

---

# AUTOMATICALLY EXTRACTING PARTIAL DIFFERENTIAL EQUATIONS FROM DATA

---

A PREPRINT

Weizhen Li

Department of Engineering  
Durham University  
Durham, UK, DH1 3HN  
weizhen.li@dur.ac.uk

Rui Carvalho

Department of Engineering  
Durham University  
Durham, UK, DH1 3HN  
rui.carvalho@dur.ac.uk

## ABSTRACT

Identifying partial differential equations (PDEs) from data is crucial for understanding the governing mechanisms of natural phenomena, yet it remains a challenging task. We present an extension to the ARGOS framework, ARGOS-RAL, which leverages sparse regression with the recurrent adaptive lasso to identify PDEs from limited prior knowledge automatically. Our method automates calculating partial derivatives, constructing a candidate library, and estimating a sparse model. We rigorously evaluate the performance of ARGOS-RAL in identifying canonical PDEs under various noise levels and sample sizes, demonstrating its robustness in handling noisy and non-uniformly distributed data. We also test the algorithm's performance on datasets consisting solely of random noise to simulate scenarios with severely compromised data quality. Our results show that ARGOS-RAL effectively and reliably identifies the underlying PDEs from data, outperforming the sequential threshold ridge regression method in most cases. We highlight the potential of combining statistical methods, machine learning, and dynamical systems theory to automatically discover governing equations from collected data, streamlining the scientific modeling process.

**Keywords** System identification · Machine learning · Sparse regression · Partial differential equations · Nonlinear dynamics

## 1 Introduction

In recent years, scientists have increasingly employed statistical and machine learning methods to uncover the governing equations of dynamical systems, particularly differential equations, from observational data [1–5]. Data-driven methods offer several advantages over traditional approaches that rely on first principles and expert knowledge. These methods can reveal patterns and relationships in the data that may not be apparent from first principles, providing new insights into complex systems [6, 7]. They are also adept at working with noisy or incomplete data commonly encountered in real-world applications, employing techniques from machine learning to enhance the robustness of discoveries [8–11]. Furthermore, by reducing the need for manual intervention and domain expertise, data-driven methods can significantly streamline the discovery process [12].

Data-driven discovery in dynamical systems has evolved from early parameter estimation using spline approximation and system reconstruction [13, 14], to leveraging statistical methods such as least squares [15–17], mixed-effects models [18, 19], and Bayesian approaches [2, 20] for parameter estimation in ODEs and PDEs. The advent of high-performance computing has further propelled symbolic regression, enabling the discovery of governing equations from data in physics and engineering [1, 21–23]. A notable development in this field is the Sparse Identification of Nonlinear Dynamics (SINDy) approach [3, 4], which constructs an extensive library of potential terms and employs the Sequential Threshold Ridge Regression (STRidge) algorithm [4] to select significant terms iteratively.

SINDy and its various enhancements [24–29] have been extensively used to discover a broad spectrum of ODEs and PDEs, describing diverse phenomena such as fluid mechanics [30], turbulence models [31], aerodynamics [32], and

biological and chemical systems [33, 34]. Recent developments have combined neural network-based techniques and SINDy, leading to innovative approaches that enhance noise tolerance in identifying PDEs [5, 23, 35–39]. Neural networks can learn complex nonlinear relationships and effectively filter out noise, complementing SINDy’s ability to identify parsimonious models. However, both neural network and SINDy methods require specific hyperparameter tuning, such as setting regularization parameters or choosing network architectures. For example, STRidge requires setting a threshold to select active terms from the candidate library [4, 29, 36, 38]. Additionally, SINDy-based methods typically approximate numerical derivatives from noisy data using the Savitzky-Golay filter, a technique for smoothing data by fitting local low-degree polynomials [40]. The parameters of this filter, such as the polynomial degree and window size, must be carefully tuned for optimal performance [4, 12]. Neural network approaches, on the other hand, require detailed decisions regarding their architecture and functioning, such as the number of neurons, the structure of hidden layers, the types of activation and loss functions, and the learning rate. In particular, using physics-informed neural networks [5, 36, 38] requires a prior understanding of the equation terms, as well as initial and boundary conditions. Consequently, using neural networks and SINDy-based methods presents a trade-off: the absence of fully automated algorithms requires users to engage in manual tuning and iterative usage of semi-automated algorithms. This scenario highlights a key challenge in the field: developing an automated algorithm to identify PDEs with minimal manual intervention, streamlining the process, and improving its applicability across diverse scientific domains.

To address the challenge of parameter tuning, Egan *et al.* [12] proposed the Automatic Regression for Governing Equations (ARGOS) algorithm, which identifies ODEs by automating the parameter tuning process. ARGOS assumes the underlying system is unknown, automates the fine-tuning of parameters for numerical differentiation, and leverages sparse regression with bootstrap confidence intervals to select active terms from the candidate library. To automatically identify PDEs, we develop ARGOS with the Recurrent Adaptive Lasso (ARGOS-RAL). This extension of the ARGOS framework employs only sparse regression to identify equations rather than engaging in large-scale bootstrapping.

We evaluate the performance of the ARGOS-RAL algorithm through a series of three numerical tests, each designed to assess its ability to identify canonical PDEs across diverse fields, including biology, neuroscience, earth science, fluid mechanics, and quantum mechanics. The first test explores the algorithm’s resilience against varying noise levels by altering the signal-to-noise ratio (SNR) in Gaussian random noise integrated into the PDE solutions, which is crucial for understanding the robustness of ARGOS-RAL under realistic noisy conditions. The second test addresses the practical challenges encountered in real-world data collection, which often results in non-uniformly distributed data points in space and time, by exploring the minimum percentage of data points necessary for the algorithm to accurately identify the underlying equation. The final evaluation assesses the algorithm’s ability to process datasets characterized by significant noise, challenging its limits and practical applicability in scenarios where data quality is compromised. Our results demonstrate that ARGOS-RAL can effectively and reliably identify the underlying PDEs from data, outperforming the STRidge method used in SINDy.

## 2 Methods

### 2.1 Overview of the ARGOS-RAL Framework

The general form of a homogeneous PDE is

$$u_t + F(x, t, u, u_x, u_{xx}, \dots) = 0 \quad (1)$$

where  $F(\cdot)$  governs the behavior of the system, with  $u = u(x, t)$  denoting its state. The notation  $u_t, u_x, u_{xx}, \dots$  represents the partial derivatives of  $u$  with respect to time and space, respectively. Equation (1) serves as a foundational representation of the dynamical system, encapsulating a wide range of phenomena through its generalized form, which can be adapted to include multiple spatial dimensions or to model systems without explicit time dependence.

To focus on data-driven modeling of spatiotemporal dynamical systems, we incorporate empirical data directly into the modeling process:

$$\mathbf{U}_t = \frac{\partial \mathbf{U}}{\partial t} = \mathbf{F}(x, \mathbf{U}, \mathbf{U}_x, \mathbf{U}_{xx}, \dots), \quad (2)$$

where  $\mathbf{U} \in \mathbb{R}^{n \times m}$  is a matrix representing the solution of the PDE as a function of  $x$  and  $t$ , and  $\mathbf{F}(\cdot)$  denotes the unknown mapping inferred from the collected data, which contains linear and nonlinear operators.

We aim to estimate the unknown mapping  $\mathbf{F}(\cdot)$  with sparse regression by constructing a comprehensive library of potential terms and assuming that only a few of them are active [3, 4, 12]. To cover a broad spectrum of possible influences on the dynamics of the system, this library includes a wide variety of functions, such as constants, monomials, interaction terms (products of variables), possibly trigonometric, and other functions, depending on the dynamical system being studied [4]. In the case of Burgers’ equation,  $u_t = -uu_x + 0.1u_{xx}$ , the true dynamics involves only two

terms: the nonlinear convection term  $uu_x$  and the linear diffusion term  $u_{xx}$ . When applying sparse regression to data generated from Burgers' equation, the method should ideally select only these two terms from the candidate library.

All features related to  $\mathbf{U}(x, t)$  in Eq. (2) are matrices. Implementing this matrix data in sparse regression leads to the creation of  $m$  distinct regression models. Each model captures the spatial dynamics of the system at a specific time point  $t_j$ , where  $j = 1, 2, \dots, m$ . To consolidate the  $m$  regression models into a single linear regression problem, we reshape the matrix  $\mathbf{U}(x, t)$  and its derivative matrices into vectors. These vectors then serve as predictors within the candidate library  $\Theta$ , which can be represented in  $\mathbb{R}^{(n \cdot m) \times p}$  or  $\mathbb{C}^{(n \cdot m) \times p}$ . By stacking the vectorized data and candidate terms, we can estimate a single sparse coefficient vector  $\beta$  that represents the governing equation across all time points rather than estimating separate models for each time point. Here,  $\mathbf{U} \in \mathbb{R}^{n \times m}$  is represented in matrix form as

$$\mathbf{U}(x, t) = \begin{pmatrix} u(x_1, t_1) & u(x_1, t_2) & \cdots & u(x_1, t_m) \\ u(x_2, t_1) & u(x_2, t_2) & \cdots & u(x_2, t_m) \\ \vdots & \vdots & \ddots & \vdots \\ u(x_n, t_1) & u(x_n, t_2) & \cdots & u(x_n, t_m) \end{pmatrix}. \quad (3)$$

Vectorizing Eq. (3) yields:

$$\mathbf{u} = \text{vec}(\mathbf{U}) = (u(x_1, t_1) \ \cdots \ u(x_n, t_1) \ \cdots \ u(x_1, t_m) \ \cdots \ u(x_n, t_m))^T. \quad (4)$$

Similarly,  $\mathbf{u}_t = \text{vec}(\mathbf{U}_t) = \text{vec}(\partial \mathbf{U} / \partial t)$ ,  $\mathbf{u}_x = \text{vec}(\mathbf{U}_x) = \text{vec}(\partial \mathbf{U} / \partial x)$ ,  $\mathbf{u}_{xx} = \text{vec}(\mathbf{U}_{xx}) = \text{vec}(\partial^2 \mathbf{U} / \partial x^2)$ ,  $\mathbf{u}^2 = \text{vec}(\mathbf{U} \odot \mathbf{U})$ , and  $\mathbf{u}\mathbf{u}_x = \text{vec}(\mathbf{U} \odot \mathbf{U}_x)$ , where  $\odot$  denotes the Hadamard product. The design matrix is given by

$$\Theta(\mathbf{u}) = \left( \begin{array}{ccccccccc} | & | & | & | & | & | & | & | & | \\ 1 & \mathbf{u} & \cdots & \mathbf{u}^d & \cdots & \mathbf{u}_x & \mathbf{u}_{xx} & \cdots & \mathbf{u}\mathbf{u}_x & \cdots & \mathbf{u}^d\mathbf{u}_{xx} & \cdots \\ | & | & | & | & | & | & | & | & | & | & | & | \end{array} \right), \quad (5)$$

where  $\mathbf{u}^d$  is a vector where all elements denote a  $d$ -th degree monomial. For example, if our data  $\mathbf{U}(x, t)$  is on a  $200 \times 100$  grid (i.e. 200 spatial measurements and 100 time-steps) and the candidate library has 30 terms, then  $\Theta \in \mathbb{R}^{20000 \times 30}$ .

After vectorization, we estimate  $\mathbf{F}(\cdot)$  by transforming Eq. (2) to a linear regression model

$$\mathbf{u}_t = \Theta(\mathbf{u})\beta + \epsilon, \quad (6)$$

where  $\beta \in \mathbb{R}^p$  is a sparse coefficient vector in which only a few values are nonzero, and  $\epsilon$  is the vector of residuals.

## 2.2 Automated Numerical Differentiation using the Savitzky-Golay Filter and the Gaussian Blur

A crucial step in constructing the candidate library in Eq. (5) is the numerical calculation of derivatives (see Fig. 1 A and B). The Savitzky-Golay filter [40] has become a favored solution in system identification for signal smoothing and differentiation [4, 41]. This method applies a least squares polynomial fit over a sliding window of data points, thereby achieving simultaneous signal smoothing and differentiation. The selection of the Savitzky-Golay filter is grounded in its proven ability to accurately maintain the original contour of the signal while significantly reducing noise and to approximate higher-order numerical derivatives with symbolic differentiation [42].

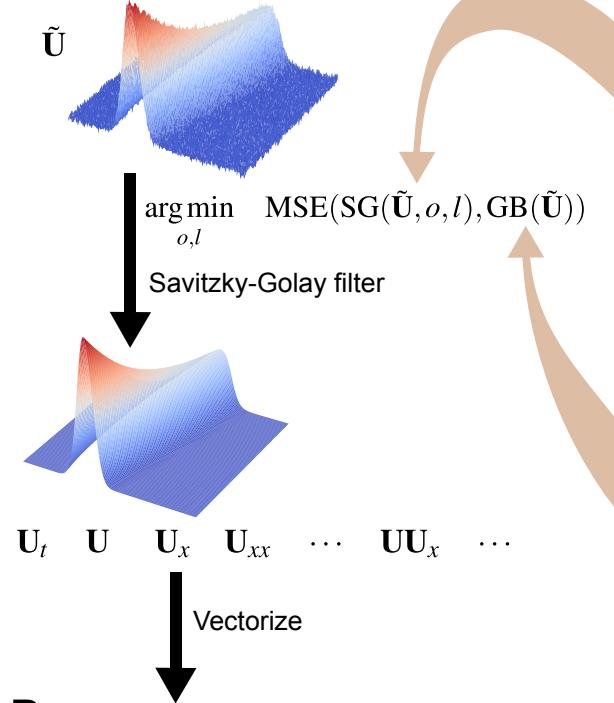
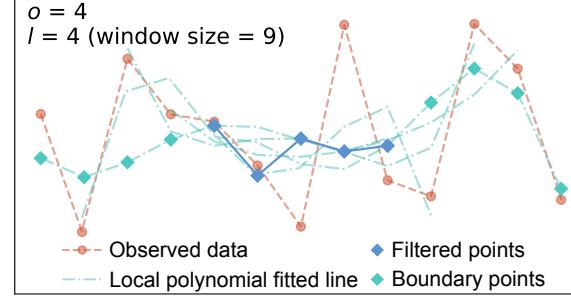
The Savitzky-Golay filter is characterized by two integer hyperparameters: the polynomial order  $o$  and the window length  $l$ , which are constrained by the conditions that  $o$  must be at least 2,  $l$  should be an odd number, and  $o + 1 + \text{mod}(o) \leq l \leq n - 1$  [42]. To automate the selection of these hyperparameters, we first apply a Gaussian blur with the kernel  $(1, 2, 1)$  to smooth the observational data (see A.1). We then treat this smoothed data, denoted as  $\text{GB}(\tilde{\mathbf{U}})$ , as the ground truth. Next, we find the optimal set of hyperparameters  $\{o^*, l^*\}$  by minimizing the mean squared error (MSE) between the Savitzky-Golay filtered data  $\text{SG}(\tilde{\mathbf{U}}, o, l)$  and the ground truth  $\text{GB}(\tilde{\mathbf{U}})$  (see Algorithm 1 in A.2). After finding the optimal set  $\{o^*, l^*\}$ , we use the Savitzky-Golay filter with these parameters to compute the smoothed data and its derivatives.

## 2.3 Sparse Regression with the Recurrent Adaptive Lasso

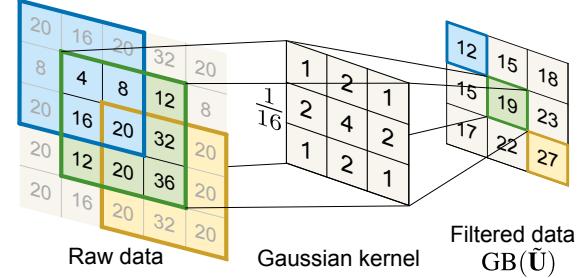
The adaptive lasso is a two-step method [12, 43]. The first step uses the ordinary least squares (OLS) to obtain unbiased estimates and derive the weights  $w$ :

$$w = |\hat{\beta}_{ols}|^{-\gamma}, \quad \gamma > 0 \quad (7)$$

## A Smoothing and Derivatives

Savitzky-Golay filter  $\text{SG}(\tilde{\mathbf{U}}, o, l)$ 

Two-dimensional Gaussian blur



## B Constructing the Candidate Library

$$\mathbf{u}_t \underbrace{\mathbf{1} \quad \mathbf{u} \quad \dots \quad \mathbf{u}^d \quad \dots \quad \mathbf{u}_x \quad \dots \quad \mathbf{u}\mathbf{u}_x \quad \dots \quad \mathbf{u}^d\mathbf{u}_{xx} \quad \dots}_{\Theta(\mathbf{u})}$$

## C The Recurrent Adaptive Lasso Algorithm

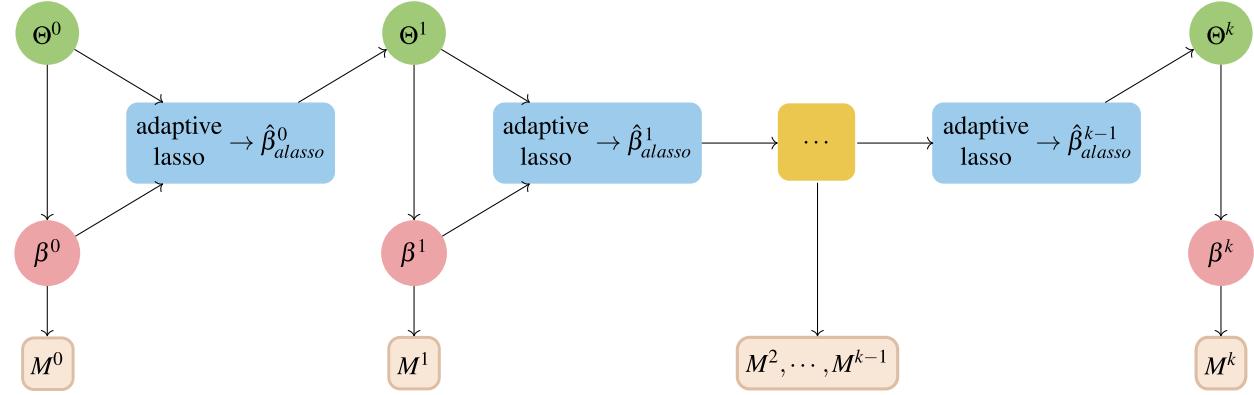


Figure 1: Process of identifying PDEs from data using ARGOS with the recurrent adaptive lasso. The identification process consists of three main steps: (A) automatic smoothing and calculation of derivatives, (B) construction of the candidate library, and (C) implementation of the recurrent adaptive lasso. We begin by collecting the data  $\tilde{\mathbf{U}}$  and applying the automatic Savitzky-Golay filter with Gaussian blur to calculate the smoothed  $\mathbf{U}$  and its partial derivatives. Next, we vectorize the smoothed data, all partial derivatives, and other related terms to construct the candidate library. Finally, we employ the recurrent adaptive lasso to identify the active features in the library, and we estimate the unbiased coefficients of the identified model using ordinary least squares regression.

where  $\hat{\beta}_{ols}$  is the OLS estimate, and  $\gamma$  is an exponent tuning the shape of the soft-thresholding function. In the second step, we obtain the estimated coefficients  $\hat{\beta}_{alasso}$  using the `glmnet` package [44] in R by solving the problem

$$\hat{\beta}_{alasso} = \arg \min_{\beta} \|u_t - \Theta \beta\|_2^2 + \lambda \sum_{j=1}^p w_j |\beta_j|, \quad (8)$$

where  $\lambda$  is a nonnegative regularization parameter controlling the amount of shrinkage applied to the coefficients of the predictors. Unlike the lasso, where the weight vector is  $w = 1$ , the adaptive lasso varies the weights in the regularization function, resulting in a stronger penalty on smaller coefficients, thus driving more of them to zero and leading to a sparser model compared to the standard lasso. The recurrent adaptive lasso applies the adaptive lasso repeatedly until convergence, resulting in a sparse model with fewer non-zero coefficients.

To balance the model's complexity against its accuracy, we determine the regularization parameter  $\lambda$  by employing the Pareto curve, which illustrates the optimal trade-off between the regularization penalty and the model residuals [45–47] (see Fig. 2). Although cross-validation is an alternative method, Cortiella *et al.* [27] have shown that it finds a  $\lambda$  optimized for prediction, potentially overfitting the true underlying equation with extra features.

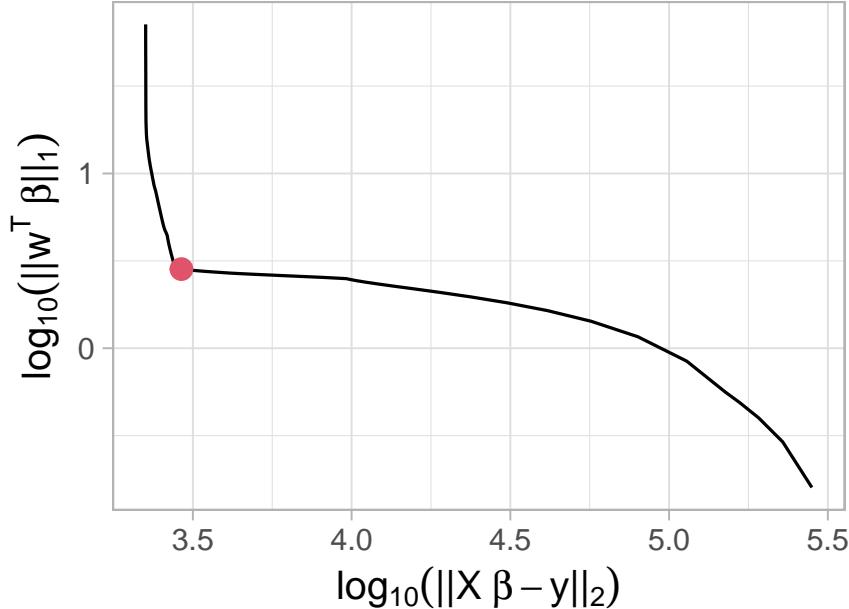


Figure 2: Pareto curve of the adaptive lasso for a sampled dataset from a Navier-Stokes system with an SNR of 36 dB. The Pareto curve balances the trade-off between sparsity and goodness-of-fit. The red point on the curve indicates the optimal value of the regularization parameter  $\lambda$  that achieves the best balance between these two competing objectives. Increasing  $\lambda$  leads to sparser solutions at the cost of a poorer fit to the data, while decreasing  $\lambda$  improves the fit but yields less sparse solutions.

The adaptive lasso regression often detects more terms than those in the true system. To improve parsimony, Egan *et al.* [12] suggested combining the adaptive lasso with bootstrap techniques to identify ODEs. Similarly, Cortiella *et al.* [27] adopted a modified version of the multi-step adaptive lasso [48] to develop a sparser model that more accurately identifies the true equations. This is achieved by iteratively adjusting the adaptive weights using previous estimates from the adaptive lasso. A significant advancement made by Cortiella *et al.* [27] is their method's ability to maintain finite weights in the adaptive lasso equation by ensuring that the estimated coefficients shrink to a small, nonzero value rather than dropping to zero. However, this approximation unintentionally introduces numerical inaccuracies as a trade-off for preventing overflow during the equation identification process.

The recurrent adaptive lasso is an iterative algorithm that estimates an initial sparse model using the adaptive lasso and subsequently refining it by trimming the candidate library (see Fig. 1 C). At each iteration, it removes terms whose coefficients the adaptive lasso penalized to zero (see A.2 Algorithm 2 step 9). It then employs least squares to re-estimate the coefficients of the remaining terms, which are used to update the adaptive weights in the next adaptive lasso iteration. This focuses the regularization on the terms that had small coefficients in the previous iteration. As

this process repeats, the recurrent adaptive lasso increasingly concentrates the  $\ell_1$ -norm shrinkage on terms that are likely irrelevant, driving their coefficients to zero [43, 49]. Meanwhile, it relaxes the regularization on terms that consistently have larger coefficients, allowing the model to retain them. The candidate set gets smaller at each iteration until the algorithm converges on a sparse model containing only the key terms. This iterative re-weighting allows the recurrent adaptive lasso to prune irrelevant terms more aggressively than the standard adaptive lasso while retaining good predictive performance. The result is a parsimonious model that identifies the true governing equation more reliably, even in the presence of many extraneous candidate terms.

Increasing the number of iterations may cause the recurrent adaptive lasso to underestimate the model. This can lead to the omission of active terms that should be included in the true underlying equation. Therefore, while iterating the candidate library  $\Theta$ , we record all candidate models and calculate the Akaike information criterion (AIC) for each model to determine the final governing equation corresponding to the lowest AIC. Given the uncertainty that the true model falls within all candidates, the AIC serves to select the model that best approximates the true model [50, 51].

### 3 Results and Discussion

#### 3.1 Evaluating the Performance of ARGOS-RAL under Varying Noise Levels and Sample Sizes

We compare the performance of ARGOS-RAL and STRidge [4] in identifying ten canonical PDEs under various SNRs and sample sizes ( $N$ ). We evaluate their performance on both noisy and noiseless data. Figure 3 demonstrates the impact of introducing increasing levels of Gaussian random noise into the solution of the Burgers' equation, effectively decreasing the SNR values.

In the evaluation of noise-contaminated data, we express the SNR as  $\text{SNR} = 20 \log_{10}(\sigma_U/\sigma_Z)$ , where  $\sigma_U$  is the standard deviation of the original data, and  $\sigma_Z$  represents the standard deviation of the added noise. We systematically vary  $\sigma_Z$  to span a broad range of noise levels, facilitating a comprehensive evaluation of the efficacy of ARGOS-RAL and STRidge in identifying various PDEs under different noise conditions. For this purpose, we generate datasets with SNRs set at  $\{0, 2, \dots, 58, 60, \infty\}$  [12], each comprising paired elements  $\{\mathbf{u}_t, \Theta(\mathbf{u})\}$ . This approach allows us to examine the robustness of each PDE identification method as it copes with varying noise levels.

In investigating sample size,  $N$ , our objective is to determine the smallest number of samples needed to reliably identify PDEs with a success rate exceeding 80%. To achieve this, we first generate a full dataset for each PDE by calculating partial derivatives and assembling a candidate library as described in Eq. (5). The size of the full dataset, denoted as  $N$ , varies depending on the specific PDE under consideration. Specifically,  $N = 10^4$  for the advection-diffusion, Burgers, and cable equations,  $N = 10^5$  for the quantum harmonic oscillator, transport, Navier-Stokes, and reaction-diffusion equations, and  $N = 10^{4.8}$  for the heat and Korteweg-De Vries (KdV) equations. Next, we randomly sample smaller subsets of size  $N$  from the full dataset, where  $N$  is chosen from a logarithmically spaced grid:  $N = 10^2, 10^{2.2}, 10^{2.4}, \dots, N$  [12] (see the blue points in Fig. 3 A). By applying the PDE identification methods to these subsets and evaluating their success rates, we can determine the smallest sample size required for reliable identification of each PDE.

#### 3.2 Quantifying Success Rates in Identifying Canonical PDEs

To evaluate the impact of different SNRs and data sizes on the method, we measure the uncertainty of model identification caused by random sampling. To do so, we create 100 unique datasets at each point on the grid, corresponding to different SNRs and  $N$  values. For each dataset, we quantify the identification accuracy with the success rate,  $\eta = \#\text{correct}/100$ , where  $\#\text{correct}$  represents the number of times the model correctly identifies all active terms. Our accuracy assessment ignores small differences between theoretical and empirical coefficients, such as a theoretical value of 0.1 compared to an estimated value of 0.098. Figure 4 illustrates these results for a selection of systems: the Burgers', Cable, Navier-Stokes, reaction-diffusion, and quantum harmonic oscillator models. We provide further analysis on additional PDEs –Transport, Heat, Advection-Diffusion, and KdV equations – in A.3 Fig. 6.

ARGOS-RAL identifies Burgers', cable, Navier-Stokes, reaction-diffusion, and advection-diffusion equations, achieving a success rate of 100% when the SNR exceeds 30 dB (see Fig. 4 and 6). However, accurately detecting specific equations requires a high SNR, particularly for the quantum harmonic oscillator, KdV, transport, and diffusion equations. The KdV equation, which involves third-order partial derivatives, presents challenges due to the significant biases in numerical approximations of these derivatives [39], resulting in datasets unsuitable for system identification with sparse regression. To implement sparse regression within the real number domain for the complex number quantum harmonic oscillator PDE, we apply the transformation shown in A.3 Eq. (17). This transformation expands the design matrix  $\Theta$  from  $nm \times p$  to  $2nm \times 2p$ , effectively quadrupling its size and potentially leading to high correlations between the variables in  $\Theta$ . The transport and diffusion equations, containing only terms  $u_x$  and  $u_{xx}$  respectively, exhibit high correlation with

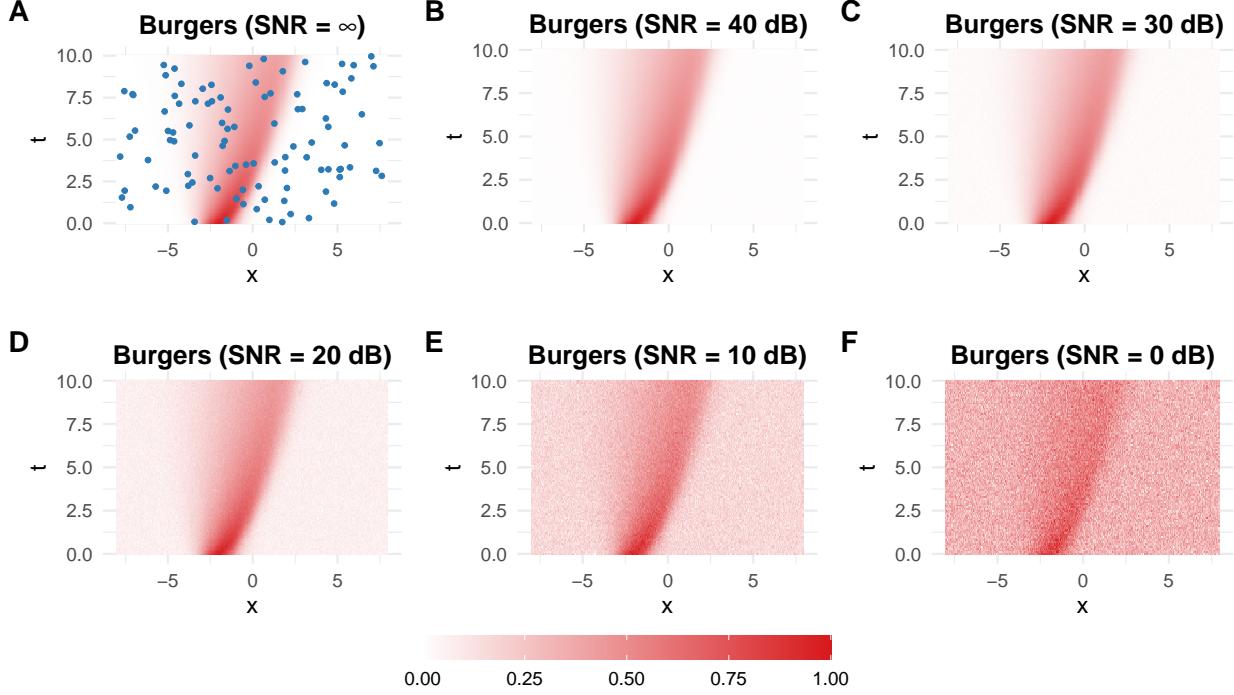


Figure 3: Influence of SNR on the Burgers' equation dataset. (A) Noiseless data points (blue) serve as a reference for evaluating the impact of sample size on PDE identification accuracy. (B-F) Noisy datasets are generated by adding Gaussian noise at SNR levels of 40 dB, 30 dB, 20 dB, 10 dB and 0 dB, respectively, to comprehensively characterize the system's behavior under varying noise conditions.

their correlated terms in the library, such as  $\{u_x, uu_x\}$  and  $\{u_{xx}, uu_{xx}\}$ , which hinders the effectiveness of  $\ell_1$ -norm shrinkage regression in identifying correct terms [43].

Figures 4 and 6 illustrate that ARGOS-RAL achieves a higher success rate than STRidge in identifying PDEs with limited data points. ARGOS-RAL consistently identifies a significant number of PDEs using as few as 1000 data points, maintaining a success rate above 80%. However, some equations, such as the reaction-diffusion and KdV equations, require larger sample sizes of approximately  $10^4$  and  $10^{3.8}$  data points, respectively, for reliable identification. We thus demonstrate ARGOS-RAL as a consistent and efficient method for PDE identification with non-uniformly sampled and noiseless datasets.

ARGOS-RAL shows a remarkable ability to identify PDEs accurately and consistently across a wide range of SNRs and sample sizes. Its success rate improves as the SNR and sample size increase, reaching 100% when both values are sufficiently large. This trend highlights the robustness of ARGOS-RAL in handling various data conditions and underscores its effectiveness in identifying PDEs, even when faced with varying levels of data quality and quantity. However, in certain scenarios, STRidge [4] with specific  $d_{tol}$  thresholds exceeds the performance of ARGOS-RAL. For instance, STRidge achieves higher success rates in identifying Navier-Stokes and reaction-diffusion equations at a 30 dB SNR, using  $d_{tol}$  settings of 2 and 10, respectively (see Fig. 4 C and D). Moreover, STRidge with  $d_{tol} = 2$  is more proficient in identifying the quantum harmonic oscillator and the transport equation with an SNR lower than 52 dB, see Fig. 4 E and 6 C, respectively. These results from the SNR and  $N$  experiments reveal that using a single fixed threshold in STRidge can lead to performance variability depending on the input data, highlighting the difficulty of selecting an appropriate  $d_{tol}$  threshold without prior knowledge of the system. This variability underscores the sensitivity of STRidge to specific threshold settings, which can impact its consistency across different datasets. Overall, STRidge surpasses ARGOS-RAL in identifying simpler PDEs, such as the transport and diffusion equations; see Fig. 6 C and D.

### 3.3 Robustness Analysis using White Gaussian Noise

To better understand the limits of identification algorithms, we designed an extreme test on a single spatial dimension. This test effectively creates a situation without valid data collection ( $\sigma_U = 0$ ), equivalent to an SNR of negative infinity,

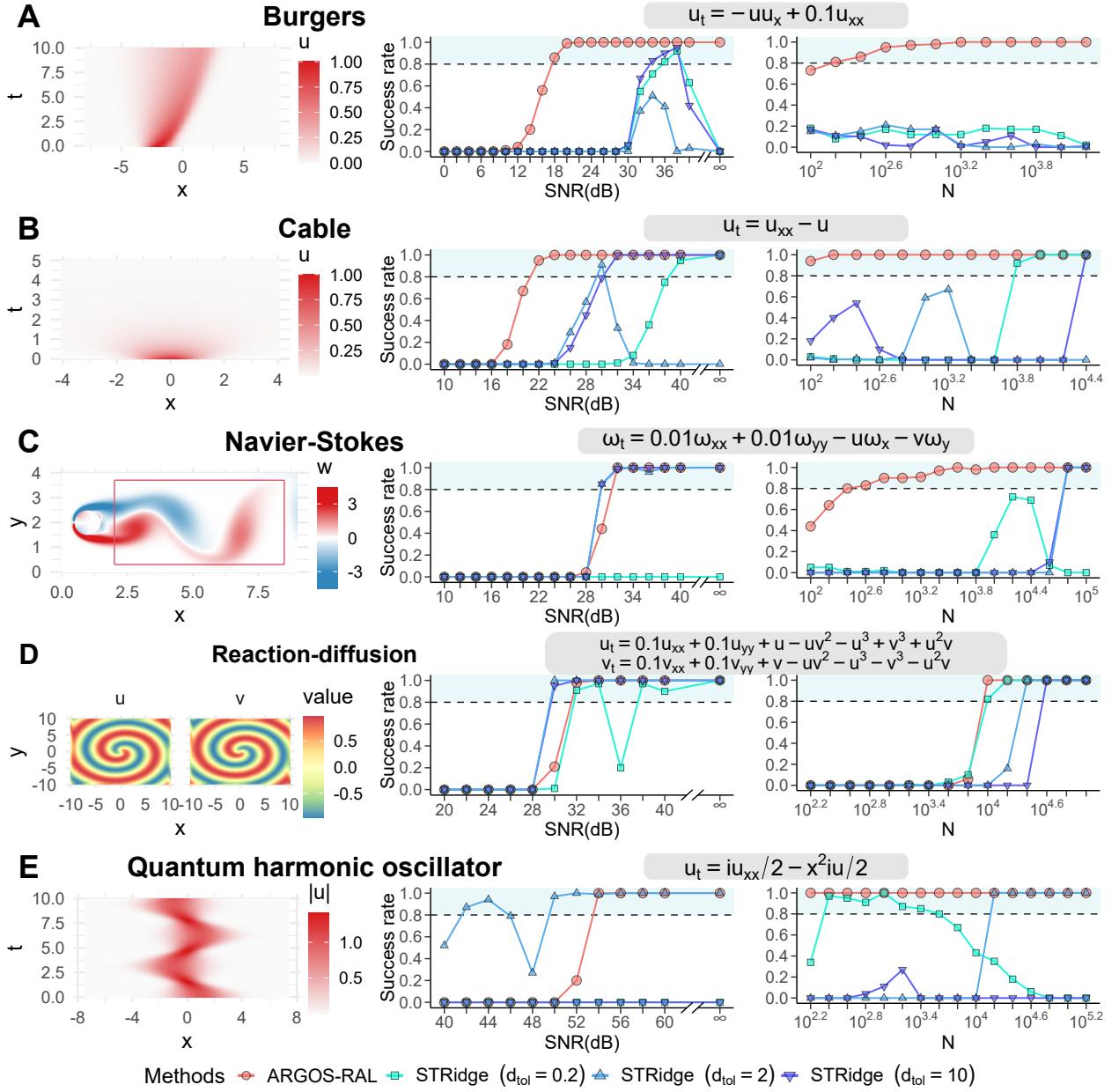


Figure 4: Success rates of ARGOS-RAL and STRidge in identifying (A) Burgers', (B) cable, (C) Navier-Stokes, (D) reaction-diffusion, and (E) quantum harmonic oscillator equations with varying SNRs and sample sizes. We analyze the noise tolerance by adding noise of different SNRs to the PDE solutions. For the sample size analysis, we randomly sample points from the set  $\{\mathbf{u}_t, \Theta(\mathbf{u})\}$  based on noiseless data. In panel (C), we use the region indicated by the red rectangle to implement both the SNR and sample size tests by sampling points within this area. PDE solution plots display time snapshots at  $t = 306$  for Navier-Stokes in panel (C) and  $t = 1$  for reaction-diffusion in panel (D). Lines connecting the points are used for visual guidance only and do not represent a fit to the data. Shaded regions represent model discovery accuracy above 80%.

representing a dataset entirely composed of random noise. This scenario sets the ultimate test stage for an algorithm: identifying dynamical systems without signal, where we expect success rates to drop to zero. When faced with this condition, an effective algorithm should identify either a null model (with no coefficients) or a dense model (with many terms from the candidate library). However, if the algorithm incorrectly identifies canonical PDEs from pure white noise data, it indicates that further improvements are needed to prevent such misidentifications and ensure the robustness of the method.

We generate 100 white Gaussian noise datasets, each consisting of 2000 spatial ( $x$ ) and 1000 temporal ( $t$ ) data points, forming a matrix in  $\mathbb{R}^{2000 \times 1000}$ . To investigate the influence of noise variance on the identification process, we use three Gaussian distributions with variances spanning three orders of magnitude:  $N(0, 0.1^2)$ ,  $N(0, 1)$ , and  $N(0, 10^2)$ . We aim to determine whether ARGOS-RAL and STRidge can identify canonical PDEs under these noise conditions. Table 1 shows the percentages of different identified models. Based on the PDEs tested by Rudy *et al.* [4] and our own study, we define parsimonious models as those having three or fewer nonzero coefficients, suggesting they may correspond to specific physical phenomena. In particular, we highlight three classic differential equations: the ODE  $u_t = c_1 u^d$ , the transport equation  $u_t = c_2 u_x$ , and the heat equation  $u_t = c_3 u_{xx}$ . In contrast, we classify models with more than three nonzero coefficients as non-parsimonious, indicating that their coefficient vectors have a dense composition.

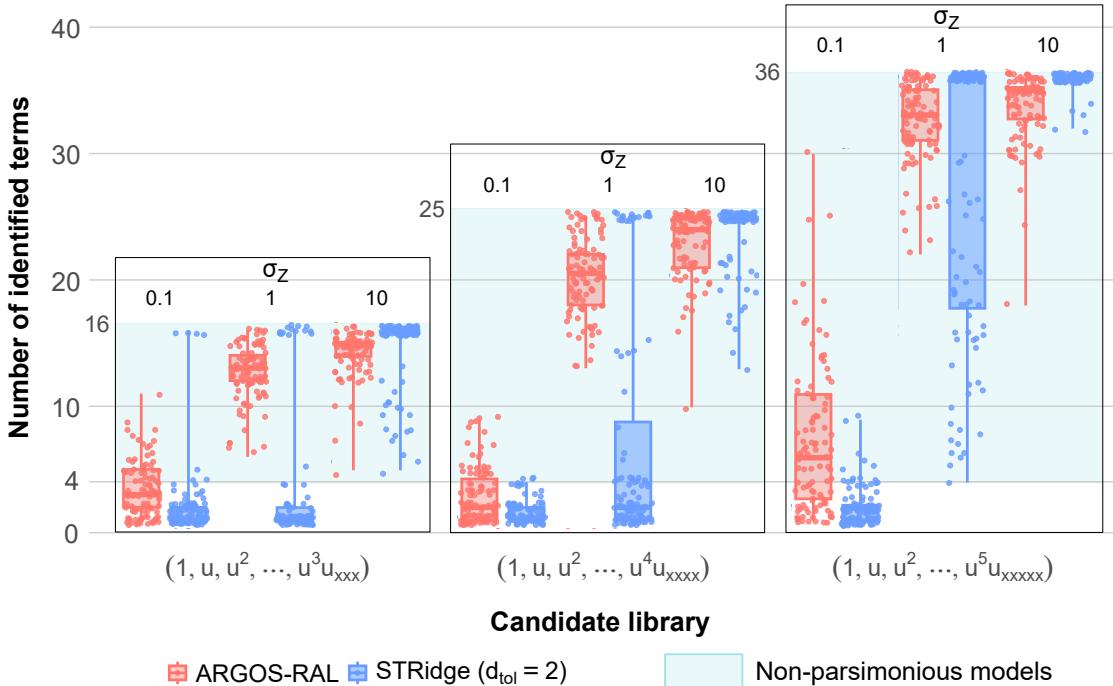


Figure 5: Number of nonzero terms identified from 100 random noise datasets using different candidate function libraries. For each case, we count the number of nonzero coefficients in the sparse regression. We display the distribution of these counts using dots for each of the 100 trials and summarize the results using box plots. Each box plot shows the median (solid horizontal line), interquartile range (box), and minimum and maximum values (whiskers) for the 100 trials. The optimal algorithm should produce boxes located either at zero, indicating a null model, or above four, representing a dense model. The box may span a wide range from four to the maximum number of terms in the library.

Table 1 and Fig. 5 demonstrate that as the standard deviation of the Gaussian noise increases, both ARGOS-RAL and STRidge tend to identify more non-parsimonious models, as indicated by the probability distributions of the number of identified terms shifting into the shaded region of Fig. 5. This is the desired behavior when the input signal is pure white noise, as we want to ensure that the algorithms do not identify parsimonious models in such cases. The difference in behavior between the two methods is most apparent when the noise level is low to moderate ( $\sigma_Z \leq 1$ ). In these cases, STRidge's distributions are more spread out and partially located in the parsimonious region, while ARGOS-RAL's distributions are more concentrated in the non-parsimonious region. This suggests that ARGOS-RAL is more effective at avoiding the identification of parsimonious models when the input signal is pure white noise with low to moderate noise levels. As the noise level increases to  $\sigma_Z = 10$ , both ARGOS-RAL and STRidge consistently

identify non-parsimonious models, as evidenced by the concentration of their distributions in the non-parsimonious region of Fig. 5. This indicates that both methods are effective at avoiding the identification of parsimonious models when the input signal is pure white noise with high noise levels.

Candidate library	Parsimonious model (%)				Non-parsimonious model (%)
	ODE $u_t = c_1 u^d$	Transport $u_t = c_2 u_x$	Heat $u_t = c_3 u_{xx}$	Others	
$\sigma_Z = 0.1$					
<b>STRidge</b> ( $d_{tol}=2$ )					
$(1, u, u^2, \dots, u^3 u_{xxx})$	4 (1)	3	2 (1)	82 (11)	9
$(1, u, u^2, \dots, u^4 u_{xxxx})$	0	0	0	91 (10)	9 (2)
$(1, u, u^2, \dots, u^5 u_{xxxxx})$	0	0	0	81 (11)	19 (1)
<b>ARGOS-RAL</b>					
$(1, u, u^2, \dots, u^3 u_{xxx})$	1 (1)	2 (1)	4 (2)	54 (29)	39 (23)
$(1, u, u^2, \dots, u^4 u_{xxxx})$	0	2 (1)	3 (2)	63 (32)	32 (23)
$(1, u, u^2, \dots, u^5 u_{xxxxx})$	0	0	0	34 (18)	66 (47)
$\sigma_Z = 1$					
<b>STRidge</b> ( $d_{tol}=2$ )					
$(1, u, u^2, \dots, u^3 u_{xxx})$	4 (1)	4	4 (2)	69 (8)	19 (1)
$(1, u, u^2, \dots, u^4 u_{xxxx})$	1	0	2	54 (7)	43 (10)
$(1, u, u^2, \dots, u^5 u_{xxxxx})$	0	0	0	0	100 (18)
<b>ARGOS-RAL</b>					
$(1, u, u^2, \dots, u^3 u_{xxx})$	0	0	0	0	100 (16)
$(1, u, u^2, \dots, u^4 u_{xxxx})$	0	0	0	0	100 (17)
$(1, u, u^2, \dots, u^5 u_{xxxxx})$	0	0	0	0	100 (13)
$\sigma_Z = 10$					
<b>STRidge</b> ( $d_{tol}=2$ )					
$(1, u, u^2, \dots, u^3 u_{xxx})$	0	0	0	0	100 (11)
$(1, u, u^2, \dots, u^4 u_{xxxx})$	0	0	0	0	100 (4)
$(1, u, u^2, \dots, u^5 u_{xxxxx})$	0	0	0	0	100 (12)
<b>ARGOS-RAL</b>					
$(1, u, u^2, \dots, u^3 u_{xxx})$	0	0	0	0	100 (12)
$(1, u, u^2, \dots, u^4 u_{xxxx})$	0	0	0	0	100 (9)
$(1, u, u^2, \dots, u^5 u_{xxxxx})$	0	0	0	0	100 (6)

Table 1: Models identified from random noise by ARGOS-RAL and STRidge. We construct the candidate library with monomials and derivatives of orders ranging from three to five. We define parsimonious models as having three or fewer nonzero coefficients. We evaluate each identified model with an F-test to determine statistical significance, with the number of significant models (p-value < 0.05) noted in parentheses. Numbers outside parentheses indicate the number of models that did not significantly differ from the null hypothesis according to the F-test.  $c_1, c_2, c_3$  are constants. For the ordinary differential equation (ODE) models, the monomial order  $d$  is a positive integer, with the maximum order corresponding to the highest order in the candidate library.

## 4 Conclusions

We designed ARGOS-RAL to automatically tune algorithm hyperparameters, enabling the identification of closed forms of PDEs directly from data. ARGOS-RAL offers several advantages over existing PDE identification methods. First, it automates the process of calculating partial derivatives and constructing the candidate library, reducing manual intervention and streamlining the modeling process. Second, the recurrent adaptive lasso employed by ARGOS-RAL provides a more robust and efficient sparse regression technique compared to the STRidge used in SINDy-based methods. This enables ARGOS-RAL to handle noisy and limited data more effectively, as demonstrated in our numerical experiments.

However, ARGOS-RAL also has some limitations. Like other library-based methods, its effectiveness depends on including the correct governing terms in the candidate library. If the true governing terms are absent, ARGOS-RAL can only approximate the PDE using the available terms, potentially leading to model misspecification. Furthermore, while ARGOS-RAL provides a more computationally efficient approach than ARGOS [12] by focusing on point estimates

rather than bootstrapping for confidence intervals, this comes at the cost of losing uncertainty quantification for the estimated coefficients.

When applying ARGOS-RAL to different scientific domains, several challenges arise. One key challenge is determining the appropriate range of candidate terms to include in the library, which often requires domain expertise. In some fields, the governing equations may involve complex nonlinearities or unconventional terms that are difficult to anticipate without prior knowledge. Another challenge is the computational cost of handling high-dimensional data, which is common in many scientific applications. As the number of variables and the complexity of the PDE increase, the size of the candidate library grows exponentially, leading to increased computational demands for sparse regression.

Despite these challenges, ARGOS-RAL offers a promising framework for automating PDE identification in various scientific domains. By leveraging sparse regression techniques and automating key steps in the modeling process, ARGOS-RAL has the potential to accelerate discovery and insight in fields ranging from physics and engineering to biology and climate science.

## Data availability

All data and codes are available at <https://github.com/Weizhenli/ARGOS-RAL>.

## References

- [1] M. Schmidt and H. Lipson, “Distilling Free-Form Natural Laws from Experimental Data,” *Science*, vol. 324, pp. 81–85, Apr. 2009.
- [2] X. Xun, J. Cao, B. Mallick, A. Maity, and R. J. Carroll, “Parameter Estimation of Partial Differential Equation Models,” *Journal of the American Statistical Association*, vol. 108, pp. 1009–1020, Sept. 2013.
- [3] S. L. Brunton, J. L. Proctor, and J. N. Kutz, “Discovering governing equations from data by sparse identification of nonlinear dynamical systems,” *Proceedings of the National Academy of Sciences*, vol. 113, no. 15, pp. 3932–3937, 2016.
- [4] S. H. Rudy, S. L. Brunton, J. L. Proctor, and J. N. Kutz, “Data-driven discovery of partial differential equations,” *Science Advances*, vol. 3, p. e1602614, Apr. 2017.
- [5] M. Raissi, P. Perdikaris, and G. Karniadakis, “Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations,” *Journal of Computational Physics*, vol. 378, pp. 686–707, Feb. 2019.
- [6] P. Y. Lu, J. Ariño Bernad, and M. Soljačić, “Discovering sparse interpretable dynamics from partial observations,” *Communications Physics*, vol. 5, p. 206, Aug. 2022.
- [7] E. Zhang, M. Dao, G. E. Karniadakis, and S. Suresh, “Analyses of internal structures and defects in materials using physics-informed neural networks,” *Science Advances*, vol. 8, p. eabk0644, Feb. 2022.
- [8] Y.-X. Jiang, X. Xiong, S. Zhang, J.-X. Wang, J.-C. Li, and L. Du, “Modeling and prediction of the transmission dynamics of COVID-19 based on the SINDy-LM method,” *Nonlinear Dynamics*, vol. 105, pp. 2775–2794, Aug. 2021.
- [9] S. Maddu, B. L. Cheeseman, I. F. Sbalzarini, and C. L. Müller, “Stability selection enables robust learning of differential equations from limited noisy data,” *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, vol. 478, p. 20210916, June 2022.
- [10] Y. Cai, X. Wang, G. Joós, and I. Kamwa, “An Online Data-Driven Method to Locate Forced Oscillation Sources From Power Plants Based on Sparse Identification of Nonlinear Dynamics (SINDy),” *IEEE Transactions on Power Systems*, vol. 38, pp. 2085–2099, May 2023.
- [11] X. Sun, J. Qian, and J. Xu, “Compressive-sensing model reconstruction of nonlinear systems with multiple attractors,” *International Journal of Mechanical Sciences*, vol. 265, p. 108905, Mar. 2024.
- [12] K. Egan, W. Li, and R. Carvalho, “Automatically discovering ordinary differential equations from data with sparse regression,” *Communications Physics*, vol. 7, pp. 1–10, Jan. 2024.
- [13] J. M. Varah, “A Spline Least Squares Method for Numerical Parameter Estimation in Differential Equations,” *SIAM Journal on Scientific and Statistical Computing*, vol. 3, pp. 28–46, Mar. 1982.
- [14] J. P. Crutchfield and B. S. McNamara, “Equations of motion from a data series,” *Complex Systems*, vol. 1, pp. 417–452, 1987.

- [15] M. Bär, R. Hegger, and H. Kantz, “Fitting partial differential equations to space-time dynamics,” *Physical Review E*, vol. 59, pp. 337–342, Jan. 1999.
- [16] T. Müller and J. Timmer, “Fitting parameters in partial differential equations from partially observed noisy data,” *Physica D: Nonlinear Phenomena*, vol. 171, pp. 1–7, Oct. 2002.
- [17] H. Liang and H. Wu, “Parameter Estimation for Differential Equation Models Using a Framework of Measurement Error in Regression Models,” *Journal of the American Statistical Association*, vol. 103, pp. 1570–1583, Dec. 2008.
- [18] H. Wu, A. A. Ding, and V. De Gruttola, “Estimation of HIV dynamic parameters,” *Statistics in Medicine*, vol. 17, pp. 2463–2485, Nov. 1998.
- [19] H. Wu and A. A. Ding, “Population HIV-1 Dynamics In Vivo: Applicable Models and Inferential Tools for Virological Data from AIDS Clinical Trials,” *Biometrics*, vol. 55, pp. 410–418, June 1999.
- [20] H. Putter, S. H. Heisterkamp, J. M. A. Lange, and F. de Wolf, “A Bayesian approach to parameter estimation in HIV dynamical models,” *Statistics in Medicine*, vol. 21, pp. 2199–2214, Aug. 2002.
- [21] J. Bongard and H. Lipson, “Automated reverse engineering of nonlinear dynamical systems,” *Proceedings of the National Academy of Sciences*, vol. 104, p. 9943, June 2007.
- [22] S.-M. Udrescu and M. Tegmark, “AI Feynman: A physics-inspired method for symbolic regression,” *Science Advances*, vol. 6, p. eaay2631, Apr. 2020.
- [23] H. Xu and D. Zhang, “Robust discovery of partial differential equations in complex situations,” *Physical Review Research*, vol. 3, p. 033270, Sept. 2021.
- [24] S. Rudy, A. Alla, S. L. Brunton, and J. N. Kutz, “Data-Driven Identification of Parametric Partial Differential Equations,” *SIAM Journal on Applied Dynamical Systems*, vol. 18, pp. 643–660, Jan. 2019.
- [25] K. Kaheman, J. N. Kutz, and S. L. Brunton, “SINDy-PI: a robust algorithm for parallel implicit sparse identification of nonlinear dynamics,” *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, vol. 476, no. 2242, p. 20200279, 2020.
- [26] D. A. Messenger and D. M. Bortz, “Weak SINDy for partial differential equations,” *Journal of Computational Physics*, vol. 443, p. 110525, Oct. 2021.
- [27] A. Cortiella, K. C. Park, and A. Doostan, “Sparse identification of nonlinear dynamical systems via reweighted  $\ell_1$ -regularized least squares,” *Computer Methods in Applied Mechanics and Engineering*, vol. 376, p. 113620, Apr. 2021.
- [28] U. Fasel, J. N. Kutz, B. W. Brunton, and S. L. Brunton, “Ensemble-SINDy: Robust sparse model discovery in the low-data, high-noise limit, with active learning and control,” *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, vol. 478, p. 20210904, Apr. 2022.
- [29] Y. Li, K. Wu, and J. Liu, “Discover governing differential equations from evolving systems,” *Physical Review Research*, vol. 5, p. 023126, May 2023.
- [30] J.-C. Loiseau, B. R. Noack, and S. L. Brunton, “Sparse reduced-order modelling: sensor-based dynamics to full-state estimation,” *Journal of Fluid Mechanics*, vol. 844, pp. 459–490, June 2018.
- [31] K. Duraisamy, G. Iaccarino, and H. Xiao, “Turbulence Modeling in the Age of Data,” *Annual Review of Fluid Mechanics*, vol. 51, pp. 357–377, Jan. 2019.
- [32] S. Li, E. Kaiser, S. Laima, H. Li, S. L. Brunton, and J. N. Kutz, “Discovering time-varying aerodynamics of a prototype bridge by sparse identification of nonlinear dynamical systems,” *Physical Review E*, vol. 100, p. 022220, Aug. 2019.
- [33] 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*, vol. 2, no. 1, pp. 52–63, 2016.
- [34] M. Hoffmann, C. Fröhner, and F. Noé, “Reactive SINDy: Discovering governing reactions from concentration data,” *The Journal of Chemical Physics*, vol. 150, p. 025101, Jan. 2019.
- [35] J. H. Lagergren, J. T. Nardini, G. Michael Lavigne, E. M. Rutter, and K. B. Flores, “Learning partial differential equations for biological transport models from noisy spatio-temporal data,” *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, vol. 476, p. 20190800, Feb. 2020.
- [36] Z. Chen, Y. Liu, and H. Sun, “Physics-informed learning of governing equations from scarce data,” *Nature Communications*, vol. 12, p. 6136, Oct. 2021.
- [37] Z. Zhang and Y. Liu, “A robust framework for identification of PDEs from noisy data,” *Journal of Computational Physics*, vol. 446, p. 110657, Dec. 2021.

- [38] P. Thanasutives, T. Morita, M. Numao, and K.-i. Fukui, “Noise-aware physics-informed machine learning for robust PDE discovery,” *Machine Learning: Science and Technology*, vol. 4, p. 015009, Mar. 2023.
- [39] D. Jia, X. Zhou, S. Li, S. Liu, and H. Shi, “Governing equation discovery based on causal graph for nonlinear dynamic systems,” *Machine Learning: Science and Technology*, vol. 4, p. 045008, Oct. 2023.
- [40] A. Savitzky and M. J. E. Golay, “Smoothing and Differentiation of Data by Simplified Least Squares Procedures.,” *Analytical Chemistry*, vol. 36, pp. 1627–1639, July 1964.
- [41] F. V. Breugel, J. N. Kutz, and B. W. Brunton, “Numerical Differentiation of Noisy Data: A Unifying Multi-Objective Optimization Framework,” *IEEE Access*, vol. 8, pp. 196865–196877, 2020.
- [42] R. W. Schafer, “What Is a Savitzky-Golay Filter? [Lecture Notes],” *IEEE Signal Processing Magazine*, vol. 28, pp. 111–117, July 2011.
- [43] H. Zou, “The Adaptive Lasso and Its Oracle Properties,” *Journal of the American Statistical Association*, vol. 101, no. 476, pp. 1418–1429, 2006.
- [44] J. H. Friedman, T. Hastie, and R. Tibshirani, “Regularization Paths for Generalized Linear Models via Coordinate Descent,” *Journal of Statistical Software, Articles*, vol. 33, no. 1, pp. 1–22, 2010.
- [45] P. C. Hansen, “Analysis of Discrete Ill-Posed Problems by Means of the L-Curve,” *SIAM Review*, vol. 34, pp. 561–580, Dec. 1992.
- [46] J. Nasehi Tehrani, A. McEwan, C. Jin, and A. van Schaik, “L1 regularization method in electrical impedance tomography by using the L1-curve (Pareto frontier curve),” *Applied Mathematical Modelling*, vol. 36, pp. 1095–1105, Mar. 2012.
- [47] A. Cultrera and L. Callegaro, “A simple algorithm to find the L-curve corner in the regularisation of ill-posed inverse problems,” *IOP SciNotes*, vol. 1, p. 025004, Aug. 2020.
- [48] P. Bühlmann and S. van de Geer, *Statistics for High-Dimensional Data*. Springer Series in Statistics, Berlin, Heidelberg: Springer Berlin Heidelberg, 2011.
- [49] R. Tibshirani, “Regression Shrinkage and Selection Via the Lasso,” *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.
- [50] Y. Yang, “Can the strengths of AIC and BIC be shared? A conflict between model identification and regression estimation,” *Biometrika*, vol. 92, pp. 937–950, Dec. 2005.
- [51] K. Aho, D. Derryberry, and T. Peterson, “Model selection for ecologists: the worldviews of AIC and BIC,” *Ecology*, vol. 95, pp. 631–636, Mar. 2014.
- [52] R. C. Gonzalez and R. E. Woods, “Smoothing (Lowpass) Spatial Filters,” in *Digital image processing*, pp. 164–175, New York, NY: Pearson, 2018.
- [53] M. C. Cross and P. C. Hohenberg, “Pattern formation outside of equilibrium,” *Reviews of Modern Physics*, vol. 65, pp. 851–1112, July 1993.
- [54] S. L. Brunton and J. N. Kutz, “Fourier and Wavelet Transforms,” in *Data-Driven Science and Engineering: Machine Learning, Dynamical Systems, and Control*, pp. 53–96, Cambridge: Cambridge University Press, 2 ed., 2022.
- [55] K. Taira and T. Colonius, “The immersed boundary method: A projection approach,” *Journal of Computational Physics*, vol. 225, no. 2, pp. 2118–2137, 2007.
- [56] T. Colonius and K. Taira, “A fast immersed boundary method using a nullspace approach and multi-domain far-field boundary conditions,” *Immersed Boundary Method and Its Extensions*, vol. 197, no. 25, pp. 2131–2146, 2008.
- [57] A. D. Polyanin and V. F. Zaitsev, “Third-Order Equations,” in *Handbook of Nonlinear Partial Differential Equations*, pp. 857–976, Chapman and Hall/CRC, 2 ed., 2012.

## A Supplementary materials

### A.1 Gaussian Blur Kernels

The Gaussian blur convolves data with a Gaussian kernel to smooth it, regardless of the data's dimensionality. This convolution method offers significant benefits for filtering out Gaussian noise, a common noise distribution encountered in data analysis [52]. For one-dimensional spatial PDEs, such as the Burgers' and cable equations, we employ the simplest 2-dimensional Gaussian kernel:

$$\frac{1}{16} \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \otimes [1 \ 2 \ 1] = \frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}. \quad (9)$$

In contrast, for two-dimensional spatial PDEs, the Navier-Stokes and reaction-diffusion equations, we use a 3-dimensional Gaussian kernel:

$$\frac{1}{64} \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \otimes [1 \ 2 \ 1] \otimes \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} = \frac{1}{64} \left[ \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix} \begin{bmatrix} 2 & 4 & 2 \\ 4 & 8 & 4 \\ 2 & 4 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix} \right]. \quad (10)$$

### A.2 Algorithms

Here, we detail the automatic Savitzky-Golay Filter and the recurrent adaptive lasso with the Pareto curve and AIC.

---

#### Algorithm 1: Automatic Savitzky-Golay Filter

---

**Input:**  $\mathbf{U} \in \mathbb{R}^{n \times m}$  or  $\mathbb{C}^{n \times m}$ ,  $dt$ ,  $dx$ .

**Output:** partial derivatives  $\mathbf{U}_t$ ,  $\mathbf{U}_x$ ,  $\mathbf{U}_{xx}$ , ... .

- 1  $\mathbf{U}_{GB} = \text{Gaussian\_Blur}(\mathbf{U}); // \text{ use Gaussian blurred data as the ground truth};$
  - 2  $(o_t^*, l_t^*) = \arg \min_{o, l} \text{MSE}(\text{Savitzky-Golay}(\tilde{\mathbf{U}}(t), o, l), \mathbf{U}_{GB});$
  - 3  $(o_x^*, l_x^*) = \arg \min_{o, l} \text{MSE}(\text{Savitzky-Golay}(\tilde{\mathbf{U}}(x), o, l), \mathbf{U}_{GB});$
  - 4  $\mathbf{U}_t = \text{Savitzky-Golay}(\mathbf{U}_{GB}, o_t^*, l_t^*, \text{derivative}=1);$
  - 5  $\mathbf{U}_x = \text{Savitzky-Golay}(\mathbf{U}_{GB}, o_x^*, l_x^*, \text{derivative}=1);$
  - 6  $\mathbf{U}_{xx} = \text{Savitzky-Golay}(\mathbf{U}_{GB}, o_x^*, l_x^*, \text{derivative}=2);$
  - 7 :
- 

### A.3 Additional PDE Test Cases

Here, we demonstrate how to solve the PDEs presented in this paper.

#### A.3.1 Burgers' equation

We can derive Burgers' equation from the Navier-Stokes equation for the velocity field by dropping the pressure gradient term. Unlike the Navier-Stokes equation, Burgers' equation does not exhibit turbulent behavior, and we can transform it to linear form via the Cole-Hopf transformation [53]:

$$u_t = -uu_x + \nu u_{xx}. \quad (11)$$

We solve Burgers' Eq. (11) using the Fourier spectral method [54] with the *ode45* function in MATLAB. We set  $\nu = 0.1$ ,  $x \in [-8, 8]$  with 256 points,  $t \in [0, 10]$  with 101 points, and the initial condition is a Gaussian function:  $\exp(-(x+2)^2)$ .

#### A.3.2 Cable equation

The cable equation, shown in Eq. (12), quantitatively describes the electrical behavior of nerve axons and other cable-like structures in biological systems. It captures the electrical circuit of current flow and voltage change both within and between neurons. The equation is derived from a circuit model of the membrane and its intracellular

**Algorithm 2:** The recurrent adaptive lasso with Pareto curve and AIC

---

**Input:**  $\Theta(u) \in \mathbb{R}^{(n \cdot m) \times p}$  or  $\mathbb{C}^{(n \cdot m) \times p}$ ,  $u_t \in \mathbb{R}^{(n \cdot m) \times 1}$  or  $\mathbb{C}^{(n \cdot m) \times 1}$ .

**Output:**  $\hat{\beta}$ 

```

1 for  $\gamma$  in 1:5 do
2    $\mathcal{J}^{(\gamma,0)} = \text{NULL}; // \text{ initialize } \mathcal{J};$ 
3    $k = 1; // \text{ iteration counter};$ 
4    $\mathcal{J}^{(\gamma,k)} = \{1, 2, \dots, p\}; // \text{ selected columns from } \Theta;$ 
5   while  $\mathcal{J}^{(\gamma,k)} \neq \mathcal{J}^{(\gamma,k-1)}$  do
6      $w^{(\gamma,k)} = \left( \arg \min_{\beta_{\mathcal{J}^{(\gamma,k)}}} \|u_t - \Theta(u)_{\mathcal{J}^{(\gamma,k)}} \beta_{\mathcal{J}^{(\gamma,k)}}\|_2^2 \right)^{-\gamma}; // \text{ ols weights};$ 
7      $\hat{\beta}^{(\gamma,k)} = \arg \min_{\beta_{\mathcal{J}^{(\gamma,k)}}} \|u_t - \Theta(u)_{\mathcal{J}^{(\gamma,k)}} \beta_{\mathcal{J}^{(\gamma,k)}}\|_2^2 + \lambda^* \sum_{j=1}^p w_j^{(\gamma,k)} |\beta_{j,\mathcal{J}^{(\gamma,k)}}|;$ 
     //  $\lambda^*$  is the optimal point on the Pareto curve;
8      $\mathcal{A}^{(\gamma,k)} = \text{AIC}(\hat{\beta}^{(\gamma,k)});$ 
9      $\mathcal{J}^{(\gamma,k)} = \left\{ j : \hat{\beta}_j^{(\gamma,k)} \neq 0 \right\}; // \text{ select active terms};$ 
10     $k = k + 1;$ 
11  end
12 end
13  $\mathcal{J}^* = \mathcal{J}^{(\gamma^*,k^*)}$  where  $(\gamma^*, k^*)$  is the index of the minimum  $\mathcal{A}$ ;
14  $\hat{\beta} = \arg \min_{\beta_{\mathcal{J}^*}} \|u_t - \Theta(u)_{\mathcal{J}^*} \beta_{\mathcal{J}^*}\|_2^2;$ 

```

---

and extracellular space. The cable equation plays a crucial role as an important PDE in biophysical studies, helping researchers understand how electrical signals change in diseases and disorders. By identifying the cable equation, researchers can diagnose these negative conditions by checking for changes in capacitances  $c_m$ , resistances  $r_m$ , and axial resistance  $r_a, r_e$ :

$$\lambda^2 \frac{\partial^2 V}{\partial x^2} = \tau \frac{\partial V}{\partial t} + V \quad \text{where} \quad \lambda = \sqrt{\frac{r_m}{r_e + r_a}} \quad \text{and} \quad \tau = r_m c_m. \quad (12)$$

We solve the cable equation using *odeint* function in Python with the Fourier spectral method. We set  $\lambda = 1$ ,  $\tau = 1$ ,  $x \in [-4, 4]$  with  $\Delta x = 0.1$ ,  $t \in [0, 5]$  with  $\Delta t = 0.01$ , and use a Gaussian function  $\exp(-x^2)$  as the initial condition.

### A.3.3 Navier-Stokes

We simulate the two-dimensional Navier-Stokes equation for fluid flow around a circular cylinder using the Immersed Boundary Projection Method [55, 56]. The two-dimensional velocity components are denoted by  $u$  and  $v$ , while  $\omega$  represents the vorticity away from the circular cylinder of diameter one and mass centre at  $(x = 1, y = 2)$ . We set the Reynolds number to 100 and aim to identify the equation

$$\omega_t = 0.01\omega_{xx} + 0.01\omega_{yy} - u\omega_x - v\omega_y. \quad (13)$$

The spatial domain spans  $x \in [0, 9]$  with  $\Delta x = 0.02$ ,  $y \in [0, 4]$  with  $\Delta y = 0.02$ , and the temporal domain covers  $t \in [300, 330]$  with  $\Delta t = 0.02$ . We save the flow data every ten snapshots. This setup generates a simulated dataset containing approximately 13.5 million points ( $449 \times 199 \times 151$ ). However, constructing the candidate library  $\Theta$  for such a large dataset poses computational challenges. To facilitate the evaluation, we randomly sample points within the red rectangular area shown in Fig. 4 C at each snapshot.

### A.3.4 Reaction-diffusion

Reaction-diffusion systems offer a versatile framework to model pattern formation in various natural phenomena in chemistry, biology, geology, physics, and ecology. These systems give rise to a rich tapestry of periodic patterns, including spots, zigzags, spiral waves, and rolls. In our analysis, we focus on a widely studied class of reaction-diffusion systems known as the  $\lambda - \omega$  systems, described by the following coupled PDEs:

$$u_t = 0.1u_{xx} + 0.1u_{yy} + u - uv^2 - u^3 + v^3 + u^2v \quad (14)$$

$$v_t = 0.1v_{xx} + 0.1v_{yy} + v - uv^2 - u^3 - v^3 - u^2v \quad (15)$$

To generate data for our analysis, we employ the simulation method described in [4]. We discretize the spatial domain  $x, y \in [-10, 10]$  using a  $512 \times 512$  grid and evolve the system over the time interval  $t \in [0, 10]$  using 201 time steps. This procedure yields a rich dataset comprising 52,690,944 spatiotemporal points on a  $512 \times 512 \times 201$  grid, providing a comprehensive characterization of the system's dynamics.

### A.3.5 Quantum harmonic oscillator

The quantum harmonic oscillator (QHO) models the parabolic potential of a harmonic oscillator in quantum mechanics. It simulates the time evolution of the wave function associated with a particle in the parabolic potential, providing the probability distribution of the particle's position at any given time by taking the squared magnitude of the wave function. The energy levels of a quantum harmonic oscillator are quantized, meaning they can only assume specific, discrete values. Furthermore, even if we form a statistical distribution from multiple experiments, it will lack information on the intricate phase of the wave function. We use the following equation:

$$u_t = \frac{1}{2}iu_{xx} - iuV = \frac{1}{2}iu_{xx} - \frac{x^2}{2}iu. \quad (16)$$

To obtain data on the QHO, we employ the operator splitting method with the Fourier transform. We consider the time domain  $t \in [0, 10]$  with  $\Delta t = 0.025$ , and the space domain  $x \in [-7.5, 7.5]$  with  $\Delta x = 15/512$ , using a Gaussian  $\exp(-((x - 1)/2)^2)$  as the initial condition. When performing a sparse regression on complex numbers, we transform the regression from complex to real numbers. For each  $y_i = y_i^R + iy_i^I$ , where the normal  $i$  is the imaginary number, the subscript  $i$  represents the  $i$ th observation, we can reform it as

$$\begin{aligned} y_i^R + iy_i^I &= \beta_0^R + i\beta_0^I + \sum_{j=1}^p [(\beta_j^R + i\beta_j^I)(x_{ij}^R + ix_{ij}^I)] + \epsilon^R + i\epsilon^I \\ &= \beta_0^R + i\beta_0^I + (\beta_1^R + i\beta_1^I)(x_{i1}^R + ix_{i1}^I) + (\beta_2^R + i\beta_2^I)(x_{i2}^R + ix_{i2}^I) + \dots + \\ &\quad (\beta_p^R + i\beta_p^I)(x_{ip}^R + ix_{ip}^I) + \epsilon^R + i\epsilon^I \\ &= \beta_0^R + i\beta_0^I + \sum_{j=1}^p (x_{ij}^R \beta_j^R - x_{ij}^I \beta_j^I) + i \sum_{j=1}^p (x_{ij}^R \beta_j^I + x_{ij}^I \beta_j^R) + \epsilon^R + i\epsilon^I \end{aligned}$$

and split it to extract two equations for both real and imaginary parts

$$y_i^R = \beta_0^R + \sum_{j=1}^p x_{ij}^R \beta_j^R - \sum_{j=1}^p x_{ij}^I \beta_j^I + \epsilon^R, \quad y_i^I = \beta_0^I + \sum_{j=1}^p x_{ij}^I \beta_j^R + \sum_{j=1}^p x_{ij}^R \beta_j^I + \epsilon^I. \quad (17)$$

Based on Eq. (17), we can organize the dataset as:

$$Y = \begin{bmatrix} y_1^R \\ y_1^I \\ y_2^R \\ y_2^I \\ \vdots \\ y_n^R \\ y_n^I \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} x_{11}^R & -x_{11}^I & x_{12}^R & -x_{12}^I & \dots & x_{1p}^R & -x_{1p}^I \\ x_{11}^I & x_{11}^R & x_{12}^I & x_{12}^R & \dots & x_{1p}^I & x_{1p}^R \\ x_{21}^R & -x_{21}^I & x_{22}^R & -x_{22}^I & \dots & x_{21}^R & -x_{21}^I \\ x_{21}^I & x_{21}^R & x_{22}^I & x_{22}^R & \dots & x_{21}^I & x_{21}^R \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{n1}^R & -x_{n1}^I & x_{n2}^R & -x_{n2}^I & \dots & x_{n1}^R & -x_{n1}^I \\ x_{n1}^I & x_{n1}^R & x_{n2}^I & x_{n2}^R & \dots & x_{n1}^I & x_{n1}^R \end{bmatrix}.$$

Finally, we implement the recurrent adaptive lasso and STRidge on the re-posed  $Y$  and  $\mathbf{X}$ .

### A.3.6 Advection-diffusion equation

The advection-diffusion equation, which combines advection and diffusion terms, describes the transport and dispersion of quantities such as temperature, substance concentration, or fluid velocity in various scientific and engineering contexts. We can express this equation as follows:

$$c_t = Dc_{xx} - uc_x \quad (18)$$

We solve the advection-diffusion equation using the Fourier spectral method and the `odeint` function in Python. We set the diffusion coefficient  $D = 1$ , the advection velocity  $u = 1$ , and consider the spatial domain  $x \in [-10, 10]$  with resolution  $\Delta x = 0.1$  and the temporal domain  $t \in [0, 10]$  with resolution  $\Delta t = 0.01$ . We use a Gaussian function of the form  $\exp(-(x + 2)^2)$  as an initial condition.

### A.3.7 The KdV equation

The Korteweg–De Vries (KdV) equation describes wave propagation on shallow water surfaces. The KdV equation solution reveals that an isolated traveling wave exhibits linear behavior, but nonlinear interactions emerge when multiple waves are present. Moreover, the dependence of wave velocity on wave amplitude ensures that any solution with multiple amplitudes will display nonlinear behavior, regardless of the interaction. Eq. (19) presents the formula for the KdV equation:

$$u_t = -6uu_x - u_{xxx}. \quad (19)$$

We employ the two-soliton solution [57] to solve the KdV equation:

$$\begin{aligned} w(x, t) &= -2 \frac{\partial^2}{\partial x^2} \ln \left( 1 + B_1 e^{\theta_1} + B_2 e^{\theta_2} + AB_1 B_2 e^{\theta_1 + \theta_2} \right) \\ \theta_1 &= a_1 x - a_1^3 t, \quad \theta_2 = a_2 x - a_2^3 t, \quad A = \left( \frac{a_1 - a_2}{a_1 + a_2} \right)^2 \end{aligned} \quad (20)$$

where  $a_1$ ,  $a_2$ ,  $B_1$ , and  $B_2$  are arbitrary constants. We set the following parameters: 201 time steps ( $n = 201$ ) with  $t \in [0, 20]$ , 512 spatial points ( $m = 512$ ) with  $x \in [-30, -30]$ , and  $a_1 = 0.5$ ,  $a_2 = 1$ ,  $B_1 = 1$ ,  $B_2 = 5$ .

### A.3.8 Transport equation

The transport equation, Eq. (21), plays a fundamental role in science and engineering, describing the spatiotemporal evolution of scalar quantities or vector fields. We solve this PDE using the analytical solution, Eq. (22), with  $c = 3$  to generate the data for our study. The spatial domain spans  $x \in [-5, 1]$  with resolution  $\Delta x = 0.01$ , while the temporal domain covers  $t \in [0, 2]$  with timestep  $\Delta t = 0.01$ .

$$u_t = cu_x, \quad c > 0 \quad (21)$$

$$u(x, t) = \exp(-(x + ct)^2) \quad (22)$$

### A.3.9 Diffusion equation

The diffusion (heat) equation, Eq. (23), plays a crucial role in many scientific and engineering fields, including solid-state physics, materials science, environmental science, and computational fluid dynamics. This equation elucidates the fundamental process of heat diffusion, enabling engineers to gain deep insights into heat conduction, thermal conductivity, and temperature-dependent phenomena in solids and other materials. Here, we use the analytic solution, Eq. (24), to generate the data. We set  $x \in [0, 5]$  with  $\Delta x = 0.01$  and  $t \in [0, 1.5]$  with  $\Delta t = 0.01$ , and choose the initial condition as  $6 \sin(\pi x/L)$ .

$$u_t = 10u_{xx} \quad (23)$$

$$u(x, t) = 6 \sin\left(\frac{\pi x}{L}\right) e^{-k\left(\frac{\pi}{L}\right)^2 t}, \quad k = 10 \quad (24)$$

The diffusion equation and its analytic solution provide a powerful framework for understanding and predicting heat transfer in various systems. By carefully selecting the spatial and temporal domains and the initial condition, we can model a wide range of real-world scenarios and gain valuable insights into the underlying physical processes.

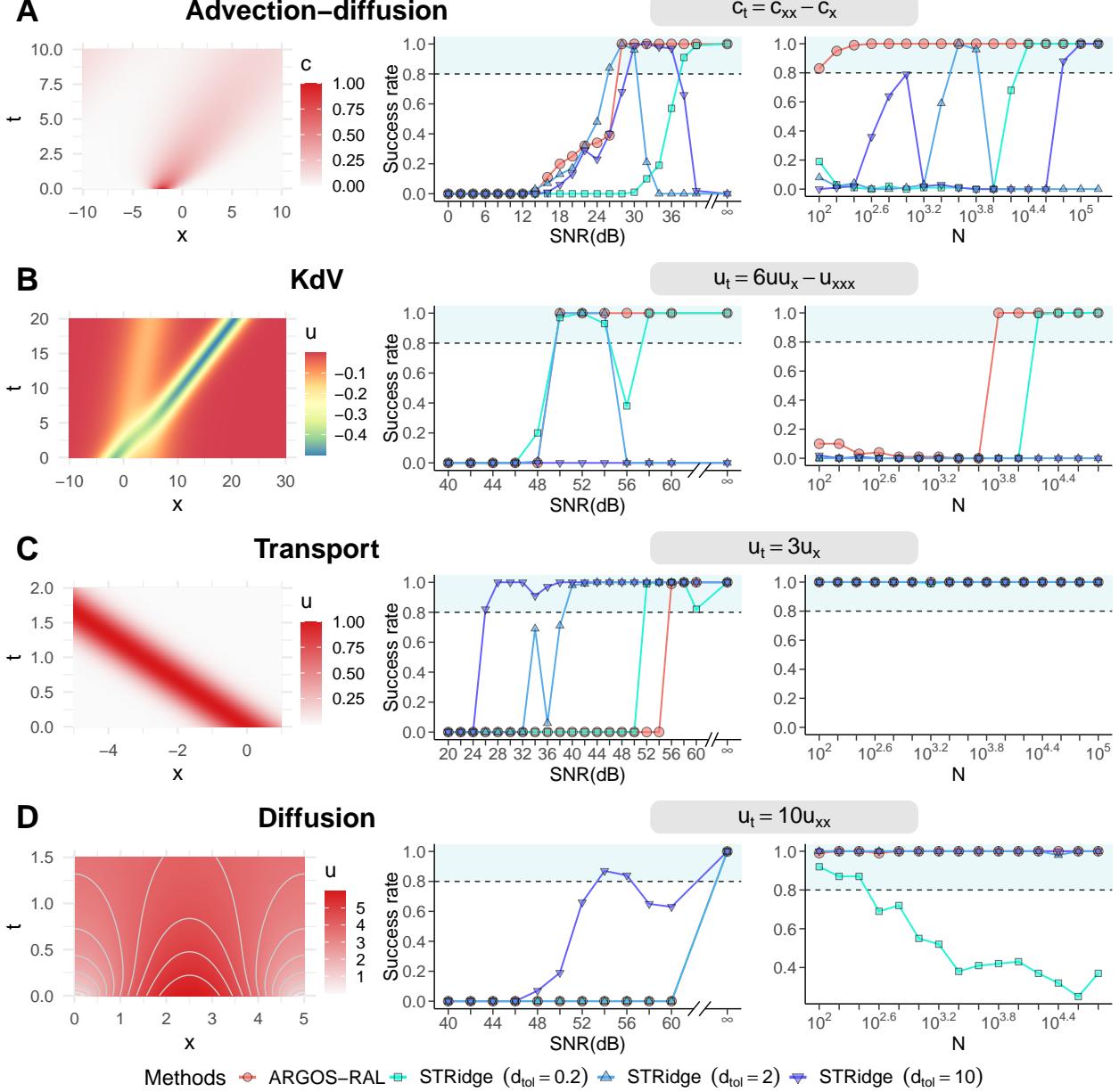


Figure 6: Success rates of ARGOS-RAL and STRidge in identifying (A) advection-diffusion, (B) KdV, (C) transport, and (D) heat equations with different SNRs and sample sizes. To analyze the noise tolerance, we added noise to the PDE solutions at different SNR levels. For the sample size analysis, we randomly sampled points from the set  $\{\mathbf{u}_t, \Theta(\mathbf{u})\}$  based on noiseless PDE solutions. For sample size analysis, we randomly sample points from  $\{\mathbf{u}_t, \Theta(\mathbf{u})\}$  set based on noiseless PDE solutions. Lines are only used to link points, not to fit points. The plots demonstrate that our method maintains high success rates in identifying the correct PDE even under significant noise and with limited sample sizes.