# DATA-DRIVEN DISCOVERY OF PDES VIA THE ADJOINT METHOD

Anonymous authors

003

010 011

012

013

014

015

016

017

018

019

021

023

025

026

Paper under double-blind review

#### ABSTRACT

In this work, we present an adjoint-based method for discovering the underlying governing partial differential equations (PDEs) given data. The idea is to consider a parameterized PDE in a general form and formulate a PDE-constrained optimization problem aimed at minimizing the error of the PDE solution from data. Using variational calculus, we obtain an evolution equation for the Lagrange multipliers (adjoint equations) allowing us to compute the gradient of the objective function with respect to the parameters of PDEs given data in a straightforward manner. In particular, we consider a family of parameterized PDEs encompassing linear, nonlinear, and spatial derivative candidate terms, and elegantly derive the corresponding adjoint equations. We show the efficacy of the proposed approach in identifying the form of the PDE up to machine accuracy, enabling the accurate discovery of PDEs from data. We also compare its performance with the famous PDE Functional Identification of Nonlinear Dynamics method known as PDE-FIND Rudy et al. (2017), on both smooth and noisy data sets. Even though the proposed adjoint method relies on forward/backward solvers, it outperforms PDE-FIND for large data sets thanks to the analytic expressions for gradients of the cost function with respect to each PDE parameter.

## 028 1 INTRODUCTION

A large portion of data-driven modelling of physical processes in literature is dedicated to deploying
 Neural Networks to obtain fast prediction given the training data set. The data-driven estimation
 methods include Physics-Informed Neural Networks Raissi et al. (2019), Pseudo-Hamiltonian neural
 networks Eidnes and Lye (2024), structure preserving Matsubara et al. (2020); Sawant et al. (2023),
 and reduced order modelling Duan and Hesthaven (2024). These methods often provide efficient
 and somewhat "accurate" predictions when tested as an interpolation method in the space of input
 or boundary parameters. Such fast estimators are beneficial when many predictions of a dynamic
 system is needed, for example in the shape optimization task in fluid dynamics.

007

However, the data-driven estimators often fail to provide accurate solution to the dynamical system when tested outside the training space, i.e. for extrapolation. Furthermore, given the 040 regression-based nature of these predictors, often they do not offer any error estimator in prediction. 041 Since we already have access to an arsenal of numerical methods in solving traditional governing 042 equations, it is attractive to learn the underlying governing equation given data instead. Once 043 the governing equation is found, one can use the standard and efficient numerical methods for 044 prediction. This way we guarantee the consistency with observed data, estimator for the numerical approximation, and interoperability. Hence, learning the underlying physics given data has motivated a new branch in the scientific machine learning for discovering the mathematical expression as the 046 governing equation given data. 047

The wide literature of data-driven discovery of dynamical systems includes equation-free modeling
Kevrekidis et al. (2003), artificial neural networks González-García et al. (1998), nonlinear regression
Voss et al. (1999), empirical dynamic modeling Sugihara et al. (2012); Ye et al. (2015), modeling
emergent behavior Roberts (2014), automated inference of dynamics Schmidt et al. (2011); Daniels
and Nemenman (2015a;b), normal form identification in climate Majda et al. (2009), nonlinear
Laplacian spectral analysis Giannakis and Majda (2012) and Koopman analysis Mezić (2013) among
others. There has been a significant advancement in this field by combining symbolic regression

with the evolutionary algorithms Bongard and Lipson (2007); Schmidt and Lipson (2009); Tohme et al. (2022), which enable the direct extraction of nonlinear dynamical system information from data. Furthermore, the concept of sparsity Tibshirani (1996) has recently been employed to efficiently and robustly deduce the underlying principles of dynamical systems Brunton et al. (2016); Mangan et al. (2016).

**Related work.** Next, we review several relevant works that have shaped the current landscape of discovering PDEs from data:

062 **PDE-FIND** Rudy et al. (2017). This method has been developed to discover underlying partial 063 differential equation by minimizing the  $L_2^2$ -norm point-wise error of the parameterized forward model 064 from the data. Estimating all the possible derivatives using Finite Difference, PDE-FIND constructs a 065 dictionary of possible terms and finds the underlying PDE by performing a sparse search using ridge regression problem with hard thresholding, also known as STRidge optimization method. Several 066 further developments in the literature has been carried out based on this idea, namely Champion et al. 067 (2019); Kaheman et al. (2020). In these methods, as the size (or dimension) of the data set increases, 068 the PDE discovery optimization problem based on point-wise error becomes extremely expensive, 069 forcing the user to arbitrarily reduce the size of data by resampling, or compress the data using proper orthogonal decomposition. Needles to say, in case of non-linear dynamics, such truncation of data 071 can introduce bias in prediction leading to finding a wrong PDE. 072

PDE-Net Long et al. (2018; 2019). In this method, the PDE is learned from data using convolution
 kernels rather than brute-force use of Finite Differences, and apply neural networks to approximate
 nonlinear responses. Similar to PDE-FIND, the loss function of PDE-Net is the point-wise error from
 data which leads to a regression task that does not scale well with the size of the data set.

Hidden Physics Models Raissi and Karniadakis (2018). This method assumes that the relevant terms of the governing PDE are already identified and finds its unknown parameters using Gaussian process regression (GPR). While GPR is an accurate interpolator which offers an estimate for the uncertainty in prediction, its training scales poorly with the size of the training data set as it requires inversion of the covariance matrix.

PINN-SR Chen et al. (2021). One of the issues with the PDE-FIND is the use of Finite Difference in estimating the derivatives. The idea of PINN-SR is to extend PDE-FIND's optimization problem to also find a PINN fit to the data in order to find smooth estimates for spatial derivatives. In particular, the training of PINN-SR combines the search for weights/biases of PINN approximation of the PDE with the sparse search in the space of possible terms to find the coefficients of the PDE given data. However, similar to PDE-FIND, the point-wise error from data is used as the loss for the regression task which does not scale well with the size of the data set.

080

**Contributions.** In this paper, we introduce a novel approach for discovering PDEs from data 090 based on the well-known adjoint method, i.e. PDE-constrained optimization method. The idea is to 091 formulate the objective (or cost) functional such that the estimate function f minimizes the  $L_2^2$ -norm 092 error from the data points  $f^*$  with the constraint that f is the solution to a parameterized PDE using the method of Lagrange multipliers. Here, we consider a parameterized PDE in a general form and 094 the task is to find all the parameters including irrelevant ones. By finding the variational extremum of the cost functional with respect to the function f, we obtain a backward-in-time evolution equation 096 for the Lagrange multipliers (adjoint equations). Next, we solve the forward parameterized PDE as 097 well as the adjoint equations numerically. Having found estimates of the Lagrange multipliers and 098 solution to the forward model f, we can numerically compute the gradient of the objective function with respect to the parameters of PDEs given data in a straightforward manner. In particular, for a family of parameterized and nonlinear PDEs, we show how the corresponding adjoint equations can 100 be elegantly derived. We note that the adjoint method has been successfully used before as an efficient 101 method for uncertainty quantification Flath et al. (2011), shape optimization and sensitivity analysis 102 method in fluid mechanics Jameson (2003); Caflisch et al. (2021) and plasma physics Antonsen et al. 103 (2019); Geraldini et al. (2021). Unlike the usual use of PDE-constrained adjoint optimization where 104 the governing equation is known, in this paper we are interested in finding the form along with the 105 coefficients of the PDE given data. 106

107 The remainder of the paper is organized as follows. First in Section 2, we introduce and derive the proposed adjoint-based method of finding the underlying system of PDEs given data. Next in Section

3, we present our results on a wide variety of PDEs and compare the solution with the celebrated PDE-FIND in terms of error and computational/training time. In Section 7, we discuss the limitations for the current version of our approach and provide concluding remarks in Section 8.

111 112

113 114

#### 2 ADJOINT METHOD FOR FINDING PDES

In this section, we introduce the problem and derive the proposed adjoint method for finding governing equations given data.

115 116

129 130

136

137

151

158 159 160

117 **Problem setup.** Assume we are given a data set on a spatial/temporal grid  $\mathcal{G} = \bigcup_{j=0}^{N_t} \mathcal{G}^{(j)}$  with 118  $\mathcal{G}^{(j)} = \{(\boldsymbol{x}^{(k)}, t^{(j)}) \mid k = 1, ..., N_{\boldsymbol{x}}\}$  for the vector of functions  $\boldsymbol{f}^*$  where k is the spatial index 120 and j the time index with  $t^{(N_t)} = T$  being the final time. Here,  $\boldsymbol{x}^{(k)} \in \Omega \subset \mathbb{R}^n$  is coordinates 121 inside the solution domain  $\Omega$ ,  $t^{(j)}$  denotes the j-th time that data is available, and output is a discrete 122 map  $\boldsymbol{f}^* : \mathcal{G} \to \mathbb{R}^N$ . The goal is to find the governing equations that accurately estimates  $\boldsymbol{f}^*$  at all 123 points on  $\mathcal{G}$ . In order to achieve this goal, we formulate the problem using the method of Lagrange 124

125 126 126 127 128 Adjoint method. For simplicity, let us first consider only the time interval  $t \in [t^{(j)}, t^{(j+1)}]$ . Consider a general a forward model  $\mathcal{L}[\cdot]$  that evolves an N-dimensional vector of continuous functions f( $x, t = t^{(j)}$ ) in  $t \in (t^{(j)}, t^{(j+1)}]$  and  $x \in \Omega$  where the *i*-th PDE is given by

$$\mathcal{L}_{i}[\boldsymbol{f}] := \partial_{t} f_{i} + \sum_{\boldsymbol{d},\boldsymbol{p}} \alpha_{i,\boldsymbol{d},\boldsymbol{p}} \nabla_{\boldsymbol{x}}^{(\boldsymbol{d})}[\boldsymbol{f}^{\boldsymbol{p}}] = 0$$
(1)

for i = 1, ..., N, resulting in a system of N-PDEs, i.e. the *i*-th PDE  $\mathcal{L}_i$  predicts  $f_i$ . Here,  $x = [x_1, x_2, ..., x_n]$  is an n-dimensional (spatial) input vector, and  $f = [f_1, f_2, ..., f_N]$  is an N-dimensional vector of functions. We use the shorthand  $f_i = f_i(x, t)$  and f = f(x, t). Furthermore,  $p = [p_1, ..., p_N]$  and  $d = [d_1, d_2, ..., d_n]$  are non-negative index vectors such that  $f^p = f_1^{p_1} f_2^{p_2} \cdots f_N^{p_N}$  and

$$\nabla_{\boldsymbol{x}}^{(\boldsymbol{d})} = \nabla_{x_1}^{(d_1)} \nabla_{x_2}^{(d_2)} \cdots \nabla_{x_n}^{(d_n)} , \qquad (2)$$

where  $\nabla_{x_i}^{(d_i)}$  for i = 1, ..., n indicates  $d_i$ -th derivative in  $x_i$  dimension, and  $\partial_t f_i$  denotes the time derivative of the *i*-th function. We denote the vector of unknown parameters by  $\boldsymbol{\alpha} = [\alpha_{i,d,p}]_{(i,d,p)\in\mathcal{D}}$ , where  $\mathcal{D}$  represents the domain of all valid combinations of *i*, *d*, and *p*.

Having written the forward model equation 1 as general as possible, the goal is to find the parameters a such that f approximates the data points of  $f^*$  at  $t = t^{(j+1)}$  given the solution  $f = f^*$  at  $t = t^{(j)}$ . To this end, we formulate a semi-discrete objective (or cost) functional that minimizes the  $L_2^2$ -norm error between what the model predicts and the data  $f^*$  on  $\mathcal{G}^{(j+1)}$ , with the constraint that f solves the forward model in Eq. (1), i.e.

$$\mathcal{C} = \sum_{i=1}^{N} \left( \sum_{k} \left( f_i^* (\boldsymbol{x}^{(k)}, t^{(j+1)}) - f_i (\boldsymbol{x}^{(k)}, t^{(j+1)}) \right)^2 + \frac{1}{\Delta \boldsymbol{x} \Delta t} \int \lambda_i(\boldsymbol{x}, t) \mathcal{L}_i[\boldsymbol{f}(\boldsymbol{x}, t)] d\boldsymbol{x} dt \right) + \epsilon_0 ||\boldsymbol{\alpha}||_2^2,$$
(3)

where  $\Delta x$  and  $\Delta t$  denotes grid spacing in  $\Omega$  and step size in t, respectively,  $||.||_2$  denotes  $L_2$ -norm, and  $\epsilon_0$  is the regularization factor. We note that PDE discovery task is ill-posed since the underlying PDE is not unique and the regularization term helps us find the PDE with the least possible coefficients.

156 Clearly, given estimates of f and Lagrange multipliers  $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_N)$ , the gradient 157 of the cost function with respect to model parameters can be simply computed via

$$\frac{\partial \mathcal{C}}{\partial \alpha_{i,d,p}} = (-1)^{|\boldsymbol{d}|} \frac{1}{\Delta \boldsymbol{x} \Delta t} \int \boldsymbol{f}^{\boldsymbol{p}} \nabla_{\boldsymbol{x}}^{(\boldsymbol{d})} [\lambda_i] d\boldsymbol{x} dt + 2\epsilon_0 \alpha_{i,d,p}$$
(4)

where i = 1, ..., N and  $|\mathbf{d}| = d_1 + ... + d_n$ , where |.| denotes  $L_1$ -norm. Here, we used integration by parts and imposed the condition that  $\lambda \to 0$  on the boundaries of  $\Omega$  at all time  $t \in [0, T]$ . The analytical expression equation 4 can be used for finding the parameters of PDE using in the gradient descent method with update rule

$$\alpha_{i,d,p} \leftarrow \alpha_{i,d,p} - \eta \frac{\partial \mathcal{C}}{\partial \alpha_{i,d,p}}$$
(5)

for i = 1, ..., N, where  $\eta = \beta \min(\Delta x)^{|d| - d_{\max}}$  is the learning rate which includes a free parameter  $\beta$  and scaling coefficient for each term of the PDE, and  $d_{\max} = \max(|d|)$  for all considered d. Let us also define  $p_{\max} = \max(|p|)$  as the highest order in the forward PDE model. We note that since the terms of the PDEs may have different scaling, the step size for the corresponding coefficient must be adjusted accordingly. This is due the fact that the gradient of the cost function is most sensitive to the highest order terms of the PDE. In Appendix E, we give a justification for our choice of the learning rate  $\eta$ .

174 However, before we can use Eq. (4) and (5), we need to find  $\lambda$ , hence the adjoint equation. This can 175 be achieved by finding the functional extremum of the cost functional C with respect to f. First, we 176 note that the semi-descrete total variation of C can be derived as

$$\delta \mathcal{C} = \sum_{i=1}^{N} \left( -\sum_{k} 2 \left( f_{i}^{*} \left( \boldsymbol{x}^{(k)}, t^{(j+1)} \right) - f_{i} \left( \boldsymbol{x}^{(k)}, t^{(j+1)} \right) \right) \delta f_{i, \boldsymbol{x}^{(k)}, t^{(j+1)}} + \frac{1}{\Delta \boldsymbol{x} \Delta t} \int \left( -\frac{\partial \lambda_{i}}{\partial t} + \sum_{\boldsymbol{d}, \boldsymbol{p}} (-1)^{|\boldsymbol{d}|} \alpha_{i, \boldsymbol{d}, \boldsymbol{p}} \nabla_{f_{i}} [\boldsymbol{f}^{\boldsymbol{p}}] \nabla_{\boldsymbol{x}}^{(\boldsymbol{d})} [\lambda_{i}] \right) \delta f_{i} d\boldsymbol{x} dt + \sum_{k} \lambda_{i} \left( \boldsymbol{x}^{(k)}, t^{(j+1)} \right) \delta f_{i, \boldsymbol{x}^{(k)}, t^{(j+1)}} \right)$$
(6)

187

188 189

190

194

196

199

213

214

181 182

177 178 179

165

166

where  $\delta f_i$  denotes variation with respect to  $f_i(\boldsymbol{x}, t)$ , and  $\delta f_{i,\boldsymbol{x}^{(k)},t^{(j+1)}}$  variation with respect to  $f_i(\boldsymbol{x} = \boldsymbol{x}^{(k)}, t = t^{(j+1)})$ . In this derivation, we descretized the last integral resulting from integration by parts in time using the same mesh as the one of data  $\mathcal{G}^{(j+1)}$ . Here again, we used integration by parts and imposed the condition that  $\boldsymbol{\lambda} \to \boldsymbol{0}$  on the boundaries of  $\Omega$  at all time  $t \in [t^{(j)}, t^{(j+1)}]$  for i = 1, ..., N. Note that  $f_i(\boldsymbol{x}, t)$  is the output of *i*-th PDE.

Next, we find the optimums of C (and hence the adjoint equations) by taking the variational derivatives with respect to  $f_i$  and  $f_{i,\boldsymbol{x}^{(k)},t^{(j+1)}}$ , i.e.

$$\frac{\delta \mathcal{C}}{\delta f_i} = 0 \implies \frac{\partial \lambda_i}{\partial t} = \sum_{\boldsymbol{d}, \boldsymbol{p}} (-1)^{|\boldsymbol{d}|} \alpha_{i, \boldsymbol{d}, \boldsymbol{p}} \nabla_{f_i} [\boldsymbol{f}^{\boldsymbol{p}}] \nabla_{\boldsymbol{x}}^{(\boldsymbol{d})} [\lambda_i]$$
(7)

197

and

$$\frac{\delta \mathcal{C}}{\delta f_{i,\boldsymbol{x}^{(k)},t^{(j+1)}}} = 0 \implies \lambda_i(\boldsymbol{x}^{(k)},t^{(j+1)}) = 2(f_i^*(\boldsymbol{x}^{(k)},t^{(j+1)}) - f_i(\boldsymbol{x}^{(k)},t^{(j+1)}))$$
(8)

200 for i = 1, ..., N and  $j = 0, ..., N_t - 1$ . We note that the adjoint equation equation 7 for the system of 201 PDEs is backward in time with the final condition at the time  $t = t^{(j+1)}$  given by Eq. (8). In order 202 to make the notation clear, we present examples for deriving the adjoint equations in Appendix F. 203 The adjoint equation is in the continuous form, while the final condition is on the discrete points, 204 i.e. on the grid  $\mathcal{G}^{(j+1)}$ . In order to obtain the Lagrange multipliers in  $t \in [t^{(j+1)}, t^{(j)})$ , a numerical 205 method appropriate for the forward equation 1 and adjoint equation equation 7 should be deployed. 206 Furthermore, the adjoint equation should have the same or coarser spatial discretization as  $\mathcal{G}^{(j+1)}$  to 207 enforce the final condition equation 8. 208

**Training with smooth data set.** The training procedure follows the standard gradient descent method. We start by taking an initial guess for parameters  $\alpha$ , e.g. here we take  $\alpha = 0$  initially. For each time interval  $t \in [t^{(j)}, t^{(j+1)}]$ , first we solve the forward model equation 1 numerically to estimate  $f(x^{(k)}, t^{(j+1)})$  given the initial condition

$$f(x^{(k)}, t^{(j)}) = f^*(x^{(k)}, t^{(j)}).$$
(9)

Then, the adjoint Eq. equation 7 is solved backwards in time with the final time condition equation 8. Finally, the estimate for parameters of the model is updated using Eq. equation 5. We repeat this for 216 all time intervals  $j = 0, ..., N_t - 1$  until convergence. In order to improve the search for coefficients 217 and enforce the PDE identification, we also deploy thresholding, i.e. set  $\alpha_{i,d,p} = 0$  if  $|\alpha_{i,d,p}| < \sigma$ 218 where  $\sigma$  is a user-defined threshold, during and at the end of training, respectively. In Algorithm 1, 219 we present a pseudocode for finding the parameters of the system of PDEs using the Adjoint method 220 (a flowchart is also shown in Fig. 5 of Appendix A). For the introduced hyperparameters, we note that  $\beta$  in the learning rate needs to be small enough to avoid unstable intermediate guessed PDEs,  $\epsilon_0$ 221 must be large enough to ensure uniqueness in cases where more than one solution may exist, and the 222 thresholding should be applied only when solution to the optimization is not improving anymore 223 up to a user-defined tolerance  $\gamma_{\rm thr}$ . For suggested default values, please see the description of the 224 algorithm. For experimental investigation on impact of these parameters for a few examples, see 225 Appendix D. 226

227 We note that the type of guessed PDE may change during the training, which adds numeri-228 cal complexity to the optimization and motivates the use of an appropriate solver for each type of 229 guessed PDE, e.g. Finite Volume method for hyperbolic and Finite Element method for Elliptic 230 PDEs. For simplicity, in this work we use the second order Finite Difference method across the board 231 to estimate the spatial and Euler for the time derivative with small enough time step sizes in solving the forward/backward equations to avoid blowups due to possible instabilities. We note that the 232 adjoint method is most effective when there is some prior knowledge of the underlying PDE type, 233 and a suitable numerical method is deployed. 234

Algorithm 1: Finding system of PDEs using Adjoint method. Default threshold  $\sigma = 10^{-3}$  applied after  $N_{\text{thr}} = 100$  iterations, with tolerances  $\gamma = 10^{-9}$  and  $\gamma_{\text{thr}} = 10^{-6}$ , and regularization factor  $\epsilon_0 = 10^{-12}$ .

238 **Input:** data  $f^*$ , learning rate  $\eta$ , tolerance  $\gamma$ , threshold  $\sigma$  applied after  $N_{\text{thr}}$ , and  $\epsilon_0$ 239 Initialize the parameters  $\alpha = 0$ ; 240 repeat 241 for  $j = 0, ..., N_t - 1$  do 242 Estimate f in  $t \in (t^{(j)}, t^{(j+1)}]$  by solving forward model (1) given initial condition (9); 243 Find  $\lambda$  in  $t \in [t^{(j)}, t^{(j+1)})$  by solving the adjoint equation in Eq. (7); 244 Compute the gradient using Eq. (4); 245 Update parameters  $\alpha$  using Eq. (5); 246 end 247 if Epochs >  $N_{\rm thr}$  or convergence in  $\alpha$  with  $\gamma_{\rm thr}$  then 248 Thresholding: set  $\alpha_i = 0$  for all *i* that  $|\alpha_i| < \sigma$ ; 249 end 250 **until** Convergence in  $\alpha$  with tolerance  $\gamma$ ; 251 Thresholding: set  $\alpha_i = 0$  for all *i* that  $|\alpha_i| < \sigma$ ; 252 Output:  $\alpha$ 

253 254

255

256

257

258

259

260

261

235

236 237

> **Training with noisy data set.** Often the data set comes with some noise. There are several preprocessing steps that can be done to reduce the noise at the expense of introducing bias, for example removing high frequencies using Fast Fourier Transform or removing small singular values from data set using Singular Value Decomposition. However, we can also reduce the sensitivity of the training algorithm to the noise by averaging the gradients before updating the parameters. Assuming that the noise is martingale, the Monte Carlo averaging gives us the unbiased estimator for the expected value of the gradient over all the data set. We adapt the training procedure by averaging gradients over all available data points and then updating the parameters (see Algorithm 2 and the flowchart in Fig. 5 of Appendix A for more details). Clearly, this will make the algorithm more robust at higher cost since the update happens only after seeing all the data.

262 263 264

### 3 RESULTS

265 266

We demonstrate the validity of our proposed adjoint-based method in discovering PDEs given
 measurements on a spatial-temporal grid. As a show case, next we show the results of the PDE
 discovery task applied to the heat equation. We also consider data collected from a variety of problems,
 including the Burgers', Kuramoto Sivashinsky, Random Walk, and Reaction Diffusion equation,

Ī	<b>nput:</b> data $f^*$ , learning rate $\eta$ , tolerance $\gamma$ , threshold $\sigma$ applied after $N_{\text{thr}}$ , and $\epsilon_0$
n	$\alpha = 0,$
1	for $i = 0, \dots, N_t - 1$ do
	Estimate $f$ in $t \in (t^{(j)}, t^{(j+1)}]$ by solving forward model (1) given initial condition (9)
	Find $\lambda$ in $t \in [t^{(j)}, t^{(j+1)})$ by solving the adjoint equation in Eq. (7):
	Compute the gradient $a^{(j)} = \partial C^{(j)} / \partial \alpha$ using Eq. (4):
	end
	Average the gradient $\mathbb{E}[\partial C/\partial \alpha] = \sum_{i} q^{(j)}/N_t$ ;
	Update parameters $\alpha$ with Eq. (5) using $\mathbb{E}[\partial C/\partial \alpha]$ ;
	if Epochs > $N_{\rm thr}$ or convergence in $\alpha$ with $\gamma_{\rm thr}$ then
	Thresholding: set $\alpha_i = 0$ for all <i>i</i> that $ \alpha_i  < \sigma$ ;
	end
u	<b>intil</b> Convergence in $\alpha$ with tolerance $\gamma$ ;
T	Thresholding: set $\alpha_i = 0$ for all <i>i</i> that $ \alpha_i  < \sigma$ ;
C	Dutput: $\alpha$

detailed in Appendix B. We have compared our approach to PDE-FIND in terms of error and time to convergence. All the results are obtained using a single core-thread of a 2.3 GHz Quad-Core Intel Core i7 CPU. In this paper, we report the execution time  $\tau$  obtained with averaging over 10 independent runs and we use error bars to show the standard deviation of the expected time, i.e.  $d_{\text{error-bar}} = \sqrt{\mathbb{E}[(\tau - \mathbb{E}[\tau])^2]}.$ 

#### 297 298 299

300

301

290 291

#### 3.1 HEAT EQUATION

As a first example, let us consider measured data collected from the solution to the heat equation, i.e.

$$\frac{\partial f}{\partial t} + D \frac{\partial^2 f}{\partial x^2} = 0, \tag{10}$$

with D = -1. The data is constructed using the Finite Difference method with initial condition  $f(x, 0) = 5 \sin(2\pi x)x(x-L)$  and a mesh with  $N_x = 100$  nodes in x covering the domain  $\Omega = [0, L]$ with L = 1 and  $N_t = 100$  steps in t with final time  $T = N_t \Delta t$  where  $\Delta t = 0.05 \Delta x^2/(1 + |D|)$  is the step size and  $\Delta x = L/N_x$  is the mesh size in x.

310 We consider a system consisting of a single PDE (i.e. N = 1, f = f, and p = p) with 311 one-dimensional input, i.e. n = 1 and  $x = x \subset \mathbb{R}$ , and  $d = d \in \mathbb{N}$ . In order to construct a general 312 forward model, here we consider derivatives and polynomials with indices  $d, p \in \{1, 2, 3\}$  as the 313 initial guess for the forward model. This leads to 9 terms with unknown coefficients  $\alpha$  that we find 314 using the proposed adjoint method (an illustrative derivation of the candidate terms can be found in 315 Appendix F.1). While we expect to recover the coefficient that corresponds to D, we expect all the 316 other coefficients (denoted by  $\alpha^*$ ) to become negligible. That is what we indeed observe in Fig. 1 317 where the error of the coefficient for each term is plotted against the number of epochs.

Next, we compare the solution obtained via the adjoint method against PDE-FIND with STRidge optimization method. Here, we test both methods in recovering the heat equation given data on the grid with discretization  $(N_t, N_x) \in \{(100, 100), (500, 100), (1000, 100)\}$ . As shown in Fig. 1, the proposed adjoint method provides more accurate results across all data sizes.

We also point out that as the size of the data set increases, PDE-FIND with STRidge regression method becomes more expensive, e.g. one order of magnitude more expensive than the adjoint method for the data on a grid size  $(N_t, N_x) = (1000, 1000)$ .



Figure 1: The estimated coefficient corresponding to D in heat equation (a) and the  $L_1$ -norm error of all considered coefficients (b) using the proposed Adjoint method with  $N_t = 100$  and  $N_x = 100$ . Also, we show  $L_1$ -norm error of the estimated coefficients (c) and the execution time (d) using the proposed Adjoint method (blue) and PDE-FIND method (red), given data on a grid with  $N_t \in \{100, 500, 1000\}$  steps in t, and  $N_x \in \{100, 1000\}$  nodes in x.

Table 1: Recovering the heat equation given sparse data set in time. Here we rounded the coefficients up to three decimals.

$N_t$	Method	Recovered PDE
100	Adjoint PDE-FIND	$\begin{aligned} f_t - f_{xx} &= 0\\ f_t - f_{xx} &= 0 \end{aligned}$
50	Adjoint PDE-FIND	$f_t - f_{xx} = 0$ $f_t - 0.999 f_{xx} + 0.177 f - 0.261 f^3 - 0.089 f_x - 0.011 f^3 f_x - 0.003 f^2 f_{xx} - 0.001 f f_{xxx} = 0$
25	Adjoint PDE-FIND	$f_t - f_{xx} = 0$ $f_t - 0.999 f_{xx} + 0.532 f - 0.778 f^3 - 0.268 f f_x - 0.035 f^3 f_x - 0.010 f^2 f_{xx} - 0.003 f f_{xxx} = 0$
12.5	Adjoint PDE-FIND	$f_t - f_{xx} = 0 \\ f_t - 0.999 f_{xx} + 1.264 f - 1.863 f^3 - 0.638 f f_x - 0.081 f^3 f_x - 0.025 f^2 f_{xx} - 0.007 f f_{xxx} - 0.001 f^3 f_{xxx} = 0$
6.25	Adjoint PDE-FIND	$\begin{aligned} f_t - f_{xx} &= 0 \\ f_t - 0.999 f_{xx} + 2.769 f - 4.051 f^3 - 1.398 f f_x - 0.185 f^3 f_x - 0.055 f^2 f_{xx} - 0.016 f f_{xxx} - 0.002 f^3 f_{xxx} &= 0 \end{aligned}$

356 357

359

360

361

367

337

338

339

340 341 342

343

#### 4 PARTIAL OBSERVATIONS IN TIME

358 Here, we investigate how the error of the discovery task increases when only a subset of the fine data set is available. Consider the heat equation presented in section 3.1 and consider a data set created by solving the exact PDE using the Finite Difference method with  $\Delta t = T/N_t$  where  $N_t = 1000$  and  $\Delta x = L/N_x$  and  $N_x = 1000$ .

362 Let us assume that we are only provided with a subset of this data set. As a test, let us take every  $\nu$  time step as the input for the PDE discovery task, where  $\nu \in \{1, 2, 4, 8, 16\}$ . This corresponds to using  $\{100, 50, 25, 12.5, 6.25\}\%$  of the total data set. By doing so, the accuracy of the Finite 364 Difference method in estimating the time derivatives using the available data deteriorates, leading to 365 large error in PDE discover task for PDE-FIND method. 366

However, the adjoint method can use a finer mesh in time compared to the data set in com-368 puting the forward and backward equations and only compare the solution to the data on the coarse 369 mesh where data is available. We use  $N_t = 1000$  for the forward and backward solvers in the adjoint 370 method, and impose the final time condition where data is available. As shown in Table 1 and Fig. 2 371 the proposed adjoint method is able to recover the exact PDE regardless of how sparse the data set is 372 in time. 373

374 We emphasize that while adjoint method can use a finer discretization in time than the one 375 for data on  $\mathcal{G}$  in solving forward and backward equations, it is bound to use similar or coarser spatial discretization as  $\mathcal{G}$ . This is due to the fact that the data points  $f^*$  are used for the initial condition 376 of the forward model eq. equation 9, and the final condition of the backward adjoint equation 377 eq. equation 8.



Figure 2: Evolution of the  $L_1$ -norm error in coefficients of all considered terms using adjoint method when only (a) 50% and (b) 6.25% (b) of the data set is available. Error and execution time of the adjoint method (blue) and PDE-FIND method (red) in finding the coefficients of true heat equation given sparse data set in time in (c) and (d).



Figure 3: Ground truth, noisy ( $\sigma = 0.1\%$ ), and denoised solution to the Burgers' equation at the initial time  $t^{(0)}$  and the final time  $t^{(N_t)}$  where  $N_t = 1000$  (a), as well as the error between noisy/denoised solution and the ground truth at these times in (b) and (c), respectively.

#### 5 SENSITIVITY TO NOISE

Here, we investigate how the error increases once noise is added to the data set. In particular, we add noise  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  to each point of the data set for  $f^*$ , where  $\mathcal{N}(0, \sigma^2)$  denotes a normal distribution with zero mean and standard deviation  $\sigma$ . As test cases, we revisit the heat (section 3.1) and Burgers' equations (section B.1) with added noise of  $\epsilon$  with  $\sigma \in \{0.001, 0.005, 0.01, 0.1\}$  %. Before searching for the PDE, we first denoise the data set using Singular Value Decomposition and drop out terms with singular value below a threshold of  $\mathcal{O}(10^{-4})$ . As an example, we show the solution to Burgers' equation in the presence of 0.1% noise in Figure 4.

As shown in Figure 4, adding noise to the data set deteriorates the accuracy in finding the correct coefficients of the underlying PDE for both the adjoint method and PDE-FIND method. We observe that the adjoint method, both with and without gradient averaging, is less susceptible to noise compared to PDE-FIND, albeit at a higher computational cost. Additionally, averaging the gradients in the adjoint method improves the accuracy around two orders of magnitude at higher computational cost.

423 424

389

390

391

392 393

394

396

397

399 400

401 402

403

404

405 406

407

#### 6 ADDRESSING ILL-POSEDNESS

There may exist more than one PDE that replicates the data set. Therefore, the PDE discovery task is
 ill-posed due to the lack of uniqueness in the solution. This is an indication that further physically
 motivated constraints are needed to narrow the search space to find the desired PDE. However, among
 all possible PDEs, which PDE is found by the Adjoint method with the loss function defined as
 equation 3?

430 To answer this question, let us consider a simple example of the wave equation

431

$$f(x,t) = \sin(x-t) \tag{11}$$

<sup>421</sup> 422



Figure 4: Error and execution time of the adjoint method with (green) and without averaging the gradients (blue), along with the PDE-FIND method (red) in finding the coefficients of the true PDE, i.e. heat equation (a)-(c) and Burgers' equation (b)-(d), given noisy data.

which is a solution to infinite PDEs. For example, one class of PDEs with solution f is

$$f_t + kf_x + (k-1)f_{xxx} + c(f_{xx} + f_{xxxx}) = 0 \quad \forall k \in \mathbb{N} \text{ and } c \in \mathbb{R},$$
(12)

defined in a domain  $x \in [0, 2\pi]$  and T = 1. We create a data set using the exact f on a grid with  $N_t = 10$  time intervals and  $N_x = 100$  spatial discretization points. Let us consider a similar setup as the heat equation example 3.1 with derivatives and polynomials indices  $d \in \{1, ..., 6\}$  and p = 1 as the initial guess for the forward model. This leads to 6 terms with unknown coefficients  $\alpha$ . Here, we enable averaging and use a finer discretization in time (100 steps for forward and backward solvers in each time interval) to cope with the instabilities of the Finite Difference solver due to the inclusion of 455 the high-order derivatives. We also disable thresholding except at the end of the algorithm.

456 The proposed Adjoint method returns the solution

$$f_t + 0.996 f_x = 0 \tag{13}$$

which is the PDE with the least number of terms compared to all possible PDEs. We note that for the same problem setting, PDE-FIND finds

$$f_t + 0.9897 f_x = 0. (14)$$

The identified form of PDE can be explained by the use of regularization term in the cost function 3, which enforces the minimization of the PDE coefficients. Clearly, the regularization term may be changed to find other possible solutions of this ill-posed problem.

7 DISCUSSION

443

444

445 446

447 448 449

450

451

452

453

454

457 458 459

460

461 462 463

464

465

466 467

468 469

470 471

472 473

474

475

476

477

478

479

480

481

482

483

Below we highlight and discuss strengths and weaknesses of the proposed adjoint method.

**Strengths.** The proposed method has several strengths:

- 1. The proposed adjoint-based method of discovering PDEs can provides coefficients of the true governing equation with significant accuracy.
- 2. Since the gradient of the cost function with respect to parameters are derived analytically, the optimization problem converges fast. In particular, the adjoint method becomes cheaper than PDE-FIND as the size of the data set increases. The adjoint method by construction finds the optimal relation between gradient of the cost function and the error in the data points. This was achieved by finding the extremum of the objective functional using variational derivative. We note that a clear difference from the point-wise loss  $||f - f^*||_2$  equipped with backprobagation used in PDE-FIND method is that the adjoint method weights the error at discrete points with the Lagrange multipliers, see Eq. equation 4 and the final condition Eq. equation 8.
- 3. Since the adjoint method uses a PDE solver to find the underlying governing equation, there 484 is a guarantee that the found PDE can be solved numerically with the same PDE solver as 485 the one used by the adjoint method.

- 4. The adjoint method can use a finer mesh in time compared to the available discretization of the data set. This allows an accurate recovery of the underlying PDE compared to the PDE-FIND, where the error in the latter increases as the data set gets coarser since it estimates derivatives directly (either with Finite Difference or a polynomial fit) using the given data set.
- 492 Weaknesses.

**knesses.** Our proposed method has some limitations:

- 1. In order to use the proposed adjoint method for discovering PDEs, a general solver of PDEs needs to be implemented. Here, we used Finite Differences which can be replaced with more advance solvers. Clearly, the proposed adjoint method is most effective when there is a prior knowledge of the underlying PDE form, and an appropriate numerical solver is deployed. We note that an inherent limitation of the proposed adjoint method is the possibility of encountering either ill-posed forward or backward equations during optimization, which limits the time step size.
  - 2. In this work, we used the same spatial discretization as the input data. If the spatial grid of input data is too coarse for the PDE solver, one has to use interpolation to estimate the data on a finer spatial grid that is more appropriate for the PDE solver.
- 3. In this work, we made the assumption that the underlying PDE can be solved numerically. This can be a limitation when there are no stable numerical methods to solve the true PDE. In this scenario, the proposed method may find another PDE that is solvable and fits to the data with a notable error.

#### 8 CONCLUSION

In this work, we introduce a novel mathematical method for the discovery of partial differential equations given data using the adjoint method. By formulating the optimization problem in the variational form using the method of Lagrange multipliers, we find an analytic expression for the gradient of the cost function with respect to the parameters of the PDE as a function of the Lagrange multipliers and the forward model estimate. Then, using variational calculus, we find a backwards-in-time evolution equation for the Lagrange multipliers which incorporates the error with a source term (the adjoint equation). Hence, we can use the same solver for both forward model and the backward Lagrange equations. Here we used Finite Differences to estimate the spatial derivatives and forward Euler for the time derivatives, which indeed can be replaced with more stable and advanced solvers.

We compared the proposed adjoint method against PDE-FIND on several test cases. While PDE-FIND seems to be faster for small size problems, we observe that the adjoint method equipped with forward/backward solvers becomes faster than PDE-FIND as the size of the data set increases. Also, the adjoint method can provide machine-accuracy in identifying and finding the coefficients when the data set is noise-free. Furthermore, in the case of discovering PDEs for PDFs given its samples, both methods seem to suffer enormously form noise/bias associated with the finite number of samples and the Finite Difference on histogram. This motivates the use of smooth and least biased density estimator in these methods such as Tohme et al. (2023) in future work. In the future work, we intend to combine the adjoint-based method for the discovery of PDE with PINNs as the solver instead of Finite Difference method. This would allow us to handle noisy and sparse data as well as deploying larger time steps in estimating the forward and backward solvers. 

# 540 REFERENCES 541

54 I	
542	Thomas Antonsen, Elizabeth J Paul, and Matt Landreman. Adjoint approach to calculating shape
543	gradients for three-dimensional magnetic confinement equilibria. Journal of Plasma Physics, 85
544	(2):905850207, 2019.
545	Josh Rongard and Hod Linson Automated reverse engineering of nonlinear dynamical systems
546	Proceedings of the National Academy of Sciences 104(24):9043_9048 2007
547	1  for evenings of the National Acta entry of Sciences, 10+(2+).77+3-77+6, 2007.
548	Steven L Brunton, Joshua L Proctor, and J Nathan Kutz. Discovering governing equations from data
549	by sparse identification of nonlinear dynamical systems. Proceedings of the national academy of
550	sciences, 113(15):3932–3937, 2016.
551	Dural Cafficale Danie Silantum and Venan Vena Adjaint dama for nonlinear holt-more constin
552	constrained optimization <i>Journal of Computational Physics</i> 430:110404, 2021
553	constrained optimization. <i>Journal of Computational Thysics</i> , 459.110404, 2021.
554	Kathleen Champion, Bethany Lusch, J Nathan Kutz, and Steven L Brunton. Data-driven discovery of
555	coordinates and governing equations. Proceedings of the National Academy of Sciences, 116(45):
556	22445–22451, 2019.
557	Theo Chan Vang Liu and Hao Sun. Dhysics informed learning of governing equations from searce
558	data Nature communications 12(1):6136-2021
559	data. Nature communications, 12(1):0150, 2021.
560	Bryan C Daniels and Ilya Nemenman. Automated adaptive inference of phenomenological dynamical
561	models. Nature communications, 6(1):8133, 2015a.
562	
563	Bryan C Dameis and Hya Nemenman. Efficient interence of parsimonious phenomenological models
564	of central dynamics using s-systems and anemating regression. <i>Plos one</i> , 10(3):e0119821, 20130.
565	Junming Duan and Jan S Hesthaven. Non-intrusive data-driven reduced-order modeling for time-
566	dependent parametrized problems. Journal of Computational Physics, 497:112621, 2024.
567	
568	Solve Eignes and Kjetil Olsen Lye. Pseudo-namiltonian neural networks for learning partial differen-
569	that equations. <i>Journal of Computational Thysics</i> , page 112758, 2024.
570	H Pearl Flath, Lucas C Wilcox, Volkan Akçelik, Judith Hill, Bart van Bloemen Waanders, and
571	Omar Ghattas. Fast algorithms for bayesian uncertainty quantification in large-scale linear inverse
572	problems based on low-rank partial hessian approximations. SIAM Journal on Scientific Computing,
573	33(1):407–432, 2011.
574	Alessandro Geraldini Matt Landreman and Elizabeth Paul. An adjoint method for determining the
575	sensitivity of island size to magnetic field variations. <i>Journal of Plasma Physics</i> 87(3):905870302
576	2021.
577	
578	Dimitrios Giannakis and Andrew J Majda. Nonlinear laplacian spectral analysis for time series with
579	intermittency and low-frequency variability. <i>Proceedings of the National Academy of Sciences</i> , 100(7):2222, 2227, 2012
580	109(7):2222-2227, 2012.
581	Raul González-García, Ramiro Rico-Martinez, and Ioannis G Kevrekidis. Identification of distributed
582	parameter systems: A neural net based approach. Computers & chemical engineering, 22:S965-
583	S968, 1998.
584	Antonio Tomanonio Anno Anno Anno anti altanza anti altanza da alta da alta da alta da alta da alta da alta da a
585	Aniony Jameson. Aerodynamic snape optimization using the adjoint method. Lectures at the Von Karman Institute, Brussels, 2003
586	Kurmun Institute, Drussets, 2005.
50C	Kadierdan Kaheman, J Nathan Kutz, and Steven L Brunton. Sindy-pi: a robust algorithm for parallel
500	implicit sparse identification of nonlinear dynamics. Proceedings of the Royal Society A, 476
509	(2242):20200279, 2020.
59U	Joannis & Kaurakidia & William Gaar Jamas M. Human Depariatio & Kaurakidia Olaf Durahara
591	Constantinos Theodoropoulos, et al. Equation free, coarse grained multiscale computation; en
592	abling microscopic simulators to perform system-level analysis Commun Math Sci 1(4):715_762
222	2003.

594 595	Zichao Long, Yiping Lu, Xianzhong Ma, and Bin Dong. Pde-net: Learning pdes from data. In <i>International conference on machine learning</i> , pages 3208–3216. PMLR, 2018.
597 598	Zichao Long, Yiping Lu, and Bin Dong. Pde-net 2.0: Learning pdes from data with a numeric- symbolic hybrid deep network. <i>Journal of Computational Physics</i> , 399:108925, 2019.
599 600	Andrew J Majda, Christian Franzke, and Daan Crommelin. Normal forms for reduced stochastic climate models. <i>Proceedings of the National Academy of Sciences</i> , 106(10):3649–3653, 2009.
601 602 603 604	Niall M Mangan, Steven L Brunton, Joshua L Proctor, and J Nathan Kutz. Inferring biological net- works by sparse identification of nonlinear dynamics. <i>IEEE Transactions on Molecular, Biological</i> <i>and Multi-Scale Communications</i> , 2(1):52–63, 2016.
605 606	Takashi Matsubara, Ai Ishikawa, and Takaharu Yaguchi. Deep energy-based modeling of discrete- time physics. <i>Advances in Neural Information Processing Systems</i> , 33:13100–13111, 2020.
607 608 609	Igor Mezić. Analysis of fluid flows via spectral properties of the koopman operator. <i>Annual review of fluid mechanics</i> , 45:357–378, 2013.
610 611	Maziar Raissi and George Em Karniadakis. Hidden physics models: Machine learning of nonlinear partial differential equations. <i>Journal of Computational Physics</i> , 357:125–141, 2018.
612 613 614 615	Maziar Raissi, Paris Perdikaris, and George E Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. <i>Journal of Computational physics</i> , 378:686–707, 2019.
616	Anthony John Roberts. Model emergent dynamics in complex systems, volume 20. SIAM, 2014.
617 618 619	Samuel H Rudy, Steven L Brunton, Joshua L Proctor, and J Nathan Kutz. Data-driven discovery of partial differential equations. <i>Science advances</i> , 3(4):e1602614, 2017.
620 621 622	Nihar Sawant, Boris Kramer, and Benjamin Peherstorfer. Physics-informed regularization and structure preservation for learning stable reduced models from data with operator inference. <i>Computer Methods in Applied Mechanics and Engineering</i> , 404:115836, 2023.
623 624 625	Michael Schmidt and Hod Lipson. Distilling free-form natural laws from experimental data. <i>science</i> , 324(5923):81–85, 2009.
626 627 628	Michael D Schmidt, Ravishankar R Vallabhajosyula, Jerry W Jenkins, Jonathan E Hood, Abhishek S Soni, John P Wikswo, and Hod Lipson. Automated refinement and inference of analytical models for metabolic networks. <i>Physical biology</i> , 8(5):055011, 2011.
629 630 631	George Sugihara, Robert May, Hao Ye, Chih-hao Hsieh, Ethan Deyle, Michael Fogarty, and Stephan Munch. Detecting causality in complex ecosystems. <i>science</i> , 338(6106):496–500, 2012.
632 633	Robert Tibshirani. Regression shrinkage and selection via the lasso. <i>Journal of the Royal Statistical Society Series B: Statistical Methodology</i> , 58(1):267–288, 1996.
634 635 636	Tony Tohme, Dehong Liu, and Kamal Youcef-Toumi. GSR: A generalized symbolic regression approach. <i>Transactions on Machine Learning Research</i> , 2022.
637 638 639	Tony Tohme, Mohsen Sadr, Kamal Youcef-Toumi, and Nicolas G Hadjiconstantinou. Messy es- timation: Maximum-entropy based stochastic and symbolic density estimation. <i>arXiv preprint</i> <i>arXiv:2306.04120</i> , 2023.
640 641 642	Henning U Voss, Paul Kolodner, Markus Abel, and Jürgen Kurths. Amplitude equations from spatiotemporal binary-fluid convection data. <i>Physical review letters</i> , 83(17):3422, 1999.
643 644 645 646	Hao Ye, Richard J Beamish, Sarah M Glaser, Sue CH Grant, Chih-hao Hsieh, Laura J Richards, Jon T Schnute, and George Sugihara. Equation-free mechanistic ecosystem forecasting using empirical dynamic modeling. <i>Proceedings of the National Academy of Sciences</i> , 112(13):E1569–E1576, 2015.

### A FLOWCHARTS OF ADJOINT METHOD

 Here, we present flowcharts to illustrate the proposed adjoint algorithms 1-2 with and without averaging the gradients in Fig. 5.



Figure 5: Training flowchart of the Adjoint method in finding PDEs (left) without and (right) with gradient averaging.

#### **B** GALLERY OF EXAMPLES

Besides the heat equation 3.1, here we show how the proposed adjoint method for discovering governing equation performs in finding Burgers', Kuramoto Sivashinsky, Random Walk, Reaction Diffusion System of Equations, and wave equation.

#### B.1 BURGERS' EQUATION

As a nonlinear test case, let us consider the data from Burgers' equation given by

$$\frac{\partial f}{\partial t} + \frac{\partial (Af^2)}{\partial x} = 0 \tag{15}$$

where A = -1. The data is obtained with similar simulation setup as for heat equation (Section 3.1) except for the time step, i.e.  $\Delta t = 0.05\Delta x/(1 + |A|)$ .

Similar to Section 3.1, we adopt a system of one PDE with one-dimensional input. We also consider derivatives and polynomials with indices  $d, p \in \{1, 2, 3\}$  in the construction of the forward model. This leads to 9 terms whose coefficients we find using the proposed adjoint method. As shown in



Figure 6: The estimated coefficient corresponding to A (left) and the  $L_1$ -norm error of all considered coefficients (right) given the discretized data of Burgers' equation during training.



Figure 7:  $L_1$ -norm error of the estimated coefficients (left) and the execution time (right) for discovering the Burgers' equation equation using the Adjoint method (blue) and PDE-FIND method (red), given data on a grid with  $N_t \in \{100, 1000\}$  steps in t, and  $N_x \in \{100, 1000\}$  nodes in x.

Fig. 6, the proposed adjoint method finds the correct coefficients, i.e.  $\alpha_{d=1,p=2}$  that corresponds to D as well as all the irrelevant ones denoted by  $\alpha^*$ , up to machine accuracy in  $\mathcal{O}(10)$  epochs.

Next, we compare the solution obtained from the adjoint method to the one from PDE-FIND using STRidge optimization method. Here, we compare the error on coefficients and computational time between the adjoint and PDE-FIND by repeating the task for the data set with increasing size, i.e.  $(N_t, N_x) \in \{(100, 100), (1000, 100), (1000, 1000)\}$ . As depicted in Fig. 7, the adjoint method provides us with more accurate solution across the different discretization sizes. Regarding the computational cost, while PDE-FIND seems faster on smaller data sets, as the size of the data grows, it becomes increasingly more expensive than the adjoint method. Similar to the heat equation, for the mesh size  $(N_t, N_x) = (1000, 1000)$  we obtain one order of magnitude speed up compared to PDE-FIND. 

#### 743 B.2 KURAMOTO SIVASHINSKY EQUATION

As a more challenging test case, let us consider the recovery of the Kuramoto-Sivashinsky (KS) equation given by

 $\frac{\partial f}{\partial t} + A \frac{\partial f^2}{\partial x} + B \frac{\partial^2 f}{\partial x^2} + C \frac{\partial^4 f}{\partial x^4} = 0$ (16)

where A = -1, B = 0.5 and C = -0.5. The data is generated similar to previous sections except for the grid  $(N_t, N_x) = (64, 256)$  and the time step size  $\Delta t = 0.01 \Delta x^4 / (1 + |C|)$ .

Here again, we adopt a system of one PDE with one-dimensional input. As a guess for the forward model, we consider terms consisting of derivatives with indices  $d \in \{1, 2, 3, 4\}$  and polynomials with indices  $b \in \{1, 2\}$ , leading to 8 terms whose coefficients we find using the proposed adjoint method. As shown in Fig. 8, the adjoint method finds the coefficient with error of  $\mathcal{O}(10^{-5})$ , yet achieving machine accuracy seems not possible.



Figure 8: The estimated coefficients corresponding to A, B, C (left) and the  $L_1$ -norm error of all considered coefficients (right) given the discretized data of the KS equation during training without (top) and with (bottom) active thresholding for Epochs >  $N_{\text{thr}} = 100$ .



Figure 9:  $L_1$ -norm error of the estimated coefficients (left) and the execution time (right) for discovering the KS equation using the Adjoint method (blue) and PDE-FIND method (red), given data on a grid with  $N_t \in \{64, 128, 256\}$  steps in  $t, N_x \in \{256, 512, 1024\}$  nodes in x.

Again, in Fig. 9 we make a comparison between the predicted PDE using the adjoint method against PDE-FIND. In particular, we consider a data set on a temporal/spatial mesh of size  $(N_t, N_x) = \{(64, 256), (128, 512), (256, 1024)\}$  and compare how the error and computational cost vary. Similar to previous sections, the error is reported by comparing the obtained coefficients against the coefficients of the exact PDE in  $L_1$ -norm. Interestingly, the PDE-FIND method has 3 to 4 orders of magnitude larger error compared to the adjoint method. Also, in terms of cost, the training time for PDE-FIND seems to grow at a higher rate than the adjoint method as the (data) mesh size increases.



Figure 10: The estimated coefficients corresponding to A and D (left) and the  $L_1$ -norm error of all considered coefficients (right) of the Fokker-Planck equation as the governing law for the PDF associated with the random walk during training.

B.3 RANDOM WALK

Next, let us consider the recovery of the governing equation on probability density function (PDF) given samples of its underlying stochastic process. As an example, we consider the Itô process

$$dX = Adt + \sqrt{2D}dW \tag{17}$$

where A = 1 is drift and D = 0.5 is the diffusion coefficient, and W denotes the standard Wiener process with  $Var(dW) = \Delta t$ . We generate the data set by simulating the random walk using Euler-Maruyama scheme starting from X(t = 0) = 0 for  $N_t = 50$  steps with a time step size of  $\Delta t = 0.01$ . We estimate the PDF using histogram with  $N_x = 100$  bins and  $N_s = 1000$  samples.

Let us denote the distribution of X by f. Itô's lemma gives us the Fokker-Planck equation

$$\frac{\partial f}{\partial t} + A \frac{\partial f}{\partial x} - D \frac{\partial^2 f}{\partial x^2} = 0.$$
(18)

Given the data set for f on a mesh of size  $(N_t, N_x)$ , we can use Finite Difference to compute the contributions from derivatives of f in the governing law. Since this is one of the challenging test cases due to noise, here we only consider three possible terms in the forward model, consisting of derivatives with indices  $d \in \{1, 2, 3\}$  and polynomial power p = 1. In Fig. 10, we show how the error of finding the correct coefficients evolves during training for the adjoint method. Clearly, the adjoint method seems to recover the true PDE with  $L_1$  error of  $\mathcal{O}(10^{-2})$  in its coefficients.

In Fig. 11, we make a comparison with PDE-FIND for the same number and order of terms as the initial guess for the PDE. We compare the two methods for a range of grid and sample sizes, i.e.  $N_t \in \{50, 100\}, N_x = 100$ , and  $N_s \in \{10^3, 10^4\}$ . It turns out that the proposed adjoint method overall provides more accurate estimate of the coefficients than PDE-FIND, though at a higher cost. In Table 2, we show the discovered PDEs for both methods across the different discretizations.

Table 2: Recovery of the Fokker-Planck equation, i.e.  $f_t + f_x - 0.5 f_{xx} = 0$ , using the proposed adjoint method against PDE-FIND method given samples of the underlying stochastic process for various discretization parameters.

	1			
$N_t$	$N_x$	$N_{ m s}$	Method	<b>Recovered PDE</b>
50	100	1000	Adjoint PDE-FIND	$ \begin{vmatrix} f_t + 1.025f_x - 0.465f_{xx} = 0 \\ f_t + 0.798f_x - 0.454f_{xx} = 0 \end{vmatrix} $
		10000	Adjoint PDE-FIND	$ \begin{aligned} f_t + 1.022f_x - 0.495f_{xx} &= 0 \\ f_t + 0.818f_x - 0.496f_{xx} &= 0 \end{aligned} $
100	100	1000	Adjoint PDE-FIND	$\begin{vmatrix} f_t + 1.010f_x - 0.543f_{xx} = 0\\ f_t + 0.863f_x - 0.560f_{xx} = 0 \end{vmatrix}$
		10000	Adjoint PDE-FIND	$\begin{cases} f_t + 1.015f_x - 0.589f_{xx} = 0\\ f_t + 0.894f_x - 0.612f_{xx} = 0 \end{cases}$



Figure 11:  $L_1$ -norm error of the estimated coefficients (left) and the execution time (right) for discovering the Fokker-Planck equation using the proposed Adjoint method (blue) and PDE-FIND method (red), given samples of its underlying stochastic process with  $N_t \in \{50, 100\}$  steps in t,  $N_x = 100$  histogram bins, and  $N_s \in \{10^3, 10^4\}$  samples.

#### **B.4** REACTION DIFFUSION SYSTEM OF EQUATIONS

In order to show scalability and accuracy of the adjoint method for a system of PDEs in a higher dimensional space, let us consider a system of PDEs given by

$$\frac{\partial u}{\partial t} + c_0^u \nabla_{x_1}^2[u] + c_1^u \nabla_{x_2}^2[u] + R^u(u, v) = 0,$$
(19)

$$\frac{\partial v}{\partial t} + c_0^v \nabla_{x_1}^2[v] + c_1^v \nabla_{x_2}^2[v] + R^v(u, v) = 0$$
(20)

where

$$R^{u}(u,v) = c_{2}^{u}u + c_{3}^{u}u^{3} + c_{4}^{u}uv^{2} + c_{5}^{u}u^{2}v + c_{6}^{u}v^{3}$$
<sup>(21)</sup>

$$R^{v}(u,v) = c_{2}^{v}v + c_{3}^{v}v^{3} + c_{4}^{v}vu^{2} + c_{5}^{v}v^{2}u + c_{6}^{v}u^{3}$$

$$(22)$$

We construct the data set by solving the system of PDEs Eqs. equation 19-equation 20 using a 2nd order Finite Difference scheme with initial values

$$u_{0} = a \sin\left(\frac{4\pi x_{1}}{L_{1}}\right) \cos\left(\frac{3\pi x_{2}}{L_{2}}\right) \left(L_{1}x_{1} - x_{1}^{2}\right) \left(L_{2}x_{2} - x_{2}^{2}\right)$$
$$v_{0} = a \cos\left(\frac{4\pi x_{1}}{L_{1}}\right) \sin\left(\frac{3\pi x_{2}}{L_{2}}\right) \left(L_{1}x_{1} - x_{1}^{2}\right) \left(L_{2}x_{2} - x_{2}^{2}\right)$$

where a = 100, and the coefficients

$$\begin{aligned} \boldsymbol{c}^{u} &= [c_{i}^{u}]_{i=0}^{6} = [-0.1, -0.2, -0.3, -0.4, 0.1, 0.2, 0.3] \\ \boldsymbol{c}^{v} &= [c_{i}^{v}]_{i=0}^{6} = [-0.4, -0.3, -0.2, -0.1, 0.3, 0.2, 0.1]. \end{aligned}$$

We generate data by solving the system of PDEs Eqs. (19)-(20) using the Finite Difference method and forward Euler scheme for  $N_t = 25$  steps with a time step size of  $\Delta t = 10^{-6}$ , and in the domain  $\Omega = [0, L_1] \times [0, L_2]$  where  $L_1 = L_2 = 1$  which is discretized using a uniform grid with  $N_{x_1} \times N_{x_2} = 50^2$  nodes leading to mesh size  $\Delta x_1 = \Delta x_2 = 0.02$ . In Fig. 12 we show the solution to the system at time  $T = N_t \Delta t$  for u and v.

We consider a system consisting of two PDEs, i.e.  $N = \dim(f) = \dim(p) = 2$ , with twodimensional input, i.e.  $n = \dim(x) = \dim(d) = 2$ . Here,  $\dim(f) = \dim(\operatorname{Ima}(f))$ , where  $\operatorname{Ima}(\cdot)$ denotes the image (or output) of a function.

In order to use the developed adjoint method, we construct a guess forward system of PDEs (or forward model) using derivatives up to 2nd order and polynomials of up to 3rd order. That is,  $d_{\text{max}} = 2$  and  $p_{\text{max}} = 3$ . This leads to 90 terms whose coefficients we find using the proposed



Figure 12: Solution to the reaction diffusion system of PDEs at time t = T for u (left) and v (right).



Figure 13:  $L_1$ -norm error in the estimated coefficients of the reaction diffusion system of PDEs during training.

adjoint method (an illustrative derivation of the candidate terms can be found in Appendix F.2). The solution to the constructed model  $f \approx [u, v]$  as well as the adjoint equation for  $\lambda$  is found using the same discretization as the data set.

As shown in Fig. 13, the adjoint method finds the correct equations with error up to  $\mathcal{O}(10^{-12})$ . Furthermore, the coefficients corresponding to the irrelevant terms  $\alpha^*$  tend to zero with error of  $\mathcal{O}(10^{-11}),$  see Fig. 14. 

Furthermore, we have compared the adjoint method against PDE-FIND for a range of grid sizes in Fig. 16. We observe that the cost of PDE-FIND grows with higher rate than adjoint method as the size of the data set increases.







Figure 16: Comparing the error and execution time of the adjoint method (blue) to PDE-FIND method (red) against the size of the data set for the tolerance of  $10^{-7}$  in the discovered coefficients.

B.5 WAVE EQUATION

Consider wave equation

$$f(x,t) = \sin(x-t) \tag{23}$$

which is a solution to infinite PDEs. For example, one class of PDEs with solution f is

$$f_t + kf_x + (k-1)f_{xxx} + c(f_{xx} + f_{xxxx}) = 0 \quad \forall k \in \mathbb{N} \text{ and } c \in \mathbb{R},$$
(24)

defined in a domain  $x \in [0, 2\pi]$  and T = 1. We create a data set using the exact f on a grid with  $N_t = 10$  time intervals and  $N_x = 100$  spatial discretization points.

 Let us consider a similar setup as the heat equation example 3.1 with derivatives and poly- **998** nomials indices  $d \in \{1, ..., 6\}$  and p = 1 as the initial guess for the forward model. This leads to **999** 6 terms with unknown coefficients  $\alpha$ . Here, we enable averaging and use a finer discretization in **1000** time (100 steps for forward and backward solvers in each time interval) to cope with the instabilities **1001** of the Finite Difference solver due to the inclusion of the high-order derivatives. We also disable **1002** thresholding except at the end of the algorithm.



f

Figure 15: Profile of f at t = 0 and t = 1 (left) and the evolution of considered coefficients during adjoint optimization (right)

As shown in Fig. 15, the proposed Adjoint method returns the solution

$$f_t + 0.996 f_x = 0$$
 (25)

which is the PDE with the least number of terms compared to all possible PDEs. We note that for thesame problem setting, PDE-FIND identifies the same form of the PDE, i.e.

$$f_t + 0.9897 f_x = 0. (26)$$

### C INCOMPLETE GUESSED PDE SPACE

1025 In this section, we investigate the outcome of the adjoint method when the exact terms are not included in the initial guessed PDE form. Here, we define the space of PDE where the exact terms



Figure 17: The adjoint method applied to the Burgers' equation for complete (top) and incomplete (bottom) space of guessed PDEs.

are included in the general forward model equation 1 as complete. If the considered general form of PDE equation 1 does not include all the terms of the exact PDE, we denote that as an incomplete guessed PDE space.

Let us take the data from the numerical solution to Burgers' equation used in section B.1 with discretization  $N_t = N_x = 100$ . For the complete forward model, we again consider derivatives and polynomials with indices  $d, p \in \{1, 2, 3\}$  in the construction of the forward model. This leads to 9 terms whose coefficients we find using the proposed adjoint method. For the incomplete space of PDE, we take derivatives and polynomials with indices as  $d \in \{1, 2, 3\}$  and  $p \in \{1, 3\}$ , leading to 6 terms. Clearly, the incomplete guessed PDE space does not include the term  $\alpha_{d=1,p=2}\partial f^2/\partial x$ . Now, we would like to see which PDE is returned by the adjoint method.

In Figure 17, we made a comparison between the evolution of coefficients and  $L_2$  norm error of the estimated forward model against the data. While the complete space monotonically converges to the exact solution up to machine accuracy, the incomplete space of PDE delivers another PDE, i.e.

1061

1046 1047

1071

$$\frac{\partial f}{\partial t} + \frac{\partial^2}{\partial x^2} (0.04f - 0.01f^3) = 0, \tag{27}$$

with the relative  $L_2$  error of  $\mathcal{O}(10^{-5})$  between forward model estimation and the data points. The fact that the  $L_2$  error between f and  $f^*$  does not decrease is an indication that the considered space of PDE is incomplete and additional terms must be included. We note that here we assumed there is no noise in the data set. However, in the presence of noise, the  $L_2$  error between f and  $f^*$  may stagnate at the noise level, which makes the analysis on the completeness of the PDE space more challenging.

#### 1069 1070 D Impact of hyperparameters on adjoint method

In this section, we study the impact of some of the hyperparameters used in the adjoint algorithm. We repeat the PDE discovery experiment for Burgers's and Kuramoto Sivashinsky equation with data on a grid with  $N_x = N_t = 100$ , as described in B.1 and B.2. We check the error in the outcome coefficients and the solution of estimated forward model compared to the data.

1076 As shown in Figure 18, by increasing the regularization factor  $\epsilon_0$ , the optimization problem seems 1077 to converge faster to a stationary solution. In case of Burgers' equation, we considered  $\epsilon_0 \in$ 1078  $\{10^{-4}, 10^{-8}, 10^{-12}, 10^{-16}\}$ , where for all values of  $\epsilon_0$  the exact solution is recovered. However, in 1079 the case of Kuramoto Sivashinsky with  $\epsilon_0 \in \{10^{-10}, 10^{-12}, 10^{-14}, 10^{-16}\}$ , the solution seems to be more sensitive to  $\epsilon_0$ . Here, we fix the other hyperparameters  $\gamma_{\text{thr}} = 10^{-16}$  and  $\beta = 2 \times 10^{-3}$ 



Figure 18: Impact of regularization factor  $\epsilon_0$  and thresholding tolerance  $\gamma_{thr}$  on the error of adjoint method for Burgers' equation (a-d) and Kuramoto Sivashinsky equation (e-h).



Figure 19: Impact of the free parameter  $\beta$  in the learning rate on the error of adjoint method for Burgers' equation (a-b) and Kuramoto Sivashinsky equation (c-d).

1101 1102 1103

1104

1105

1106

1107

1108

1109

1110

1111 1112

1117

for the Burgers's equation and  $\beta = 20$  for Kuramoto Sivashinsky equation. We observe that high regularization factor deteriorates the accuracy, while stabilizing the regression problem.

1120 Next, we investigate how the error changes with the thresholding tolerance where  $\gamma_{\text{thr}} \in \{10^{-4}, 10^{-8}, 10^{-12}, 10^{-16}\}$ . Here, we fix the other hyperparameters  $\epsilon_0 = 10^{-16}$  and  $\beta = 2 \times 10^{-3}$ 1122 for the Burgers's equation and  $\beta = 20$  for Kuramoto Sivashinsky equation. Although using smaller 1123  $\gamma_{\text{thr}}$  allows faster convergence to a stationary solution almost in all cases, we remind the reader that 1124  $\gamma_{\text{thr}}$  should be large enough to allow enough training of the coefficients before truncating terms. In 1125 other words, the user should avoid trivial scenarios where the initial guess for coefficients  $\alpha$  are zero 1126 and the thresholding is applied from very beginning of the training.

Finally, we show the impact of the free parameter  $\beta$  in the learning rate on the resulting PDE discovered by the adjoint method. We compared the solution of adjoint method using  $\beta \in \{10^{-3}, 2 \times 10^{-3}, 3 \times 10^{-3}, 4 \times 10^{-3}\}$  for the Burgers' equation, and  $\beta \in \{2, 5, 10, 20\}$  for the Kuramoto Sivashinsky equation. Also, we fix the other hyperparameters  $\gamma_{\text{thr}} = \epsilon_0 = 10^{-16}$ . As shown in Figure 19, regardless of the value of  $\beta$ , adjoint method delivers the same solution. However, larger values of  $\beta$  lead to faster convergence to the solution, if the numerical solver does not become unstable. The upper bound of  $\beta$  is limited by the stability of the guessed PDE, and can be found with try-and-error.

### <sup>1134</sup> E JUSTIFICATION FOR THE CHOICE OF THE LEARNING RATE

<sup>1136</sup> In the proposed adjoint method, we considered the update rule

$$\alpha_{i,d,p} \leftarrow \alpha_{i,d,p} - \eta \frac{\partial \mathcal{C}}{\partial \alpha_{i,d,p}}$$
(28)

(29)

for i = 1, ..., N. Here, we given a justification for our choice of the learning parameter  $\eta$ .

From the expression for the gradient of cost function with respect to parameters equation 4, i.e.

 $\frac{\partial \mathcal{C}}{\partial \alpha_{i,\boldsymbol{d},\boldsymbol{p}}} = (-1)^{|\boldsymbol{d}|} \frac{1}{\Delta \boldsymbol{x} \Delta t} \int \boldsymbol{f}^{\boldsymbol{p}} \nabla_{\boldsymbol{x}}^{(\boldsymbol{d})} [\lambda_i] d\boldsymbol{x} dt + 2\epsilon_0 \alpha_{i,\boldsymbol{d},\boldsymbol{p}},$ 

1145

1142

1138 1139 1140

1151 1152

1153

1154

1149 we can see that

$$\left|\frac{\partial \mathcal{C}}{\partial \alpha_{i,d,p}}\right| = \mathcal{O}(\nabla_{\boldsymbol{x}}^{(\boldsymbol{d})}) \tag{30}$$

$$\leq \mathcal{O}(h^{-|\boldsymbol{d}|}) \tag{31}$$

1155 where  $h = \min(\Delta x)$ . So, the magnitude of the gradient scales exponentially with the order of the 1156 derivative d. The highest order terms, i.e. the terms with  $d = d_{\max} = \max(|d|)$ , have the largest 1157 magnitude for their gradients. This means that by taking a constant learning rate  $\eta$ , the adjoint method 1158 would find the coefficients of the highest order terms first. This effect leads to the non-uniform 1159 convergence of the adjoint method.

1160 In order to enforce uniform convergence on all PDE parameters, in this paper we consider 1161

$$\eta = \beta \min(\Delta \boldsymbol{x})^{|\boldsymbol{d}| - d_{\max}} \tag{32}$$

as the learning rate which encodes the scaling with respect to the order of derivative for each PDE term. With this choice of learning rate, we have

1166 1167

1168 1169 1170

1162

1163

$$\eta \Big| \frac{\partial \mathcal{C}}{\partial \alpha_{i,d,p}} \Big| \le \mathcal{O}(h^{|d| - d_{\max}}) \mathcal{O}(h^{-|d|})$$
(33)

$$\leq \mathcal{O}(h^{-d_{\max}}) , \tag{34}$$

for all i, d, p. Hence, our choice of  $\eta$ , i.e. Eq. equation 32, enforces uniform convergence on all PDE parameters.

# F ILLUSTRATIVE DERIVATION OF THE ADJOINT EQUATIONS FOR THECONSIDERED CASES.

1177

1174

Although the proposed method and its algorithm can be and has been computed in an automated fashion, here we show two detailed illustrative examples for 1-dimensional and 2-dimensional cases presented in Section 3 for the sake of better understanding the used notation and how the library of candidate terms looks like.

1182

1183 F.1 HEAT AND BURGERS' EQUATIONS 1184

As mentioned in Sections 3.1 and B.1, for these two cases, we consider a system consisting of a single PDE, i.e.  $N = \dim(f) = \dim(p) = 1$  where f = f and p = p, in a one-dimensional input space, i.e.  $n = \dim(x) = \dim(d) = 1$  where x = x, and d = d. In addition, we consider candidate terms consisting of derivatives with indices  $d \in \{1, 2, 3\}$  and polynomials with indices  $p \in \{1, 2, 3\}$ . In other words,  $d_{\text{max}} = 3$  and  $p_{\text{max}} = 3$ . The resulting forward model in Eq. 1 takes the form 

$$\mathcal{L}[f] = \frac{\partial f}{\partial t} + \sum_{d=1}^{3} \sum_{p=1}^{3} \alpha_{d,p} \frac{\partial^d (f^p)}{\partial x^d}$$

$$= \frac{\partial f}{\partial t} + \alpha_{1,1} \frac{\partial f}{\partial x} + \alpha_{1,2} \frac{\partial (f^2)}{\partial x} + \alpha_{1,3} \frac{\partial (f^3)}{\partial x}$$

 $+\alpha_{3,1}\frac{\partial^3 f}{\partial x^3} + \alpha_{3,2}\frac{\partial^3 (f^2)}{\partial x^3} + \alpha_{3,3}\frac{\partial^3 (f^3)}{\partial x^3}$ (35) 

 $+ \alpha_{2,1} \frac{\partial^2 f}{\partial x^2} + \alpha_{2,2} \frac{\partial^2 (f^2)}{\partial x^2} + \alpha_{2,3} \frac{\partial^2 (f^3)}{\partial x^2}$ 

where  $\alpha_{d,p}$  denotes the parameter corresponding to the term with d-th derivative and p-th polynomial order. As we can observe, we have 9 terms with unknown coefficients  $\alpha = [\alpha_{d,p}]_{d \in \{1,2,3\}, p \in \{1,2,3\}}$ that we aim to find using the proposed adjoint method. 

The cost functional in this case is simply

$$\mathcal{C} = \sum_{j,k} \left( f^*(x^{(k)}, t^{(j)}) - f(x^{(k)}, t^{(j)}) \right)^2 + \frac{1}{\Delta x \Delta t} \int \lambda(x, t) \mathcal{L}[f(x, t)] dx dt + \epsilon_0 ||\boldsymbol{\alpha}||_2^2 \,.$$
(36)

Letting variational derivatives of C with respect to f to be zero, and using integration by parts, the corresponding adjoint equation can be obtained as 

$$\frac{\partial \lambda}{\partial t} = \sum_{d=1}^{3} \sum_{p=1}^{3} (-1)^{d} \alpha_{d,p} \frac{\partial \left(f^{p}\right)}{\partial f} \frac{\partial^{d} \lambda}{\partial x^{d}}$$

1215  
1216  
1217  
1218  

$$= -\alpha_{1,1}\frac{\partial\lambda}{\partial x} - \alpha_{1,2}(2f)\frac{\partial\lambda}{\partial x} - \alpha_{1,3}(3f^2)\frac{\partial\lambda}{\partial x}$$

$$= -\alpha_{1,1}\frac{\partial\lambda}{\partial x} - \alpha_{1,2}(2f)\frac{\partial\lambda}{\partial x} - \alpha_{1,3}(3f^2)\frac{\partial\lambda}{\partial x}$$

$$= -\alpha_{1,1}\frac{\partial\lambda}{\partial x} - \alpha_{1,2}(2f)\frac{\partial\lambda}{\partial x} - \alpha_{1,3}(3f^2)\frac{\partial\lambda}{\partial x}$$

$$+ \alpha_{2,1} \frac{\partial^2 x^2}{\partial x^2} + \alpha_{2,2} (2f) \frac{\partial^2 x^2}{\partial x^2} + \alpha_{2,3} (3f^2) \frac{\partial^2 x^2}{\partial x^2}$$

$$+ \alpha_{2,1} \frac{\partial^3 \lambda}{\partial x^3} - \alpha_{3,2} (2f) \frac{\partial^3 \lambda}{\partial x^3} - \alpha_{3,3} (3f^2) \frac{\partial^3 \lambda}{\partial x^3}$$

$$+ \alpha_{2,1} \frac{\partial^3 \lambda}{\partial x^3} - \alpha_{3,2} (2f) \frac{\partial^3 \lambda}{\partial x^3} - \alpha_{3,3} (3f^2) \frac{\partial^3 \lambda}{\partial x^3}$$

with final condition  $\lambda(x^{(k)}, t^{(j+1)}) = 2(f^*(x^{(k)}, t^{(j+1)}) - f(x^{(k)}, t^{(j+1)}))$  for all j, k. The parame-ters  $\alpha$  are then found using the gradient descent method with update rule 

$$\alpha_{d,p} \leftarrow \alpha_{d,p} - \eta \, \frac{\partial \mathcal{C}}{\partial \alpha_{d,p}} \tag{38}$$

(37)

where 
$$\eta = \beta \min(\Delta x)^{d-d_{\max}}$$
 and  $\frac{\partial \mathcal{C}}{\partial \alpha_{d,p}} = (-1)^d \frac{1}{\Delta x \Delta t} \int f^p \frac{\partial^d \lambda}{\partial x^d} dx dt + 2\epsilon_0 \alpha_{d,p}$ . (39)

This leads to the update rule for each coefficient, for example 

#### $\alpha_1$

$$\alpha_{1,1} \leftarrow \alpha_{1,1} - \frac{\beta}{\min(\Delta x)^2} \frac{1}{\Delta x \Delta t} \int f \frac{\partial \lambda}{\partial x} dx dt - 2\beta \epsilon_0 \alpha_{1,1}$$
  
$$\alpha_{1,2} \leftarrow \alpha_{1,2} - \frac{\beta}{\frac{1}{\Delta x} (\Delta x)^2} \frac{1}{\Delta x \Delta t} \int f^2 \frac{\partial \lambda}{\partial x} dx dt - 2\beta \epsilon_0 \alpha_{1,2}$$

$$\alpha_{1,3} \leftarrow \alpha_{1,3} - \frac{\beta}{\min(\Delta x)^2} \frac{1}{\Delta x \Delta t} \int f^3 \frac{\partial \lambda}{\partial x} dx dt - 2\beta \epsilon_0 \alpha_{1,3} \,.$$

#### F.2 REACTION DIFFUSION SYSTEM OF EQUATIONS

As mentioned in Section B.4, for this case, we consider a system consisting of two PDEs, i.e.  $N = \dim(f) = \dim(p) = 2$  where  $f = [f_1, f_2]$  and  $p = [p_1, p_2]$ , in a two-dimensional input

space, i.e.  $n = \dim(\mathbf{x}) = \dim(\mathbf{d}) = 2$  where  $\mathbf{x} = [x_1, x_2]$ , and  $\mathbf{d} = [d_1, d_2]$ . In addition, we consider candidate terms with derivatives such that  $d \in \mathcal{D}_d = \{[0,0], [1,0], [0,1], [2,0], [0,2]\}$  and polynomials such that  $p \in \mathcal{D}_p = \{[1,0], [0,1], [1,1], [2,0], [0,2], [2,1], [1,2], [3,0], [0,3]\}$ . In other words,  $d_{\text{max}} = 2$  and  $p_{\text{max}} = 3$ . The resulting forward model in Eq. 1 takes the form  $\mathcal{L}_i[\boldsymbol{f}] = \partial_t f_i + \sum_{\boldsymbol{d},\boldsymbol{r}} \alpha_{i,\boldsymbol{d},\boldsymbol{p}} \nabla_{\boldsymbol{x}}^{(\boldsymbol{d})}[\boldsymbol{f}^{\boldsymbol{p}}]$ (40)where  $i \in \{1, 2\}$ ,  $f^{p} = f_{1}^{p_{1}} f_{2}^{p_{2}}$  and  $\nabla_{x}^{(d)} = \nabla_{x_{1}}^{(d_{1})} \nabla_{x_{2}}^{(d_{2})}$ . This is equivalent to  $\mathcal{L}_{i}[f_{1}, f_{2}] = \frac{\partial f_{i}}{\partial t} + \sum_{[d_{1}, d_{2}] \in \mathcal{D}_{d}} \sum_{[p_{1}, p_{2}] \in \mathcal{D}_{p}} \alpha_{i, [d_{1}, d_{2}], [p_{1}, p_{2}]} \frac{\partial^{d_{1}+d_{2}} \left(f_{1}^{p_{1}} f_{2}^{p_{2}}\right)}{\partial x_{1}^{d_{1}} \partial x_{2}^{d_{2}}}$  $= \frac{\partial f_i}{\partial t} + \alpha_{i,[0,0],[1,0]} f_1 + \alpha_{i,[0,0],[0,1]} f_2 + \alpha_{i,[0,0],[1,1]} f_1 f_2 + \ldots + \alpha_{i,[0,0],[0,3]} f_2^3$  $+\alpha_{i,[1,0],[1,0]}\frac{\partial f_1}{\partial x_1} + \alpha_{i,[1,0],[0,1]}\frac{\partial f_2}{\partial x_1} + \alpha_{i,[1,0],[1,1]}\frac{\partial (f_1 f_2)}{\partial x_1} + \dots + \alpha_{i,[1,0],[0,3]}\frac{\partial (f_2^3)}{\partial x_1}$  $+\alpha_{i,[0,2],[1,0]}\frac{\partial^2 f_1}{\partial x_2^2} + \alpha_{i,[0,2],[0,1]}\frac{\partial^2 f_2}{\partial x_2^2} + \alpha_{i,[0,2],[1,1]}\frac{\partial^2 (f_1 f_2)}{\partial x_2^2} + \dots + \alpha_{i,[0,2],[0,3]}\frac{\partial^2 (f_2^3)}{\partial x_2^2}$ where  $i \in \{1, 2\}$ . As we can observe, we have  $|\mathcal{D}_d| \times |\mathcal{D}_p| = 5 \times 9 = 45$  terms with unknown coefficients  $\alpha_i = [\alpha_{i,d,p}]_{d \in \mathcal{D}_d, p \in \mathcal{D}_p}$  for the *i*-th PDE, i.e. a total of 90 terms for the considered system, that we aim to find using the proposed adjoint method. The cost functional in this case is simply  $\mathcal{C} = \sum_{i=1}^{2} \left( \sum_{i=1}^{2} \left( f_{i}^{*}(\boldsymbol{x}^{(k)}, t^{(j+1)}) - f_{i}(\boldsymbol{x}^{(k)}, t^{(j+1)}) \right)^{2} + \frac{1}{\Delta \boldsymbol{x} \Delta t} \int \lambda_{i}(\boldsymbol{x}, t) \mathcal{L}_{i}[\boldsymbol{f}(\boldsymbol{x}, t)] d\boldsymbol{x} dt \right) + \epsilon_{0} ||\boldsymbol{\alpha}||_{2}^{2}.$ 

The corresponding adjoint equation is given by

$$= \sum_{[d_1,d_2]\in\mathcal{D}_{\boldsymbol{d}}} \sum_{[p_1,p_2]\in\mathcal{D}_{\boldsymbol{p}}} (-1)^{d_1+d_2} \alpha_{i,[d_1,d_2],[p_1,p_2]} \frac{\partial \left(f_1^{p_1} f_2^{p_2}\right)}{\partial f_i} \frac{\partial^{d_1+d_2} \lambda_i}{\partial x_1^{d_1} \partial x_2^{d_2}}$$
(43)

and 
$$\lambda_i(\boldsymbol{x}^{(k)}, t^{(j+1)}) = 2(f_i^*(\boldsymbol{x}^{(k)}, t^{(j+1)}) - f_i(\boldsymbol{x}^{(k)}, t^{(j+1)}))$$
 for all  $j, k$  and where  $i \in \{1, 2\}$ .

Assume, without loss of generality, that i = 1. Then, we can write

 $\frac{\partial \lambda_i}{\partial t} = \sum (-1)^{|\boldsymbol{d}|} \alpha_{i,\boldsymbol{d},\boldsymbol{p}} \nabla_{f_i} [\boldsymbol{f}^{\boldsymbol{p}}] \nabla_{\boldsymbol{x}}^{(\boldsymbol{d})} [\lambda_i]$ 

$$\frac{1282}{1283} \qquad \frac{\partial \lambda_1}{\partial t} = + \alpha_{1,[0,0],[1,0]} \lambda_1 + \alpha_{1,[0,0],[1,1]} f_2 \lambda_1 + \alpha_{1,[0,0],[2,0]} (2f_1) \lambda_1 + \ldots + \alpha_{1,[0,0],[3,0]} (3f_1^2) \lambda_1 \\
\frac{1284}{1285} \qquad - \alpha_{1,[1,0],[1,0]} \frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[1,1]} f_2 \frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[2,0]} (2f_1) \frac{\partial \lambda_1}{\partial x_1} - \ldots - \alpha_{1,[1,0],[3,0]} (3f_1^2) \frac{\partial \lambda_1}{\partial x_1} \\
\frac{1286}{1286} \qquad - \alpha_{1,[1,0],[1,0]} \frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[1,1]} f_2 \frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[2,0]} (2f_1) \frac{\partial \lambda_1}{\partial x_1} - \ldots - \alpha_{1,[1,0],[3,0]} (3f_1^2) \frac{\partial \lambda_1}{\partial x_1} \\
\frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[1,0]} \frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[1,1]} f_2 \frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[2,0]} (2f_1) \frac{\partial \lambda_1}{\partial x_1} - \ldots - \alpha_{1,[1,0],[3,0]} (3f_1^2) \frac{\partial \lambda_1}{\partial x_1} \\
\frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[1,0]} \frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[1,1]} f_2 \frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[2,0]} (2f_1) \frac{\partial \lambda_1}{\partial x_1} - \ldots - \alpha_{1,[1,0],[3,0]} (3f_1^2) \frac{\partial \lambda_1}{\partial x_1} \\
\frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[1,0]} \frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[1,1]} f_2 \frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[2,0]} (2f_1) \frac{\partial \lambda_1}{\partial x_1} - \ldots - \alpha_{1,[1,0],[3,0]} (3f_1^2) \frac{\partial \lambda_1}{\partial x_1} \\
\frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[1,0]} \frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[1,1]} f_2 \frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[2,0]} (2f_1) \frac{\partial \lambda_1}{\partial x_1} - \ldots - \alpha_{1,[1,0],[3,0]} (3f_1^2) \frac{\partial \lambda_1}{\partial x_1} \\
\frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[1,0]} \frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[1,1]} \frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[2,0]} (2f_1) \frac{\partial \lambda_1}{\partial x_1} - \ldots - \alpha_{1,[1,0],[3,0]} (3f_1^2) \frac{\partial \lambda_1}{\partial x_1} \\
\frac{\partial \lambda_1}{\partial x_1} - \alpha_{1,[1,0],[1,0]} \frac{\partial \lambda_1}{\partial x_1}$$

+

. . .

$$+ \alpha_{1,[0,2],[1,0]} \frac{\partial^2 \lambda_1}{\partial x_2^2} + \alpha_{1,[0,2],[1,1]} f_2 \frac{\partial^2 \lambda_1}{\partial x_2^2} + \alpha_{1,[0,2],[2,0]} (2f_1) \frac{\partial^2 \lambda_1}{\partial x_2^2} + \dots + \alpha_{1,[0,2],[3,0]} (3f_1^2) \frac{\partial^2 \lambda_1}{\partial x_2^2}$$
(44)

and  $\lambda_1(\boldsymbol{x}^{(k)}, t^{(j+1)}) = 2(f_1^*(\boldsymbol{x}^{(k)}, t^{(j+1)}) - f_1(\boldsymbol{x}^{(k)}, t^{(j+1)}))$  for all j, k. We can follow the same procedure for i = 2. The parameters  $\alpha_i$  are then found using the gradient descent method with update rule

 $\alpha$ 

$$_{i,d,p} \leftarrow \alpha_{i,d,p} - \eta \frac{\partial \mathcal{C}}{\partial \alpha_{i,d,p}}$$

$$\tag{45}$$

(42)

where  $\eta = \beta \min(\Delta \boldsymbol{x})^{|\boldsymbol{d}| - d_{\max}} \quad \text{and} \quad \frac{\partial \mathcal{C}}{\partial \alpha_{i,\boldsymbol{d},\boldsymbol{p}}} = (-1)^{|\boldsymbol{d}|} \frac{1}{\Delta \boldsymbol{x} \Delta t} \int \boldsymbol{f}^{\boldsymbol{p}} \nabla_{\boldsymbol{x}}^{(\boldsymbol{d})} [\lambda_i] d\boldsymbol{x} dt + 2\epsilon_0 \alpha_{i,\boldsymbol{d},\boldsymbol{p}} \quad (46)$ with  $\Delta x = \Delta x_1 \Delta x_2$ , leading to the update rule for each coefficient, for example 
$$\begin{split} &\alpha_{i,[0,0],[1,0]} \leftarrow \alpha_{i,[0,0],[1,0]} - \frac{\beta}{\min(\Delta \boldsymbol{x})^2} \frac{1}{\Delta \boldsymbol{x} \Delta t} \int f_1 \lambda_i d\boldsymbol{x} dt - 2\beta \epsilon_0 \alpha_{i,[0,0],[1,0]} \\ &\alpha_{i,[1,0],[1,0]} \leftarrow \alpha_{i,[1,0],[1,0]} - \frac{\beta}{\min(\Delta \boldsymbol{x})} \frac{1}{\Delta \boldsymbol{x} \Delta t} \int f_1 \frac{\partial \lambda_i}{\partial x_1} d\boldsymbol{x} dt - 2\beta \epsilon_0 \alpha_{i,[1,0],[1,0]} \,. \end{split}$$