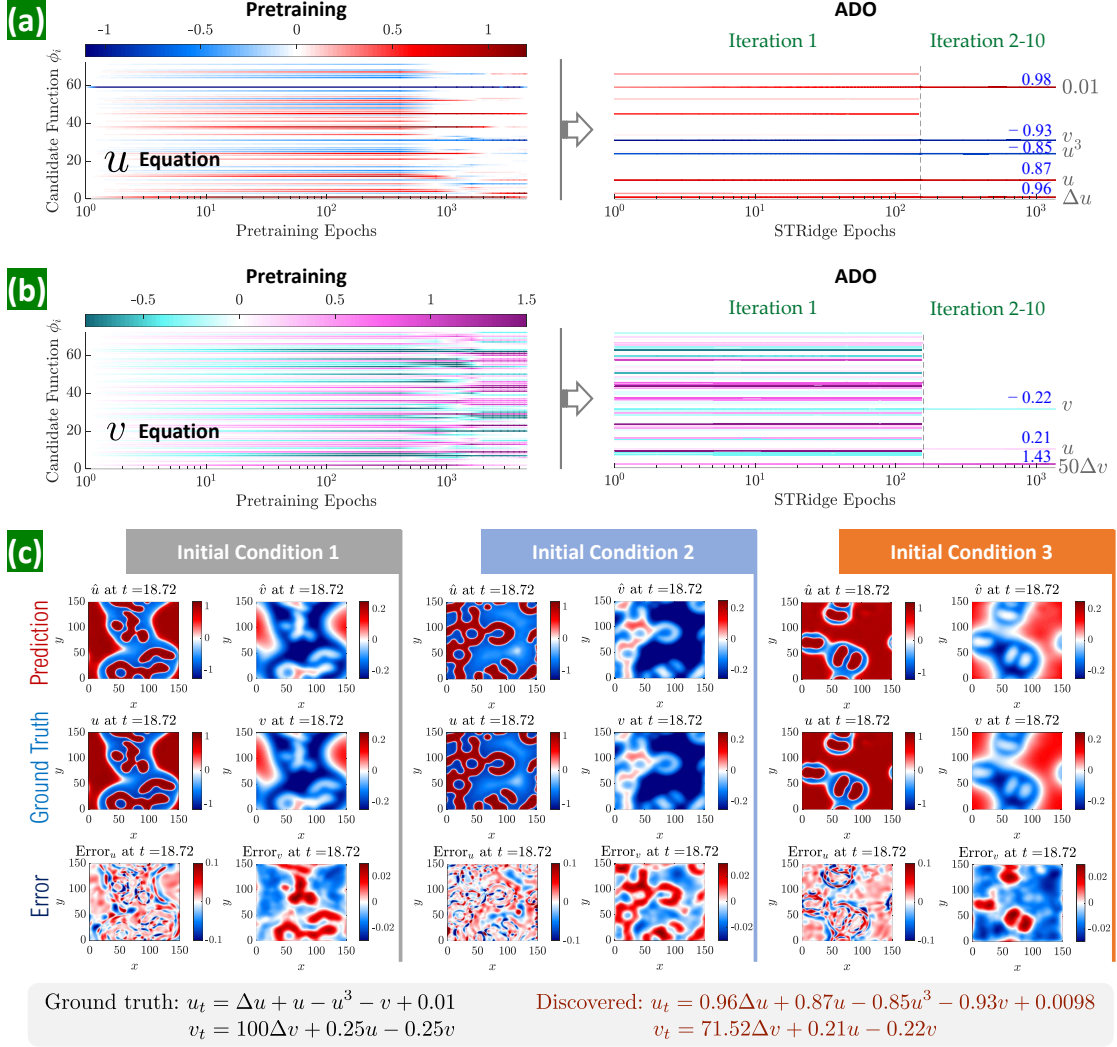


Fig. 5: Discovered Fitzhugh-Nagumo equations based on data sampled under three initial conditions (ICs) with 10% noise. (a) Evolution of the sparse coefficients $\lambda_u \in \mathbb{R}^{72 \times 1}$ for 72 candidate functions used to construct the first PDE (u -equation), where the color represents the coefficient value. (b) Evolution of the sparse coefficients $\lambda_v \in \mathbb{R}^{72 \times 1}$ for the second PDE (v -equation). For visualization purpose, we re-scale the identified coefficients of the constant stimulus term “1” in the u -equation by multiplying 100 and of the diffusion term Δv in the v -equation by dividing 50. (c) Snapshots of predicted response, ground truth and error distributions for all three ICs at an unmeasured time instance ($t = 18.72$). The relative ℓ_2 error for the predicted full-field response (stacked u and v) is 5.04%.

naturally alleviated by increasing the spatiotemporal measurement resolution (even still under fairly large noise pollution, e.g., 10%). Nevertheless, the exact form of the PDEs is successfully discovered in this challenging example, which is deemed more important since the coefficients can be further tuned/calibrated when additional data arrives. Fig. 5c shows typical snapshots of the predicted u and v components, the ground truth reference and the error distributions for one unmeasured time instance ($t = 18.72$). The stacked full-field ℓ_2 error is 5.04%.

Experimental Discovery of Cell Migration and Proliferation

The last example is placed to demonstrate the proposed approach for discovering a governing PDE that describes cell migration and proliferation, based on the sparse and noisy experimental

data collected from *in vitro* cell migration (scratch) assays [44]. The 1D cell density distributions at different time instants (0h, 12h, 24h, 36h, 48h) were extracted from high-resolution imaging via image segmentation and cell counting. A series of assays were performed under different initial cell densities (e.g., the total number of cells spans from 10,000 to 20,000 following the designated initial distribution in the test well shown in Extended Data Fig. 3a at $t = 0$ h). More detailed description of the experiment setup and datasets can be found in [44]. Our objective herein is to uncover a parsimonious PDE for modeling the dynamics of cell density $\rho(x, t)$. Here, we consider four scenarios with the initial number of cells ranging from 14,000, 16,000, 18,000 to 20,000. We take the mean of the test data from three identically-prepared experimental replicates for each scenario (see Extended Data Fig. 3b-e) to train our model shown in Fig. 1a for PDE discovery. Given our prior knowledge that the cell dynamics can be described by a diffusion (migration) and reaction (proliferation) process, we assume the PDE holds the form of $\rho_t = \gamma\rho_{xx} + \mathcal{F}(\rho)$, where γ is the unknown diffusion coefficient and \mathcal{F} denotes the underlying nonlinear reaction functional. We use 8 additional candidate terms (e.g., $\{1, \rho, \rho^2, \rho^3, \rho_x, \rho\rho_x, \rho^2\rho_x, \rho^3\rho_x\}$) to reconstruct \mathcal{F} , whose coefficients are sparse. Hence, the total number of trainable coefficients remains 9 (e.g., $\mathbf{\Lambda} \in \mathbb{R}^{9 \times 1}$). Other details on the PiDL model setting and training can be found in [Supplementary Note B.4](#).

Fig. 6a shows the evolution of 9 coefficients for the example case of 18,000 cells, where redundant candidate terms are pruned right after the first ADO iteration via hard thresholding of the corresponding coefficients to zero. The next ADO iterations followed by post-tuning refine the coefficients of active terms for final reconstruction of the PDE. Fig. 6b depicts the identified active term coefficients and the corresponding PDEs for different quantities of cells, sharing a unified form of $\rho_t = \gamma\rho_{xx} + \lambda_1\rho + \lambda_2\rho^2$ which exactly matches the famous Fisher-Kolmogorov model [45]. The rates of migration (diffusion) and proliferation (reaction) generally increase along with the number of cells, as seen from the identified coefficients in Fig. 6b. With the discovered PDEs, we simulate/predict the evolution of cell densities at different time instants (12h, 24h, 36h and 48h) presented in Fig. 6c-f, where the measurement at 0h is used as the initial condition while $\rho_x(x = 0, t) = \rho_x(x = 1900, t) = 0$ is employed as the Neumann boundary condition. The satisfactory agreement between the prediction and the measurement provides a clear validation of our discovered PDEs. It is noted that the extremely scarce and noisy experimental datasets unfortunately pose intractable challenge for any other existing methods (e.g., SINDy [5, 6]) to produce a reasonable discovery. This experimental example further demonstrates the strength and capacity of the proposed methodology in regard to handling high level of data scarcity and noise for PDE discovery.

DISCUSSION

In summary, we have presented a novel interpretable deep learning method for discovering physical laws, in particular parsimonious closed-form PDE(s), from scarce and noisy data (commonly seen in scientific investigations and real-world applications) for multi-dimensional nonlinear spatiotemporal systems. This approach combines the strengths of DNNs for rich representation learning of nonlinear functions, automatic differentiation for accurate derivative calculation as well as ℓ_0 sparse regression to tackle the fundamental limitation faced by existing sparsity-promoting methods that scale poorly with respect to data noise and scarcity. The use of collocation points (having no correlation with the measurement data) can render the proposed framework tolerable to scarce and noisy measurements, making the DNN for PDE solution approximation generalizable (see [Supplementary Note C2](#)). The special network architecture design is able to account for multiple independent datasets sampled under different initial/boundary conditions. An alternating direction optimization

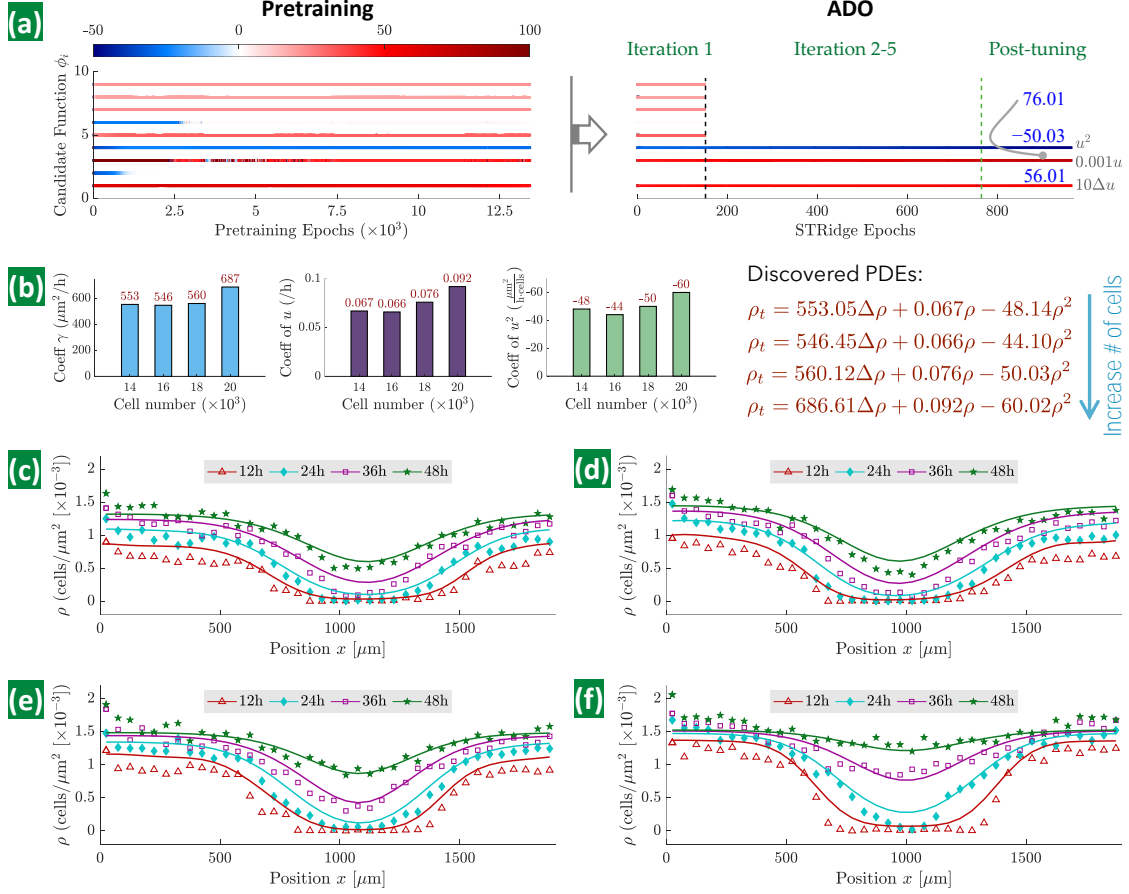


Fig. 6: Discovery result for cell migration and proliferation. (a) Example evolution of the sparse coefficients $\mathbf{\Lambda} \in \mathbb{R}^{9 \times 1}$ for 9 candidate functions used to construct the underlying PDE for the case of 18,000 cells. The diffusion and reaction coefficients for Δu and u are re-scaled for visualization purpose. (b) Discovered active terms $\{\Delta\rho, \rho, \rho^2\}$, their coefficients and the corresponding PDEs for 14,000, 16,000, 18,000 and 20,000 cells, respectively. (c)-(f) Simulated cell densities at different time instants based on the discovered PDEs for 14,000, 16,000, 18,000 and 20,000 cells, respectively, where the measurement at 0h is used as the initial condition while $\rho_x(x=0, t) = \rho_x(x=1900, t) = 0$ is employed as the Neumann boundary condition. The simulation result is represented by solid curves while the markers denote the measurement data.

strategy is proposed to simultaneously train the DNN and determine the optimal sparse coefficients of selected candidate terms for reconstructing the PDE(s). The synergy of interpretable DNN and sparse PDE representation results in the following outcome: the DNN provides accurate modeling of the solution and its derivatives as a basis for constructing the governing equation(s), while the sparsely represented PDE(s) in turn informs and constraints the DNN which makes it generalizable and further enhances the discovery. The overall approach is rooted in a comprehensive integration of bottom-up (data-driven) and top-down (physics-informed) processes for scientific discovery, with fusion of physics-informed deep learning, sparse regression and optimization. We demonstrate this method on a number of dynamical systems exhibiting nonlinear spatiotemporal behaviors (e.g., chaotic, shock, propagating front, etc.) governed by multi-dimensional PDEs based on either single or multiple datasets, numerically or experimentally. Results highlight that the approach is capable of accurately discovering the exact form of the governing equation(s), even in an information-poor space where the multi-dimensional measurements are scarce and noisy.

There still remain some potential limitations associated with the present PiDL framework for

physical law discovery. For example, although the fully connected DNN used in this work has advantage of analytical approximation of the PDE derivatives via automatic differentiation, directly applying it to model the solution of higher dimensional systems (such as long/short-term response evolution in a 3D domain) results in computational bottleneck and optimization challenges, e.g., due to the need for a vast number of collocation points to maintain satisfactory accuracy. Advances in discrete DNNs with spatiotemporal discretization (e.g., the convolutional long short-term memory network (ConvLSTM) [46] or similar) have the potential to help resolve this challenge, which will be demonstrated in our future work. In addition, the “root-branch” scheme might suffer from scalability issues when a large number of independent datasets sampled under various IBCs are available, resulting in many branches of the network for PDE solution approximation. The number of DNN trainable variables, the requirement of collocation points for retaining solution accuracy, and thus the computing memory, will grow in general linearly with the number of independent datasets (e.g., $\mathcal{O}(r)$). Nevertheless, this issue can be potentially well resolved by multi-GPU parallelization. Ideally, if the IBCs are known *a priori* and can be parameterized, a parametric DNN learning scheme could be developed into the proposed PiDL for parametric PDE solution approximation that accounts for different IBCs [40]. Several other aspects, such as the design of library of candidate functions and discovery with unknown source terms, are further discussed in [Supplementary Note C1, C3, C4](#).

METHOD

The innovations of this work are built upon seamless integration of the strengths of deep neural networks for rich representation learning, physics embedding, automatic differentiation and sparse regression to (1) approximate the solution of system variables, (2) compute essential derivatives, as well as (3) identify the key derivative terms and parameters that form the structure and explicit expression of the PDE(s). The resulting approach is able to deal with scarce/sparse and highly noisy measurement data while accounting for different initial/boundary conditions. The key method components are discussed below.

Network Architecture

The proposed network architectures of PiDL with sparse regression are shown in Figs. 1a and 1b that respectively deal with single-IBC dataset and multiple-IBC (r) independent datasets. The latent solution \mathbf{u} is interpreted by a dense (fully connected) DNN shown in Fig. 1a, namely, $\mathbf{u}^\theta = \mathbf{u}(\mathbf{x}, t; \theta)$, for the case of single dataset, while a “root-branch” dense DNN depicted in Fig. 1b is designed to approximate the latent solutions \mathbf{u}_i ($i = 1, \dots, r$) corresponding to different IBCs, viz., $\mathbf{u}_i^\theta = \mathbf{u}(\mathbf{x}, t; \theta^{(0)}, \theta^{(i)})$, for multiple independent datasets. Here, θ ’s denote the DNN trainable parameters. The DNNs take the spatiotemporal domain coordinates $\{\mathbf{x}, t\}$ as input followed by multiple fully-connected feedforward hidden layers (each layer has dozens of nodes). We use the hyperbolic tangent (tanh) or sine (sin) as the universal activation function thanks to their strength for high-order differentiation and unbiased estimation for both positive and negative values. The sin function is used when the system response exhibits periodic patterns. The output later is based on linear activation for universal magnitude mapping. When multiple datasets are available, e.g. sampled from different IBCs, domain coordinates are input to the “root” net (shared hidden layers), followed by r “branch” nets (individual hidden layers) that predict system response corresponding to each IBC/dataset. The “root” learns the common patterns across all datasets while the “branches” learn specific details determined by each IBC for each independent dataset. Such an architecture integrates information from different measurements at the expense of larger computational efforts

and produces solution approximations satisfying a unified physics (e.g., governing PDE(s)). The DNNs essentially play a role as a nonlinear functional to approximate the latent solution.

The DNN is connected to the physical law (reconstruction of PDE(s)) through a graph-based automatic differentiator where derivatives on \mathbf{u} 's are evaluated at machine precision. The library of candidate functions ϕ^θ can be computed from the DNNs. For the case of multiple independent datasets, the libraries $\phi^{(i)}$ resulted from the ‘‘branch’’ nets are concatenated to build one unified ϕ^θ . If there is unknown source input, the candidate functions for \mathbf{p} can also be incorporated into the library for discovery. The sparse representation of the reconstructed PDE(s) is then expressed in a residual form: $\mathcal{R}^\theta := \mathbf{u}_t^\theta - \phi^\theta \mathbf{\Lambda} \rightarrow \mathbf{0}$ s.t. $\mathbf{\Lambda} \in \mathcal{S}$, where $\mathcal{R}^\theta \in \mathbb{R}^{1 \times n}$ denotes the PDE residuals, \mathcal{S} represents the sparsity constraint set, and n is the dimension of the system variable (e.g., $\mathbf{u} \in \mathbb{R}^{1 \times n}$). Thus, the overall network architecture consists of heterogeneous trainable variables, namely, DNN parameters $\theta \in \mathbb{R}^{n_\theta \times 1}$ and PDE coefficients $\mathbf{\Lambda} \in \mathcal{S} \subset \mathbb{R}^{s \times n}$, where n_θ denotes the number of DNN trainable parameters and $n_\theta \gg sn$.

Physics-constrained Sparsity-regularized Loss Function

The physics-constrained sparsity-regularized loss function, expressed in Eq. (3), is composed of three components, the data loss \mathcal{L}_d , the residual physics loss \mathcal{L}_p and a sparsity regularization term imposed on $\mathbf{\Lambda}$. The data loss function reads

$$\mathcal{L}_d(\theta; \mathcal{D}_u) = \frac{1}{N_m} \|\mathbf{u}^\theta - \mathbf{u}^m\|_2^2 \quad (4)$$

where \mathbf{u}^m is the measurement data, \mathbf{u}^θ is the corresponding DNN-approximated solution, N_m is the total number of data points, and $\|\cdot\|_2$ denotes the Frobenius norm. The responses are stacked when multiple datasets are available, e.g., $\mathbf{u}^m = \{\mathbf{u}_1^m, \dots, \mathbf{u}_r^m\}$ and $\mathbf{u}^\theta = \{\mathbf{u}_1^\theta, \dots, \mathbf{u}_r^\theta\}$, where $r \geq 2$, as shown in Fig. 1b. The PDE residuals \mathcal{R}^θ are evaluated on a large number of randomly sampled collocation points \mathcal{D}_c , and used to form the residual physics loss function given by

$$\mathcal{L}_p(\theta, \mathbf{\Lambda}; \mathcal{D}_c) = \frac{1}{N_c} \|\dot{\mathbf{U}}(\theta) - \mathbf{\Phi}(\theta)\mathbf{\Lambda}\|_2^2 \quad (5)$$

where $\dot{\mathbf{U}}$ and $\mathbf{\Phi}$ denote respectively the discretization of the first-order time derivative term and the library of candidate functions evaluated on the collocation points; N_c is the total number of spatiotemporal collocation points. For the case of multiple datasets, $\dot{\mathbf{U}}$ and $\mathbf{\Phi}$ are concatenated over the index of different IBCs to ensure the identical physical law (in particular, the governing PDE(s)) is imposed, as depicted in Fig. 1b. Note that \mathcal{L}_d ensures that the DNN accurately interpret the latent solution of the PDE(s) via fitting the data, while \mathcal{L}_p generalizes and provides constraints for the DNN through reconstructing the closed form of the PDE(s). The ℓ_0 regularization term in Eq. (3) promotes the sparsity of the coefficients $\mathbf{\Lambda}$ for sparse representation of the PDE(s).

Alternating Direction Optimization

The total loss function in Eq. (3) has an implicit complex form, and thus, directly solving the optimization problem is highly intractable since the ℓ_0 regularization makes this problem np -hard. Though relaxation of the ℓ_0 term by the less rigorous ℓ_1 regularization improves the well-posedness and enables the optimization in a continuous space, false positive identification occurs [42, 43]. To address this challenge, we present an alternating direction optimization (ADO) algorithm that divides the overall optimization problem into a set of tractable subproblems to sequentially optimize

θ and Λ within a few alternating iterations (denoted by k), namely,

$$\Lambda_{k+1}^* := \arg \min_{\Lambda} \left[\left\| \dot{\mathbf{U}}(\theta_k^*) - \Phi(\theta_k^*)\Lambda \right\|_2^2 + \beta \|\Lambda\|_0 \right] \quad (6a)$$

$$\theta_{k+1}^* := \arg \min_{\theta} \left[\mathcal{L}_d(\theta; \mathcal{D}_u) + \alpha \mathcal{L}_p(\theta, \Lambda_{k+1}^*; \mathcal{D}_c) \right] \quad (6b)$$

The fundamental concept of the ADO algorithm shares similarity with the alternating direction methods of multipliers [47]. In each alternating iteration $k + 1$, the sparse PDE coefficients Λ in Eq. (6a) are updated (denoted by Λ_{k+1}^*) via STRidge (a sequential thresholding regression process that serves as a proxy for ℓ_0 regularization [5, 6]), based on the DNN parameters from the previous iteration (e.g., θ_k^*). The DNN parameters θ in the current iteration are then updated (denoted by θ_{k+1}^*) through a standard neural network training algorithm (in particular, the combined Adam [48] + L-BFGS [49] optimizer), taking Λ_{k+1}^* as known. The alternations between the sub-optimal solutions will lead to a high-quality optimization solution. It is noteworthy that the Adam optimizer plays a role for global search while the L-BFGS optimizer takes responsibility of fine tuning in a local solution region. The learning rate of Adam ranges from 10^{-5} to 10^{-3} in the test examples. The algorithm design of ADO, the choice of hyperparameters (e.g., the relative weighting of the loss functions, α and β), as well as the implementation details and specifications are given in [Supplementary Algorithm 1](#) and [Algorithm 2](#).

Pre-training of PiDL is conducted before running the ADO algorithm for discovery, by simply replacing $\|\Lambda\|_0$ in Eq. (3) with $\|\Lambda\|_1$ where brute-force gradient-based optimization (e.g., Adam + L-BFGS) for both θ and Λ becomes applicable, namely,

$$\{\theta^*, \Lambda^*\} = \arg \min_{\{\theta, \Lambda\}} \{ \mathcal{L}_d(\theta; \mathcal{D}_u) + \alpha \mathcal{L}_p(\theta, \Lambda; \mathcal{D}_c) + \beta \|\Lambda\|_1 \} \quad (7)$$

The ℓ_1 -regularized pre-training can accelerate the convergence of ADO by providing an admissible “initial guess”. Post-training (or post-tuning) is also applicable, which can be applied after the closed form of the PDE(s) is uncovered. This can be done by training the DNN along with the identification of the discovered non-zero coefficients, viz.,

$$\{\theta^*, \Lambda^*\} = \arg \min_{\{\theta, \Lambda\}} \{ \mathcal{L}_d(\theta; \mathcal{D}_u) + \alpha \mathcal{L}_p(\theta, \Lambda; \mathcal{D}_c) \} \quad (8)$$

where the initialization of the unknown parameters $\{\theta, \Lambda\}$ can be inherited from the ADO result. The post-training step is completely optional since the ADO method can already provides a high-quality solution as shown in the test examples. Nevertheless, the post-training could add additional discovery accuracy through fine tuning.

Data availability

All the used datasets in this study are available on GitHub at <https://github.com/isds-neu/EQDiscovery> upon final publication.

Code availability

All the source codes to reproduce the results in this study are available on GitHub at <https://github.com/isds-neu/EQDiscovery> upon final publication.

References

- [1] Josh Bongard and Hod Lipson. Automated reverse engineering of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*, 104(24):9943–9948, 2007.
- [2] Michael D. Schmidt and Hod Lipson. Distilling free-form natural laws from experimental data. *Science*, 324 5923:81–5, 2009.
- [3] Hayden Schaeffer, Russel Caffisch, Cory D. Hauck, and Stanley Osher. Sparse dynamics for partial differential equations. *Proceedings of the National Academy of Sciences*, 110(17):6634–6639, 2013.
- [4] Bryan C. Daniels and Ilya Nemenman. Automated adaptive inference of phenomenological dynamical models. *Nature Communications*, 6:8133, 2015.
- [5] Steven L. Brunton, Joshua L. Proctor, and J. Nathan Kutz. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*, 113(15):3932–3937, 2016.
- [6] Samuel H. Rudy, Steven L. Brunton, Joshua L. Proctor, and J. Nathan Kutz. Data-driven discovery of partial differential equations. *Science Advances*, 3(4):e1602614, 2017.
- [7] Hayden Schaeffer. Learning partial differential equations via data discovery and sparse optimization. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 473(2197):20160446, 2017.
- [8] Bethany Lusch, J Nathan Kutz, and Steven L Brunton. Deep learning for universal linear embeddings of nonlinear dynamics. *Nature Communications*, 9(1):1–10, 2018.
- [9] Z. Wang, X. Huan, and K. Garikipati. Variational system identification of the partial differential equations governing the physics of pattern-formation: Inference under varying fidelity and noise. *Computer Methods in Applied Mechanics and Engineering*, 356:44 – 74, 2019.
- [10] Kathleen Champion, Bethany Lusch, J. Nathan Kutz, and Steven L. Brunton. Data-driven discovery of coordinates and governing equations. *Proceedings of the National Academy of Sciences*, 116(45):22445–22451, 2019.
- [11] Niklas Pfister, Stefan Bauer, and Jonas Peters. Learning stable and predictive structures in kinetic systems. *Proceedings of the National Academy of Sciences*, 116(51):25405–25411, 2019.
- [12] Zhilong Huang, Yanping Tian, Chunjiang Li, Guang Lin, Lingling Wu, Yong Wang, and Hanqing Jiang. Data-driven automated discovery of variational laws hidden in physical systems. *Journal of the Mechanics and Physics of Solids*, 137:103871, 2020.
- [13] Jean-Christophe Loiseau and Steven L Brunton. Constrained sparse galerkin regression. *Journal of Fluid Mechanics*, 838:42–67, 2018.
- [14] Jean-Christophe Loiseau, Bernd R Noack, and Steven L Brunton. Sparse reduced-order modelling: sensor-based dynamics to full-state estimation. *Journal of Fluid Mechanics*, 844:459–490, 2018.
- [15] Zhilu Lai and Satish Nagarajaiah. Sparse structural system identification method for nonlinear dynamic systems with hysteresis/inelastic behavior. *Mechanical Systems and Signal Processing*, 117:813 – 842, 2019.

- [16] Shanwu Li, Eurika Kaiser, Shujin Laima, Hui Li, Steven L. Brunton, and J. Nathan Kutz. Discovering time-varying aerodynamics of a prototype bridge by sparse identification of nonlinear dynamical systems. *Physics Review E*, 100:022220, Aug 2019.
- [17] N. M. Mangan, S. L. Brunton, J. L. Proctor, and J. N. Kutz. Inferring biological networks by sparse identification of nonlinear dynamics. *IEEE Transactions on Molecular, Biological and Multi-Scale Communications*, 2(1):52–63, 2016.
- [18] Moritz Hoffmann, Christoph Fröhner, and Frank Noé. Reactive SINDy: Discovering governing reactions from concentration data. *The Journal of chemical physics*, 150(2):025101, 2019.
- [19] Bhavana Bhadriraju, Abhinav Narasingam, and Joseph Sang-II Kwon. Machine learning-based adaptive model identification of systems: Application to a chemical process. *Chemical Engineering Research and Design*, 152:372–383, 2019.
- [20] Frank Cichos, Kristian Gustavsson, Bernhard Mehlig, and Giovanni Volpe. Machine learning for active matter. *Nature Machine Intelligence*, 2(2):94–103, 2020.
- [21] E. Kaiser, J. N. Kutz, and S. L. Brunton. Sparse identification of nonlinear dynamics for model predictive control in the low-data limit. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 474(2219):20180335, 2018.
- [22] Kathleen P. Champion, Steven L. Brunton, and J. Nathan Kutz. Discovery of nonlinear multiscale systems: Sampling strategies and embeddings. *SIAM Journal on Applied Dynamical Systems*, 18(1):312–333, 2019.
- [23] Magnus Dam, Morten Brøns, Jens Juul Rasmussen, Volker Naulin, and Jan S. Hesthaven. Sparse identification of a predator-prey system from simulation data of a convection model. *Physics of Plasmas*, 24(2):022310, 2017.
- [24] Lorenzo Boninsegna, Feliks Nuske, and Cecilia Clementi. Sparse learning of stochastic dynamical equations. *The Journal of Chemical Physics*, 148(24):241723, 2018.
- [25] Kadierdan Kaheman, J Nathan Kutz, and Steven L Brunton. SINDy-PI: A robust algorithm for parallel implicit sparse identification of nonlinear dynamics. *arXiv preprint arXiv:2004.02322*, 2020.
- [26] Hayden Schaeffer, Giang Tran, and Rachel Ward. Extracting sparse high-dimensional dynamics from limited data. *SIAM Journal on Applied Mathematics*, 78(6):3279–3295, 2018.
- [27] Linan Zhang and Hayden Schaeffer. On the convergence of the SINDy algorithm. *Multiscale Modeling & Simulation*, 17(3):948–972, 2019.
- [28] Samuel Rudy, Alessandro Alla, Steven L. Brunton, and J. Nathan Kutz. Data-driven identification of parametric partial differential equations. *SIAM Journal on Applied Dynamical Systems*, 18(2):643–660, 2019.
- [29] Sheng Zhang and Guang Lin. Robust data-driven discovery of governing physical laws with error bars. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 474(2217):20180305, 2018.
- [30] Harsha Vaddireddy, Adil Rasheed, Anne E. Staples, and Omer San. Feature engineering and symbolic regression methods for detecting hidden physics from sparse sensor observation data. *Physics of Fluids*, 32(1):015113, 2020.

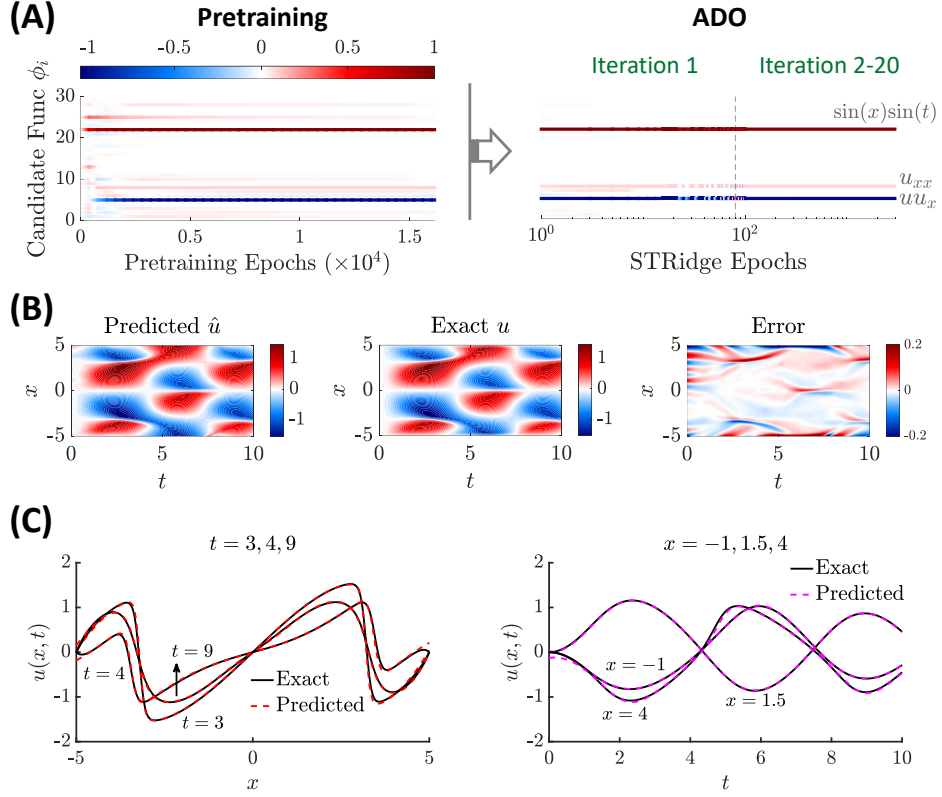
- [31] Jun Zhang and Wenjun Ma. Data-driven discovery of governing equations for fluid dynamics based on molecular simulation. *Journal of Fluid Mechanics*, 892:A5, 2020.
- [32] John H Lagergren, John T Nardini, G Michael Lavigne, Erica M Rutter, and Kevin B Flores. Learning partial differential equations for biological transport models from noisy spatio-temporal data. *Proceedings of the Royal Society A*, 476(2234):20190800, 2020.
- [33] Daniel R. Gurevich, Patrick A. K. Reinbold, and Roman O. Grigoriev. Robust and optimal sparse regression for nonlinear PDE models. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 29(10):103113, 2019.
- [34] Atılım Günes Baydin, Barak A Pearlmutter, Alexey Andreyevich Radul, and Jeffrey Mark Siskind. Automatic differentiation in machine learning: a survey. *The Journal of Machine Learning Research*, 18(1):5595–5637, 2017.
- [35] Jiequn Han, Arnulf Jentzen, and Weinan E. Solving high-dimensional partial differential equations using deep learning. *Proceedings of the National Academy of Sciences*, 115(34):8505–8510, 2018.
- [36] Justin Sirignano and Konstantinos Spiliopoulos. DGM: A deep learning algorithm for solving partial differential equations. *Journal of Computational Physics*, 375:1339–1364, 2018.
- [37] M. Raissi, P. Perdikaris, and G.E. Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- [38] Yibo Yang and Paris Perdikaris. Adversarial uncertainty quantification in physics-informed neural networks. *Journal of Computational Physics*, 394:136–152, 2019.
- [39] Yohai Bar-Sinai, Stephan Hoyer, Jason Hickey, and Michael P. Brenner. Learning data-driven discretizations for partial differential equations. *Proceedings of the National Academy of Sciences*, 116(31):15344–15349, 2019.
- [40] Luning Sun, Han Gao, Shaowu Pan, and Jian-Xun Wang. Surrogate modeling for fluid flows based on physics-constrained deep learning without simulation data. *Computer Methods in Applied Mechanics and Engineering*, 361:112732, 2020.
- [41] Maziar Raissi, Alireza Yazdani, and George Em Karniadakis. Hidden fluid mechanics: Learning velocity and pressure fields from flow visualizations. *Science*, 367(6481):1026–1030, 2020.
- [42] Jens Berg and Kaj Nyström. Data-driven discovery of pdes in complex datasets. *Journal of Computational Physics*, 384:239–252, 2019.
- [43] Gert-Jan Both, Subham Choudhury, Pierre Sens, and Remy Kusters. Deepmod: Deep learning for model discovery in noisy data. *Journal of Computational Physics*, page 109985, 2020.
- [44] Wang Jin, Esha T Shah, Catherine J Penington, Scott W McCue, Lisa K Chopin, and Matthew J Simpson. Reproducibility of scratch assays is affected by the initial degree of confluence: experiments, modelling and model selection. *Journal of Theoretical Biology*, 390:136–145, 2016.
- [45] Philip K Maini, DL Sean McElwain, and David I Leavesley. Traveling wave model to interpret a wound-healing cell migration assay for human peritoneal mesothelial cells. *Tissue Engineering*, 10(3-4):475–482, 2004.

- [46] Shi Xingjian, Zhouong Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *Advances in Neural Information Processing Systems*, pages 802–810, 2015.
- [47] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122, 2011.
- [48] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations (ICLR)*, 2015.
- [49] R. Byrd, P. Lu, J. Nocedal, and C. Zhu. A limited memory algorithm for bound constrained optimization. *SIAM Journal on Scientific Computing*, 16(5):1190–1208, 1995.

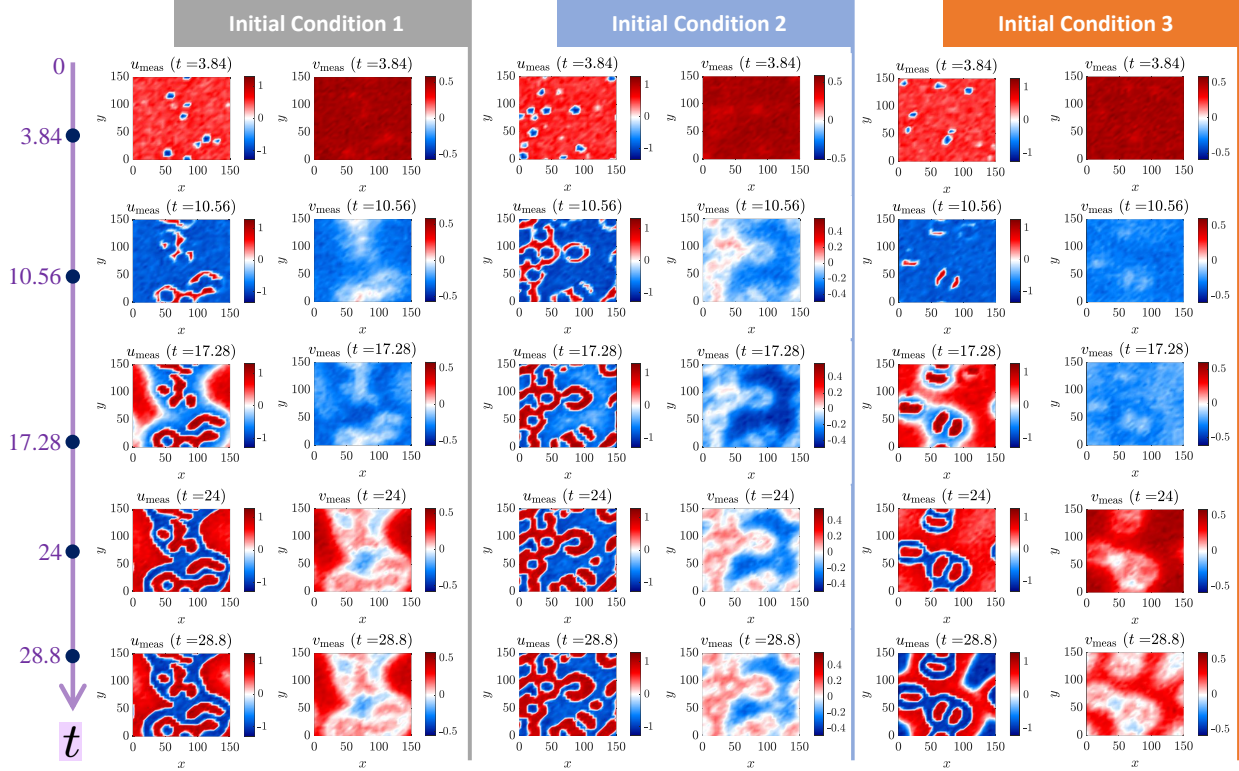
Acknowledgement: We acknowledge the support by the Engineering for Civil Infrastructure program at National Science Foundation under grant CMMI-2013067, the research award from MathWorks, and the Tier 1 Seed Grant Program at Northeastern University.

Competing interests: The authors declare no competing interests.

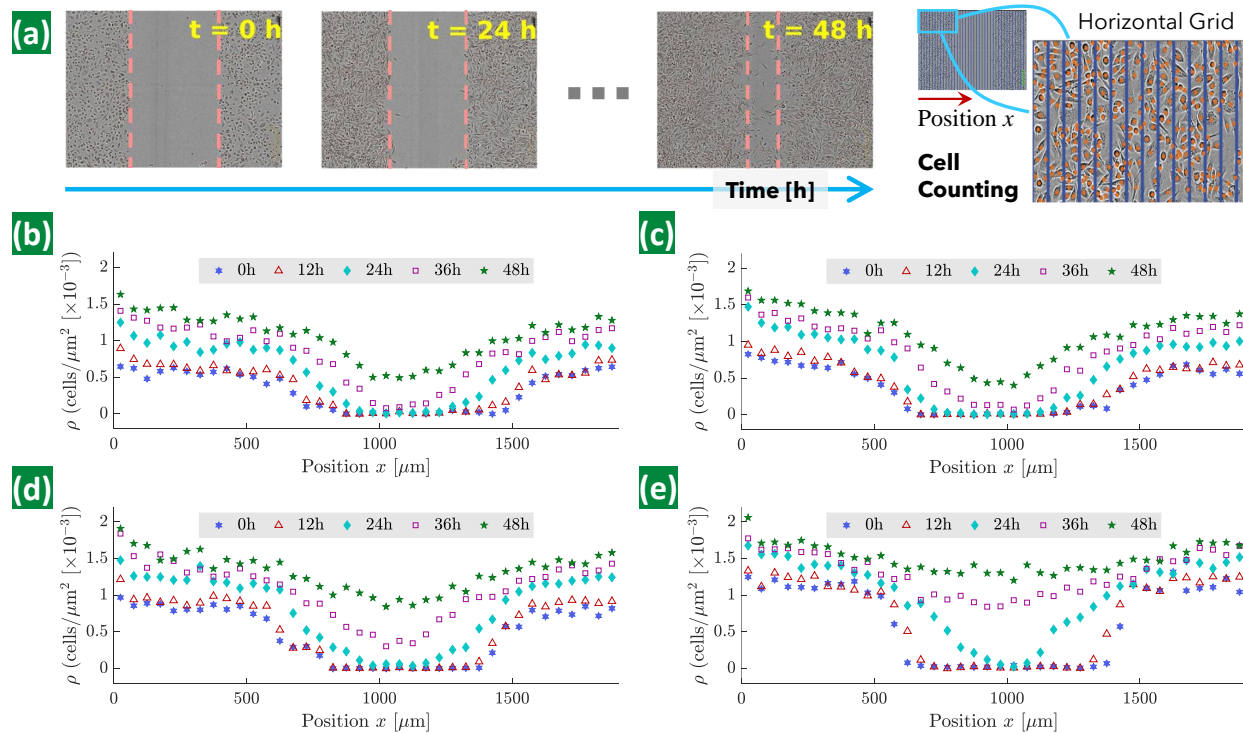
Supplementary information: The supplementary information is attached.



Extended Data Fig. 1: Discovered Burgers' equation and source term for measurement data with 10% noise. (a) Evolution of the sparse coefficients $\Lambda \in \mathbb{R}^{30 \times 1}$ for 30 candidate functions $\phi \in \mathbb{R}^{1 \times 30}$ used to form the PDE and the unknown source term, where the color represents the coefficient value. (b) The predicted response in comparison with the exact solution with the prediction error. The relative full-field ℓ_2 error of the prediction is 13.8%. The major errors are mostly distributed close to the boundaries due to scarce training data.



Extended Data Fig. 2: A few typical snapshots of low-resolution noisy measurements (10% noise) sampled from the system response under three different initial conditions (ICs) for discovering Fitzhugh-Nagumo equations. Note that the measurement data consists of 31 low-resolution noisy snapshots (with a grid size of 31×31) for each IC uniformly sampled within the time range of $[0, 28.8]$.



Extended Data Fig. 3: Measurement datasets of cell densities, ρ , based on scratch assays [44]. (a) Example scratch assay imaging of 16,000 cells in the test well with a width of $1,900 \mu\text{m}$ (the images are reproduced from Jin *et al.* [44]). The images are taken at different time instants (0h, 12h, 24h, 36h, 48h). The dashed lines show the approximate location of the positions of the leading edge. These images are then evenly divided into 38 segments ($50 \mu\text{m}$ each) horizontally, where the cells are counted in each segment to determine the horizontal cell densities. (b)-(e) the cell densities at different time instants for 14,000, 16,000, 18,000 and 20,000 cells, respectively.

Supplementary Information for:

Physics-informed learning of governing equations from scarce data

Zhao Chen¹, Yang Liu^{2,*}, and Hao Sun^{1,3,‡}

¹Department of Civil and Environmental Engineering, Northeastern University, Boston, MA 02115, USA

²Department of Mechanical and Industrial Engineering, Northeastern University, Boston, MA 02115, USA

³Department of Civil and Environmental Engineering, MIT, Cambridge, MA 02139, USA

*Corresponding author. E-mail: yang1.liu@northeastern.edu

‡Corresponding author. E-mail: h.sun@northeastern.edu

Contents

A Alternating Direction Optimization (ADO): Algorithm	1
B Examples	3
B.1 Discovery of Benchmark PDEs with Single Dataset	3
B.1.1 Burgers' equation	3
B.1.2 Kuramoto-Sivashinsky equation	4
B.1.3 Nonlinear Schrödinger equation	5
B.1.4 Navier-Stokes equation	7
B.1.5 λ - ω type Reaction-Diffusion equations	7
B.2 Comparison with SINDy	11
B.3 Discovery of PDEs with Multiple Independent Datasets	12
B.3.1 Burgers' equation with shock behavior	12
B.3.2 Fitzhugh-Nagumo type of Reaction-Diffusion equations	13
B.4 Experimental Discovery of Cell Migration and Proliferation	15
C Discussion	17
C.1 Selection of candidate functions	17
C.2 Noisy measurements and collocation points	19
C.3 Simultaneous identification of unknown source term	21
C.4 Other network architecture	22

This supplementary document provides a detailed description of the proposed algorithm, examples, and discussion of technical challenges for discovering closed-form partial differential equations (PDEs) from scarce and noisy data.

A Alternating Direction Optimization (ADO): Algorithm

The proposed ADO algorithm for training the network of PiDL with sparse regression is outlined in Algorithm 1, where the STRidge sub-function (a sequential thresholding regression process that serves as a proxy for ℓ_0 regularization [1, 2]) is given in Algorithm 2.

Algorithm 1 The proposed ADO for network training: $[\boldsymbol{\theta}_{\text{best}}, \mathbf{\Lambda}_{\text{best}}] = \text{ADO}(\mathcal{D}_u, \mathcal{D}_c, \Delta\delta, n_{\text{max}}, n_{\text{str}})$

- 1: **Input:** Measurement data \mathcal{D}_u , collocation points $\mathcal{D}_c = \{\mathbf{x}_i, t_i\}_{i=1,2,\dots,N_c}$, relative weighting of loss functions α and β , threshold tolerance increment $\Delta\delta$ for STRidge, maximum number of alternating iterations n_{max} , and maximum number of STRidge iterations n_{str} .
 # we take a 2D system in a 2D domain as an example: $\mathbf{u} = \{u, v\}$ and $\mathbf{x} = \{x, y\}$
 # α typically takes 1 (or 10 if the magnitude of \mathcal{L}_p is at least one order lower compared with the data loss)
 # β takes a small positive value, e.g., $10^{-6}\sigma_u$ where σ_u is the standard variation of the measurement data
 - 2: Split measurement data \mathcal{D}_u into training-validation sets ($n_{\text{tr}}/n_{\text{va}} = 80/20$): $\mathcal{D}_u^{\text{tr}} \in \mathbb{R}^{n_{\text{tr}} \times 2}$ and $\mathcal{D}_u^{\text{va}} \in \mathbb{R}^{n_{\text{va}} \times 2}$. #
 $N_m = n_{\text{tr}} + n_{\text{va}}$
 - 3: Split collocation points \mathcal{D}_c into training-validation sets ($m_{\text{tr}}/m_{\text{va}} = 80/20$): $\mathcal{D}_c^{\text{tr}} \in \mathbb{R}^{m_{\text{tr}} \times 3}$ and $\mathcal{D}_c^{\text{va}} \in \mathbb{R}^{m_{\text{va}} \times 3}$. #
 $N_c = m_{\text{tr}} + m_{\text{va}}$
 - 4: Initialize the *Tensor Graph* for the entire network.
 - 5: Pre-train the network via combined Adam and L-BFGS with $\{\mathcal{D}_u^{\text{tr}}, \mathcal{D}_c^{\text{tr}}\}$, and validate the trained model with $\{\mathcal{D}_u^{\text{va}}, \mathcal{D}_c^{\text{va}}\}$, namely,

$$\{\hat{\boldsymbol{\theta}}_0, \hat{\mathbf{\Lambda}}_0\} = \arg \min_{\{\boldsymbol{\theta}, \mathbf{\Lambda}\}} \{\mathcal{L}_d(\boldsymbol{\theta}; \mathcal{D}_u) + \alpha \mathcal{L}_p(\boldsymbol{\theta}, \mathbf{\Lambda}; \mathcal{D}_c) + \beta \|\mathbf{\Lambda}\|_1\}. \quad \# \text{ pre-train the network; } \hat{\mathbf{\Lambda}}_0 = \{\hat{\boldsymbol{\lambda}}_0^u, \hat{\boldsymbol{\lambda}}_0^v\}$$
 - 6: **for** $k = 1, 2, \dots, n_{\text{max}}$ **do**
 - 7: Assemble the system states over the collocation points $\mathcal{D}_c^{\text{tr}}$ and $\mathcal{D}_c^{\text{va}}$:

$$\begin{aligned} \dot{\mathbf{U}}_u^{\text{tr}} &= \bigcup_{i=1}^{N_c^{\text{tr}}} u_t(\hat{\boldsymbol{\theta}}_{k-1}; \mathbf{x}_i^{\text{tr}}, t_i^{\text{tr}}) \quad \text{and} \quad \dot{\mathbf{U}}_u^{\text{va}} = \bigcup_{i=1}^{N_c^{\text{tr}}} u_t(\hat{\boldsymbol{\theta}}_{k-1}; \mathbf{x}_i^{\text{va}}, t_i^{\text{va}}) \\ \dot{\mathbf{U}}_v^{\text{tr}} &= \bigcup_{i=1}^{N_c^{\text{va}}} v_t(\hat{\boldsymbol{\theta}}_{k-1}; \mathbf{x}_i^{\text{tr}}, t_i^{\text{tr}}) \quad \text{and} \quad \dot{\mathbf{U}}_v^{\text{va}} = \bigcup_{i=1}^{N_c^{\text{tr}}} v_t(\hat{\boldsymbol{\theta}}_{k-1}; \mathbf{x}_i^{\text{va}}, t_i^{\text{va}}). \end{aligned}$$
 - 8: Assemble the candidate library matrices over the collocation points \mathcal{D}_c , $\mathcal{D}_c^{\text{tr}}$ and $\mathcal{D}_c^{\text{va}}$:

$$\tilde{\boldsymbol{\Phi}} = \bigcup_{i=1}^{N_c} \boldsymbol{\phi}(\hat{\boldsymbol{\theta}}_{k-1}; \mathbf{x}_i, t_i), \quad \tilde{\boldsymbol{\Phi}}^{\text{tr}} = \bigcup_{i=1}^{N_c^{\text{tr}}} \boldsymbol{\phi}(\hat{\boldsymbol{\theta}}_{k-1}; \mathbf{x}_i^{\text{tr}}, t_i^{\text{tr}}) \quad \text{and} \quad \tilde{\boldsymbol{\Phi}}^{\text{va}} = \bigcup_{i=1}^{N_c^{\text{va}}} \boldsymbol{\phi}(\hat{\boldsymbol{\theta}}_{k-1}; \mathbf{x}_i^{\text{va}}, t_i^{\text{va}}).$$
 - 9: Normalize candidate library matrices $\tilde{\boldsymbol{\Phi}}$, $\tilde{\boldsymbol{\Phi}}^{\text{tr}}$ and $\tilde{\boldsymbol{\Phi}}^{\text{va}}$ column-wisely ($j = 1, \dots, s$) to improve matrix condition:

$$\boldsymbol{\Phi}_{:,j} = \tilde{\boldsymbol{\Phi}}_{:,j} / \|\tilde{\boldsymbol{\Phi}}_{:,j}\|_2, \quad \boldsymbol{\Phi}_{:,j}^{\text{tr}} = \tilde{\boldsymbol{\Phi}}_{:,j}^{\text{tr}} / \|\tilde{\boldsymbol{\Phi}}_{:,j}^{\text{tr}}\|_2 \quad \text{and} \quad \boldsymbol{\Phi}_{:,j}^{\text{va}} = \tilde{\boldsymbol{\Phi}}_{:,j}^{\text{va}} / \|\tilde{\boldsymbol{\Phi}}_{:,j}^{\text{va}}\|_2.$$
 - 10: Determine ℓ_0 regularization parameter $\gamma = 0.001\kappa(\boldsymbol{\Phi})$. # $\kappa(\cdot)$ denotes the matrix condition number;
 - 11: Initialize the error indices: $\hat{\epsilon}^u = \|\boldsymbol{\Phi}^{\text{va}} \hat{\boldsymbol{\lambda}}_{k-1}^u - \dot{\mathbf{U}}_u^{\text{va}}\|_2^2 + \gamma \|\hat{\boldsymbol{\lambda}}_{k-1}^u\|_0$ and $\hat{\epsilon}^v = \|\boldsymbol{\Phi}^{\text{va}} \hat{\boldsymbol{\lambda}}_{k-1}^v - \dot{\mathbf{U}}_v^{\text{va}}\|_2^2 + \gamma \|\hat{\boldsymbol{\lambda}}_{k-1}^v\|_0$.
 - 12: Set the initial threshold tolerance $\delta_1 = \Delta\delta$.
 - 13: **for** $iter = 1, 2, \dots, n_{\text{str}}$ **do**
 - 14: Run STRidge as shown in Algorithm 2 to determine:

$$\tilde{\boldsymbol{\lambda}}^u = \text{STRidge}(\dot{\mathbf{U}}_u^{\text{tr}}, \boldsymbol{\Phi}^{\text{tr}}, \delta_{iter}) \quad \text{and} \quad \tilde{\boldsymbol{\lambda}}^v = \text{STRidge}(\dot{\mathbf{U}}_v^{\text{tr}}, \boldsymbol{\Phi}^{\text{tr}}, \delta_{iter}).$$
 - 15: Update the error indices: $\epsilon^u = \|\boldsymbol{\Phi}^{\text{va}} \tilde{\boldsymbol{\lambda}}^u - \dot{\mathbf{U}}_u^{\text{va}}\|_2^2 + \gamma \|\tilde{\boldsymbol{\lambda}}^u\|_0$ and $\epsilon^v = \|\boldsymbol{\Phi}^{\text{va}} \tilde{\boldsymbol{\lambda}}^v - \dot{\mathbf{U}}_v^{\text{va}}\|_2^2 + \gamma \|\tilde{\boldsymbol{\lambda}}^v\|_0$.
 - 16: **if** $\epsilon^u \leq \hat{\epsilon}^u$ or $\epsilon^v \leq \hat{\epsilon}^v$ (run in parallel) **then**
 - 17: Increase threshold tolerance with increment: $\delta_{iter+1} = \delta_{iter} + \Delta\delta$.
 - 18: **else**
 - 19: Decrease threshold tolerance increment $\Delta\delta = \Delta\delta/1.618$.
 - 20: Update threshold tolerance with the new increment $\delta_{iter+1} = \max\{\delta_{iter} - 2\Delta\delta, 0\} + \Delta\delta$.
 - 21: **end if**
 - 22: **end for**
 - 23: Return and re-scale the current best solution from STRidge iterations: $\hat{\mathbf{\Lambda}}_k = \{\tilde{\boldsymbol{\lambda}}^u, \tilde{\boldsymbol{\lambda}}^v\}$. # re-scaling due to normalization of $\boldsymbol{\Phi}$
 - 24: Train the DNN via combined Adam and L-BFGS with $\{\mathcal{D}_u^{\text{tr}}, \mathcal{D}_c^{\text{tr}}\}$, and validate the trained model with $\{\mathcal{D}_u^{\text{va}}, \mathcal{D}_c^{\text{va}}\}$, namely,

$$\hat{\boldsymbol{\theta}}_k = \arg \min_{\boldsymbol{\theta}} \{\mathcal{L}_d(\boldsymbol{\theta}; \mathcal{D}_u) + \alpha \mathcal{L}_p(\boldsymbol{\theta}, \hat{\mathbf{\Lambda}}_k; \mathcal{D}_c)\}. \quad \# \text{ train DNN given } \hat{\mathbf{\Lambda}}_k \text{ as known}$$
 - 25: **end for**
 - 26: **Output:** the best solution $\boldsymbol{\theta}_{\text{best}} = \hat{\boldsymbol{\theta}}_{n_{\text{max}}}$ and $\mathbf{\Lambda}_{\text{best}} = \hat{\mathbf{\Lambda}}_{n_{\text{max}}}$
-

Algorithm 2 Sequential threshold ridge regression (STRidge): $\hat{\lambda} = \text{STRidge}(\dot{\mathbf{U}}, \Phi, \delta)$

1: **Input:** Time derivative vector $\dot{\mathbf{U}}$, candidate function library matrix Φ , and threshold tolerance δ .

2: Inherit coefficients $\hat{\lambda}$ from the DNN pre-training or the previous update.

3: **repeat**

4: Determine indices of coefficients in $\hat{\lambda}$ falling below or above the sparsity threshold δ :

$$\mathcal{I} = \{i \in \mathcal{I} : |\hat{\lambda}_i| < \delta\} \text{ and } \mathcal{J} = \{j \in \mathcal{J} : |\hat{\lambda}_j| \geq \delta\}.$$

5: Enforce sparsity to small values by setting them to zero: $\hat{\lambda}_{\mathcal{I}} = \mathbf{0}$.

6: Update remaining non-zero values with ridge regression:

$$\hat{\lambda}_{\mathcal{J}} = \arg \min_{\lambda_{\mathcal{J}}} \{\|\Phi_{\mathcal{J}} \lambda_{\mathcal{J}} - \dot{\mathbf{U}}\|_2^2 + 1 \times 10^{-5} \|\lambda_{\mathcal{J}}\|_2^2\}. \quad \# \text{ the parameter } 1 \times 10^{-5} \text{ is tunable}$$

7: **until** maximum number of iterations reached.

8: **Output:** The best solution $\hat{\lambda} = \hat{\lambda}_{\mathcal{I}} \cup \hat{\lambda}_{\mathcal{J}}$

B Examples

We observe the efficacy and robustness of our methodology on a group of canonical PDEs used to represent a wide range of physical systems with nonlinear, periodic and/or chaotic behaviors. In particular, we discover the closed forms of Burgers', Kuramoto-Sivashinsky (KS), nonlinear Schrödinger, Navier-Stokes (NS), and λ - ω Reaction-Diffusion (RD) equations from scarce and noisy time-series measurements recorded by a number of sensors at fixed locations from a single IBC. Gaussian white noise is added to the synthetic response with the noise level defined as the root-mean-square ratio between the noise and the exact solution. To demonstrate the “root-branch” network for discovery of PDE(s) based on multiple independent datasets sampled under different IBCs, we consider (1) the 1D Burgers' equation with light viscosity that exhibits a shock behavior, and (2) a 2D Fitzhugh-Nagumo (FN) type reaction-diffusion system that describes activator-inhibitor neuron activities excited by external stimulus. At last, we test our framework on the experimental data of cell migration and proliferation. Our method is also compared with SINDy (the PDE-FIND approach presented in [2]) which is also presented herein. The identification error is defined as the average relative error of the identified non-zero PDE coefficients with respect to the ground truth, which is used to evaluate the accuracy of the discovered PDEs for the following examples. Simulations in this paper are performed on a workstation with 28 Intel Core i9-7940X CPUs and 2 NVIDIA GTX 1080Ti GPU cards.

B.1 Discovery of Benchmark PDEs with Single Dataset

B.1.1 Burgers' equation

We first consider a dissipative system with the dynamics governed by a 1D viscous Burgers' equation expressed as

$$u_t = -uu_x + \nu u_{xx}$$

where u is a field variable, x and t are the spatial and temporal coordinates, and ν denotes the diffusion coefficient. The equation describes the decaying stationary viscous shock of a system after a finite period of time, commonly found in simplified fluid mechanics, nonlinear acoustics, gas dynamics and traffic flow. In this work, solution for the Burgers' Equation is from an open dataset [2], in which the diffusion coefficient ν is assumed to be 0.1 and u is discretized into 256 spatial grid points for 101 time steps with a Gaussian initial condition. In particular, 5 sensors are randomly placed at fixed locations among the 256 spatial grid points to record the wave response for 101 time steps, leading to 1.95% of the dataset used in [2]. A total number of 1.6×10^5 collocation points,

e.g., in the pair of $\{x, t\}$, are sampled by the Sobol sequence [3]. A group of 16 candidate functions ($\phi \in \mathbb{R}^{1 \times 16}$) are used to reconstruct the PDE, consisting of polynomial terms (u, u^2, u^3), derivatives (u_x, u_{xx}, u_{xxx}) and their multiplications. The fully connected DNN has 8 hidden layers and a width of 20 neuron nodes in each layer. The training efforts are performed via 1×10^3 epochs of L-BFGS for pre-training and 20 ADO iterations. In each ADO iteration, we use the combination of 100 epochs of Adam and 1×10^3 (or less, depending on the relative loss decay) epochs of L-BFGS to train the DNN for alternation with STRidge. The discovered equation for the dataset with 10 % noise reads:

$$u_t = -0.997uu_x + 0.098u_{xx}$$

where the aggregated relative identification error for all non-zero elements in $\mathbf{\Lambda}$ is $1.15 \pm 1.20\%$. The discovery result is summarized in Fig. S.1. The evolution of the PDE coefficients for all candidate functions is illustrated in Fig. S.1A. While all coefficients are initialized as zeros, they vary evidently in pre-training process. In the ADO stage, the adaptive sparsity threshold in STRidge gradually prunes redundant components right after the first alternating iteration. Despite that only 1.95% subsampled responses are measured, the PiDL approach can accurately extrapolate the full-field solution with a ℓ_2 error of 2.02% (see Fig. S.1B). Fig. S.1C shows the comparison of spatial and temporal snapshots between the predicted and the exact solutions which agree extremely well.

B.1.2 Kuramoto-Sivashinsky equation

Another dissipative system with intrinsic instabilities is considered, governed by the 1D Kuramoto-Sivashinsky (KS) equation:

$$u_t = -uu_x - u_{xx} - u_{xxxx}$$

where the reverse diffusion term $-u_{xx}$ leads to the blowup behavior while the fourth-order derivative u_{xxxx} introduces chaotic patterns as shown in Fig. S.2B, making an ideal test problem for equation discovery. Starting with a smooth initial condition, the KS system evolves to an unstable laminar status due to the highly nonlinear terms including the high-order derivative. The KS equation is widely used to model the instabilities in laminar flame fronts and dissipative trapped-ion modes among others. We subsample the open dataset [2] by randomly choosing 320 points from the 1024 spatial grid nodes as fixed sensors and record the wave response for 101 time steps, occupying about 12.3% of the original dataset. A set of 2×10^4 collocation points, sampled using the Sobol sequence in the spatiotemporal domain, are employed to evaluate the residual physics loss. A library of 36 candidate functions are used to construct the PDE, consisting of polynomials (u, u^2, u^3, u^4, u^5), derivatives ($u_x, u_{xx}, u_{xxx}, u_{xxxx}, u_{xxxxx}$) and their multiplications. The DNN architecture and hyperparameters are same as those in the Burgers' equation example.

It is notable that the chaotic behavior poses significant challenges in approximating the full-field spatiotemporal derivatives, especially the high-order u_{xxxx} , from poorly measured data for discovery of such a PDE. Existing methods (for example the family of SINDy methods [2, 4]) eventually fail in this case given very coarse and noisy measurements. Nevertheless, the PiDL approach successfully distills the closed form of the KS equation from subsampled sparse data even with 10% noise:

$$u_t = -0.992uu_x - 0.990u_{xx} - 0.991u_{xxxx}$$

where the coefficients have an average relative error, for all non-zero elements in $\mathbf{\Lambda}$, of $0.71 \pm 0.06\%$. The evolution of the coefficients $\mathbf{\Lambda} \in \mathbb{R}^{36 \times 1}$ in Fig. S.2A illustrates that both the candidate terms and the corresponding coefficients are correctly identified (close to the original parameters) within a small number of ADO iterations. Although the available measurement data are sparsely sampled in the spatiotemporal domain under a high-level noise corruption, the predicted full-field wave by

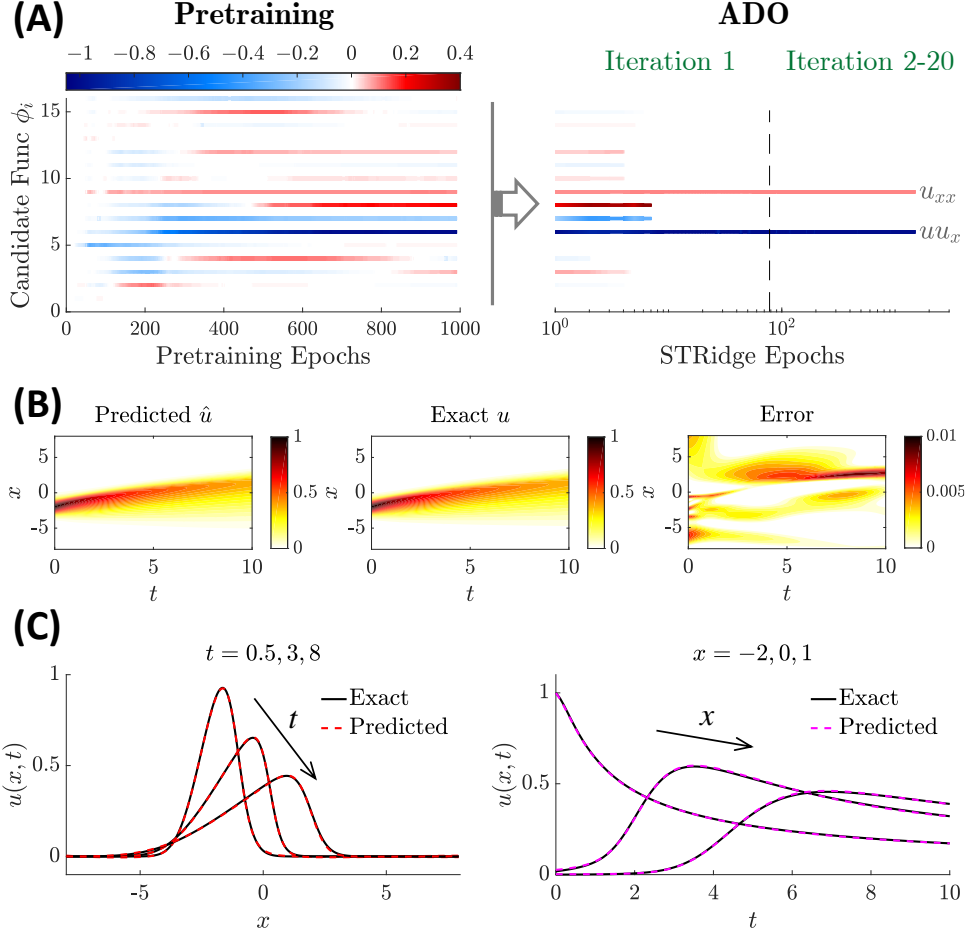


Fig. S.1: Discovered Burgers' equation for data with 10% noise. (A) Evolution of the sparse coefficients $\Lambda \in \mathbb{R}^{16 \times 1}$ for 16 candidate functions $\phi \in \mathbb{R}^{1 \times 16}$ used to form the PDE, where the color represents the coefficient value. (B) The predicted response in comparison with the exact solution with the prediction error. (C) Comparison of spatial and temporal snapshots between the predicted and the exact solutions. The relative full-field ℓ_2 error of the prediction is 2.02%.

the trained PiDL also agrees well with the exact solution with a relative ℓ_2 error of 1.87% (Fig. S.2B). The spatial and temporal snapshots of the predicted response match seamlessly the ground truth as shown in Fig. S.2C.

B.1.3 Nonlinear Schrödinger equation

In the third example, we discover the nonlinear Schrödinger equation, originated as a classical wave equation, given by

$$iu_t = -0.5u_{xx} - |u|^2u$$

where u is a complex field variable. This well-known equation is widely used in modeling the propagation of light in nonlinear optical fibers, Bose-Einstein condensates, Langmuir waves in hot plasmas, and so on. The solution to this Schrödinger equation is simulated based on a Gaussian initial condition with the problem domain meshed into 512 spatial points and 501 temporal steps, while the measurements are taken from 256 randomly chosen spatial “sensors” for 375 time instants, resulting in 37.5% data used in [2] for uncovering the closed form of the equation. A library of

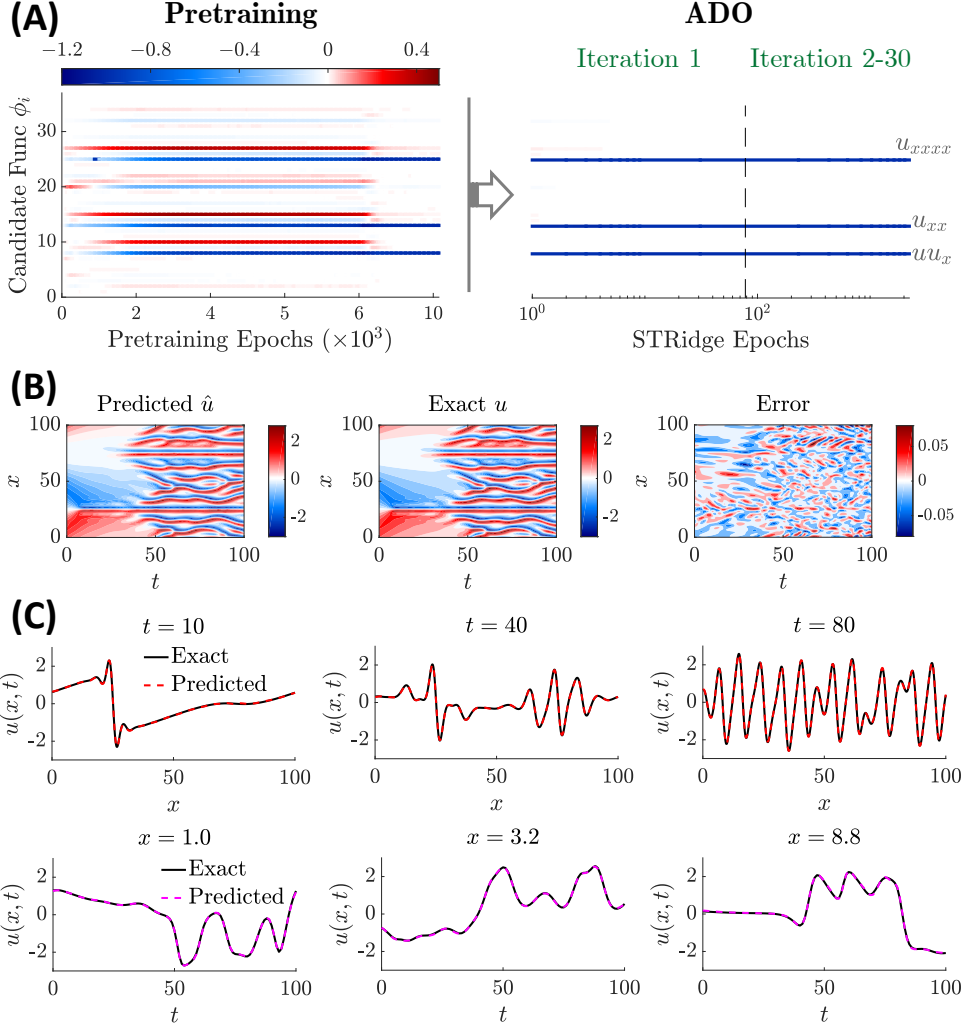


Fig. S.2: Discovered the KS equation for data with 10% noise. (A) Evolution of the sparse coefficients $\Lambda \in \mathbb{R}^{36 \times 1}$ for 36 candidate functions $\phi \in \mathbb{R}^{1 \times 36}$ used to reconstruct the PDE. (B) The predicted response compared with the exact solution. (C) Comparison of spatial and temporal snapshots between the predicted and the exact solutions. The relative full-field ℓ_2 error of the prediction is 1.87%.

40 candidate functions are used for constructing the PDE, varying among polynomial functions (u, u^2, u^3), absolute values ($|u|, |u|^2, |u|^3$), derivatives (u_x, u_{xx}, u_{xxx}) and their combination. Since the function is complex valued, we model separately the real part (u_R) and the imaginary part (u_I) of the solution in the output of the DNN, assemble them to obtain the complex solution $u = u_R + iu_I$, and construct the complex valued candidate functions for discovery. To avoid complex gradients in optimization, we use the modulus (magnitude, $|u|$), instead of the ℓ_2 norm, for the residual physics loss \mathcal{L}_p . The fully connected DNN has 8 hidden layers and a width of 40 neuron nodes in each layer. The pre-training takes 1.6×10^5 epochs of Adam (with additional L-BFGS tuning) followed by 30 ADO iterations. In each ADO iteration, we use 1×10^3 Adam epochs and up to 4×10^3 (depending on the relative loss decay) L-BFGS epochs to train the DNN for alternation with STRidge.

The discovered equation under 10 % noise is written as

$$iu_t = -0.490u_{xx} - 0.974|u|^2u$$

where the average relative error for non-zero coefficients is $2.31 \pm 0.28\%$. The evolution history of

the sparse coefficients $\mathbf{\Lambda} \in \mathbb{R}^{40 \times 1}$ clearly shows the convergence to the actual values (Fig. S.3A) resulting in accurate closed-form identification of the PDE. Even though the evolution of $\mathbf{\Lambda}$ is quite intense in the pre-training stage, the most dominant components remain after the first few ADO iterations. The predicted full-field response, for both real and imaginary parts, matches well the exact solution with a slight relative ℓ_2 error of about 1% (Fig. S.3B and C). The comparison of spatiotemporal snapshots between the predicted and the exact solutions for the real part (Fig. S.3D) and imaginary part (Fig. S.3E) also shows almost perfect agreement.

B.1.4 Navier-Stokes equation

We consider a 2D fluid flow passing a circular cylinder with the local rotation dynamics (see Fig. S.4). For incompressible and isotropic fluids which also have conservative body forces, the well-known Navier-Stokes vorticity equation reads

$$w_t = -uw_x - vw_y + 0.01w_{xx} + 0.01w_{yy}$$

where w is the spatiotemporally variant vorticity, $\mathbf{u} = \{u, v\}$ denotes the fluid velocities at Reynolds number 100, ∇ is the gradient, and ∇^2 is the Laplace operator. The full-field solution to the NS vorticity equation is obtained using the immersed boundary projection method [5]. The dimensionless domain is discretized into a 499×199 spatial grid and 151 time steps. The cylinder has a unit diameter and the input flow from the left side has a unit velocity. Measurements of velocities $\{u, v\}$ and vorticity w are taken at 500 random spatial locations lasting 60 time steps in the boxed area behind the cylinder as shown in Fig. S.4, namely 0.22% subsamples from the total dataset and 1/10 of the data used in [2]. The residual physics loss is evaluated on 6×10^4 collocation points randomly sampled in the spatiotemporal domain using the Sobol sequence [3]. The library of candidate functions consists of 60 components including polynomial terms ($u, v, w, uv, uw, vw, u^2, v^2, w^2$), derivatives ($w_x, w_y, w_{xx}, w_{xy}, w_{yy}$) and their combination. The latent output in the DNN contains u, v and w . The DNN has 8 fully connected hidden layers and a width of 60 nodes in each layer. The pre-training takes 5×10^3 epochs of Adam (with additional L-BFGS tuning up to 1×10^4 epochs) followed by 6 ADO iterations. In each ADO iteration, we use the Adam optimizer with 500 epochs and the L-BFGS with up to 1×10^3 epochs to train the DNN for each alternation within STRidge.

The discovered NS vorticity equation for the case of 10% noise is given as follows

$$w_t = -0.999uw_x - 0.994vw_y + 0.010w_{xx} + 0.010w_{yy}$$

where the aggregated relative identification error for all non-zero elements in $\mathbf{\Lambda}$ is $1.40 \pm 1.83\%$. It is encouraging that the uncovered vorticity equation is almost identical to the ground truth, for both the derivative terms and their coefficients, even under 10% noise corruption. The coefficients $\mathbf{\Lambda} \in \mathbb{R}^{60 \times 1}$, corresponding to 60 candidate functions $\phi \in \mathbb{R}^{1 \times 60}$, converge very quickly to the correct values with precise sparsity right after the first ADO iteration (Fig. S.5A). The vorticity patterns and magnitudes are also well predicted as indicated by multiple spatial snapshots at different time instants ($t = 0.2, 7.6, 15, 22.4, 29.8$) shown in Fig. S.5B in comparison with the exact solution (Fig. S.5C, with small errors as depicted in Fig. S.5D). Note that the response in these snapshots is not used in training the net work. The ℓ_2 error of the predicted full-field vorticity response is about 2.57%. This example provides a compelling test case for the proposed PiDL approach which is capable of discovering the closed-form NS equation with scarce and noisy data.

B.1.5 λ - ω type Reaction-Diffusion equations

The examples discussed previously are low-dimensional (1D) models with limited complexity. We herein consider a λ - ω reaction-diffusion (RD) system in a 2D domain with the pattern forming

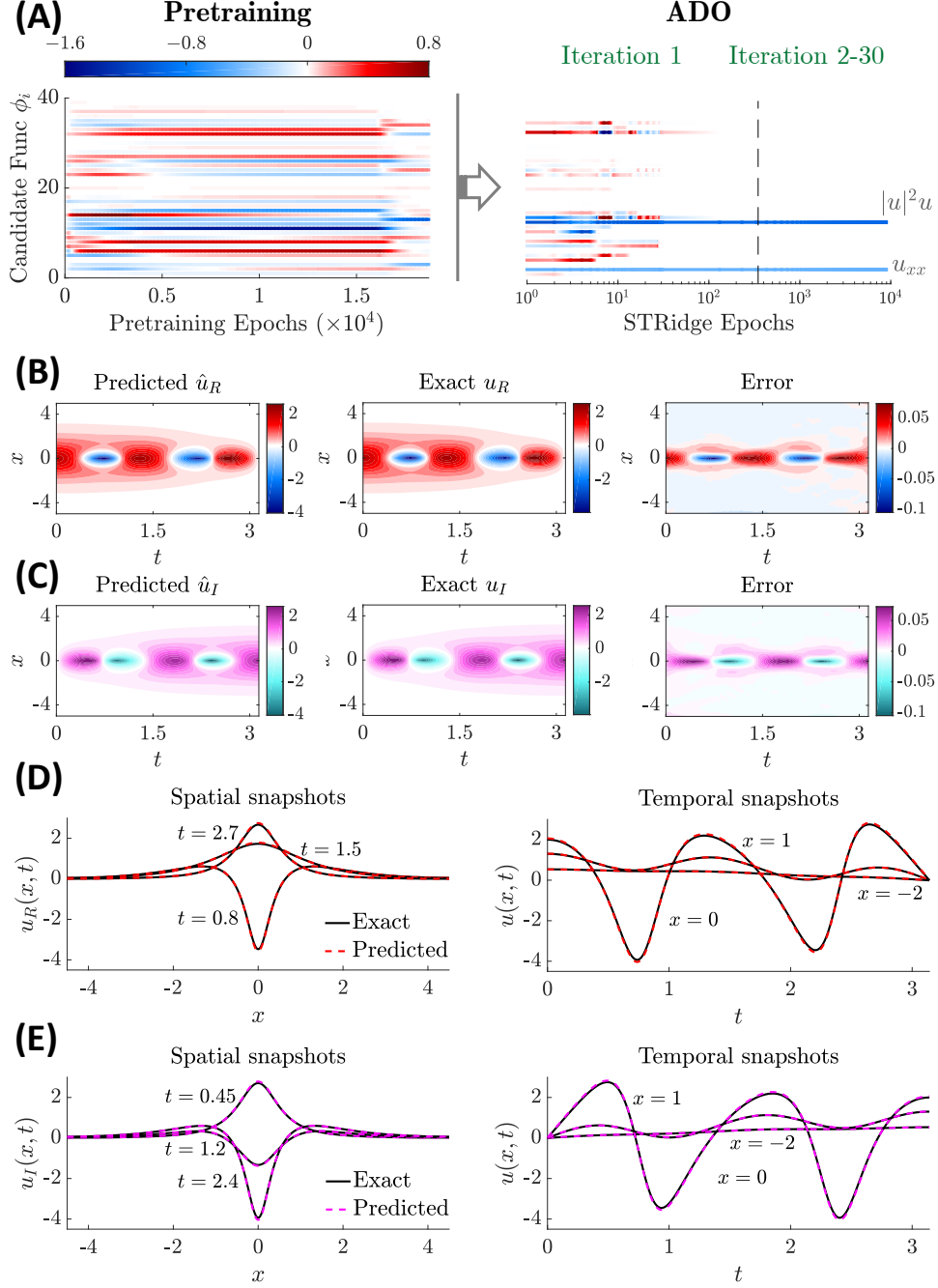


Fig. S.3: Discovered nonlinear Schrödinger equation for a dataset with 10% noise. (A) Evolution of the sparse coefficients $\Lambda \in \mathbb{R}^{40 \times 1}$ for the candidate functions $\phi \in \mathbb{R}^{1 \times 40}$ used to reconstruct the PDE. (B and C) The predicted real-part (B) and imaginary-part (C) responses compared with the exact solution. (D and E) Comparison of spatial and temporal snapshots between the predicted and the exact solutions for the real part (D) and imaginary part (E). The relative full-field ℓ_2 error of the prediction is about 1%.

behavior governed by two coupled PDEs:

$$\begin{cases} u_t = 0.1u_{xx} + 0.1u_{yy} - uv^2 - u^3 + v^3 + u^2v + u \\ v_t = 0.1v_{xx} + 0.1v_{yy} - uv^2 - u^3 - v^3 - u^2v + v \end{cases}$$

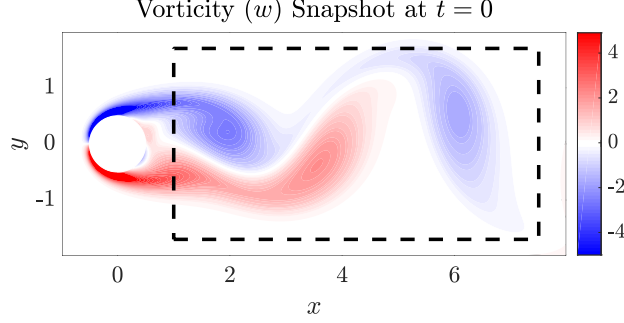


Fig. S.4: Vorticity field $w(\mathbf{x}, t)$ at $t = 0$ for a steady flow passing a cylinder. Measurements are sampled from the the boxed area surrounded by the dashed line.

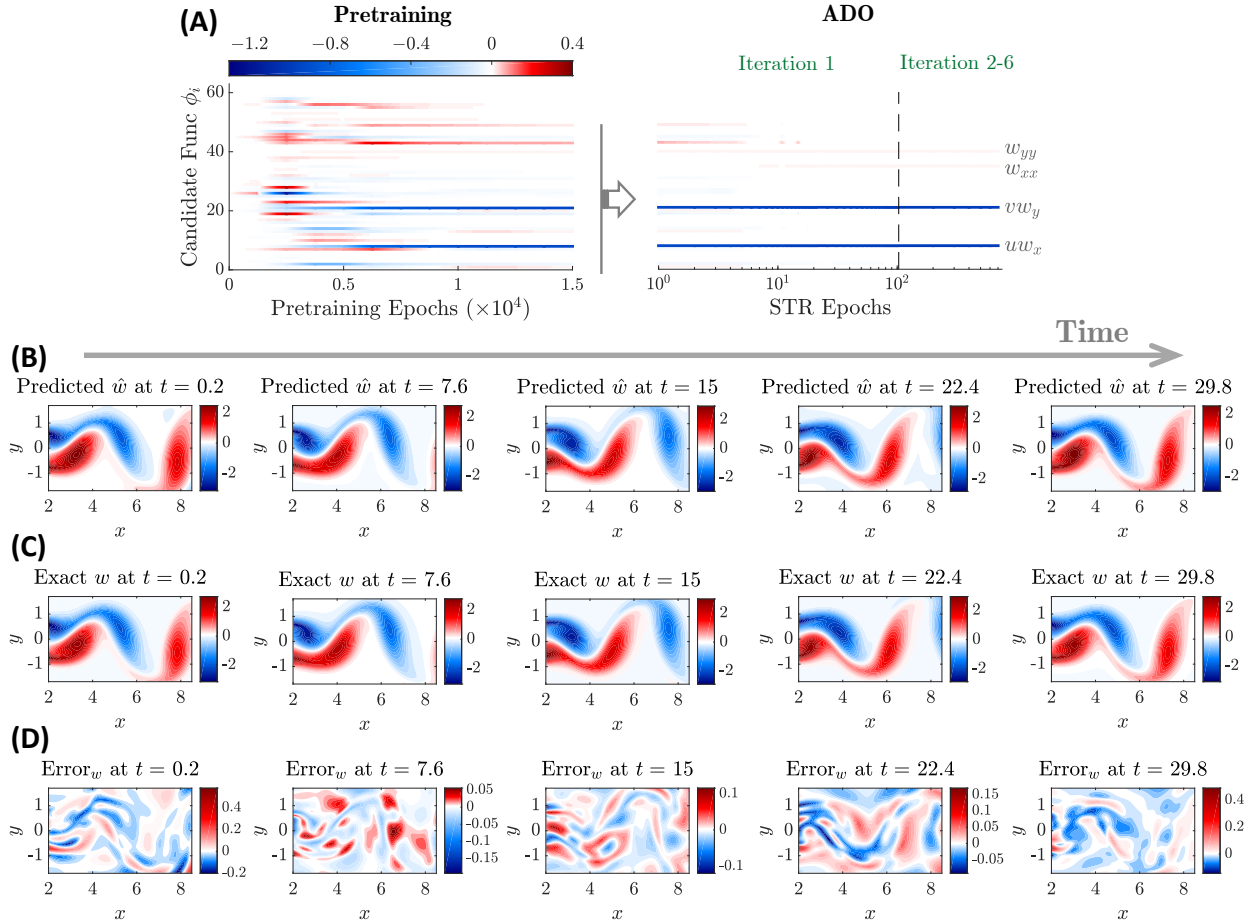


Fig. S.5: Discovered NS equation for data with 10% noise. (A) Evolution of the sparse coefficients $\Lambda \in \mathbb{R}^{60 \times 1}$ for 60 candidate functions $\phi \in \mathbb{R}^{1 \times 60}$ used to form the vorticity equation. (B-D) Vorticity snapshots at different time instants ($t = 0.2, 7.6, 15, 22.4, 29.8$) for the prediction (B), the exact solution (C) and the prediction error (D). Note that response at these time instants are not included in dataset for training the PiDL model and equation discovery. The relative full-field ℓ_2 error of the prediction is about 2.57%.

where u and v are two field variables. The λ - ω equations are typically used to describe the multi-scale phenomenon of local reactive transformation and the global diffusion in chemical reactions, with wide applications in pattern formation [6], biological morphogenesis [7], and ecological inva-

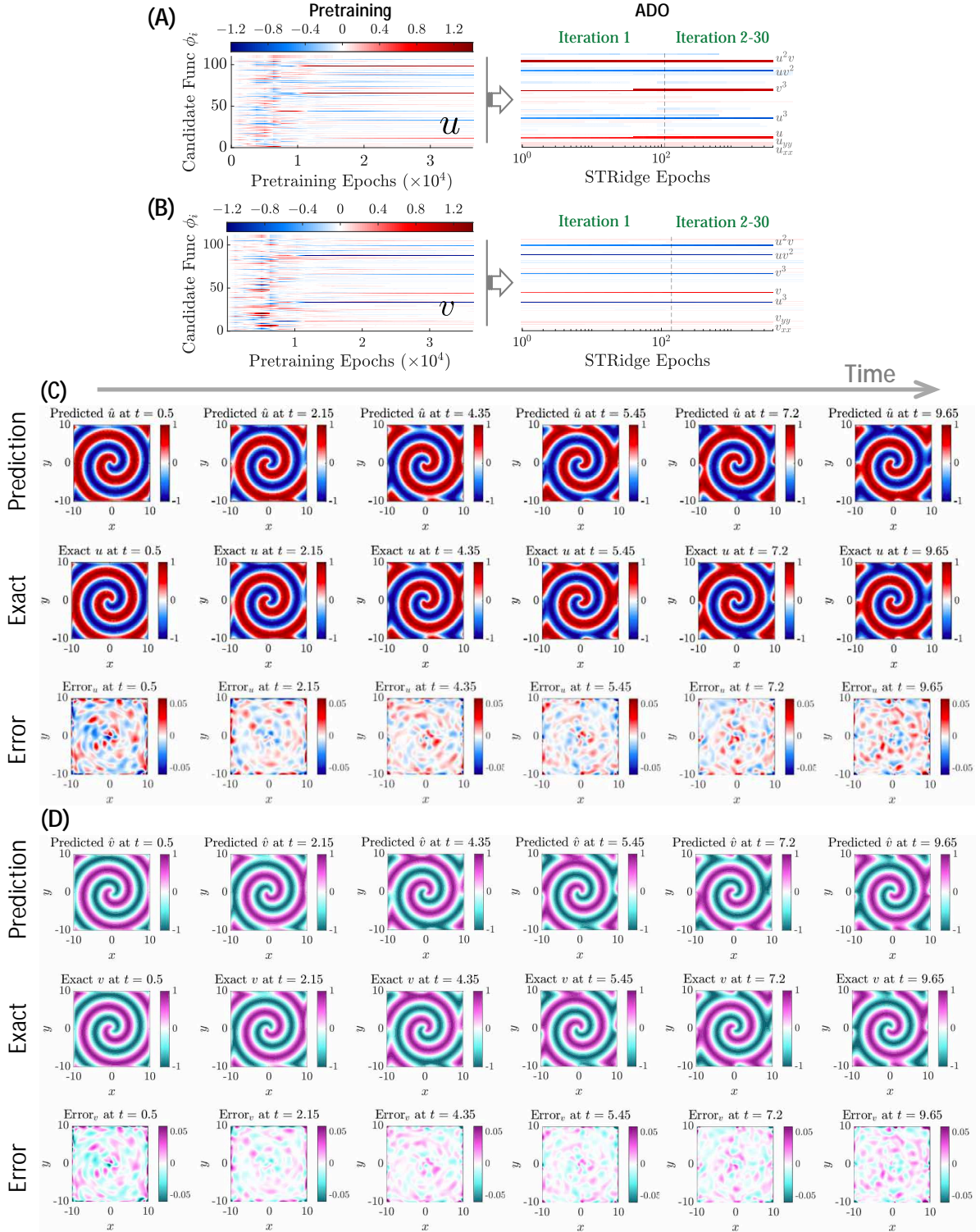


Fig. S.6: Discovered λ - ω equations for a dataset with 10% noise. (A and B) Evolution of the sparse coefficients $\lambda^u \in \mathbb{R}^{110 \times 1}$ (A) and $\lambda^v \in \mathbb{R}^{110 \times 1}$ (B) for 110 candidate functions $\phi \in \mathbb{R}^{1 \times 110}$ used to reconstruct the u -equation and the v -equation, respectively. (C and D) The response snapshots u (C) and v (D) at different time instants ($t = 0.5, 2.15, 4.35, 5.45, 7.2, 9.65$), showing the predictions and the exact solutions, as well as the prediction error maps. The relative full-field ℓ_2 error of the prediction is about 2.21%.

sions [8], among others. The λ - ω equations exhibit a wide range of behaviors including wave-like phenomena and self-organized patterns found in chemical and biological systems. The binomial system is also called an activator-inhibitor system because one state variable encourages the increase of both states while the other state component inhibits their growth. The particular λ - ω equations in this test example display spiral waves subjected to periodic boundary conditions. The domain for generating the solution is divided into 65,536 (256×256) spatial points with 201 time steps. We take randomly 2,500 spatial points as fixed sensors recording the wave response for 15 randomly sampled time steps, leading to 1/4 of the subsampled dataset used in [2] and 0.29% of the total data. We sample 1×10^5 collocation points using the Sobol sequence [3] to evaluate the residual physics loss. A total of 110 candidate functions are employed, including polynomials up to the 3rd order ($u, v, u^2, v^2, uv, u^3, u^2v, uv^2, v^3$), derivatives up to the 2nd order ($u_x, u_y, v_x, v_y, u_{xx}, u_{xy}, u_{yy}, v_{xx}, v_{xy}, v_{yy}$) and their combination, for the sparse discovery of the two PDEs. Since the system dimension is relatively high, we enhance the discovery by post-training (post-tuning) of the DNN and the uncovered non-zero PDE coefficients, after the ADO stage, resulting in refined/improved discovery. The DNN has 8 fully connected hidden layers and a width of 60 nodes in each layer. The pre-training takes 1×10^4 epochs of Adam (with additional L-BFGS tuning up to 1×10^4 epochs) followed by 30 ADO iterations. In each ADO iteration, we use 1×10^3 Adam epochs and up to 4×10^3 (depending on the relative loss decay) L-BFGS epochs to train the DNN for each alternation within STRidge.

The reconstructed equations for the case of 10% noise are given by

$$\begin{cases} u_t = 0.091u_{xx} + 0.092u_{yy} - 0.907uv^2 - 0.918u^3 + 0.997v^3 + 0.997u^2v + 0.917u \\ v_t = 0.099v_{xx} + 0.099v_{yy} - 1.006uv^2 - 1.003u^3 - 0.930v^3 - 0.935u^2v + 0.934u \end{cases}$$

where the the average relative error for all non-zero coefficients is $4.78 \pm 3.66\%$. The evolution process illustrates that the sparse patterns are iteratively recovered out of a mixture of more than 100 candidate functions. At the end, both the sparse terms and the associated coefficients are precisely identified (as depicted in Fig. S.6A and B). Due to the complexity of the PDEs and the high dimension, slightly more epochs are required in ADO to retain reliable convergence. The predicted response snapshots by the trained PiDL at different time instants, e.g., $t = \{0.5, 2.15, 4.35, 5.45, 7.2, 9.65\}$, are shown in in Fig. S.6C and D, which are very close to the ground truth (the errors are distributed within a small range). This example shows especially the great ability and robustness of our method for discovering governing PDEs for high-dimensional systems from highly noisy data. The relative full-field ℓ_2 error of the prediction is about 2.21%.

B.2 Comparison with SINDy

We have performed the comparison study between the proposed PiDL approach and the state-of-the-art PDE-FIND method (an extended version of SINDy) [2], in the context of different levels of data size and noise. We test the five PDEs described previously and summarize the discovery errors for the sparse coefficients in Table S.1. The error is defined as the average relative error of the identified non-zero PDE coefficients with respect to the ground truth. If the terms in the PDEs are discovered incorrectly, we mark it as “NA” (not applicable). It is seen from Table S.1 that the proposed PiDL approach is capable of correctly uncovering the closed-form PDEs for all cases, regardless of the varying levels of data size and noise, which demonstrates excellent robustness. Although PDE-FIND shows great success in PDE discovery with negligible error for large and clean (or approximately noise-free) measurement data, this method eventually fails when the level of data scarcity and/or noise increases. In general, PDE-FIND relies on the strict requirement of measurement quality and quantity. However, PiDL is able to alleviate and resolve this limitation thanks to

Table S.1: Summary of the PiDL discovery results in comparison with PDE-FIND [2] for canonical models.

PDE name	Method	Error (noise 0%)	Error (noise 1%)	Error (noise 10%)	# of Measurement points
Burgers'	PiDL	0.01±0.01%	0.19±0.11%	1.15±1.20%	~505
	PDE-FIND	NA	NA	NA	~505
		0.15±0.06%	0.80±0.60%	NA	~26K
KS	PiDL	0.07±0.01%	0.61±0.04%	0.71±0.06%	~32K
	PDE-FIND	35.75±16.30%	NA	NA	~32K
		1.30±1.30%	52.00±1.40%	NA	~257K
Schrödinger	PiDL	0.09±0.04%	0.65±0.29%	2.31±0.28%	~96K
	PDE-FIND	NA	NA	NA	~96K
		0.05±0.01%	3.00±1.00%	NA	~257K
NS	PiDL	0.66±0.72%	0.86±0.63%	1.40±1.83%	~30K
	PDE-FIND	NA	NA	NA	~30K
		1.00±0.20%	7.00±6.00%	NA	~300K
λ - ω RD	PiDL	0.07±0.08%	0.25±0.30%	4.78±3.66%	~37.5K
	PDE-FIND	NA	NA	NA	~37.5K
		0.02±0.02%	NA	NA	~150K

Note: In the table, KS, NS and RD refer to the Kuramoto-Sivashinsky, Navier-Stokes and the λ - ω Reaction-Diffusion PDEs. Gaussian white noise is added to the synthetic response with the noise level defined as the root-mean-square ratio between the noise and exact solution. NA denotes “not applicable” (e.g., failure in correct identification of the sparse PDE coefficients). The identification error is defined as the average relative error of the identified non-zero PDE coefficients with respect to the ground truth.

the combination of the strengths of DNNs for rich representation learning of nonlinear functions, automatic differentiation for accurate derivative calculation as well as ℓ_0 sparse regression. In addition, the use of collocation points introduces additional “pseudo datasets”, compensates indirectly the scarcity of measurement data, and enriches the constraint for constructing the closed form of PDEs. Nonetheless, we have to mention that the proposal PiDL approach is much more computationally costly compared to PDE-PIND, primarily due to the training of DNNs. Fortunately, this issue can be well managed through parallel computing on a powerful GPU platform and remains a less important concern compared to the aim for successful discovery of correct underlying PDEs.

B.3 Discovery of PDEs with Multiple Independent Datasets

B.3.1 Burgers' equation with shock behavior

We consider to discover the previously discussed Burgers' equation (see Section B.1.1) with a small diffusion/viscosity parameter, expressed as

$$u_t = -uu_x + \frac{0.01}{\pi}u_{xx}$$

based on datasets generated by imposing three different IBCs. The small diffusion coefficient $0.01/\pi \approx 0.0032$ creates shock formation in a compact area with sharp gradient and poses notorious difficulty for many numerical methods to resolve, which could challenge the DNN's ap-

proximation ability and thus affect the discovery. The three IBCs used for data generation include:

$$\text{IBC 1: } u(x, 0) = -\sin(\pi x), u(-1, t) = u(1, t) = 0$$

$$\text{IBC 2: } u(x, 0) = \mathcal{G}(x), u(-1, t) = u(1, t) = 0$$

$$\text{IBC 3: } u(x, 0) = -x^3, u(-1, t) = 1, u(1, t) = -1$$

where \mathcal{G} denotes a Gaussian function. The ground truth solution is simulated by MATLAB function `pdede` in a spatiotemporal domain $\Omega \times [0, T] = [-1, 1]_{d=200} \times [0, 1]_{d=1000}$. For all IBCs, we assume that there are 30 sensors randomly deployed in space measuring the wave traveling (e.g., u) for 500 time instants (7.5% of the total grid points). A denser sensor grid is needed herein, compared with the previous Burgers' example, in order to capture the shock behaviors. All measurements are polluted with 10% Gaussian noise. The noisy measurements are depicted in Fig. S.7A for the three datasets. For visualization purpose, we only draw a handful of signals out of a total of 30 time series for each IBC.

We design a ‘‘root-branch’’ DNN: the root takes the spatiotemporal coordinates $\{x, t\}$ as input followed by 4 hidden layers of 20 nodes, while each of the three branches is separately connected to the last hidden layer of the root followed by 4 hidden layers of 30 nodes before the output layer. The motivation for this design is that the branch nets can capture the solution difference due to different IBCs while the shared root net learns the common response patterns that obey the unique Burgers' equation. Note that although we have three distinctive solution approximations, we stack them into one candidate library followed by a unified form of PDE. Therefore, we can combine the information from three datasets to discover one physics equation. A group of 4.5×10^4 collocation points are generated by the Latin hypercube sampling strategy [9] for determining the residual physics loss. The PDE library consists of 16 candidate functions, exactly the same as the Burgers' case in Section B.1.1. The training efforts include 3×10^4 epochs pretraining by L-BFGS followed by 6 ADO iterations. In each ADO iteration, we use 1×10^3 Adam epochs (with an initial learning rate of 1e-4 and drops by 50% every 1×10^3 Adam epochs) and up to 3×10^4 L-BFGS epochs (depending on relative loss decay) to train the DNN in synergy with STRidge. The discovered PDE is given by

$$u_t = -1.006uu_x + 0.0039u_{xx}$$

which shows great agreement with the ground truth. Fig. S.7B depicts the evolution of the coefficients ($\mathbf{\Lambda} \in \mathbb{R}^{16 \times 1}$) of candidate functions, where the correct terms in the library (uu_x and u_{xx}) are successfully distilled while other redundant terms are eliminated (e.g., hardly thresholded to zero) by ADO. The coefficients of the active terms are accurately identified as well (in particular the small viscosity parameter that leads to shock formation, e.g., 0.0039). While the ℓ_1 penalty lays the foundation for sparsity in pretraining stage, while ADO finds the correct sparsity pattern in after the first iteration and refines the coefficients in the subsequent iterations. Fig. ??c-d show the predicted responses and errors for three IBC cases, with a stacked full-field ℓ_2 error of 2.24%. The trained ‘‘root-branch’’ network can accurately reproduce distinctive system responses even with limited measurements under 10% noise, giving a full-field ℓ_2 error of 2.24%, as shown in Fig. S.7C-E.

B.3.2 Fitzhugh-Nagumo type of Reaction-Diffusion equations

We consider the Fitzhugh-Nagumo (FN) type reaction-diffusion system, in a 2D domain $\Omega = [0, 150] \times [0, 150]$ with periodic boundary conditions, whose governing equations are expressed by two coupled PDEs [10, 11]:

$$\begin{aligned} u_t &= \gamma_u \Delta u + u - u^3 - v + \alpha \\ v_t &= \gamma_v \Delta v + \beta(u - v) \end{aligned}$$

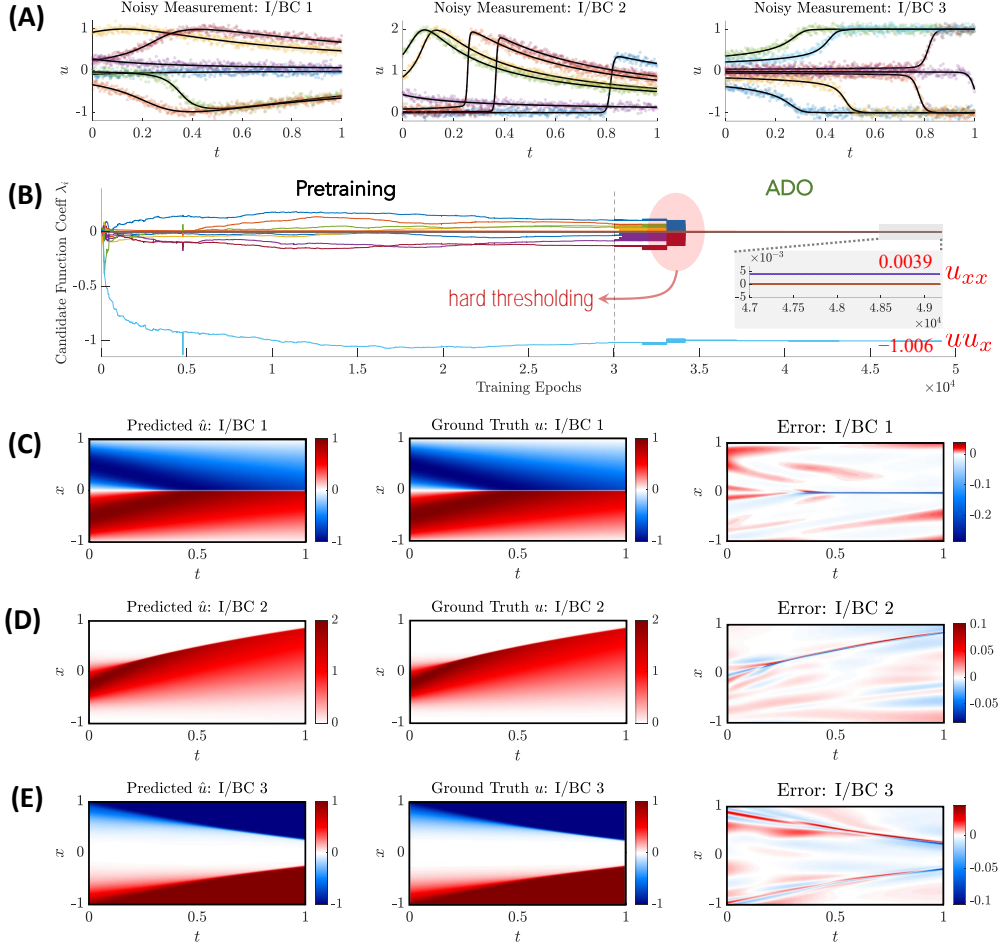


Fig. S.7: Discovered Burgers’ equation with small viscosity based on datasets sampled under three IBCs with 10% noise. (A) Visualization of noisy measurements for the three datasets. Note that there are 30 sensors and only a few are illustrated in this figure. (B) Evolution of the sparse coefficients $\Lambda \in \mathbb{R}^{16 \times 1}$ for 16 candidate functions $\phi \in \mathbb{R}^{1 \times 16}$ used to construct the PDE, where the color represents the coefficient value. The correct terms (uu_x and u_{xx}) and their coefficients are successfully identified while other redundant terms are eliminated by ADO. (C-E) The predicted response in comparison with the exact solution for three IBCs. The relative full-field ℓ_2 error of all the stacked predictions is 2.24%.

where u and v represent two interactive components/matters (e.g., biological), $\gamma_u = 1$ and $\gamma_v = 100$ are diffusion coefficients, $\alpha = 0.01$ and $\beta = 0.25$ are the coefficients for reaction terms, and Δ is the Laplacian operator. The FN equations are commonly used to describe biological neuron activities excited by external stimulus (α), which exhibit an activator-inhibitor system because one equation boosts the production of both components while the other equation dissipates their new growth. The ground truth data is generated by the finite difference method ($dx = dy = 0.5$ and $dt = 0.0002$) for the time period of $[0, T] = [0, 36]$, with three random fields as initial conditions. Three measurement datasets are then generated, each of which consists of 31 low-resolution snapshots (projected into a 31×31 grid) uniformly down-sampled from full-field synthetic data during the period of $[7.18, 36]$ under a 10% noise condition. Similar to the previous example in Section B.3.1, we design a “root-branch” DNN with three branches: the root net has 2 hidden layers of 60 nodes while each of the three branch nets has 3 hidden layers of 60 nodes. We sample 5×10^4 spatiotemporal collocation points using Latin hypercube sampling [9] to construct the physics residuals.

We assume the diffusion terms (Δu and Δv) are known in the PDEs, whose coefficients (γ_u and γ_v) yet need to be identified. We employ the bounds to these two positive coefficients to speed up the convergence, namely, $\gamma_u \in [0, 5]$ and $\gamma_v \in [0, 150]$. We design 70 candidate functions, composed of up to third-order polynomials (including the constant term “1” as the zero order), derivatives $\{u_x, u_y, u_{xy}, v_x, v_y, v_{xy}\}$ and their mutual multiplication, to reconstruct the nonlinear reaction terms in the PDEs. Hence, the final library has 72 candidate terms. To account for the small stimulus term (e.g., 0.01 in the first equation), we increase the sensitivity of the constant candidate “1” in the library by down-scaling its magnitude to the order of 10^{-5} and 5×10^{-4} for u and v equations respectively. The training efforts include the pretraining stage with 6×10^3 Adam epochs and 4×10^4 L-BFGS epochs, 10 ADO iterations, and an extra post-training with 1×10^5 Adam epochs. In each ADO iteration, we use 1×10^4 Adam epochs and up to 1×10^4 L-BFGS epochs in synergy with STRidge. To deal with the aforementioned bounds for γ_u and γ_v in an unconstrained optimization process, we set $\gamma_u = 5\sigma(\tilde{\gamma}_u)$ and $\gamma_v = 150\sigma(\tilde{\gamma}_v)$ and take $\{\tilde{\gamma}_u, \tilde{\gamma}_v\}$ as trainable variables, where $\sigma(\cdot)$ denotes the Sigmoid function. The discovered equations under 10 % noise is

$$\begin{aligned} u_t &= 0.962\Delta u + 0.874u - 0.847u^3 - 0.931v + 0.0098 \\ v_t &= 71.515\Delta v + 0.214u - 0.224v \end{aligned}$$

It is seen that the form of the PDEs is precisely uncovered with all correct active terms (including the unknown external stimulus in the first equation). The corresponding identified coefficients are generally close to the ground truth (error of non-zero coefficients: $11.72 \pm 8.34\%$) except the diffusion coefficient for v (i.e., γ_v) which seems to be a less sensitive parameter according to our test. It should be noted that, given very scarce and noisy measurement datasets in this example, the “root-branch” DNN is faced with challenges to accurately model the solutions with sharp propagating fronts (see Fig. S.8C-D). The less accurate solution approximation by DNN then affects the discovery precision. This issue can be naturally alleviated by increasing the spatiotemporal measurement resolution (even still under fairly large noise pollution, e.g., 10%). Nevertheless, the exact form of the PDEs is successfully discovered in this challenging example, which is deemed more important since the coefficients can be further tuned/calibrated when additional data arrives. The evolution of the PDE coefficients corresponding to 72 candidate functions for \hat{u} and \hat{v} is illustrated in Fig. S.8A and B, respectively. Note that, for visualization purpose, we re-scale the identified coefficients of the constant stimulus term “1” in the u -equation by multiplying 100 in Fig. S.8A and the diffusion term Δv in the v -equation by dividing 50 in Fig. S.8B. The trained network is finally used to predict the full-field responses under three IBCs (see the snapshots in Fig. S.8C-D at two unmeasured time instants). The stacked full-field ℓ_2 error is 5.04%.

B.4 Experimental Discovery of Cell Migration and Proliferation

In this example, we consider to discover a biological system based on scratch assay experiments [12] investigating the cell migration and proliferation process. The 1D cell density distributions at different time instants (0h, 12h, 24h, 36h, 48h) were extracted and simplified from high-resolution imaging via image segmentation and cell counting (see Extended Data Fig. 3 in the Main Text). A series of assays were performed under different initial cell densities (e.g., the total number of cells spans from 10,000 to 20,000 following the designated initial distribution in the test well. More detailed description of the experiment setup and datasets can be found in [12].

Our objective herein is to uncover a parsimonious PDE for modeling the dynamics of cell density $\rho(x, t)$. Here, we consider four scenarios with the initial number of cells ranging from 14,000, 16,000, 18,000 to 20,000. We take the mean of the test data from three identically-prepared experimental

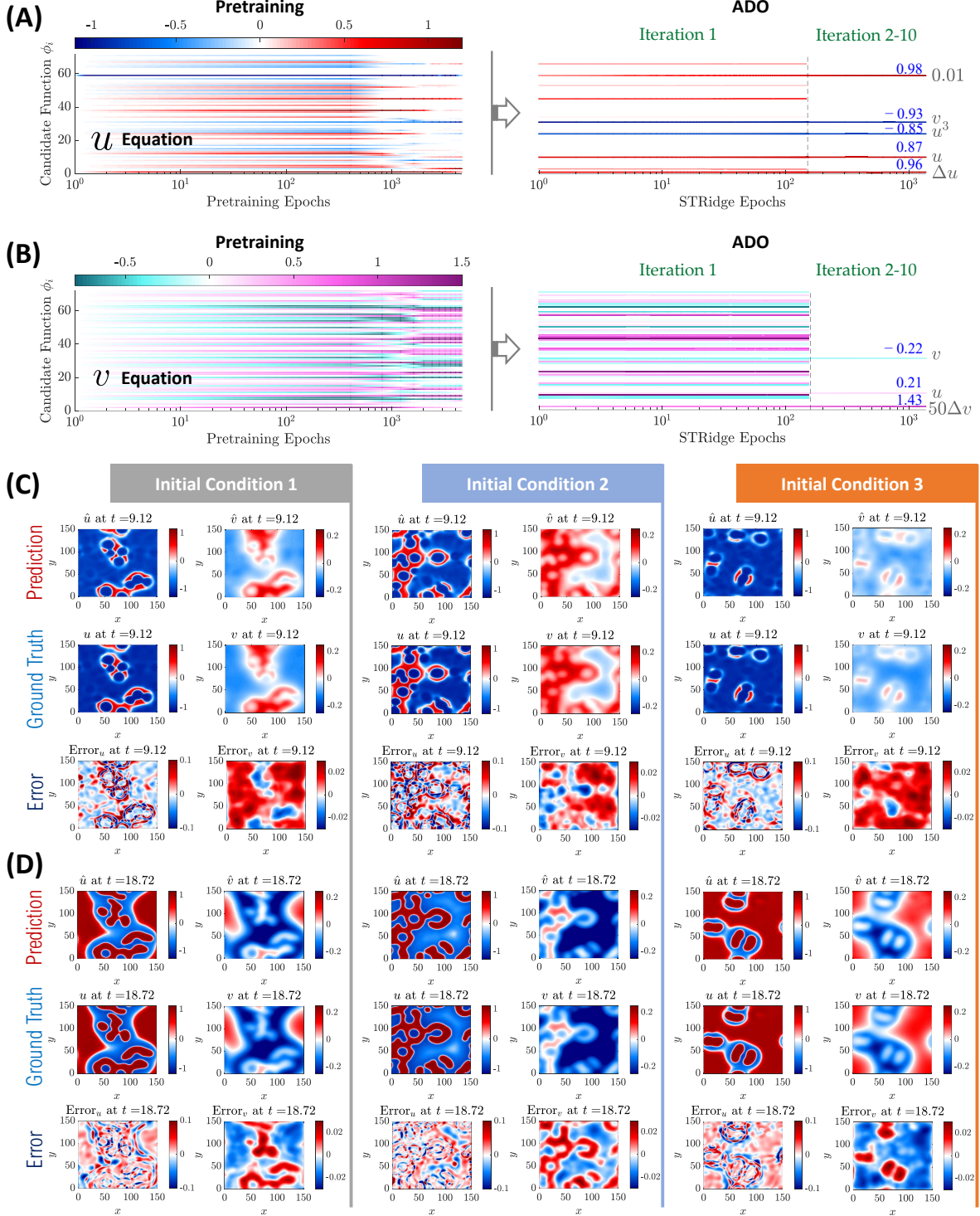


Fig. S.8: Discovered Fitzhugh-Nagumo equations based on measurements sampled under three initial conditions (ICs) with 10% noise. (A) Evolution of the sparse coefficients $\lambda_u \in \mathbb{R}^{72 \times 1}$ for 72 candidate functions used to construct the first PDE (u -equation), where the color represents the coefficient value. (B) Evolution of the sparse coefficients $\lambda_v \in \mathbb{R}^{72 \times 1}$ for the second PDE (v -equation). Note that, for visualization purpose, we re-scale the identified coefficients of the constant stimulus term “1” in the u -equation by multiplying 100 in A and the diffusion term Δv in the v -equation by dividing 50 in B. (C-D) Snapshots of predicted response, ground truth and error distributions for all three ICs at two unmeasured time instances ($t = 9.12$ and $t = 18.72$). The relative ℓ_2 error for the predicted full-field response (stacked u and v) is 5.04%.

replicates for each scenario for PDE discovery. Each mean dataset has a total of 38×5 measurement points for five time instances. Given our prior knowledge that the cell dynamics can be described by a diffusion (migration) and reaction (proliferation) process, we assume the PDE holds the form of $\rho_t = \gamma\rho_{xx} + \mathcal{F}(\rho)$, where γ is the unknown diffusion coefficient and \mathcal{F} denotes the underlying nonlinear reaction functional. We use 8 additional candidate terms (e.g., $\{1, \rho, \rho^2, \rho^3, \rho_x, \rho\rho_x, \rho^2\rho_x, \rho^3\rho_x\}$) to reconstruct \mathcal{F} , whose coefficients are sparse. Hence, the total number of trainable coefficients remains 9 (e.g., $\mathbf{\Lambda} \in \mathbb{R}^{9 \times 1}$).

We sample 1×10^4 collocation pairs using Latin hypercube sampling [9] in the spatiotemporal domain of $\Omega \times [0, T] = [0, 1900]\mu\text{m} \times [0, 48]\text{h}$. The DNN has 3 hidden layers of 30 nodes activated by the tanh function (see Fig. 1 in the Main Text). Considering that the cell density is constantly positive, we impose a *softplus* function (e.g., $\ln(1 + e^z)$) in the output layer to curb the final output of ρ . To account for potential large magnitude variation of the candidate term coefficients, we apply the sigmoid and tanh functions to squash magnitude gaps. Specifically, we set $\gamma = 1000\text{sig}(\tilde{\gamma})$ and $\boldsymbol{\lambda} = 50\text{tanh}(\tilde{\boldsymbol{\lambda}})$, where $\tilde{\gamma}$ and $\tilde{\boldsymbol{\lambda}}$ are trainable ‘‘proxies’’ for diffusion coefficient γ and other 8 coefficients $\boldsymbol{\lambda}$ (note: $\mathbf{\Lambda} = \{\gamma, \boldsymbol{\lambda}\} \in \mathbb{R}^{9 \times 1}$). The training efforts include the pretraining stage with 8×10^3 Adam epochs and 8×10^3 L-BFGS epochs, 5 ADO iterations, and extra post-training with 1×10^5 Adam epochs. In each ADO iteration, we use 2×10^3 Adam epochs in synergy with STRidge. Fig. S.9A shows the evolution of 9 coefficients for the example case of 18,000 cells, where redundant candidate terms are pruned right after the first ADO iteration via hard thresholding of the corresponding coefficients to zero. The next ADO iterations followed by post-tuning refine the coefficients of active terms for final reconstruction of the PDE. The discovered underlying PDEs under different initial cell states are given as follows:

$$\begin{aligned} 14\text{k cells: } \rho_t &= 553.05\rho_{xx} + 0.067\rho - 48.14\rho^2 \\ 16\text{k cells: } \rho_t &= 546.45\rho_{xx} + 0.066\rho - 44.10\rho^2 \\ 18\text{k cells: } \rho_t &= 560.12\rho_{xx} + 0.076\rho - 50.03\rho^2 \\ 20\text{k cells: } \rho_t &= 686.61\rho_{xx} + 0.092\rho - 60.02\rho^2 \end{aligned}$$

which share a unified form of $\rho_t = \gamma\rho_{xx} + \lambda_1\rho + \lambda_2\rho^2$ which exactly matches the famous Fisher-Kolmogorov model [13, 14]. The rates of migration (diffusion) and proliferation (reaction) generally increase along with the number of cells, as seen from the identified coefficients. Fig. S.9B-E depict the learned cell density profiles by the trained DNN, which capture the critical patterns of the measurement while showing little evidence of overfitting. With the discovered PDEs, we simulate/predict the evolution of cell densities at different time instants (12h, 24h, 36h and 48h) presented in Fig. S.9F-I, where the measurement at 0h is used as the initial condition while $\rho_x(x = 0, t) = \rho_x(x = 1900, t) = 0$ is employed as the Neumann boundary condition. The satisfactory agreement between the prediction and the measurement provides a clear validation of our discovered PDEs.

C Discussion

In this section, we discuss several other features, influence factors and limitations of the proposed PiDL method for data-driven discovery of PDEs, and highlight the potential future work.

C.1 Selection of candidate functions

The library of candidate functions is a significant component in PiDL, similar to the SINDy framework. On one hand, we prefer to make the candidate library as diverse as possible. On

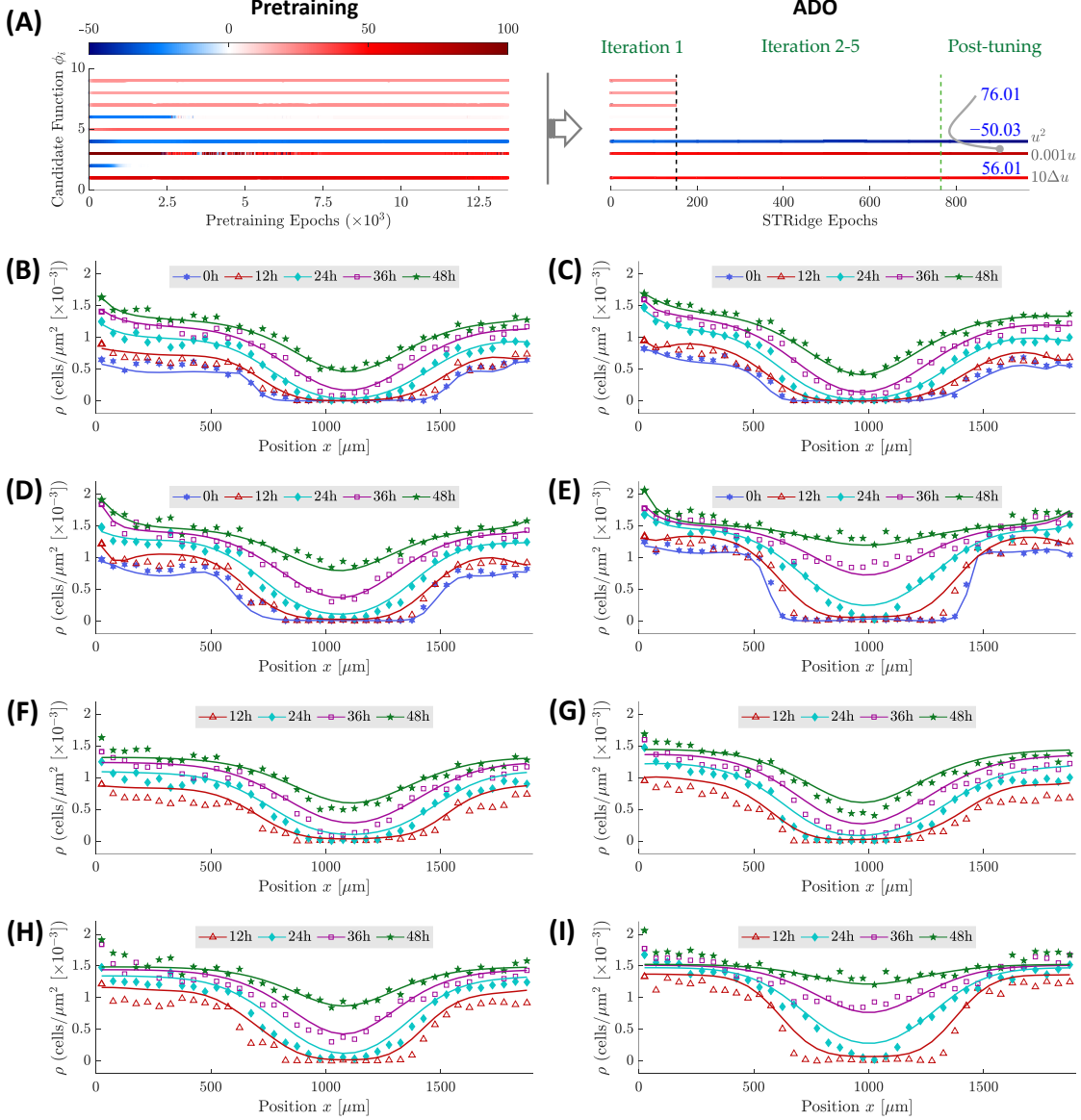


Fig. S.9: Experimental discovery of cell migration and proliferation. (A) Example evolution of the sparse coefficients $\Lambda \in \mathbb{R}^{9 \times 1}$ for 9 candidate functions used to construct the underlying PDE for the case of 18,000 cells. The diffusion and reaction coefficients for Δu and u are re-scaled for visualization purpose. (B)-(E) Predicted cell densities (represented by solid curves) by the trained DNNs in comparison with the measurement data (denoted by markers) for 14,000, 16,000, 18,000 and 20,000 cells, respectively. (F)-(I) Simulated cell densities, represented by solid curves, at different time instants based on the discovered PDEs for 14,000, 16,000, 18,000 and 20,000 cells, respectively, where the measurement at 0h is used as the initial condition while $\rho_x(x=0, t) = \rho_x(x=1900, t) = 0$ is employed as the Neumann boundary condition. The simulation result is represented by solid curves while the markers denote the measurement data.

the other hand, balancing the increasing theoretical and computational complexity is crucial for applications. We believe that a specialized library hinged by our domain-specific knowledge and statistical experience can constrain the search space and reduce the complexity of PDE discovery. Although the higher the dimension of the library is, the more likely the exact terms will be uncovered from data. Nevertheless, a highly large-scale library (e.g., the number of components on the order of magnitude of $\geq 10^3$), essentially approximated by the DNN, is very likely to be rank deficient

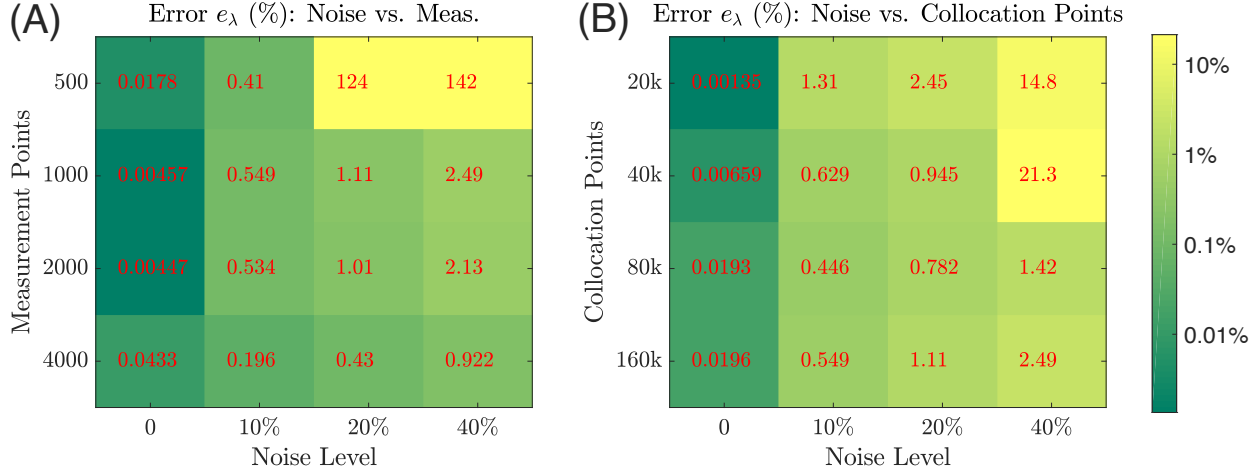


Fig. S.10: Error e_λ for discovery of Burgers’ equation under different measurement points, collocation points and noise levels. Numbers in each cell denote the percentage error of e_λ for a specific condition, which is the relative ℓ_2 norm error between the identified coefficients $\hat{\mathbf{A}}$ and the ground truth \mathbf{A}_{true} . The color also indicates the error level. The collocation points are fixed at 1.6×10^5 in (A), while the measurement points are always 1×10^3 in (B).

and have poor conditioning, in addition to the growing theoretical complexity and computational burden. Balancing these concerns and finding mathematical principles based on domain-specific knowledge to establish an efficient candidate library remain an open problem. Noteworthy, failing to include essential candidate functions will lead to false positive discovery of parsimonious closed form of PDEs, despite that a “best-of-fit” form can be found. Alternatively, we can first include rich polynomial terms and discover the governing PDEs in an approximate form, followed by Taylor series analysis [15] and power-law classes [16] to infer a more parsimonious form. This will also help inform the redesign and enrichment of the library of candidate functions for potentially improved discovery.

C.2 Noisy measurements and collocation points

The total loss function is evaluated on the measurement data (for \mathcal{L}_d) and the collocation points (for \mathcal{L}_p). Therefore, the availability of noisy measurement data and the number of collocation points sampled from the spatiotemporal space will affect the convergence of the PiDL model and thus the PDE discovery accuracy. We herein study the sensitivity of PiDL to these factors in the context of discovery accuracy based on the Burgers’ equation example. In particular, we use the relative ℓ_2 -norm error to reflect the global accuracy of the identified sparse coefficients, defined as $e_\lambda = \|\hat{\mathbf{A}} - \mathbf{A}_{\text{true}}\|_2 / \|\mathbf{A}_{\text{true}}\|_2$ where $\hat{\mathbf{A}}$ denotes the identified coefficients and \mathbf{A}_{true} is the ground truth. Fig. S.10 shows the error metrics for discovering the Burgers’ equation under different quantities of measurement points and collocation points and noise levels. Increasing the number of data points in the measurement (e.g., recorded by more sensors) can well compensate the noise effect as shown in Fig. S.10A (the number of collocation points is fixed at 1.6×10^5), which agrees with our common sense. Although optimal sensor placement might alleviate the need of large datasets [17], this is out of the scope of this work. The use of more collocation points can mitigate the noise effect and improve the discovery accuracy as illustrated in Fig. S.10B (the number of measurement points is fixed at 1×10^3). For this specific case, 2×10^4 (or more) collocation points are able to maintain a satisfactory discovery accuracy for measurements under noise corruption at a realistic level (e.g., $\leq 20\%$). When the data are sampled under a very noisy condition (e.g., 40% noise level),

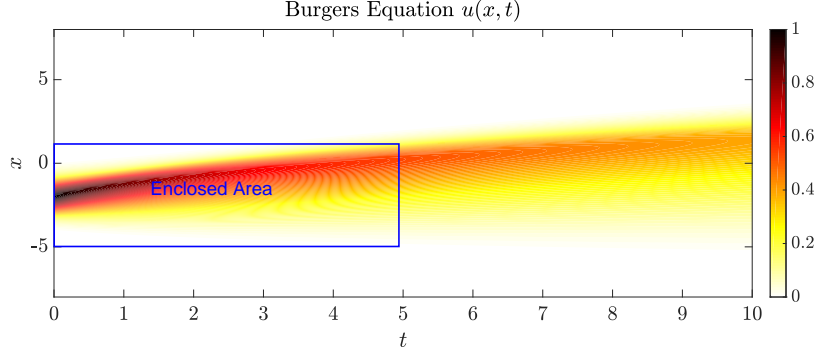


Fig. S.11: Parametric study on the effect of collocation points for discovering the Burgers’ equation. The measurements are only taken from the enclosed area, while the collocation points are sampled across the full field.

Table S.2: On the extrapolation (generalization) ability of PiDL

Case	Meas. points	Collocation points	Noise level	Training error	Validation error	Full-field error	ℓ_2 error of Λ
1	1.5×10^3	8×10^4	0	0.03%	0.04%	2.41%	0.02%
2	1.5×10^3	8×10^4	10%	5.73%	5.95%	4.50%	0.79%
3	1.5×10^3	0	0	58.79%	60.19%	144.99%	142.66%
4	3×10^3	0	0	0.10%	0.10%	14.36%	0.38%

Note: The training error, validation error and full-field error are calculated in the form of $\|\hat{\mathbf{u}} - \mathbf{u}_{\text{true}}\|_2 / \|\mathbf{u}_{\text{true}}\|_2$, where $\hat{\mathbf{u}}$ denotes the DNN-predicted response and \mathbf{u}_{true} is the reference ground truth solution.

the proposed method is still robust if a larger number of collocation points are used (e.g., $\geq 8 \times 10^4$).

It is noteworthy that the collocation points require no correlation with the measurement data. In particular, we use the Sobol sequence [3] (or Latin hypercube sampling [9] which is also applicable) to simulate a finer uniform partitions of the problem domain, making the random sampling of collocation points more representative. Intuitively, the more the collocation points are used, the more generalizable the trained network will be and the more accurate the discovered PDE is. However, a large number of collocation points also impose heavy computational burden, limited by hardware resources. A fairly large amount of collocation samples (e.g., on the order of magnitude of $> 10^4$), comparable to the complexity and dimension of the discovery problem, are suggested in practical applications meanwhile considering the memory of the computing machine.

We further conduct a comparative study on the role of collocation points and seek for numerical understanding of how much they can help for extrapolation. Taking the Burgers’ equation for instance, we define an enclosed area, part of the full-field response, as shown in Fig. S.11, and sample the measurement data only within such an area. We intend to reconstruct the full-field response beyond the enclosed area and discover the PDE taking advantage of collocation points. More specifically, the enclosed area is meshed by 100×50 spatiotemporal points. We take 30 randomly selected spatial locations as fixed sensors recording the dynamic response of the system, resulting in 1.5×10^3 data points. Additionally, we sample 8×10^4 collocation points from the full spatiotemporal field for evaluating the residual physics loss during model training. Four cases are considered to demonstrate the function of collocation points with measurements sampled in the enclosed area (see Table S.2).

Provided with clean measurements from the enclosed area and global collocation points, PiDL does an impressive job on both full-field response reconstruction and sparse coefficients identification (see Case 1 in Table S.2). When the measurements become noisy (e.g., 10% noise level), despite the response prediction errors increase, the PDE can still be accurately discovered (see Case 2 in Table S.2). If we consider removing all collocation points and only train the network with clean

measurements, the response prediction errors (even during the training and validation stage) all remain over 50%, meanwhile the PDE is also completely misidentified (see Case 3 in Table S.2). Once we double the clean measurement points to 3×10^3 , the trained DNN has strong interpolation and discovery abilities; however, the trained network does a poor job in extrapolating the full-field response (see Case 4 in Table S.2). Concluding from this parametric test, we can see that the collocation points can render PiDL tolerable to scarce and noisy measurements, making the DNN generalizable.

C.3 Simultaneous identification of unknown source term

In practical applications, the physical system might be subjected to spatiotemporal source input (\mathbf{p}) which is unknown and can be only sparsely measured. When discovering the underlying governing equation for such a system, the source should be considered and reconstructed concurrently. In this case, we incorporate the source candidate functions into the library ϕ for simultaneous discovery of the PDE and reconstruction of the unknown source. Thus, the sparse representation of the PDE(s) can be written as

$$\mathbf{u}_t = [\phi^u \ \phi^p][\Lambda^u \ \Lambda^p]^T$$

where ϕ^u and ϕ^p denote the libraries of candidate functions, while Λ^u and Λ^p are the corresponding sparse coefficients, for the field variable \mathbf{u} and the source \mathbf{p} , respectively. To demonstrate this concept, we test the Burgers' equation driven by a source term, expressed as

$$u_t + uu_x - 0.1u_{xx} = \sin(x) \sin(t).$$

To generate the solution, the problem domain is meshed into 201 spatial grid points for $x \in [-5, 5]$ and 101 time steps for $t \in [0, 10]$. We use 20 fixed sensors randomly selected from the spatial grid points to record the wave response (u) for 50 time steps, polluted with 10% noise. Note that the source is not measured and regarded as unknown.

The following libraries of candidate function are used to reconstruct the PDE and the source:

$$\phi^u = \{1, u, u^2, u^3, u_x, uu_x, u^2u_x, u^3u_x, u_{xx}, uu_{xx}, u^2u_{xx}, u^3u_{xx}, u_{xxx}, uu_{xxx}, u^2u_{xxx}, u^3u_{xxx}\}$$

$$\phi^p = \{a, b, c, d, a^2, b^2, c^2, d^2, ac, ab, ad, bc, bd, cd\}$$

where $a = \sin(t)$, $b = \sin(x)$, $c = \cos(t)$ and $d = \cos(x)$. The hyperparameters for the PiDL network are similar those used in the previous Burgers' example. The pre-training takes up to 15×10^3 epochs of Adam and about 1×10^3 epochs of L-BFGS, followed by 20 ADO iterations. In each ADO iteration, we use the Adam optimizer with 1×10^3 epochs and the L-BFGS with up to 1×10^3 epochs to train the DNN for each alternation within STRidge. The discovered PDE along with the uncovered source term is given by

$$u_t + 1.002uu_x - 0.088u_{xx} = 0.995 \sin(x) \sin(t).$$

It can be seen by comparing the above two equations that both the sparse terms and the corresponding coefficients are accurately identified, despite only scarce and noisy measurement of the system response is supplied. The discovery result is summarized in Fig. S.12. The evolution of the sparse coefficients for both the PDE and the source term shows robust and quick convergence to the ground truth (Fig. S.12A), with the average relative error for non-zero coefficients of $4.39\% \pm 7.03\%$. Although only 4.9% subsampled responses are measured while the source information is completely unknown, the PiDL approach can reasonably well extrapolate the full-field solution with a ℓ_2 error of 13.8% (see Fig. S.12B). The major errors are mostly distributed close to the

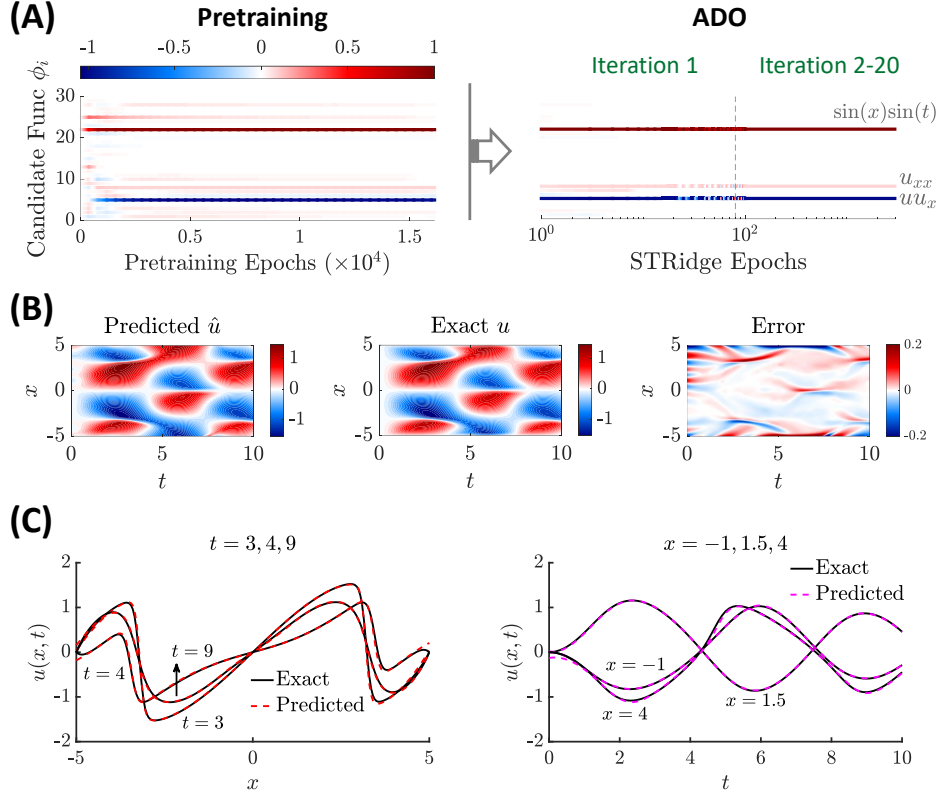


Fig. S.12: Discovered Burgers' equation and source term for measurement data with 10% noise. (A) Evolution of the sparse coefficients $\Lambda \in \mathbb{R}^{30 \times 1}$ for 30 candidate functions $[\phi^u \ \phi^p] \in \mathbb{R}^{1 \times 30}$ used to form the PDE and the unknown source term, where the color represents the coefficient value. (B) The predicted response in comparison with the exact solution with the prediction error. (C) Comparison of spatial and temporal snapshots between the predicted and the exact solutions. The relative full-field ℓ_2 error of the prediction is 13.8%. The major errors are mostly distributed close to the boundaries due to scarce training data.

boundaries due to scarce training data. Fig. S.12C shows the comparison of spatial and temporal snapshots between the predicted and the exact solutions which match well with each other.

Nevertheless, if the source is very complex with its general expression or form completely unknown, distinct challenges arise when designing the source library of candidate functions ϕ^p . This may require an extraordinarily large-space library to retain diversifying representations, and thus pose additional computational complexity for accurate discovery of the PDEs. In some specific cases, the unknown source term can probably be approximated by the combination of continuous basis functions such the Fourier series, instead of finding its closed form. These open questions will be addressed in our future work.

C.4 Other network architecture

There still remain some potential limitations associated with the present PiDL framework for physical law discovery. For example, although the fully connected DNN used in this work has advantage of analytical approximation of the PDE derivatives via automatic differentiation, directly applying it to model the solution of higher dimensional systems (e.g., long-short term response evolution in a 3D domain) results in computational bottleneck and optimization challenges. Advances in discrete DNNs with spatiotemporal discretization (e.g., the convolutional long-short term memory network (ConvLSTM) [18] or similar) have the potential to help resolve this challenge,

which will be demonstrated in our future work. However, a careful design of the spatiotemporal differentiator is required for the discrete DNNs (e.g., high-order finite difference filter for derivative approximation, accounting for domain irregularity, etc.).

References

- [1] Steven L. Brunton, Joshua L. Proctor, and J. Nathan Kutz. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*, 113(15):3932–3937, 2016.
- [2] Samuel H. Rudy, Steven L. Brunton, Joshua L. Proctor, and J. Nathan Kutz. Data-driven discovery of partial differential equations. *Science Advances*, 3(4):e1602614, 2017.
- [3] I.M Sobol. On the distribution of points in a cube and the approximate evaluation of integrals. *USSR Computational Mathematics and Mathematical Physics*, 7(4):86–112, 1967.
- [4] Hayden Schaeffer. Learning partial differential equations via data discovery and sparse optimization. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 473(2197):20160446, 2017.
- [5] Kunihiro Taira and Tim Colonius. The immersed boundary method: A projection approach. *Journal of Computational Physics*, 225(2):2118–2137, 2007.
- [6] Ankur Gupta and Saikat Chakraborty. Linear stability analysis of high- and low-dimensional models for describing mixing-limited pattern formation in homogeneous autocatalytic reactors. *Chemical Engineering Journal*, 145(3):399–411, 2009.
- [7] Lionel G Harrison. Kinetic theory of living pattern. *Endeavour*, 18(4):130–136, 1994.
- [8] Elizabeth E Holmes, Mark A Lewis, JE Banks, and RR Veit. Partial differential equations in ecology: spatial interactions and population dynamics. *Ecology*, 75(1):17–29, 1994.
- [9] M. D. McKay, R. J. Beckman, and W. J. Conover. A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics*, 21(2):239–245, 1979.
- [10] Richard FitzHugh. Impulses and physiological states in theoretical models of nerve membrane. *Biophysical journal*, 1(6):445, 1961.
- [11] Jinichi Nagumo, Suguru Arimoto, and Shuji Yoshizawa. An active pulse transmission line simulating nerve axon. *Proceedings of the IRE*, 50(10):2061–2070, 1962.
- [12] Wang Jin, Esha T Shah, Catherine J Penington, Scott W McCue, Lisa K Chopin, and Matthew J Simpson. Reproducibility of scratch assays is affected by the initial degree of confluence: experiments, modelling and model selection. *Journal of theoretical biology*, 390:136–145, 2016.
- [13] Ronald Aylmer Fisher. The wave of advance of advantageous genes. *Annals of Eugenics*, 7(4):355–369, 1937.
- [14] Philip K Maini, DL Sean McElwain, and David I Leavesley. Traveling wave model to interpret a wound-healing cell migration assay for human peritoneal mesothelial cells. *Tissue Engineering*, 10(3-4):475–482, 2004.
- [15] Kailiang Wu and Dongbin Xiu. Numerical aspects for approximating governing equations using data. *Journal of Computational Physics*, 384:200–221, 2019.
- [16] Bryan C. Daniels and Ilya Nemenman. Automated adaptive inference of phenomenological dynamical models. *Nature Communications*, 6:8133, 2015.
- [17] Krithika Manohar, J. Nathan Kutz, and Steven L. Brunton. Optimal Sensor and Actuator Selection using Balanced Model Reduction. *arXiv e-prints*, page arXiv:1812.01574, December 2018.
- [18] Shi Xingjian, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *Advances in neural information processing systems*, pages 802–810, 2015.