

Solving Inverse-PDE Problems with Physics-Aware Neural Networks

Samira Pakravan^{a,1,*}, Pouria A. Mistani^{a,1}, Miguel A. Aragon-Calvo^c, Frederic Gibou^{a,b}

^a*Department of Mechanical Engineering, University of California, Santa Barbara, CA 93106-5070*

^b*Department of Computer Science, University of California, Santa Barbara, CA 93106-5110*

^c*Instituto de Astronomía, UNAM, Apdo. Postal 106, Ensenada 22800, B.C., México*

Abstract

We propose a novel composite framework to find unknown fields in the context of inverse problems for partial differential equations (PDEs). We blend the high expressibility of deep neural networks as universal function estimators with the accuracy and reliability of existing numerical algorithms for partial differential equations as custom layers in semantic autoencoders. Our design brings together techniques of computational mathematics, machine learning and pattern recognition under one umbrella to incorporate domain-specific knowledge and physical constraints to discover the underlying hidden fields. The network is explicitly aware of the governing physics through a hard-coded PDE solver layer in contrast to most existing methods that incorporate the governing equations in the loss function or rely on trainable convolutional layers to discover proper discretizations from data. This subsequently focuses the computational load to only the discovery of the hidden fields and therefore is more data efficient. We call this architecture Blended inverse-PDE networks (hereby dubbed BiPDE networks) and demonstrate its applicability for recovering the variable diffusion coefficient in Poisson problems in one and two spatial dimensions, as well as the diffusion coefficient in the time-dependent and nonlinear Burgers' equation in one dimension. We also show that this approach is robust to noise.

Keywords: inverse problems, differential equations, deep learning, scientific machine learning, numerical methods

1. Introduction

Inverse differential problems, where given a set of measurements one seeks a set of optimal parameters in a governing differential equation, arise in numerous scientific and technological domains. Some well-known applications include X-ray tomography [20, 55], ultrasound [83], MRI imaging [36], and transport in porous media [35]. Moreover, modeling and control of dynamic complex systems is a common problem in a broad range of scientific and engineering domains, with examples ranging from understanding the motion of bacteria colonies in low Reynolds number flows [70], to the control of spinning rotorcrafts in high speed flights [30, 31]. Other applications in medicine, navigation, manufacturing, *etc.* need estimation of the unknown parameters in *real-time*, *e.g.* in electroporation [91, 53] the pulse optimizer has to be informed about tissue parameters in microsecond time. On the other hand, high resolution data-sets describing spatiotemporal evolution of complex systems are becoming increasingly available by advanced multi-scale numerical simulations (see *e.g.* [53, 52]). These advances have become possible partly due to recent developments in discretization techniques for nonlinear partial differential equations with sharp boundaries (see *e.g.* the reviews [25, 24]). However, solving these inverse problems poses substantial computational and mathematical challenges that makes it difficult to infer reliable parameters from limited data and in real-time.

*Corresponding author: spakravan@ucsb.edu

¹Equal contribution.

The problem can be mathematically formulated as follows. Let the values of $u = u(t, x_1, \dots, x_n)$ be given by a set of measurements, which may include noise. Knowing that u satisfies the partial differential equation:

$$\frac{\partial u}{\partial t} = f \left(t, x_1, \dots, x_n; u, \frac{\partial u}{\partial x_1}, \dots, \frac{\partial u}{\partial x_n}; \frac{\partial^2 u}{\partial x_1 \partial x_1}, \dots, \frac{\partial^2 u}{\partial x_1 \partial x_n}; \dots; \mathbf{c} \right),$$

find the hidden fields stored in \mathbf{c} , where the hidden fields can be constant or variable coefficients (scalars, vectors or tensors).

Deep neural networks have, rather recently, attracted considerable attention for data modeling in a vast range of scientific domains, in part due to freely available modern deep learning libraries (in particular `TensorFlow` [1]). For example, deep neural networks have shown astonishing success in emulating sophisticated simulations [29, 90, 89, 10, 75], discovering governing differential equations from data [67, 7, 47, 72], as well as potential applications to study and improve simulations of multiphase flows [25]. We refer the reader to [60, 61] for a comprehensive survey of interplays between numerical approximation, statistical inference and learning. However, these architectures require massive datasets and extensive computations to train numerous hidden weights and biases. Therefore, reducing complexity of deep neural network architectures for inverse problems poses a significant practical challenge for many applications in physical sciences, especially when the collection of large datasets is a prohibitive task [66]. One remedy to reduce the network size is to embed the knowledge from existing mathematical models [79] or known physical laws within a neural network architecture [45, 23]. Along these lines, semantic autoencoders were recently proposed by Aragon-Calvo [2], where they replaced the decoder stage of an autoencoder architecture with a given physical law that can reproduce the provided input data given a physically meaningful set of parameters. The encoder is then constrained to discover optimal values for these parameters, which can be extracted from the bottleneck of the network after training. We shall emphasize that this approach reduces the size of the unknown model parameters, and that the encoder can be used independently to infer hidden parameters in real time, while adding interpretability to deep learning frameworks. Inspired by their work, we propose to blend traditional numerical solver algorithms with custom deep neural network architectures to solve inverse PDE problems more efficiently, and with higher accuracy.

1.1. Existing works

Recently, the most widely used approach for solving forward and inverse partial differential equations using neural networks has been the constrained optimization technique. These algorithms augment the cost function with terms that describe the PDE, its boundary and its initial conditions, while the neural network acts as a surrogate for the solution field. Depending on how the derivatives in the PDEs are computed, there may be two general classes of methods that we review in the next paragraph.

In the first class, spatial differentiations in the PDE are performed exclusively using automatic differentiation, while temporal differentiation may be handled using the traditional Runge-Kutta schemes (called *discrete time models*) or using automatic differentiations (called *continuous time models*) [68]. In these methods, automatic differentiation computes gradients of the output of a neural network with respect to its input variables. Hence, the input must always be the independent variables, *i.e.* the input coordinates \mathbf{x} , time and the free parameters. In this regard, network optimization aims to calibrate the weights and biases such that the neural network outputs the closest approximation of the solution of a PDE; this is enforced through a regularized loss function. An old idea that was first proposed by Lagaris *et al.* (1998) [43]. In 2015, the general framework of solving differential equations as a learning problem was proposed by Owhadi [57, 58, 59] which revived interest in using neural networks for solving differential equations in recent years. Raissi *et al.* (2017) [68, 69] presented the regularized loss function framework under the name *physics informed neural networks* or PINNs and applied it to time-dependent PDEs. Ever since, other authors have mostly adopted PINNs, see *e.g.* [76, 4]. The second class of constrained optimization methods was proposed by Xu and Darve [88] who examined the possibility of directly using pre-existing finite discretization schemes within the loss function.

An alternative approach for solving PDE systems is through explicit embedding of the governing equations inside the architecture of deep neural networks via convolutional layers, activation functions or augmented neural networks. Below we review some of these methods:

- A famous approach is PDE-Net [47, 46] which relies on the idea of numerical approximation of differential operators by convolutions. Therefore, PDE-Nets use convolution layers with trainable and constrained kernels that mimic differential operators (such as U_x, U_y, U_{xx}, \dots) whose outputs are fed to a (symbolic) multilayer neural network that models the nonlinear response function in the PDE system, *i.e.* the right hand side in $U_t = F(U, U_x, U_y, U_{xx}, \dots)$. Importantly, PDE-Nets can only support *explicit* time integration methods, such as the forward Euler method [47]. Moreover, because the differential operators are being learned from data samples, these methods have hundreds of thousands of trainable parameters that demand hundreds of data samples; *e.g.* see section 3.1 in [47] that uses 20 δt -blocks with 17,000 parameters in each block, and use 560 data samples for training.
- Berg and Nyström [6] (hereby BN17) proposed an augmented design by using neural networks to estimate PDE parameters whose output is fed into a forward finite element PDE solver, while the adjoint PDE problem is employed to compute gradients of the loss function with respect to weights and biases of the network using automatic differentiation. Even though their loss function is a simple L_2 -norm functional, the physics is not localized in the structure of the neural network as the adjoint PDE problem is also employed for the optimization process. It is important to recognize that in their approach the numerical solver is a separate computational object than the neural network, therefore computing gradients of error functional with respect to the network parameters has to be done explicitly through the adjoint PDE problem. Moreover, their design can not naturally handle trainable parameters in the numerical discretization itself, a feature that is useful for some meshless numerical schemes. In contrast, *in BiPDEs the numerical solver is a computational layer added in the neural network architecture and naturally supports trainable parameters in the numerical scheme.* For example in the meshless method developed in section 4 we leverage this unique feature of BiPDEs to also train for shape parameters and interpolation seed locations of the numerical scheme besides the unknown diffusion coefficient.
- Dal Santos *et al.* [16] proposed an embedding of a reduced basis solver as *activation function* in the last layer of a neural network. Their architecture resembles an autoencoder in which the decoder is the reduced basis solver and the parameters at the bottleneck “are the values of the physical parameters themselves or the affine decomposition coefficients of the differential operators” [16].
- Lu *et al.* [48] proposed an unsupervised learning technique using variational autoencoders to extract physical parameters (not inhomogeneous spatial fields) from noisy spatiotemporal data. Again the encoder extracts physical parameters and the decoder propagates an initial condition forward in time given the extracted parameters. These authors use convolutional layers both in the encoder to extract features as well as in the decoder with recurrent loops to propagate solutions in time; *i.e.* the decoder leverages the idea of estimating differential operators with convolutions. Similar to PDE-Nets, this architecture is also a “PDE-integrator with explicit time stepping”, and also they need as few as 10 samples in the case of Kuramoto-Sivashinsky problem.

In these methods, a recurring idea is treating latent space variables of autoencoders as physical parameters passed to a physical model decoder. This basic idea pre-dates the literature on solving PDE problems and has been used in many different domains. Examples include Aragon-Calvo [2] who developed a galaxy model fitting algorithm using *semantic autoencoders*, or Google Tensorflow Graphics [82] which is a well-known application of this idea for scene reconstruction.

1.2. Present work

Basic criteria of developing numerical schemes for solving partial differential equations are *consistency* and *convergence* of the method, *i.e.* increasing resolution of data should yield better results.

Not only there is no guarantee that approximating differential operators through learning convolution kernels or performing automatic differentiations provide a consistent or even stable numerical method, but also the learning of convolution kernels to approximate differential operators requires more data and therefore yield less data-efficient methods. Therefore it seems reasonable to explore the idea of blending classic numerical discretization methods in neural network architectures, hence informing the neural network about proper discretization methods. This is the focus of the present manuscript.

In the present work, we discard the framework of constrained optimization altogether and instead choose to explicitly blend fully traditional finite discretization schemes as the decoder layer in semantic autoencoder architectures. In our approach, the loss function is only composed of the difference between the actual data and the predictions of the solver layer, but contrary to BN17 [6] we do not consider the adjoint PDE problem to compute gradients of the error functional with respect to network parameters. This is due to the fact that in our design the numerical solver is a custom layer inside the neural network through which backpropagation occurs naturally. This is also in contrast to PINNs where the entire PDE, its boundary and its initial conditions are reproduced by the output of a neural network by adding them to the loss function. Importantly, the encoder learns an approximation of the inverse transform in a *self-supervised* fashion that can be used to evaluate the hidden fields underlying unseen data without any further optimization. Moreover, the proposed framework is versatile as it allows for straightforward consideration of other domain-specific knowledge such as symmetries or constraints on the hidden field. In this work, we develop this idea for stationary and time-dependent PDEs on structured and unstructured grids and on noisy data using mesh-based and mesh-less numerical discretization methods.

1.3. Novelties and features of BiPDEs

A full PDE solver is implemented as a *custom layer inside the architecture of semantic autoencoders* to solve inverse-PDE problems in a self-supervised fashion. Technically this is different than other works that implement a propagator decoder by manipulating activation functions or kernels/biases of convolutional layers, or those that feed the output of a neural network to a separate numerical solver such as in BN17 which requires the burden of considering the adjoint problem in order to compute partial differentiations. The novelties and features of this framework are summarized below:

1. **General discretizations.** We do not limit numerical discretization of differential equations to only finite differences that are emulated by convolution operations, our approach is more general and permits employing more sophisticated numerical schemes such as meshless discretizations. It is a more general framework that admits any existing discretization method directly in a decoder stage.
2. **Introducing solver layers.** All the information about the PDE system is *only* localized in a solver layer; *i.e.* we do not inform the optimizer or the loss function with the adjoint PDE problem, or engineer regularizers or impose extra constraints on the kernels of convolutions, or define exotic activation functions as reviewed above. In other words, PDE solvers are treated as custom layers similar to convolution operations that are implemented in convolutional layers. An important aspect is the ability to employ any of the usual loss functions used in deep learning, for example we arbitrarily used mean absolute error or mean squared error in our examples.
3. **Blending meshless methods with trainable parameters.** Another unique proposal made in this work is the use of Radial Basis Function (RBF) based PDE solver layers as a natural choice to blend with deep neural networks. Contrary to other works, the neural network is not only used as an estimator for the unknown field but also it is tasked to optimize the shape parameters and interpolation points of the RBF scheme. In fact, our meshless decoder is not free of trainable parameters similar to reviewed works, instead shape parameters and seed locations are trainable parameters that define the RBF discretization, this is analogous to convolutional layers with trainable weights/biases that are used in machine learning domain. In fact this presents an example of neural networks complementing numerical

discretization schemes. Choosing optimal shape parameters or seed locations is an open question in the field of RBF-based PDE solvers and here we show neural networks can be used to optimally define these discretization parameters.

4. **Explicit/implicit schemes.** Most of the existing frameworks only accept explicit numerical discretizations in time, however our design naturally admits implicit methods as well. Using implicit methods allows taking bigger timesteps for stiff problems such as the diffusion problem, hence not only providing faster inverse-PDE solvers, but also present more robust/stable inverse PDE solvers.
5. **Data efficient.** Our design lowers the computational cost as a result of reusing classical numerical algorithms for PDEs during the learning process, which focuses provided data to infer the actual unknowns in the problem, *i.e.* reduces the load of learning a discretization scheme from scratch.
6. **Physics informed.** Domain-specific knowledge about the unknown fields, such as symmetries or specialized basis functions, can be directly employed within our design.
7. **Inverse transform.** After training, the encoder can be used independently as a real-time estimator for unknown fields, *i.e.* without further optimization. In other words, the network can be pre-trained and then used to infer unknown fields in real-time applications.

2. Blended inverse-PDE network (BiPDE-Net)

The basic idea is to embed a numerical solver into a deep learning architecture to recover unknown functions in inverse-PDE problems, and all the information about the governing PDE system is only encoded inside the DNN architecture as a solver layer. In this section we describe our proposed architectures for inverse problems in one and two spatial dimensions.

2.1. Deep neural networks (DNN)

The simplest neural network is a single layer of perceptron that mathematically performs a linear operation followed by a nonlinear composition applied to its input space,

$$\mathcal{N} = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b}), \quad (1)$$

where σ is called the *activation function*. Deep neural networks are multiple layers stacked together within some architecture. The simplest example is a set of layers connected in series without any recurrent loops, known as feedforward neural networks (FNN). In a densely connected FNN, the action of the network is simply the successive compositions of previous layer outputs with the next layers, *i.e.*,

$$\mathcal{N}_l = \sigma(\mathbf{W}_l \mathcal{N}_{l-1}(\mathbf{x}) + \mathbf{b}_l), \quad (2)$$

where l indicates the index of a layer. This compositional nature of NNs is the basis of their vast potential as universal function estimators of any arbitrary function on the input space \mathbf{x} , see e.g. [80, 15, 13]. Another important feature of NNs is that they can effectively express certain high dimensional problems with only a few layers, for example Darbon *et al.* [17] have used NNs to overcome the curse of dimensionality for some Hamilton-Jacobi PDE problems (also see [27, 76]).

Most machine learning models are reducible to composition of simpler layers which allows for more abstract operations at a higher level. Common layers include dense layers as described above, convolutional layers in convolutional neural networks (CNNs) [44, 42], Long-short term memory networks (LSTM) [33], Dropout layers [77] and many more. In the present work, we pay particular attention to CNNs owing to their ability to extract complicated spatial features from high dimensional input datasets. Furthermore, we define custom PDE solver layers as new member of the family of pre-existing layers by directly implementing numerical discretization schemes inside the architecture of deep neural networks.

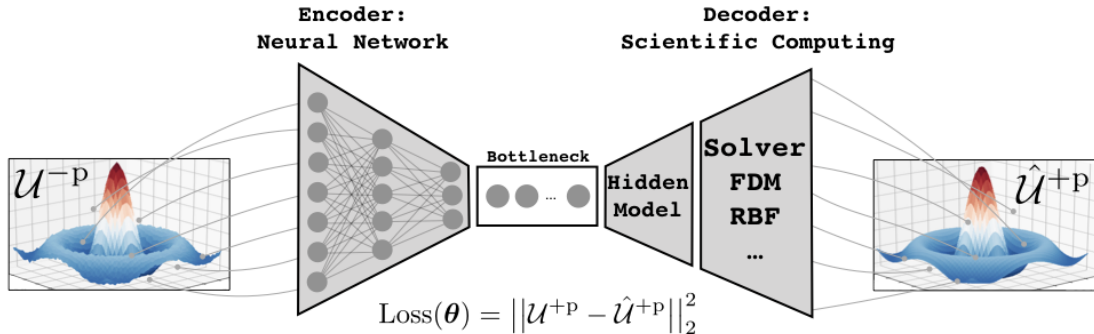


Figure 1: Architecture of the BiPDE to infer unknown parameters of hidden fields. Here the loss function is the mean squared error between data and output of the autoencoder, however other choices for loss function may be used depending on the nature of data.

2.2. Custom solver layers

A *layer* is a high level abstraction that plays a central role in existing deep learning frameworks such as TensorFlow² [1], Keras API [12], PyTorch [62], *etc.* Each Layer encapsulates a state, *i.e.* trainable parameters such as weights/biases, and a transformation of inputs to outputs. States in a layer could also be non-trainable parameters in which case they will be excluded from backpropagation during training.

We implement different explicit or implicit numerical discretization methods as custom layers that transform an unknown field, initial data and boundary conditions to outputs in the solution space. Solver layers encapsulate numerical discretization schemes with trainable (*e.g.* shape parameters and seeds in meshless methods) or non-trainable (*e.g.* the finite difference methods) state parameters. Interestingly, solver layers with trainable parameters are new computational objects analogous to pre-existing convolutional layers with trainable kernel parameters.

An important aspect of layer objects is that they can be composed with other layers in any order. Particularly, this offers an interesting approach for solving inverse problems given by systems of partial differential equations with several unknown fields that can be modeled with neural layers. We will explore this avenue in future work. In the remainder of this manuscript we will only focus on different inverse-PDE examples given by a single PDE equation and one unknown field.

2.3. Blended neural network architectures

BiPDE is a two-stage architecture, with the first stage responsible for learning the unknown coefficients and the second stage performing numerical operations as in traditional numerical solvers (see figure 1). To achieve higher performance, it is essential to use GPU-parallelism. We leverage the capability provided by the publicly available library TensorFlow [1] by implementing our PDE-solver as a *custom layer* into our network using the Keras API [12]. Details of this includes vectorized operations to build the linear system associated by the PDE discretization.

We propose a semantic autoencoder architecture as proposed by Aragon-Calvo (2019) [2] with hidden parameters being represented at the bottleneck of the autoencoder. Figure 1 illustrates the architecture for the proposed semantic autoencoder. Depending on static or time dependent nature of the governing PDE, one may train this network over pairs of input-output solutions that are shifted p steps in time, such that for a static PDE we have $p = 0$ while dynamic PDEs correspond to $p \geq 1$. We call this parameter the *shift parameter*, which will control the accuracy of the method (*cf.* see section 4).

²For example see TensorFlow manual page at https://www.tensorflow.org/guide/keras/custom_layers_and_models

An important aspect is that the input to BiPDE is the solution data itself. In other words the neural network in a BiPDE is learning the *inverse transform*,

$$\mathcal{NN} : \{u\} \rightarrow \text{hidden field}, \quad (3)$$

where $\{u\}$ indicates an ensemble of solutions, *e.g.* solutions obtained with different boundary conditions or with different hidden fields. Note that in other competing methods such as PINNs the input is sanctioned to be the coordinates in order for automatic differentiation to compute spatial and temporal derivatives; as a consequence PINNs can only be viewed as *surrogates* for the solution of the differential problem defined on the space of coordinates. However, we emphasize that semantic autoencoders are capable to approximate the inverse transformation from the space of solutions to the space of hidden fields, a feature that we exploit in section 3.1.2.

Essentially different numerical schemes can be implemented in the decoder stage. We will blend examples of both mesh-based and mesh-less numerical discretizations and present numerical results and comparisons with PINNs. We will show how BiPDEs can handle data on unstructured grids and data with added noise. In section 3, we demonstrate performance of mesh-based BiPDEs on inverse problems in two spatial dimensions by using a finite difference discretization and Zernike expansion of the non-homogeneous hidden field, we will consider both stationary and dynamic PDE problems in this section. Then in section 4, we develop a mesh-less BiPDE and consider a dynamic nonlinear inverse partial differential problem.

3. Mesh-based BiPDE: Finite Differences

We consider a variable coefficient Poisson problem in one and two spatial dimensions as well as the one dimensional nonlinear Burger’s equation as an example of a nonlinear dynamic PDE problem with a scalar unknown parameter.

3.1. Stationary Poisson problem

We consider the governing equation for diffusion dominated processes in heterogeneous media:

$$\nabla \cdot (D(\mathbf{x})\nabla u) = -f(\mathbf{x}), \quad \mathbf{x} \in \Omega \quad (4)$$

$$u(\mathbf{x}) = u_0(\mathbf{x}), \quad \mathbf{x} \in \partial\Omega \quad (5)$$

Here we consider a rectangular domain with Dirichlet boundary conditions.

Discretization. In our architecture, we use the standard 5-point stencil finite difference discretization of the Poisson equation in the solver layer, *i.e.*

$$\frac{D_{i-1/2,j}u_{i-1,j} - (D_{i-1/2,j} + D_{i+1/2,j})u_{i,j} + D_{i+1/2,j}u_{i+1,j}}{\Delta x^2} + \frac{D_{i,j-1/2}u_{i,j-1} - (D_{i,j-1/2} + D_{i,j+1/2})u_{i,j} + D_{i,j+1/2}u_{i,j+1}}{\Delta y^2} + f_{i,j} = 0,$$

and we use the linear algebra solver implemented in `TensorFlow` to solve for the solution field, *i.e.* we used `tf.linalg.solve` method that is a dense linear system solver. Of course, this can be improved by implementing a sparse linear solver.

Hidden Model. We decompose the hidden field into a finite number of eigenfunctions and search for their optimal coefficients. This is also advantageous from a physics point of view, because domain’s knowledge of hidden fields can be naturally formulated in terms of basis functions into this framework. One such family of series expansions are the moment-based methods that have been largely exploited in image reconstruction [40, 5, 64, 3]. In particular, Zernike moments [84] provide a linearly independent set of polynomials defined on the unit circle/sphere in two/three spatial dimensions. Zernike moments are well-suited for such a task and are commonly used for representing optical aberration in astronomy and atmospheric sciences [65], for image reconstruction and for enhanced ultrasound focusing in biomedical imaging [19, 50, 39].

n	 m 	R_{nm}	Z_{nm}^o	Z_{nm}^e	Aberration/Pattern
0	0	1	0	1	Piston
1	1	ρ	$\rho \sin(\theta)$	$\rho \cos(\theta)$	Tilt
2	0	$2\rho^2 - 1$	0	$2\rho^2 - 1$	Defocus
	2	ρ^2	$\rho^2 \sin(2\theta)$	$\rho^2 \cos(2\theta)$	Oblique/Vertical Astigmatism
3	1	$3\rho^3 - 2\rho$	$(3\rho^3 - 2\rho) \sin(\theta)$	$(3\rho^3 - 2\rho) \cos(\theta)$	Vertical/Horizontal Coma
	3	ρ^3	$\rho^3 \sin(3\theta)$	$\rho^3 \cos(3\theta)$	Vertical/Oblique Trefoil
4	0	$6\rho^4 - 6\rho^2 + 1$	0	$6\rho^4 - 6\rho^2 + 1$	Primary Spherical
	2	$4\rho^4 - 3\rho^2$	$(4\rho^4 - 3\rho^2) \sin(2\theta)$	$(4\rho^4 - 3\rho^2) \cos(2\theta)$	Oblique/Vertical Secondary Astigmatism
	4	ρ^4	$\rho^4 \sin(4\theta)$	$\rho^4 \cos(4\theta)$	Oblique/Vertical Quadrafoil

Table 1: First 15 odd and even Zernike polynomials according to Noll’s nomenclature. Here, the ordering is determined by ordering polynomial with lower radial order first, cf. [86].

Zernike moments are advantageous over regular moments in that they intrinsically provide rotational invariance, higher accuracy for irregular patterns, and are orthogonal, which reduces information redundancy in the different coefficients. Zernike polynomials capture deviations from zero mean as a function of radius and azimuthal angle. Furthermore, the complete set of orthogonal bases provided by Zernike moments can be obtained with lower computational precision from input data, which enhances the robustness of the reconstruction procedure.

Odd and even Zernike polynomials are given as a function of the azimuthal angle θ and the radial distance ρ between 0 and 1 measured from the center of image,

$$\begin{bmatrix} Z_{nm}^o(\rho, \theta) \\ Z_{nm}^e(\rho, \theta) \end{bmatrix} = R_{nm}(\rho) \begin{bmatrix} \sin(m\theta) \\ \cos(m\theta) \end{bmatrix},$$

with

$$R_{nm}(\rho) = \begin{cases} \sum_{l=0}^{(n-|m|)/2} \frac{(-1)^l (n-l)!}{l! [(n+|m|)/2-l]! [(n-|m|)/2-l]!} \rho^{n-2l} & \text{for } n-m \text{ even,} \\ 0 & \text{for } n-m \text{ odd,} \end{cases}$$

where n and m are integers with $n \geq |m|$. A list of radial components is given in table 1 (from [85]). For an extensive list of Zernike polynomials in both two and three spatial dimensions, we refer the interested reader to [51].

Furthermore, each Zernike moment is defined by projection of the hidden field $f(x, y)$ on the orthogonal basis,

$$\begin{bmatrix} A_{nm} \\ B_{nm} \end{bmatrix} = \frac{n+1}{\epsilon_{mn}^2 \pi} \int_x \int_y f(x, y) \begin{bmatrix} Z_{nm}^o(x, y) \\ Z_{nm}^e(x, y) \end{bmatrix} dx dy, \quad x^2 + y^2 \leq 1,$$

where for $m = 0$, $n \neq 0$ we defined $\epsilon_{0n} = 1/\sqrt{2}$ and $\epsilon_{mn} = 1$ otherwise. Finally, superposition of these moments expands the hidden field in terms of Zernike moments:

$$\hat{f}(x, y) = \sum_{n=0}^{N_{max}} \sum_{|m|=0}^n [A_{nm} Z_{nm}^o(r, \theta) + B_{nm} Z_{nm}^e(r, \theta)]. \quad (6)$$

In order to identify the coefficients in the Zernike expansion (6) for hidden fields, we use a semantic autoencoder architecture with Zernike moments being represented by the code at the bottleneck of the autoencoder. Figure 2 illustrates the architecture for the proposed semantic autoencoder.

Architecture. Even though a shallow neural network with as few neurons as the number of considered Zernike terms suffices to estimate values of the unknown Zernike moments in each of

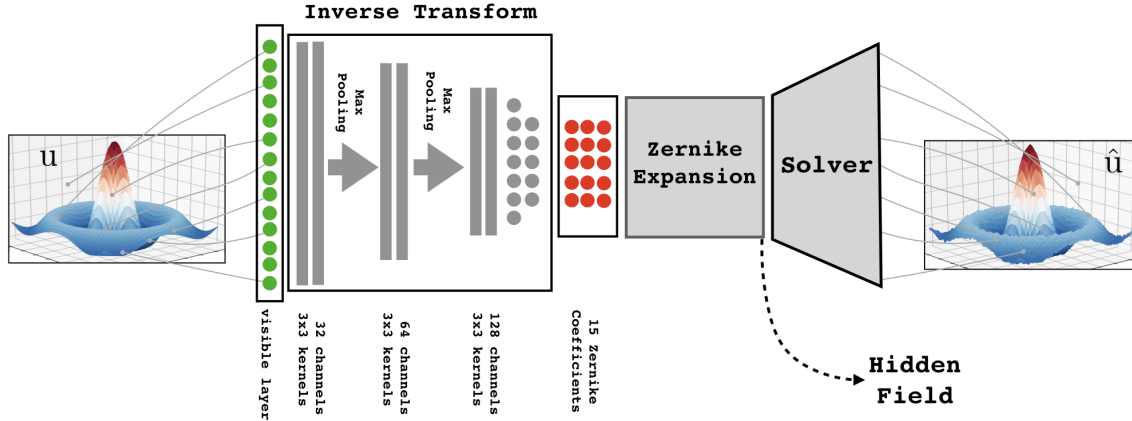


Figure 2: Architecture of the semantic autoencoder to infer hidden fields. Zernike moments are discovered at the bottleneck of the architecture.

the problems considered in this section, however we will use a deep convolutional neural network (detailed below) in order to achieve our ultimate goal of approximating the inverse transform for the Poisson problem in a broad range of diffusion coefficient fields. Therefore we design one deep neural network and uniformly apply it to several problems in this section.

In the training process of a CNN, the kernels are trained at each layer such that several feature maps are extracted at each layer from input data. The CNN is composed of 3 convolutional blocks with 32, 64, 128 channels respectively and kernel size 3×3 . Moreover, we use the **MaxPooling** filter with kernel size (2, 2) after each convolutional block to downsample the feature maps by calculating the maximum values of each patch within these maps. We use the **ReLU** activation function [26], *i.e.* a piecewise linear function that only outputs positive values: $\text{ReLU}(x) = \max(0, x)$, in the convolutional layers followed by a **Sigmoid** activation in dense layers and a scaled **Sigmoid** activation at the final layer,

$$\tilde{\sigma}(x) = D_{\min} + (D_{\max} - D_{\min})\sigma(x), \quad (7)$$

such that the actual values of the diffusion coefficient are within the range (D_{\min}, D_{\max}) , known from domain specific knowledge. After each dense layer, we apply **Dropout** layers with a rate of 0.2 to prevent overfitting [32, 77] (a feature that is most useful in estimating the inverse transform operator) and avoid low quality local minima during training.

3.1.1. Test cases.

Case I. A tilted plane. In the first example we consider a linear diffusion model given by

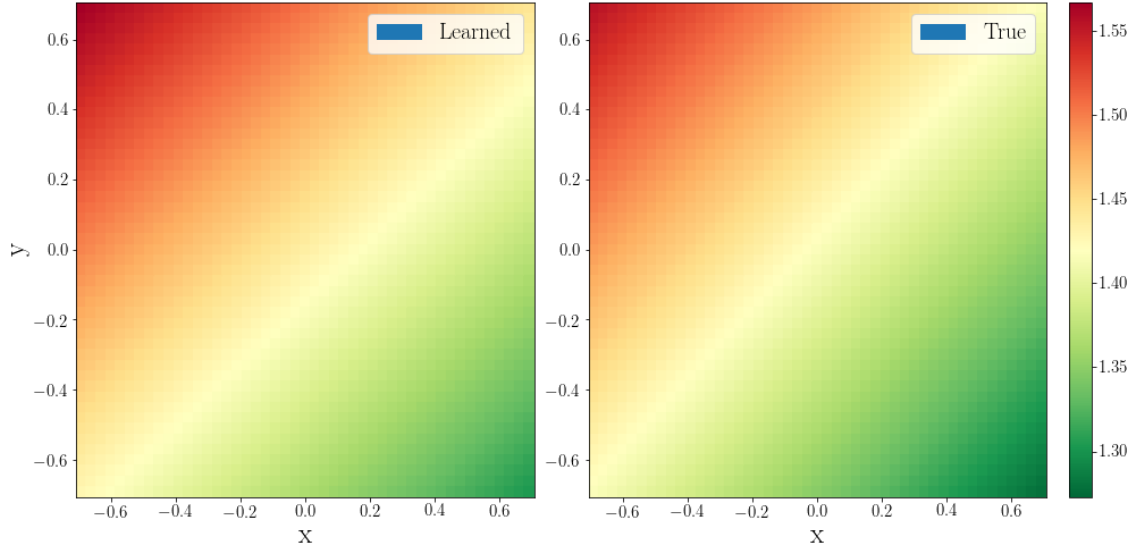
$$D(x, y) = \sqrt{2} + 0.1(y - x)$$

where the boundary condition function u_{BC} and the source field f are given by

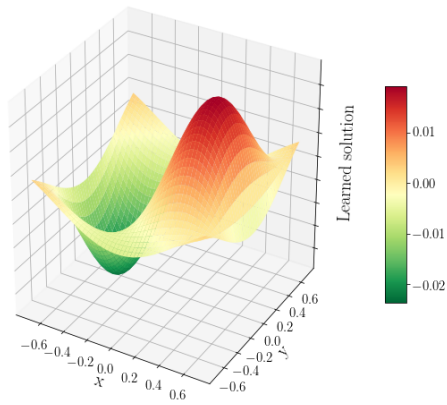
$$u_{BC}(x, y) = 0.01 \cos(\pi x) \cos(\pi y) \quad \text{and} \quad f(x, y) = \sin(\pi x) \cos(\pi y)$$

In this experiment we only use a *single solution field* for training. Even though in our experiments the method succeeded to approximate the hidden field even with a *single grid point* to compute the loss function, here we consider all the grid points in the domain to obtain improved accuracy in the results. We trained the network for 30 epochs using an **Adam** optimizer [41] that takes 170 seconds on a Tesla T4 GPU available on a free Google Colaboratory account³. Figure 3 depicts the results

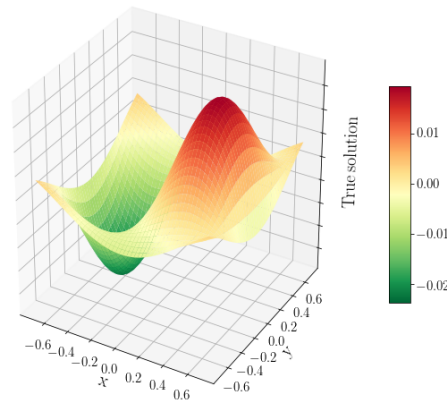
³<https://colab.research.google.com/>



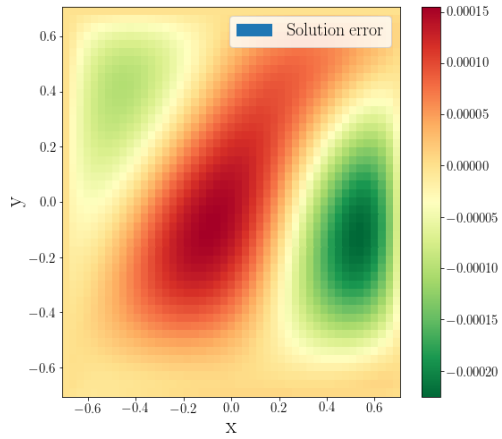
(a) Comparison of learned (left) versus true diffusion coefficient (right).



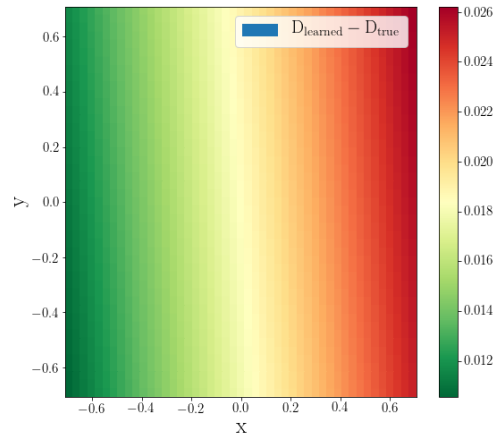
(b) Learned solution.



(c) True solution.



(d) Error in learned solution $u - \hat{u}$.



(e) Error in learned diffusion coefficient.

Figure 3: Results for the two dimensional tilted plane (case I).

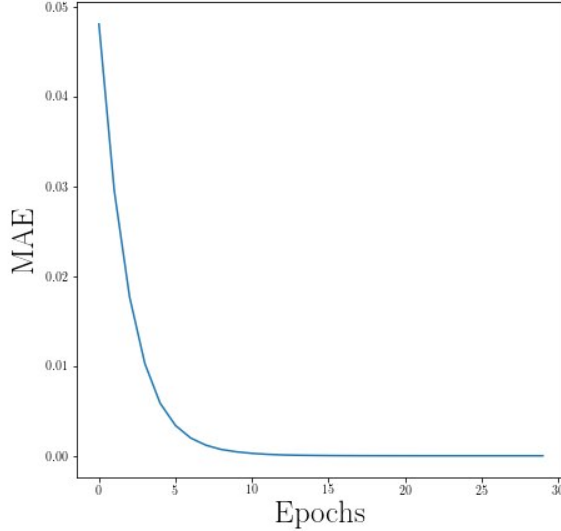


Figure 4: Mean absolute error in solution vs. epochs for the two dimensional tilted plane (case I).

obtained by the proposed scheme. The diffusion map is discovered with a maximum relative error of only 2%, while the error in the solution field is 1%. It is noteworthy to mention that the accuracy of the results in this architecture are influenced by the accuracy of the discretizations used in the solver layer. While we used a second-order accurate finite difference discretization, it is possible to improve these results by using higher order discretizations instead. We leave such optimizations to future work.

T	# params	C(32)	C(32)	C(64)	C(64)	C(128)	C(128)	D(64)	D(32)	MAE _D	L _D [∞]
1	1,468,323	Y	Y	Y	Y	Y	Y	Y	Y	0.0144207	0.0294252
2	1,459,075	Y	-	Y	Y	Y	Y	Y	Y	0.0193128	0.0267854
3	1,422,147	Y	-	Y	-	Y	Y	Y	Y	0.0226252	0.0527432
4	1,274,563	Y	-	Y	-	Y	-	Y	Y	0.0199361	0.0272122
5	682,627	Y	-	Y	-	Y	-	-	Y	0.0141946	0.0243868
6	313,859	Y	-	Y	-	-	-	-	Y	0.0301841	0.0544990
7	46,467	Y	-	Y	-	-	-	-	-	0.0190432	0.0264254
8	6,915	-	-	-	-	-	-	-	-	0.0183808	0.0267156

Table 2: Influence of architecture of the decoder stage on mean absolute error $\text{MAE}_D \equiv \sum |D(\mathbf{x}) - \hat{D}(\mathbf{x})|/N$ and maximum error L_D^∞ in the discovered hidden field in case I. Double vertical lines correspond to `MaxPooling2D()` layers and triple vertical lines correspond to `Flatten()` layer. C(o) and D(o) stand for `conv2D(filters)` and `Dense(neurons)` layers respectively. There are 3 neurons at the bottleneck not shown in the table.

Influence of architecture. Table 2 tabulates the mean absolute error in the discovered tilted plane diffusion coefficient for different architectures of the encoder stage. No significant improvement is observed for deeper or shallower encoder network for the example considered here.

Case II. superimposed Zernike polynomials. We consider a more complicated hidden diffusion field given by

$$D(x, y) = 4 + a_0 + 2a_1x + 2a_2y + \sqrt{3}a_3(2x^2 + 2y^2 - 1).$$

The boundary condition function u_{BC} and the source field f are given by

$$u_{BC}(x, y) = \cos(\pi x) \cos(\pi y) \quad \text{and} \quad f(x, y) = x + y.$$

Figure 5 illustrates the performance of the proposed Zernike-based network using a mean absolute error measure for the loss function. We trained the network for 100 epochs using an Adam optimizer [41].

Resilience to noise. We also assess the performance of our scheme on noisy datasets. We consider a zero-mean Gaussian noise with standard deviation 0.025 superimposed on the solution field. Figure 6 depicts the solution learned from a noisy input image. The network succeeds in discovering the diffusion field with comparable accuracy as in the noise-free case. Note that this architecture naturally removes the added noise from the learned solution, a feature that is similar to applying a low-pass filter on noisy images.

3.1.2. Learning the inverse transform

In the previous sections, we have applied BiPDE to find the variable diffusion coefficient from a single input image. Another interesting feature of the proposed semantic autoencoder architecture is its ability to train neural networks in order to discover the inverse transform for the underlying hidden fields *in a self-supervised fashion*. In this scenario, the trained encoder learns the inverse transform function that approximates the hidden parameters given a solution field to its input. Note that even though the same task could be accomplished by supervised learning of the hidden fields, *i.e.* by explicitly defining loss on the hidden fields without considering the governing equations, BiPDEs substitute the data labels with a governing PDE and offer comparable prediction accuracy. In this section we train BiPDEs over ensembles of solution fields to estimate hidden Zernike moments of diffusion coefficients underlying unseen data.

One dimensional inverse transform

We build a one dimensional semantic autoencoder using 3 layers with 100, 40, and 2 neurons respectively. We used the ReLU activation function for the first two layers and a Sigmoid activation function for the last layer representing the hidden parameters. A linear solver is then stacked with this encoder that uses the second order accurate finite difference discretization, *i.e.*

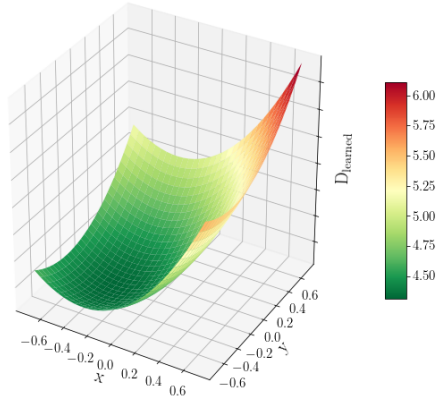
$$\frac{D_{i-1/2}u_{i-1} - (D_{i-1/2} + D_{i+1/2})u_i + D_{i+1/2}u_{i+1}}{\Delta x^2} + f_i = 0, \quad D_{i+1/2} = \frac{D_i + D_{i+1}}{2}$$

However, the diffusion map is internally reconstructed using the hidden parameters before feeding the output of the encoder to the solver. As a test problem, we consider the one dimensional Poisson problem with a generic linear form for the diffusion coefficient,

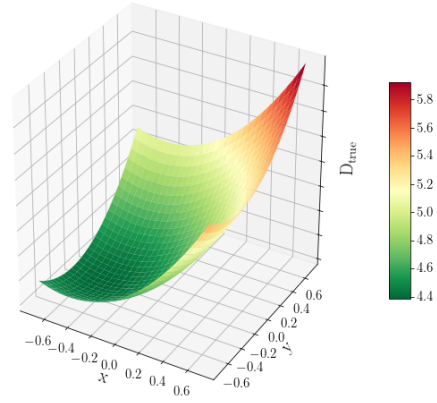
$$D(x) = 1 + a_0 + a_1x.$$

We consider identical left and right Dirichlet boundary conditions of 0.2 for all images and let the source term be $f(x) = \sin(\pi x)$. We consider random diffusion coefficients a_0 and a_1 with a uniform distribution in $[0.25, 0.75]$ and we generate 1000 solutions over the domain $x \in [-1, 1]$. We train BiPDE over 900 images from this dataset and validate its performance over the remaining 100 images using a mean squared error loss function for 1000 epochs. Each image is generated on a uniform grid with $N_x = 160$ grid points. We used a batch size of 100 in these experiments using the Adam optimizer. Figure 7(c) shows loss versus epochs in this experiment. Figure 8 compares learned and true coefficients over two independent test samples containing 1000 solutions, with and without a zero-mean Gaussian noise with standard deviation 0.025, *i.e.* amounting to $\sim 13\%$ added noise over the images.

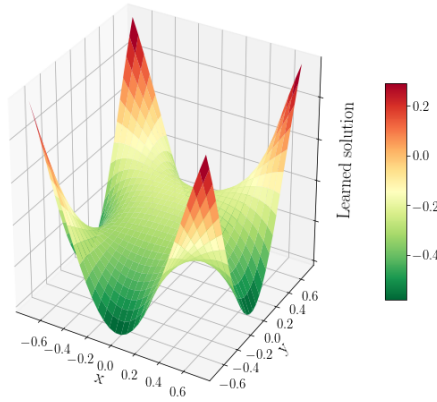
In figure 8, we expanded the range of unknown parameters $a_0, a_1 \in [0.15, 0.85]$ in our test sample to assess performance of trained encoder over unseen data that are outside the range of training data (as a measure of generalizability of BiPDEs). In this figure blue points correspond to new images whose true unknown parameters fall inside the training range, and red data points correspond to those outside the training range. We observe that the encoder is able to predict the unknown parameters even outside of its training range, although its accuracy gradually diminishes



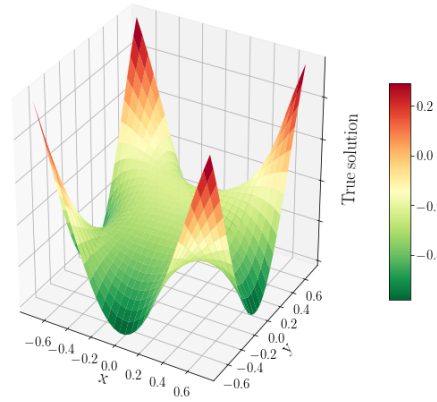
(a) Learned diffusion.



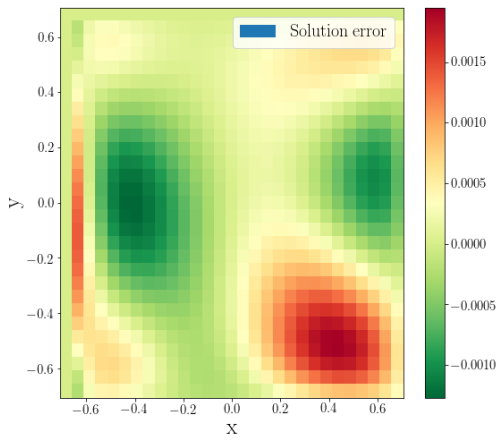
(b) True diffusion.



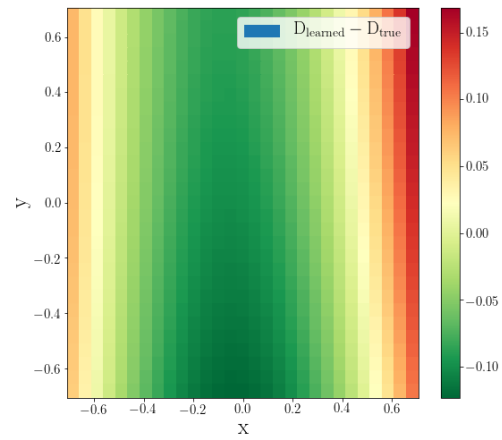
(c) Learned solution.



(d) True solution.

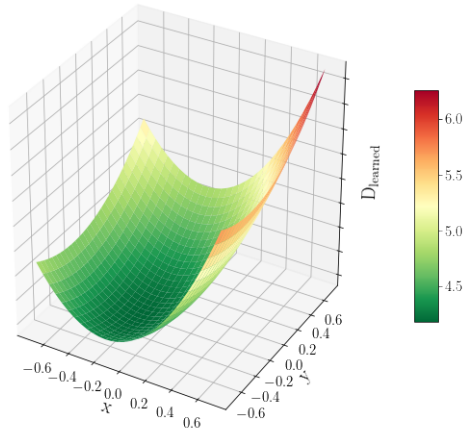


(e) Error in learned solution $u - \hat{u}$.

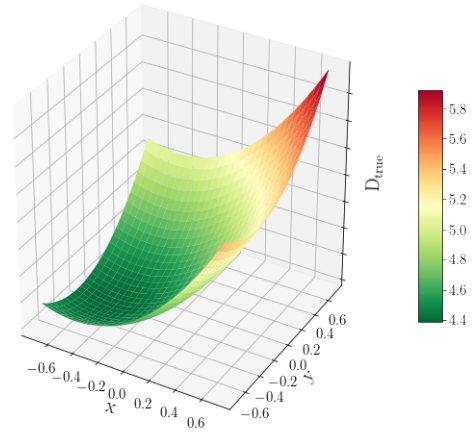


(f) Error in learned diffusion coefficient.

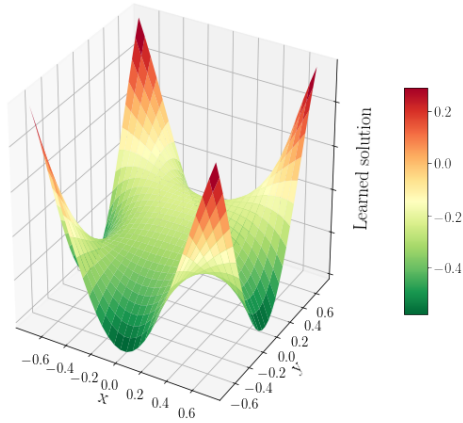
Figure 5: Results in the two dimensional parabolic case.



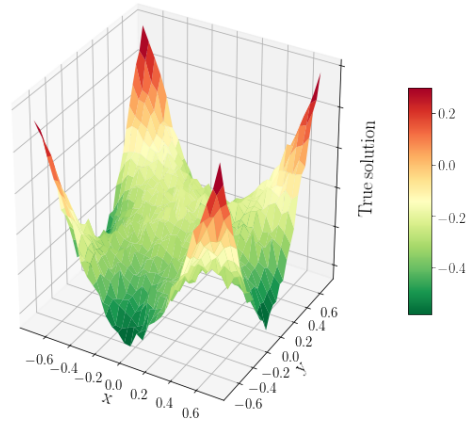
(a) Learned diffusion.



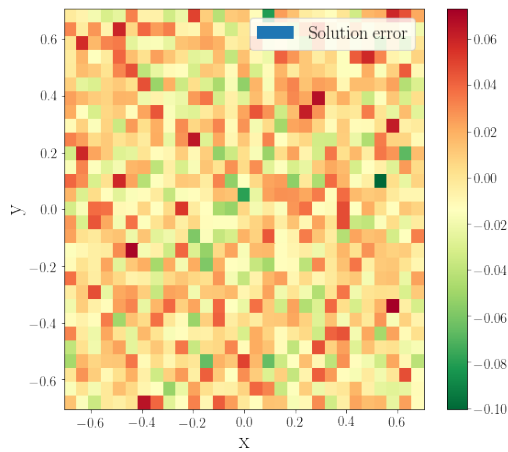
(b) True diffusion.



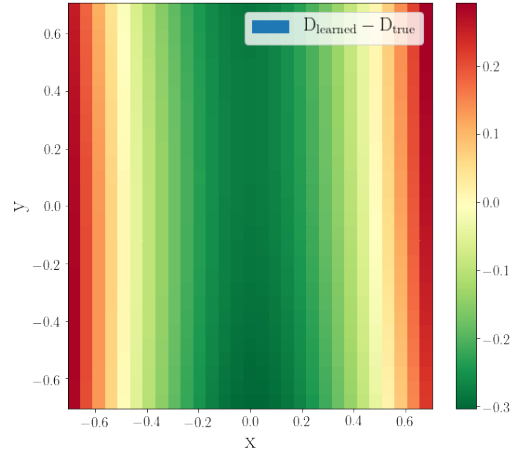
(c) Learned solution.



(d) Noisy input solution.

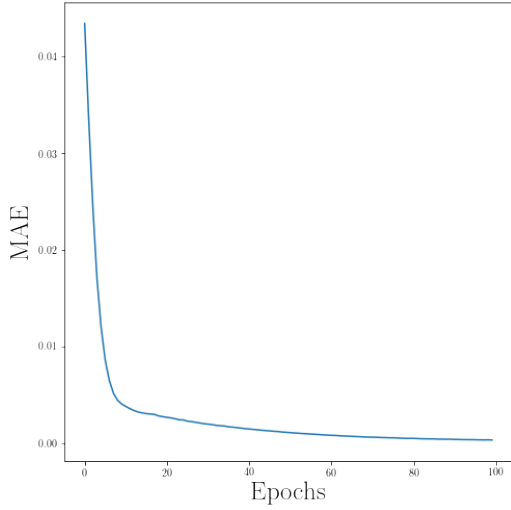


(e) Error in learned solution $u - \hat{u}$.

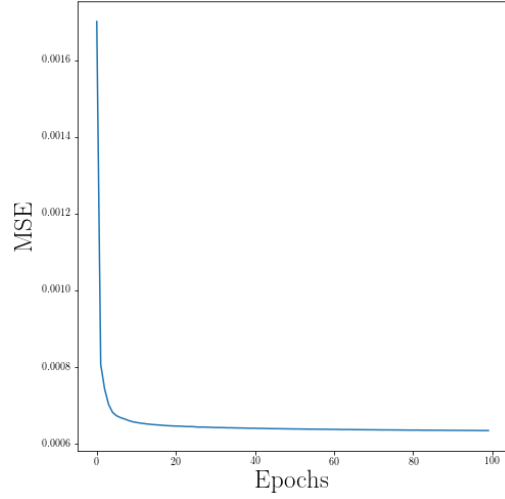


(f) Error in learned diffusion coefficient.

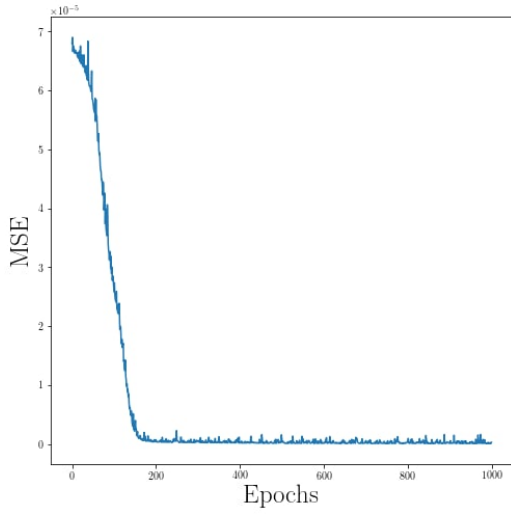
Figure 6: Results in the two dimensional case with added noise. After 300 epochs the network discovers the hidden diffusion field with a maximum relative error of 5%. Interestingly the learned solution is resilient to added noise and the network approximates a noise-free solution.



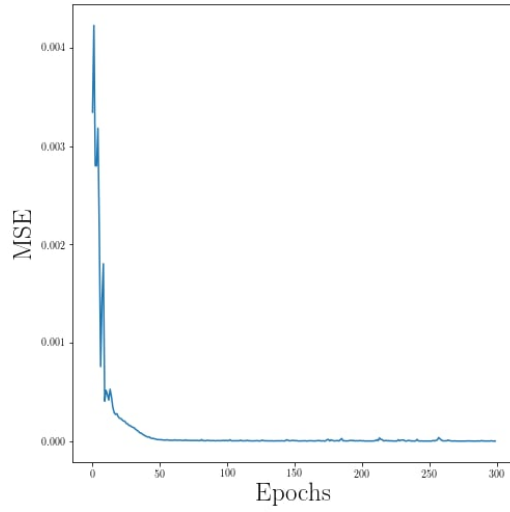
(a) L_1 loss vs. epoch for case II without added noise.



(b) L_2 loss vs. epoch for case II with added noise.

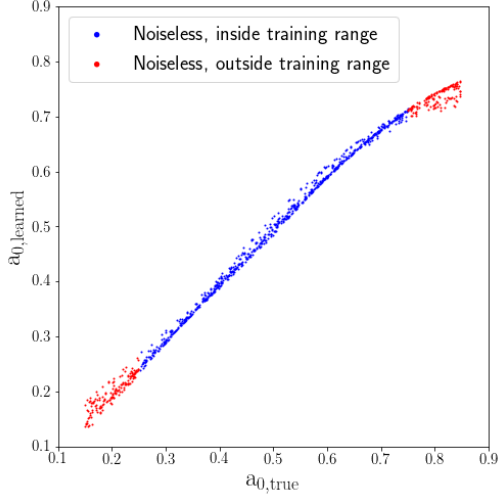


(c) L_2 loss vs. epoch for 1D inverse transform.

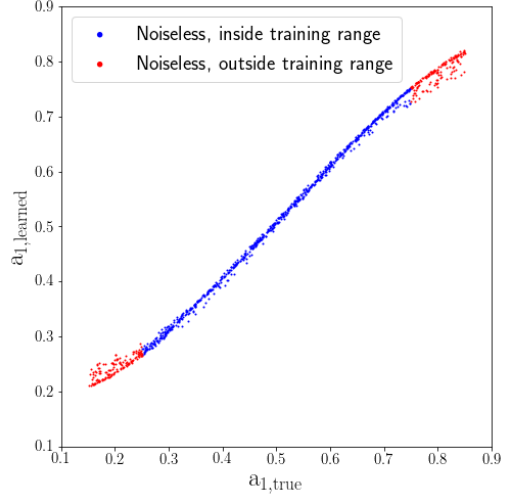


(d) L_2 loss vs. epoch for 2D inverse transform.

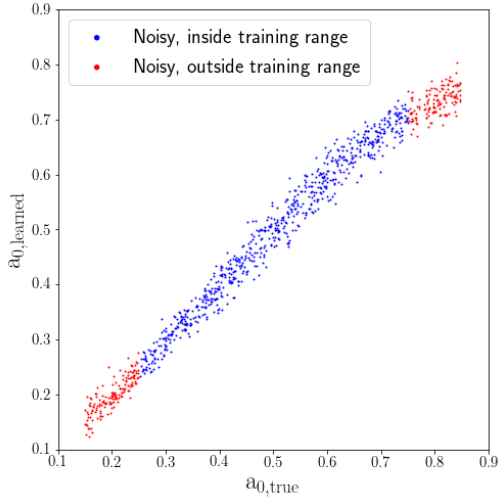
Figure 7: Mean absolute/square error vs. epochs for (top panel) the two dimensional parabolic experiment (case II) with and without added Gaussian noise of section 3.1.1, and (bottom panel) the inverse transform for 1D and 2D experiments of section 3.1.2.



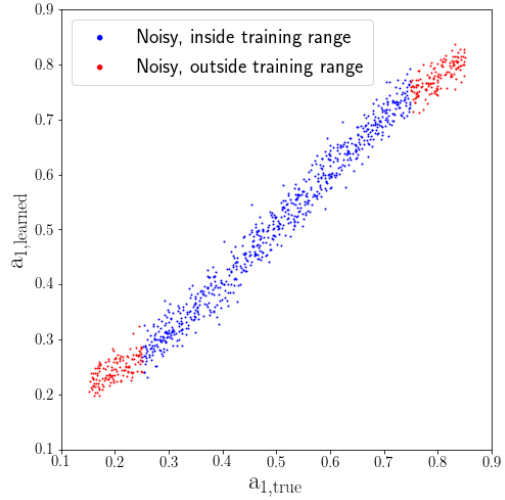
(a) Regression quality is $R^2 = 0.9906891$.



(b) Regression quality is $R^2 = 0.9953392$.



(c) Regression quality is $R^2 = 0.9796781$.



(d) Regression quality is $R^2 = 0.9834912$.

Figure 8: (Top, bottom) panel shows performance of BiPDE over 1000 randomly chosen one-dimensional images with $N_x = 160$ grid points after 1000 epochs (with,without) added zero-mean Gaussian noise with standard deviation 0.025 to the test sample. The hidden diffusion coefficient is $D(x) = 1 + a_0 + a_1x$. In each case the R^2 coefficient is reported for the blue data points, where unknown parameters fall within the training range $[0.25, 0.75]$. Red data points show predictions outside of training range. Network has 20, 222 trainable parameters, and training takes ~ 2 seconds per epoch on a Tesla T4 GPU available on a free Google Colaboratory account.

1D inverse transform	Noiseless		Noisy (13% relative noise)	
Sample Size, $N_x = 100$	a_0	a_1	a_0	a_1
$N_{\text{data}} = 250$	0.9953634	0.9977753	0.9609166	0.9570264
$N_{\text{data}} = 500$	0.9979478	0.9988417	0.9644154	0.9640230
$N_{\text{data}} = 1000$	0.9990417	0.9992921	0.9600430	0.9586783
$N_{\text{data}} = 2000$	0.9995410	0.9997107	0.9599427	0.9652383
$N_{\text{data}} = 4000$	0.9994279	0.9994974	0.9603496	0.9661519
$N_{\text{data}} = 8000$	0.9998054	0.9998115	0.9614795	0.9619859
Grid Points, $N_{\text{data}} = 1000$	a_0	a_1	a_0	a_1
$N_x = 20$	0.9900532	0.9987560	0.8623348	0.8822680
$N_x = 40$	0.9932568	0.9975166	0.9161125	0.9081806
$N_x = 80$	0.9986574	0.9993274	0.9509870	0.9511483
$N_x = 160$	0.9991550	0.9990234	0.9747287	0.9762977
$N_x = 320$	0.9985649	0.9987451	0.9861375	0.9860783
$N_x = 640$	0.9991842	0.9991606	0.9920950	0.9922520

Table 3: R^2 coefficient for predicted Zernike coefficients of the one dimensional Poisson problem at different training sample size and number of grid points.

far away from the training range. Note that the predicted values for a_0 and a_1 exhibit a systematic error towards the lower and upper bounds of the **Sigmoid** activation function, indicative of the influence of the **Sigmoid** activation function used in the last layer. This emphasizes the significance of properly designing activation functions at the bottleneck.

Using the R^2 statistical coefficient as a measure of accuracy for the trained encoder, we assess effects of sample size and grid points on the performance of BiPDEs and report the results in table 3.

1. *Effect of sample size:* First, we fix number of grid points and vary sample size. We find that increasing sample size improves accuracy of the predictions in the case of clean data, however in the case of noisy data the accuracy does not show significant improvement by enlarging sample size.
2. *Effect of grid points:* Second, we fix the sample size and gradually increase number of grid points. We find that accuracy of predictions on noisy data is strongly correlated with number of grid points, however this dependence is weaker for clean data.

Two dimensional inverse transform

We consider an example of variable diffusion coefficients parameterized as $D(x, y) = 4 + 2a_2y + \sqrt{3}a_3(2x^2 + 2y^2 - 1)$, with unknown coefficients randomly chosen in range $a_2, a_3 \in [0.25, 0.75]$. The equations are solved on a square domain $\Omega \in [-\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}]^2$ governed by the Poisson equation:

$$\begin{aligned} \nabla \cdot ([4 + 2a_2y + \sqrt{3}a_3(2x^2 + 2y^2 - 1)]\nabla u) + x + y &= 0, & (x, y) \in \Omega, \\ u_{BC} &= \cos(\pi x) \cos(\pi y), & (x, y) \in \partial\Omega. \end{aligned}$$

The encoder is composed of two convolutional layers with 32 and 64 channels followed by a 2×2 average pooling layer and a dense layer with 128 neurons, at the bottleneck there are 2 neurons representing the two unknowns. All activation functions are **ReLU** except at the bottleneck that has **Sigmoid** functions. An **Adam** optimizer is used on a mean squared error loss function.

We trained BiPDE over 900 generated solutions with randomly chosen parameters a_2, a_3 and validated its performance on 100 independent solution fields for 300 epochs, evolution of loss function is shown in figure 7(d). Then we tested the trained model over another set of 1000 images with randomly generated diffusion maps independent of the training dataset. Furthermore, we repeated this exercise over 1000 images with added zero-mean Gaussian noise with standard deviation 0.025. In each case, the learned parameters are in good agreement with the true values, as illustrated in

2D inverse transform	Noiseless		Noisy (13% relative noise)	
	Sample Size	a_2	a_3	a_2
$N_{\text{data}} = 250$	0.9897018	0.9958963	0.9872887	0.9892064
$N_{\text{data}} = 500$	0.9917211	0.9977917	0.9910183	0.9900091
$N_{\text{data}} = 1000$	0.9915683	0.9986852	0.9896654	0.9915149
$N_{\text{data}} = 2000$	0.9940470	0.9993891	0.9909640	0.9883151
$N_{\text{data}} = 4000$	0.9938268	0.9997119	0.9919061	0.9898697

Table 4: R^2 coefficient for predicted Zernike coefficients of the two dimensional Poisson problem by increasing training sample size. Number of grid points are fixed at 30×30 .

figure 9. Moreover, we performed a sensitivity analysis on the accuracy of the predictions with respect to sample size. We measured quality of fit by the R^2 statistical coefficient. Results are tabulated in table 4 and indicate training over more samples leads to more accurate predictions when applied to clean data, while noisy data do not show a strong improvement by increasing sample size.

3.2. Dynamic Burger’s problem

In this section, we demonstrate the applicability of BiPDEs on time-dependent nonlinear partial differential equations, and we use those results to illustrate the consistency and accuracy of the proposed framework. Similar to previous works [69], we consider the nonlinear Burgers’ equation in one spatial dimension,

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} = \nu \frac{\partial^2 u}{\partial x^2} \quad x \in [-1, 1], \quad t \in [0, 1) \quad (8)$$

$$u(-1, t) = u(1, t) = 0 \quad u(x, 0) = -\sin(\pi x) \quad (9)$$

where $\nu = 1/Re$ with Re being the Reynolds number. Notably, Burgers’ equation is of great practical significance for understanding evolution equations as it is nonlinear. Burgers’ equation has been used as a model equation for the Navier-Stokes equation and by itself can be used to describe shallow water waves [18], turbulence [9], traffic flow [54], and many more.

- **Discretization.** In our design we adopted the 6th-order compact finite difference scheme proposed by Sari and Gurarslan (2009) [71] for its simplicity of implementation, its high accuracy and because it leads to a linear system with narrow band and subsequently ensures computational efficiency. This scheme combines a tridiagonal⁴ sixth-order compact finite difference scheme (CFD6) in space with a low-storage third-order accurate total variation diminishing Runge-Kutta scheme (TVD-RK3) for its time evolution ([74]). In particular, the high-order accuracy associated with this discretization provides highly accurate results on coarse grids. This is an important aspect of BiPDEs as a data-efficient inverse solver, which stems from their capacity to seamlessly blend highly accurate and sophisticated discretization methods with deep neural networks.

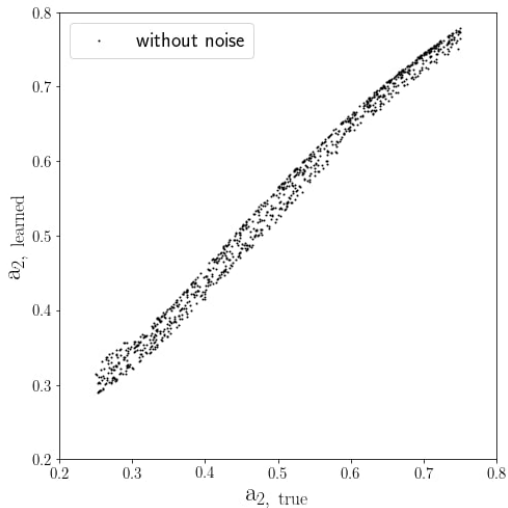
The first-order spatial derivatives are given at intermediate points by

$$\alpha u'_{i-1} + u'_i + \alpha u'_{i+1} = b \frac{u_{i+2} - u_{i-2}}{4h} + a \frac{u_{i+1} - u_{i-1}}{2h}, \quad i = 3, \dots, N-2 \quad (10)$$

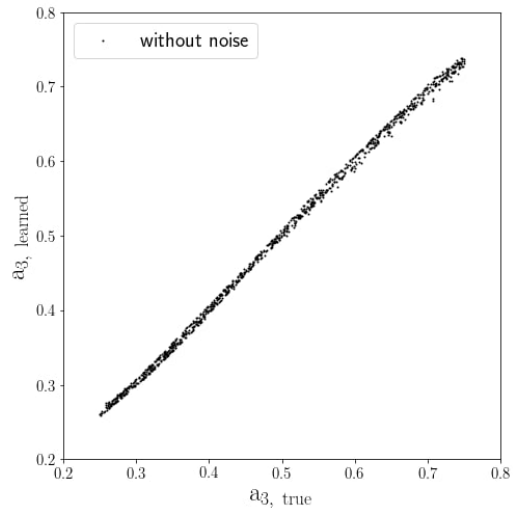
where

$$a = \frac{2}{3}(\alpha + 2), \quad b = \frac{1}{3}(4\alpha - 1), \quad (11)$$

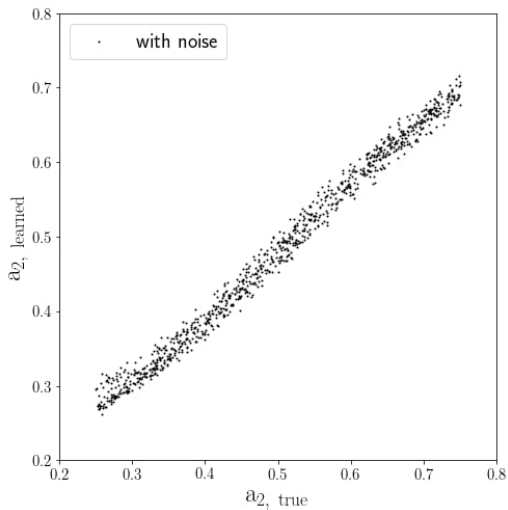
⁴Tridiagonal systems of equations can be obtained in $\mathcal{O}(N)$ operations.



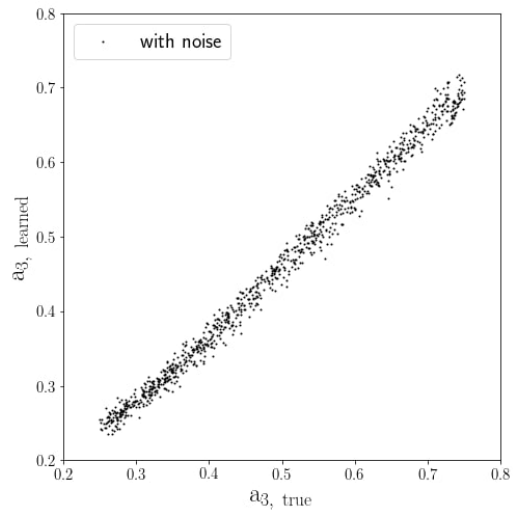
(a) Regression quality is $R^2 = 0.9915683$.



(b) Regression quality is $R^2 = 0.9986852$.



(c) Regression quality is $R^2 = 0.9896654$.



(d) Regression quality is $R^2 = 0.9915149$.

Figure 9: Top row shows performance of BiPDE over 1000 randomly chosen clean 2D images after 1000 epochs, and the bottom panel shows performance of the same network on noisy images given a zero-mean Gaussian noise with standard deviation 0.025. Network has 1,852,000 trainable parameters and training takes ~ 11 seconds on a Tesla T4 GPU available on a free Google Colaboratory account.

and $h = x_{i+1} - x_i$ is the mesh size, with grid points identified by the index $i = 1, 2, \dots, N$. For $\alpha = 1/3$ we obtain a sixth order accurate tridiagonal scheme. The boundary points (for non-periodic boundaries) are treated by using the formulas [22, 71],

$$\begin{aligned} u'_1 + 5u'_2 &= \frac{1}{h} \left[-\frac{197}{60}u_1 - \frac{5}{12}u_2 + 5u_3 - \frac{5}{3}u_4 + \frac{5}{12}u_5 - \frac{1}{20}u_6 \right] \\ \frac{2}{11}u'_1 + u'_2 + \frac{2}{11}u'_3 &= \frac{1}{h} \left[-\frac{20}{33}u_1 - \frac{35}{132}u_2 + \frac{34}{33}u_3 - \frac{7}{33}u_4 + \frac{2}{33}u_5 - \frac{1}{132}u_6 \right] \\ \frac{2}{11}u'_{N-2} + u'_{N-1} + \frac{2}{11}u'_N &= \frac{1}{h} \left[\frac{20}{33}u_N + \frac{35}{132}u_{N-1} - \frac{34}{33}u_{N-2} + \frac{7}{33}u_{N-3} - \frac{2}{33}u_{N-4} + \frac{1}{132}u_{N-5} \right] \\ 5u'_{N-1} + u'_N &= \frac{1}{h} \left[\frac{197}{60}u_N + \frac{5}{12}u_{N-1} - 5u_{N-2} + \frac{5}{3}u_{N-3} - \frac{5}{12}u_{N-4} + \frac{1}{20}u_{N-5} \right] \end{aligned}$$

This can be easily cast in the matrix form

$$BU' = AU \quad (12)$$

where $U = [u_1, u_2, \dots, u_N]^T$ is the vector of solution values at grid points. Furthermore, second order derivatives are computed by applying the first-order derivatives twice⁵,

$$BU'' = AU'' \quad (13)$$

Burgers' equation are thus discretized as:

$$\frac{dU}{dt} = \mathcal{L}U, \quad \mathcal{L}U = \nu U'' - U \otimes U', \quad (14)$$

where \otimes is the element-wise multiplication operator and \mathcal{L} is a *nonlinear* operator. We use a low storage TVD-RK3 method to update the solution field from time-step k to $k+1$,

$$U^{(1)} = U_k + \Delta t \mathcal{L}U_k \quad (15)$$

$$U^{(2)} = \frac{3}{4}U_k + \frac{1}{4}U^{(1)} + \frac{1}{4}\Delta t \mathcal{L}U^{(1)} \quad (16)$$

$$U_{k+1} = \frac{1}{3}U_k + \frac{2}{3}U^{(2)} + \frac{2}{3}\Delta t \mathcal{L}U^{(2)} \quad (17)$$

with a CFL coefficient of 1. Note that this method only requires two storage units per grid point, which is useful for large scale scientific simulations in higher dimensions.

- **Training protocol.** For training, we first solve Burgers' equation for M time-steps, then we construct two shifted solution matrices that are separated by a single time-step, *i.e.*,

$$\mathcal{U}^{-1} = [U^1 \quad U^2 \quad \dots \quad U^{M-1}] \quad \mathcal{U}^{+1} = [U^2 \quad U^3 \quad \dots \quad U^M] \quad (18)$$

Basically, one step of TVD-RK3 maps a column of \mathcal{U}^{-1} to its corresponding column in \mathcal{U}^{+1} given an accurate prediction for the hidden parameter. Hence, a semantic BiPDE is trained with \mathcal{U}^{-1} and \mathcal{U}^{+1} as its input and output respectively. The unknown diffusion coefficient is discovered by the code at the bottleneck of the architecture.

- **Numerical setup.** To enable direct comparison with PINNs, we also consider a second parameter γ in Burger's equation. In these current experiments we train for 2 unknown parameters (ν, γ) in the Burger's equation given by

$$u_t + \gamma uu_x - \nu u_{xx} = 0, \quad t \in [0, 1], \quad x \in [-1, 1]$$

⁵From implementation point of view this is a very useful feature of this scheme, because A and B are constant matrices that do not change through training it is possible to pre-compute them using `numpy`'s [56] basic data structures, and then simply import the derivative operators into `TensorFlow`'s custom solver layer using `tf.convert_to_tensor()` command.

# epochs = 200	$\Delta t = 0.001$		$\Delta t = 0.0005$		$\Delta t = 0.00025$	
grid size	ν	γ	ν	γ	ν	γ
$N_x = 20$	0.0028751	0.9500087	0.0028731	0.9500685	0.0028828	0.9499334
$N_x = 40$	0.0030294	0.9750050	0.0030341	0.9750042	0.0030391	0.9750047
$N_x = 80$	0.0031067	0.9875077	0.0031101	0.9875285	0.0031167	0.9875455
$N_x = 160$	0.0031455	0.9937580	0.0031443	0.9937674	0.0031519	0.9937985
$N_x = 320$	0.0031659	0.9968843	0.0031679	0.9968919	0.0031738	0.9969027
$N_x = 640$	0.0031775	0.9984500	0.0031797	0.9984597	0.0031841	0.9984711
$N_x = 700$	0.0031773	0.9985866	0.0031779	0.9985945	0.0031865	0.9986123

Table 5: Discovering two unknown values of ν and γ in Burger’s equation. These values are plotted in figure 10.

Similar to Raissi *et al.* [69] we consider $\nu = 0.01/\pi$ and $\gamma = 1.0$. For completeness we also recall the loss function used in PINN that encodes Burger’s equation as a regularization,

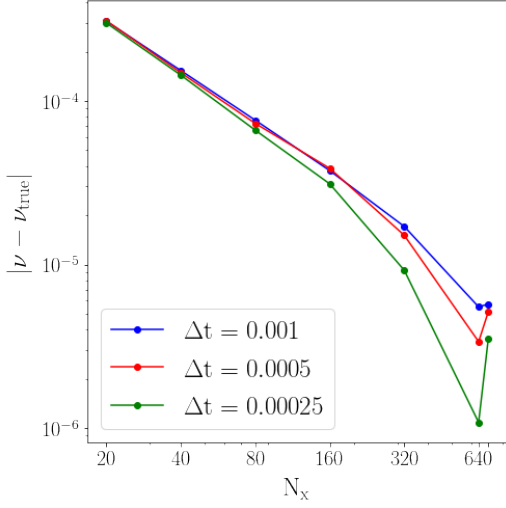
$$MSE = \frac{1}{N} \sum_{i=1}^N \left| u(t_u^i, x_u^i) - u^i \right|^2 + \frac{1}{N} \sum_{i=1}^N \left| u_t(t_u^i, x_u^i) + \gamma u(t_u^i, x_u^i) u_x(t_u^i, x_u^i) - \nu u_{xx}(t_u^i, x_u^i) \right|^2$$

where (t_u^i, x_u^i, u^i) constitute training data with $N = 2000$ observation points in the spatio-temporal domain. In this experiment PINN is composed of 9 layers with 20 neurons per hidden layer. It is worth mentioning that we are reporting BiPDE results by considering solutions in a narrower time span $t \in [0, 0.2]$.

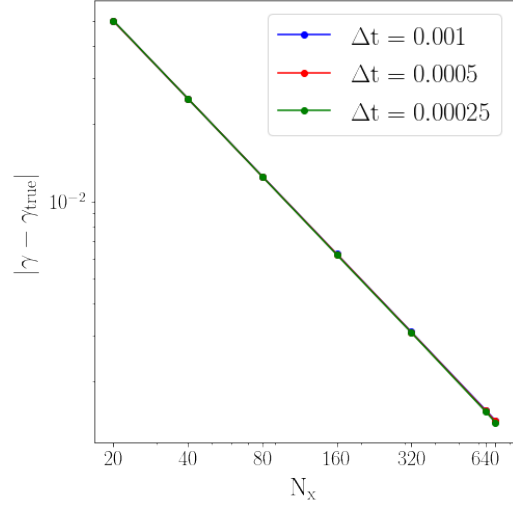
- Architecture.** Obviously, one can choose a single neuron to represent an unknown parameter ν or γ and in a few iterations an approximate value can be achieved. However, our goal is to train a general purpose encoder that is capable of predicting the unknown value from an input solution pair with arbitrary values of ν and γ and without training on new observations (*cf.* see part 4.1). Therefore, we consider a BiPDE that is composed of a `conv1D` layer with 128 filters and a kernel size of 10 with `tanh` activation function, followed by `AveragePooling1D` with a pool size of 2 that after flattening is fed to two `Dense` layers with 20 and 2 neurons respectively that apply `Sigmoid` activation function. We used the `Adam` optimizer to minimize the mean absolute error measure for the loss function.
- Accuracy test.** First, we train for two unknowns in Burger’s equation, namely ν and γ . We perform a sensitivity analysis for 200 epochs with different numbers of spatial grid points, as well as different time-steps. In each case, we measure the error between the learned values of ν and γ with their true value $\nu_{\text{true}} = 0.01/\pi$ and $\gamma_{\text{true}} = 1.0$. Convergence results of this experiment are tabulated in table 5 and shown in figure 10. We find that increasing the number of grid points (*i.e.* the resolution) improves the accuracy up to almost 700 grid points before accuracy in ν (but not γ) starts to deteriorate. Note the decrease in time-step size does not have a significant effect on accuracy unless large number of grid points $N_x > 160$ are considered where decreasing time-step clearly improves results.

For comparison purposes we report numerical results from table 1 of Raissi *et al.* (2017) [69] in our figures 10(c)–10(d). Here we only presented noise-less results of BiPDE, therefore only the 0% added noise case of PINN is comparable, *i.e.* the solid red line in figures 10(c)–10(d). Even though the two test cases have significant differences and much care should be taken to directly compare the two methods, however BiPDEs have a clear advantage by exhibiting *convergence* in the unknown values, *i.e.* more data means better results.

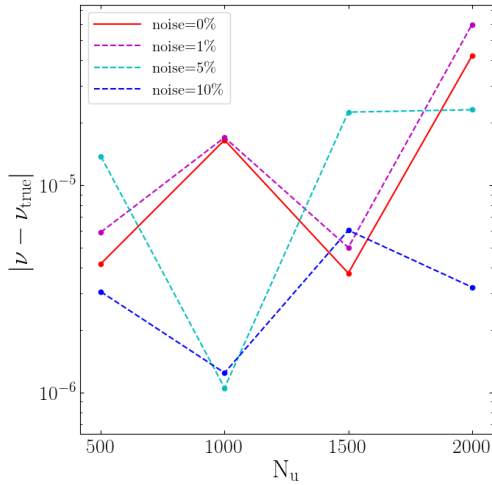
In a second experiment, we fix the value of $\gamma = 1.0$ and only train for the unknown diffusion coefficient ν . Similar to previous test we trained the network for 200 epochs and figure 11 shows the error in the discovered value of ν at different levels of resolution. In this case



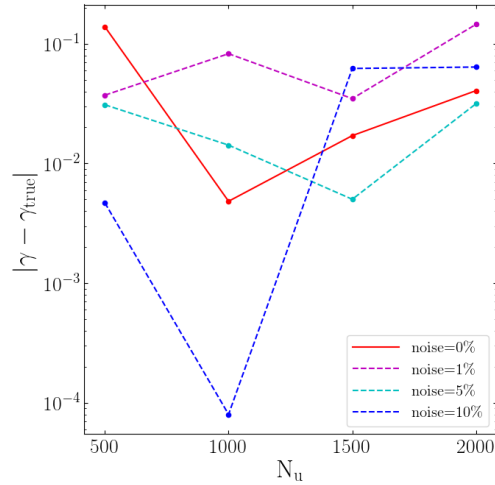
(a) Error in ν - BiPDE with finite difference method.



(b) Error in γ - BiPDE with finite difference method.



(c) Error in ν - PINN.



(d) Error in γ - PINN.

Figure 10: Sensitivity analysis in training both parameters γ and ν with BiPDE (a,b), also results from table 1 of Raissi *et al.* (2017) [69] are shown for comparison (c,d) - note only the solid red line may be compared to BiPDE results where no noise is considered on the solution data. True values are $\nu_{\text{true}} = 0.01/\pi$ and $\gamma = 1.0$. In figure (a) the data points at the right end of N_x axis correspond to $N_x = 700$ grid points where the accuracy in the discovered ν value deteriorates.

decreasing time-step size does not seem to have a significant effect on accuracy. A curious observation is the absolute value of error for ν is two orders of magnitude more precise when the network is trained for both parameters ν and γ than when only tuning for ν . Again, convergence in unknown parameter is retained in this experiment.

4. Mesh-less BiPDE: Multi-Quadratic Radial Basis Functions

Not only are direct computations of partial derivatives from noisy data extremely challenging, in many real world applications, measurements can only be made on scattered point clouds. Tikhonov regularization type approaches have been devised to avoid difficulties arising from high sensitivity of

$\langle \hat{\nu} \rangle$	$\Delta t = 0.001$	$\Delta t = 0.0005$	$\Delta t = 0.00025$
$N_x = 20$	0.0064739	0.0065189	0.0065514
$N_x = 40$	0.0048452	0.0048200	0.0047086
$N_x = 80$	0.0040260	0.0040324	0.0039963
$N_x = 160$	0.0036042	0.0036011	0.0036310
$N_x = 320$	0.0033958	0.0034144	0.0033827
$N_x = 640$	0.0032919	0.0032895	0.0032916
$N_x = 700$	0.0032829	0.0032816	0.0032906

Table 6: Discovering one unknown parameter in Burger’s equation, values for ν .

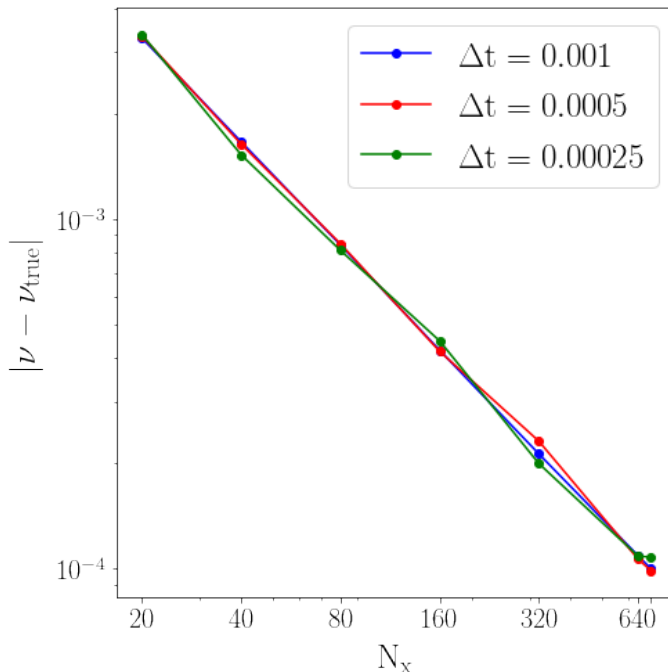


Figure 11: Sensitivity analysis in training only one parameter ν . True value of $\nu_{\text{true}} = 0.01/\pi$ is sought in Burgers’ equation at different levels of resolution. Rightmost data points correspond to $N_x = 700$ grid points.

differencing operations on noisy data [14, 11, 78]; for neural network based approaches, see [49, 73]. Recently, Trask *et al.* [81] have proposed an efficient framework for learning from unstructured data that is based on the Generalized Moving Least Squares (GMLS) technique. They show performance of GMLS-Nets to identify differential operators and to regress quantities of interest from unstructured data fields. Another interesting approach had been put forth in the late 80s by [8, 63] that designed neural networks based on Radial Basis Functions (RBF) to perform functional interpolation tasks. In these networks, the activation function is defined as the radial basis function of interest and the training aims to discover the weights of this network, which interestingly coincide with the coefficients in the corresponding radial basis expansion.

Since the early 70s, RBFs have been used for highly accurate interpolation from scattered data. In particular, Hardy [28] introduced a special kind of RBF called the *multiquadratic* series expansion, that provides superior performance in terms of accuracy, stability, efficiency, simplicity and memory usage [21]. Kansa (1990) [37, 38] pioneered the use of radial basis functions to solve time dependent partial differential equations through deriving a modified multi-quadratic scheme. In 1998, Hon and Mao [34] applied multiquadratics as a spatial discretization method for the nonlinear Burgers’ equation and solved it for a wide range of Reynolds numbers (from 0.1 to 10,000). Their scheme was

later enhanced to second-order accuracy in time by Xie and Li (2013) [87] via introducing a compact second-order accurate time discretization. Interestingly, the accuracy of these mesh-free methods can be improved by fine-tuning distributions of collocation points or their *shape parameters*. For example, Hon and Mao devised an adaptive point to chase the peak of shock waves, which improved their results. Fortunately, such fine-tuning of parameters can be automated using BiPDE networks; we demonstrate this in this section.

- **Discretization.** We chose to blend the second-order accurate method of Xie and Li, briefly described next and we leave further details to their original paper [87]. Initially, one can represent a distribution $u(\mathbf{x})$ in terms of a linear combination of radial basis functions,

$$u(\mathbf{x}) \approx \sum_{j=0}^{N_s} \lambda_j \phi_j(\mathbf{x}) + \psi(\mathbf{x}), \quad \mathbf{x} \in \Omega \subset \mathcal{R}^{dim}, \quad (19)$$

where $\phi(\mathbf{x})$ is the radial basis function that we adopt,

$$\phi_j(\mathbf{x}) = \sqrt{r_j^2 + c_j^2}, \quad r_j^2 = \|\mathbf{x} - \mathbf{x}_j\|_2^2, \quad (20)$$

and c_j is the *shape parameter* that has been experimentally shown to follow $c_j = Mj + b$ with $j = 0, 1, \dots, N_s$ (N_s is number of seed points). Moreover M and b are tuning parameters. In equation (19), $\psi(\mathbf{x})$ is a polynomial to ensure solvability of the resulting system when ϕ_j is only conditionally positive definite. To solve PDEs, one only needs to represent the solution field with an appropriate form of equation (19). In the case of Burgers' equation the solution at any time-step n can be represented by

$$u^n(x) \approx \sum_{j=0}^{N_s} \lambda_j^n \phi_j(x) + \lambda_{N_s+1}^n x + \lambda_{N_s+2}^n \quad (21)$$

over a set of reference points for the basis functions that are given by $x_j = j/N_s$, $j = 0, 1, \dots, N_s$. Xie and Li derived the following compact second-order accurate system of equations

$$\left[1 + \frac{\Delta t}{2} u_x^n(\hat{x}_j)\right] u^{n+1}(\hat{x}_j) + \frac{\Delta t}{2} u^n(\hat{x}_j) u_x^{n+1}(\hat{x}_j) - \frac{\nu \Delta t}{2} u_{xx}^{n+1}(\hat{x}_j) = u^n(\hat{x}_j) + \frac{\nu \Delta t}{2} u_{xx}^n(\hat{x}_j) \quad (22)$$

over a set of $N_d + 1$ distinct collocation points $\hat{x}_j = (1 + j)/(N_d + 2)$ with $j = 0, 1, \dots, N_d$. Two more equations are obtained by considering the left and right boundary conditions $u^{n+1}(x_L) = u^{n+1}(x_R) = 0$. Note that spatial derivatives are directly computed by applying derivative operator over equation (21). At every time-step, one solves for the $N + 3$ coefficients $\lambda_0^n, \dots, \lambda_{N+2}^n$, while the spatial components of the equations remain intact (as long as points are not moving). The solution is obtained over the initial conditions given by $u^0(\hat{x}_j)$.

For implementation purposes, we represent the system of equations (22) in a matrix notation that is suitable for tensorial operations in `TensorFlow`. To this end, we first write equation (21) as

$$U^n(\hat{x}) = A(\hat{x})\Lambda^n, \quad (23)$$

where

$$U_{(N_d+1) \times 1}^n = \begin{bmatrix} u^n(\hat{x}_0) \\ u^n(\hat{x}_1) \\ \vdots \\ u^n(\hat{x}_{N_d}) \end{bmatrix}, \quad \Lambda_{(N_s+1) \times 1}^n = \begin{bmatrix} \lambda_0^n \\ \lambda_1^n \\ \vdots \\ \lambda_{N_s}^n \end{bmatrix}, \quad (24)$$

$$\left[A_{ij}(\hat{\mathbf{x}}) \right]_{(N_d+1) \times (N_s+1)} = \left[\phi_j(\hat{x}_i) - \phi_j(x_L) - \frac{\phi_j(x_R) - \phi_j(x_L)}{x_R - x_L} (\hat{x}_i - x_L) \right], \quad (25)$$

with $i = 0, 1, \dots, N_d$ and $j = 0, 1, \dots, N_s$. Note that we already injected the homogeneous boundary conditions into equation (23). Therefore, equation (22) can be written as,

$$\left[A + (\mathbf{g}_x \mathbf{1}^T) \otimes A + (\mathbf{g} \mathbf{1}^T) \otimes A_x - \frac{\nu \Delta t}{2} A_{xx} \right] \Lambda^{n+1} = \left[A + \frac{\nu \Delta t}{2} A_{xx} \right] \Lambda^n, \quad (26)$$

where $\mathbf{1}^T = [1, 1, \dots, 1]_{1 \times (N_s+1)}$, \otimes is component-wise multiplication, and

$$\mathbf{g} = \frac{\Delta t}{2} A \Lambda^n, \quad \mathbf{g}_x = \frac{\Delta t}{2} A_x \Lambda^n, \quad (27)$$

$$(A_x)_{ij} = \phi'_j(\hat{x}_i) - \frac{\phi_j(x_R) - \phi_j(x_L)}{x_R - x_L}, \quad (A_{xx})_{ij} = \phi''_j(\hat{x}_i). \quad (28)$$

Note that in case of training for two parameters (ν, γ) , expression for \mathbf{g} in equation 26 needs to be modified by letting $\mathbf{g} = \frac{\gamma \Delta t}{2} A \Lambda^n$.

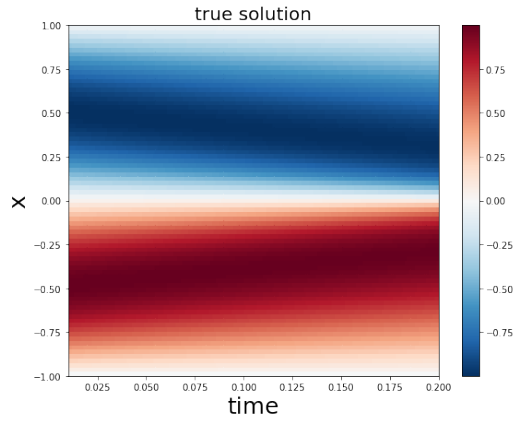
- **Architecture.** Note that both the collocation points and the interpolation seed points can be any random set of points within the domain and not necessarily a uniform set of points as we chose above. In fact, *during training we allow BiPDE to find a suitable set of interpolation points as well as the shape parameters on its own.* The input data is calculated using aforementioned finite difference method over uniform grids and later interpolated on a random point cloud to produce another sample of solutions on unstructured grids for training. Thus, in our architecture the last layer of the encoder has $2N_s + 1$ neurons with `sigmoid` activation functions representing the $2N_s$ shape parameters and seed points, as well as another neuron for the unknown diffusion coefficient. Note that for points to the left of origin, in the range $x \in [-1, 0]$, we simply multiplied the output of N_s activation functions by “ -1 ” within the solver layer (because output of `Sigmoid` function is always positive). We use the mean squared error between data and predicted solution at time-step $n+p$ as the loss function. We used the `Adam` optimizer to minimize the loss function.
- **Training protocol.** As in the previous case, we apply successive steps of MQ-RBF scheme to march the input data forward to a future time-step. Not surprisingly, we observed that taking higher number of steps improves the results because erroneous guess of the diffusion coefficient leads to more pronounced discrepancy after longer periods of time. Hence, we map the data \mathcal{U}^{-p} to \mathcal{U}^{+p} , which is p time-steps shifted in time,

$$\mathcal{U}^{-p} = [U^1, U^2, \dots, U^{M-p}] \quad \mathcal{U}^{+p} = [U^{1+p}, U^{2+p}, \dots, U^M] \quad (29)$$

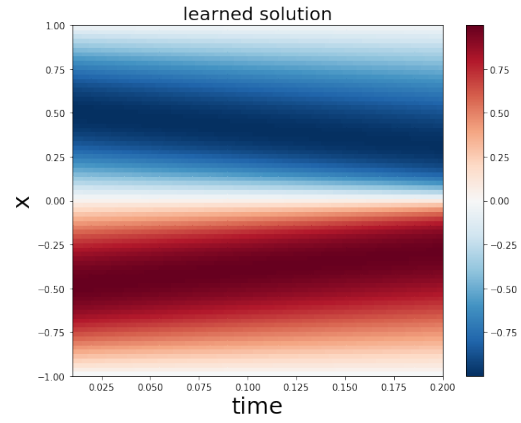
In our experiments a value of $p = 10$ was sufficient to get satisfactory results at the absence of noise. However, at the presence of Gaussian noise and for smaller values of the diffusion coefficient (such as for $\nu_{\text{true}} = 0.01/\pi$) we had to increase the shift parameter to $p = 100$.

- **Numerical setup.** Once again, we let $\nu_{\text{true}} = 0.01/\pi \approx 0.00318$ and integrate Burgers’ equation up to $t_f = 0.2$ with a fixed time-step of $\Delta t = 0.001$. We use the finite difference method of the previous section to generate the datasets. We then interpolate the solution on 80 data points, uniformly distributed in the range $(-1, 1)$ with 20 interpolation seed points. For this experiment, we set the batch size to 1. We trained the network using `Adam` optimizer. The results after 50 epochs are given in figure 12.

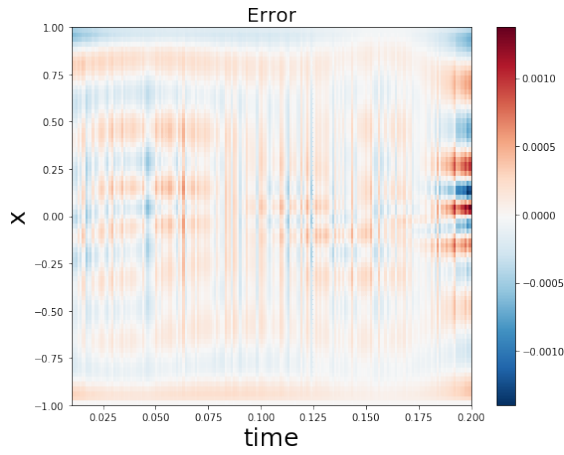
Interestingly, for every pair of input-output, the network discovers a distinct value for the diffusion coefficient that provides a measure of uncertainty for the unknown value. We report the average value of all diffusion coefficients as well as the probability density of these values. We observe that for all pairs of solutions, the predicted value for the diffusion coefficient is



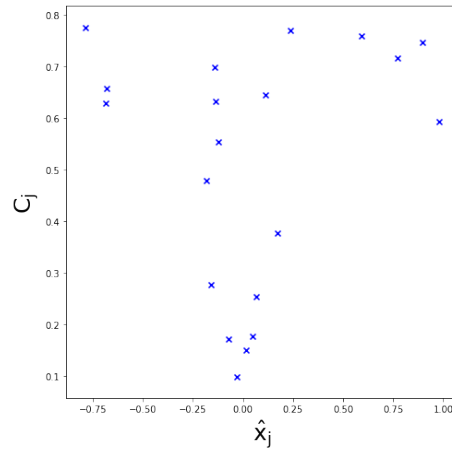
(a) True solution generated by finite differences (input data).



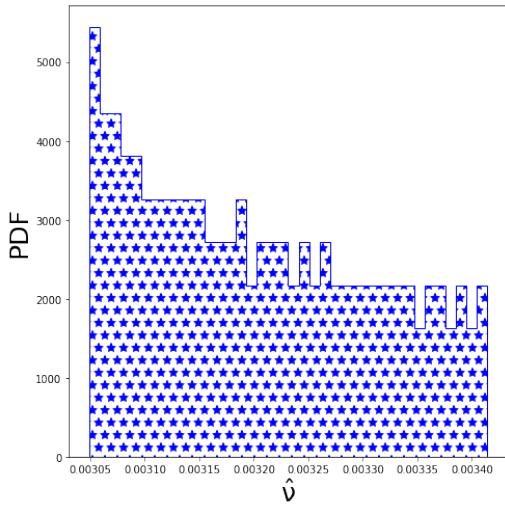
(b) Learned solution generated by MQ-RBF BiPDE (output data).



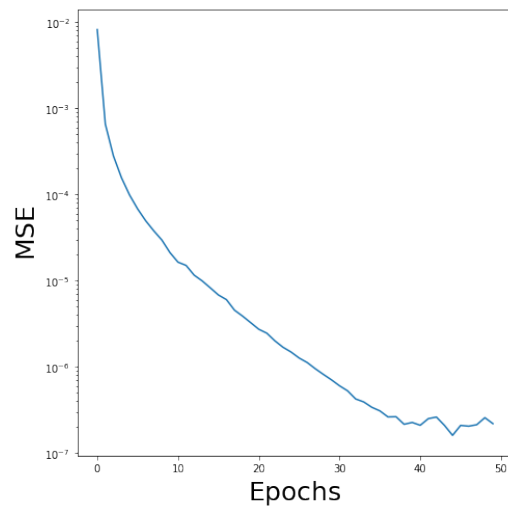
(c) Error in solution.



(d) Discovered seeds and shape parameters.



(e) Distribution of diffusion coefficients.



(f) Evolution of mean squared error during training.

Figure 12: Results of applying the RBF-BiPDE to Burgers' equation with a true diffusion coefficient of $\nu_{\text{true}} = 0.003183$. The average value of the predicted diffusion coefficients is $\hat{\nu} = 0.00320$.

distributed in the range $0.00305 \leq \hat{\nu} \leq 0.00340$ with an average value of $\langle \hat{\nu} \rangle = 0.00320$, which is in great agreement with the true value, indeed with 0.6% relative error. Interestingly, we observe that the BiPDE network has learned to concentrate its interpolation seed points around the origin where the solution field varies more rapidly. Furthermore, around $x = \pm 0.5$, the interpolation points are more sparse, which is in agreement with the smooth behavior of the solution field at these coordinates. Therefore, this network may be used as an automatic shock tracing method to improve numerical solutions of hyperbolic problems with shocks and discontinuities as was shown by Hon and Mao.

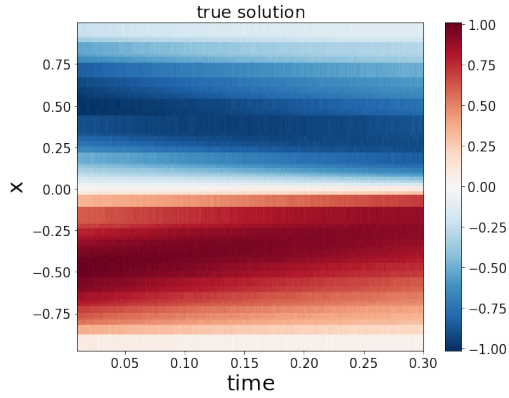
- **Resilience to noise on unstructured grids.** We consider several cases to assess robustness to noise. In each case, we pick 80 *randomly* distributed points along the x -axis and linearly interpolate the solution field on this set of points. Then, we add a Gaussian noise with a given standard deviation. This noisy and unstructured data field is then fed into the MQ-RBF based BiPDE of this section. We use a batch size of 10, with 10% of each sample for validation during training. A summary of our results follows:

1. Let $\nu_{\text{true}} = 0.1/\pi$, $p = 10$, $N_d = 80$, $N_s = 20$, $\Delta t = 0.001$, and consider a Gaussian noise with a standard deviation of 1%. After 100 epochs, we obtain the results in figure 13.
2. Let $\nu_{\text{true}} = 0.1/\pi$, $p = 100$, $N_d = 200$, $N_s = 20$, $\Delta t = 0.001$, and consider a Gaussian noise with a standard deviation of 5%. After 150 epochs, we obtain the results in figure 14.
3. Let $\nu_{\text{true}} = 0.01/\pi$, $p = 100$, $N_d = 80$, $N_s = 20$, $\Delta t = 0.001$, and consider a Gaussian noise with a standard deviation of 1%. After 200 epochs, we obtain the results in figure 15.
4. Let $\nu_{\text{true}} = 0.01/\pi$, $p = 100$, $N_d = 200$, $N_s = 20$, $\Delta t = 0.001$, and consider a Gaussian noise with a standard deviation of 5%. After 150 epochs, we obtain the results in figure 16.

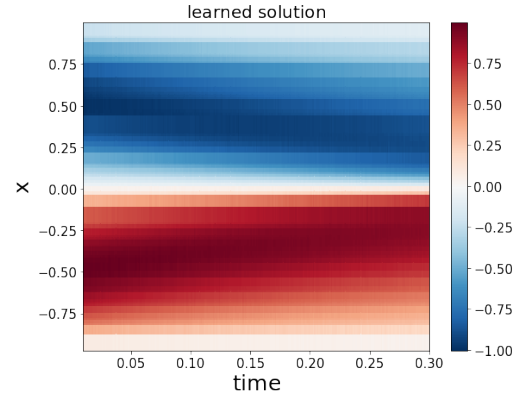
We observe that this architecture is generally robust to noise. However, at higher noise values require more tuning of hyperparameters, as well as longer training.

- **Accuracy tests.** We report the values of the discovered diffusion coefficients in the Burgers' equation for different grid sizes and different time-steps. We use the same setting as that detailed in the numerical setup part in this section. Particularly, the interpolation seeds are determined by the network and the training data is on a uniformly distributed set of points computed by the finite difference method of the previous section. We consider three different time-steps, $\Delta t = 0.001$, 0.0005 , 0.00025 , and two diffusion coefficients of $\nu_{\text{true}} = 0.01/\pi$, $0.1/\pi$ over grids of size $N_x = 80, 160$. At each time-step, for all experiments with different grid sizes, we stop the training when the mean squared error in the solution field converges to a constant value and does not improve by more epochs; this roughly corresponds to 50, 25, 12 epochs for each of the training time-steps, respectively. This indicates that by choosing smaller time steps less number of epochs are needed to obtain the same level of accuracy in the unknown parameter. Furthermore, we use an Adam optimizer with a learning rate of 0.001.

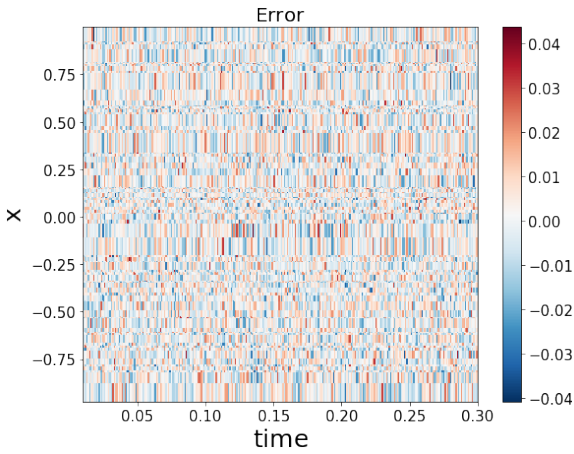
The results of the accuracy tests are tabulated in tables 7–8. We observe, for all experiments, that the discovered coefficient is in great agreement with the true values. Due to adaptivity of the interpolation seed points and their shape parameters for different experiments, the observed error values do not seem to follow the trend of traditional finite difference methods, as depicted in previous sections. This could also be due to lower order of accuracy of the MQ-RBF method, *i.e.* being a second-order accurate method, compared to the higher-order accurate finite difference method used in the previous section.



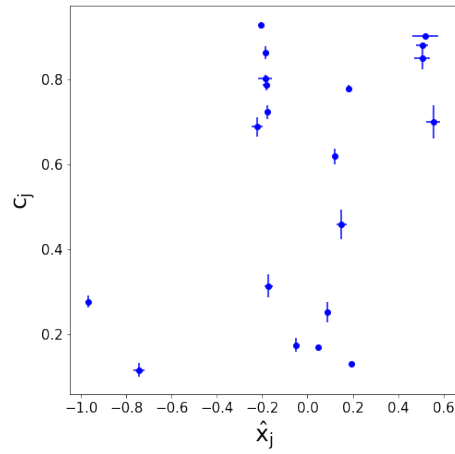
(a) True solution generated by finite differences and with added noise. Solution is interpolated on a random grid.



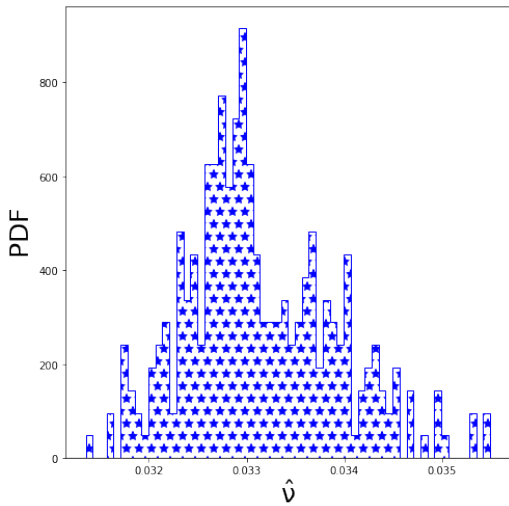
(b) Learned solution generated by MQ-RBF BiPDE (output data).



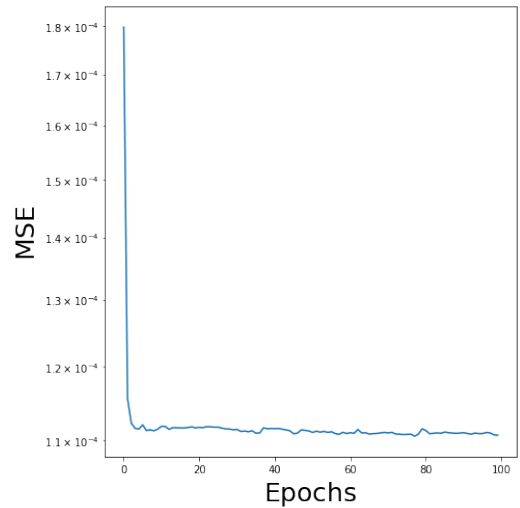
(c) Error in solution.



(d) Discovered seeds and shape parameters. Error bars indicate one standard deviation.

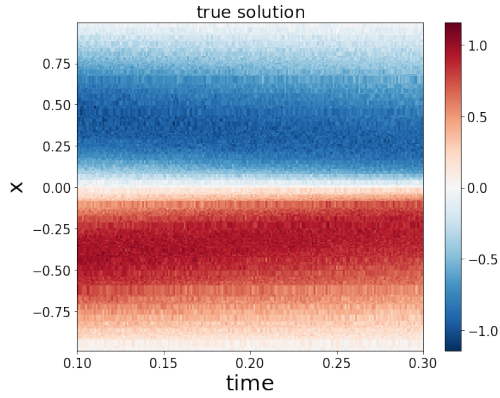


(e) Probability density of diffusion coefficients.

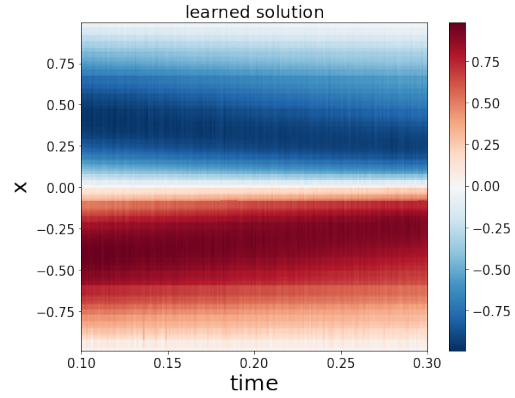


(f) Evolution of mean squared error versus number of epochs.

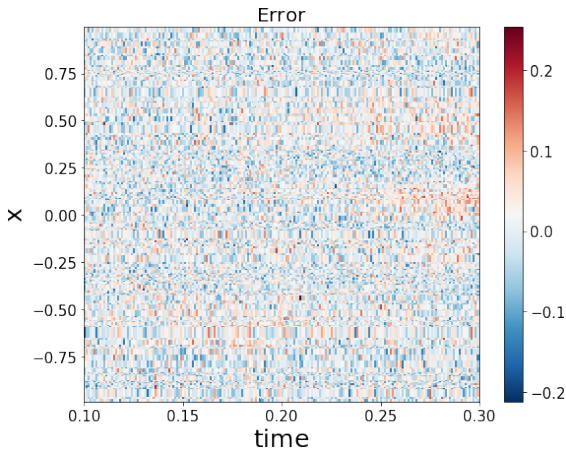
Figure 13: Results of applying the RBF-BiPDE to Burgers' equation with a true diffusion coefficient of $\nu_{\text{true}} = 0.03183$. The average value of the predicted diffusion coefficients is $\hat{\nu} = 0.0331$. The data is provided on a scattered point cloud with added Gaussian noise with 1% standard deviation.



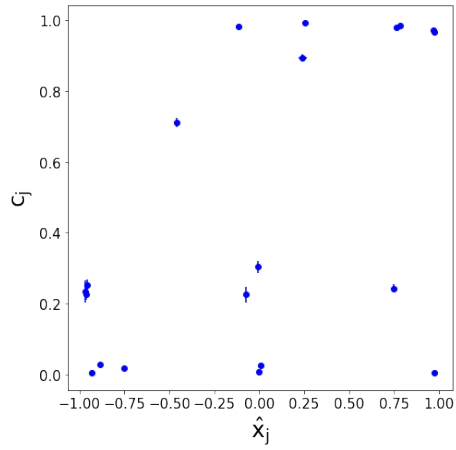
(a) True solution generated by finite differences and with added noise. Solution is interpolated on a random grid.



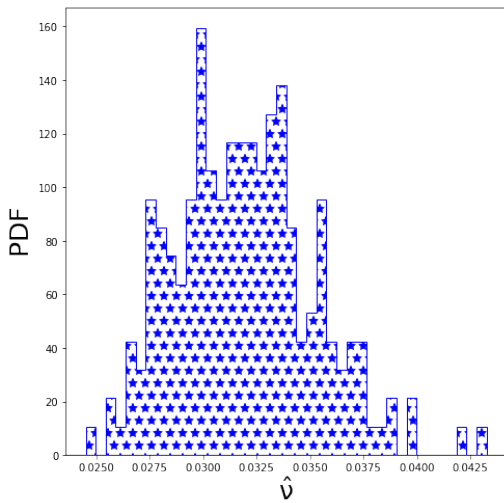
(b) Learned solution generated by MQ-RBF BiPDE (output data).



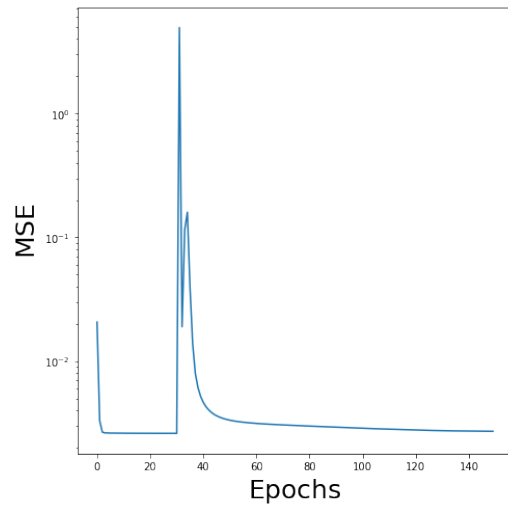
(c) Error in solution.



(d) Discovered seeds and shape parameters. Error bars indicate one standard deviation.

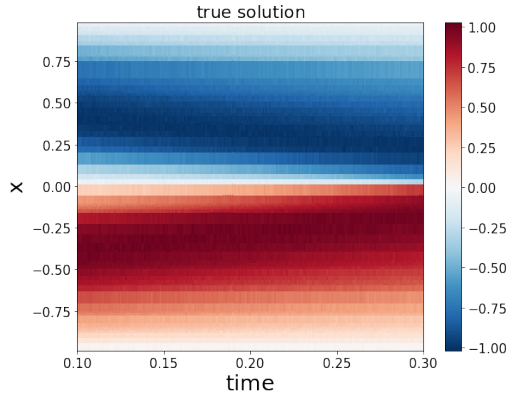


(e) Probability density of diffusion coefficients.

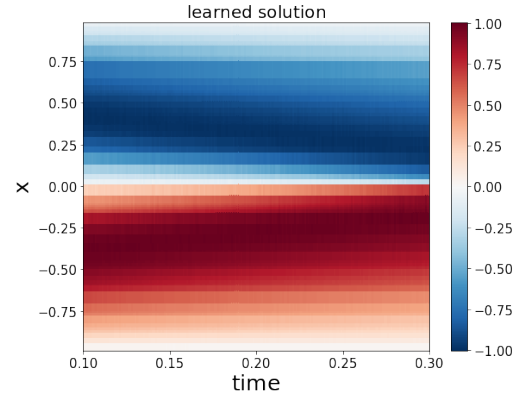


(f) Evolution of mean squared error versus number of epochs.

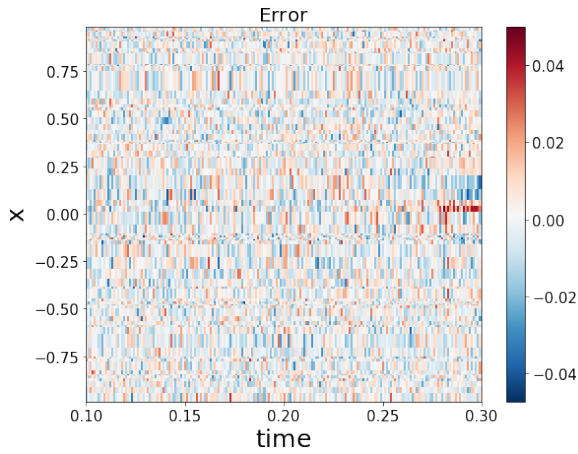
Figure 14: Results of applying the RBF-BiPDE to Burgers' equation with a true diffusion coefficient of $\nu_{\text{true}} = 0.03183$. The average value of the predicted diffusion coefficients is $\hat{\nu} = 0.03160$. The data is provided on a scattered point cloud with added Gaussian noise with 5% standard deviation.



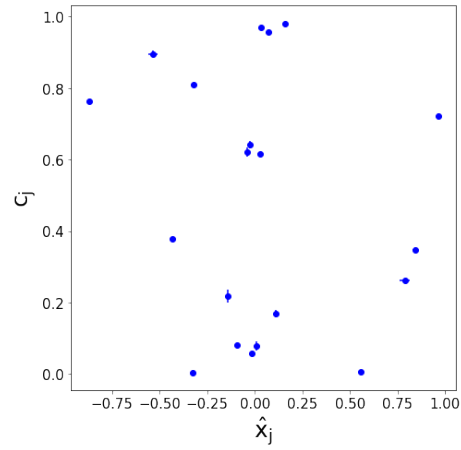
(a) True solution generated by finite differences and with added noise. Solution is interpolated on a random grid.



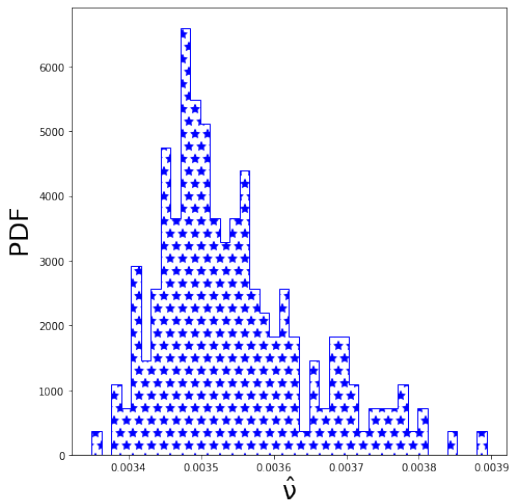
(b) Learned solution generated by MQ-RBF BiPDE (output data).



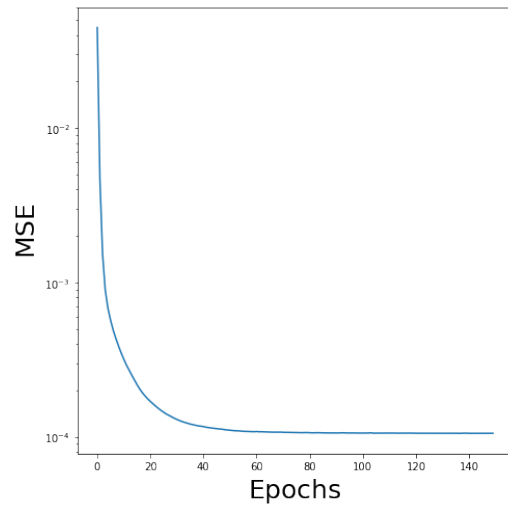
(c) Error in solution.



(d) Discovered seeds and shape parameters. Error bars indicate one standard deviation.

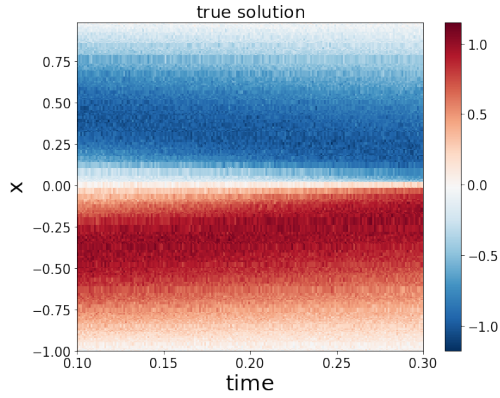


(e) Probability density of diffusion coefficients.

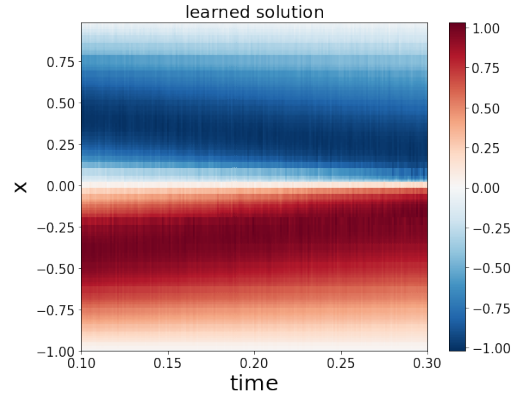


(f) Probability density of diffusion coefficients.

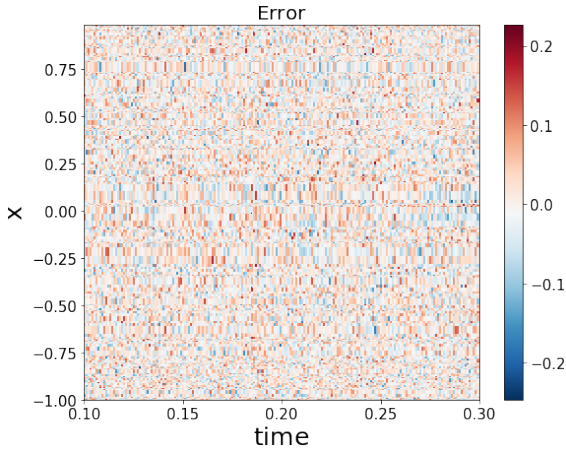
Figure 15: Results of applying the RBF-BiPDE to Burgers' equation with a true diffusion coefficient of $\nu_{\text{true}} = 0.003183$. The average value of the predicted diffusion coefficients is $\hat{\nu} = 0.00352$. The data is provided on a scattered point cloud with added Gaussian noise with 1% standard deviation.



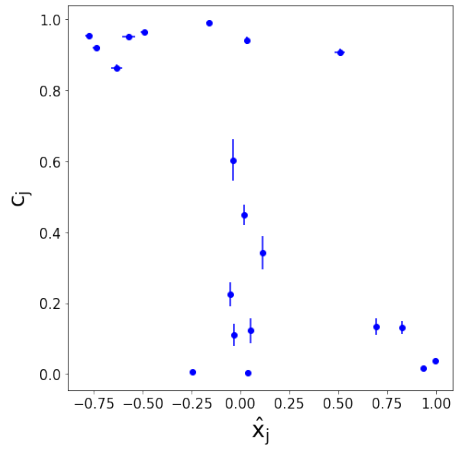
(a) True solution generated by finite differences and with added noise. Solution is interpolated on a random grid.



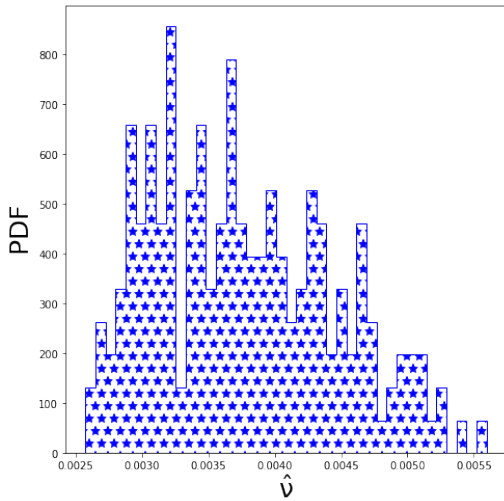
(b) Learned solution generated by MQ-RBF BiPDE (output data).



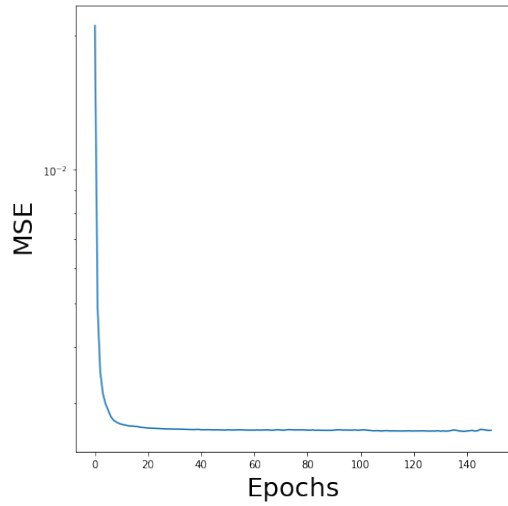
(c) Error in solution.



(d) Discovered seeds and shape parameters. Error bars indicate one standard deviation.



(e) Probability density of diffusion coefficients.



(f) Probability density of diffusion coefficients.

Figure 16: Results of applying the RBF-BiPDE to Burgers' equation with a true diffusion coefficient of $\nu_{\text{true}} = 0.003183$. The average value of the predicted diffusion coefficients is $\hat{\nu} = 0.003677$. The data is provided on a scattered point cloud with added Gaussian noise with 5% standard deviation.

$\langle \hat{\nu} \rangle$	$\Delta t = 0.001$	$\Delta t = 0.0005$	$\Delta t = 0.00025$
# epochs	50	25	12
$N_x = 80$	$0.03173 \pm 3.4 \times 10^{-4}$	$0.03196 \pm 4.2 \times 10^{-4}$	$0.03188 \pm 2.8 \times 10^{-4}$
$N_x = 160$	$0.03186 \pm 5.8 \times 10^{-5}$	$0.03191 \pm 3.6 \times 10^{-4}$	$0.03137 \pm 1.2 \times 10^{-4}$

Table 7: Discovered values of the diffusion coefficient for $\nu_{\text{true}} = 0.03183$ at different time-steps and grid sizes.

$\langle \hat{\nu} \rangle$	$\Delta t = 0.001$	$\Delta t = 0.0005$	$\Delta t = 0.00025$
# epochs	50	25	12
$N_x = 80$	$0.003326 \pm 5.1 \times 10^{-5}$	$0.003162 \pm 2.2 \times 10^{-4}$	$0.003155 \pm 1.2 \times 10^{-4}$
$N_x = 160$	$0.003264 \pm 1.0 \times 10^{-4}$	$0.003151 \pm 1.3 \times 10^{-4}$	$0.003192 \pm 1.2 \times 10^{-4}$

Table 8: Discovered values of the diffusion coefficient for $\nu_{\text{true}} = 0.003183$ obtained with different time-steps and grid sizes.

4.1. Learning the inverse transform

As we emphasized before, a feature of BiPDE is to produce self-supervised pre-trained encoder models for inverse differential problems that are applicable in numerous applications where hidden values should be estimated in real-time. We train an encoder over a range of values $\nu \in [0.1/\pi, 1/\pi]$ and assess the performance of the trained model on new data with arbitrarily chosen ν values. We choose 50 diffusion coefficients that are distributed uniformly in this range, then integrate the corresponding Burgers' equation up to $t_f = 0.2$ with a constant time-step of $\Delta t = 0.0005$ on a grid with $N_x = 100$ grid points using the aforementioned finite difference method. There are 4000 time-steps in each of the 50 different realizations of Burgers' equation. For a fixed value of $p = 20$, we draw 10 solution pairs for each value of ν at uniformly distributed time instances and discard the first two instances to improve convergence of the network. Hence, the training data uniformly samples the space of solutions over a 8×50 grid of (t, ν) , as illustrated in figure 17. We use the resulting 400 pairs in training of a semantic BiPDE, with 320 pairs used for training and 80 pairs for validation.

Architecture. Given an arbitrary input, the signature of the hidden physical parameters will be imprinted on the data in terms of complex patterns spread in space and time. We use a CNN layer as a front end unit to transform the input pixels to internal image representations. The CNN unit has 32 filters with kernel size of 5. The CNN is followed by max pooling with pool size of 2, which is then stacked with another CNN layer of 16 filters and kernel size of 5 along with another max pooling layer. The CNN block is stacked with two dense layers with 100 and 41 neurons, respectively. CNN and dense layers have **ReLU** and **Sigmoid** activation functions, respectively. Overall, there are 42,209 trainable parameters in the network. Conceptually, the CNN extracts features on every snapshot that characterizes the evolution of the solution field through time-steps with a proper physical parameter. This parameter is enforced to be the diffusion coefficient through the PDE solver decoder stage. We train this network for 500 epochs using an **Adam** optimizer.

Resilience to noise. Even though the encoder is trained on ideal datasets, we demonstrate a semantic BiPDE still provides accurate results on noisy datasets. In contrast to other methods, we pre-train the network in a self-supervised fashion on clean data and later we apply the trained encoder on unseen noisy data⁶. In figure 18, we provide the performance of this network on training as well as on unseen clean/noisy data-sets. Furthermore, the network determines optimal parameters of the MQ-RBF method by evaluating interpolation seed points as well as their corresponding shape parameters to obtain the best approximation over *all* input data.

⁶Note that the network could also be trained on noisy data as we showed before; however training would take longer in that case.

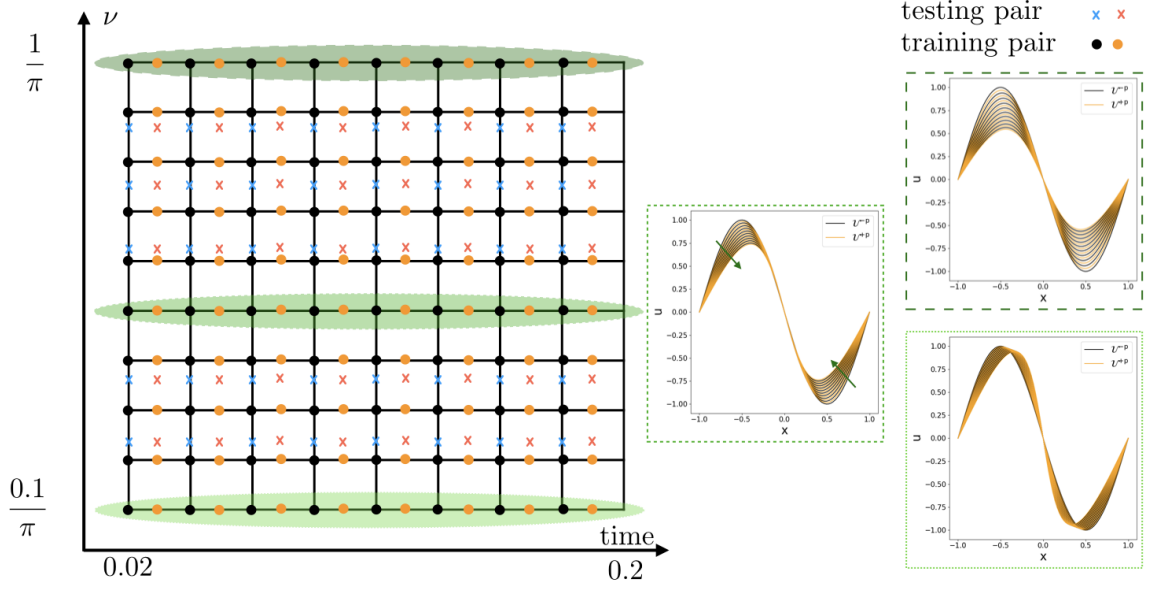


Figure 17: Topology of data points for training and testing of the semantic BiPDE. Along the ν dimension, we depict 10 (out of 50) of the selected data points, while along the time dimension we illustrate the actual 8 data points. Training pairs of \mathcal{U}^{-p} and \mathcal{U}^{+p} are color coded by black and orange dots, respectively; testing pairs are depicted by blue and red crosses. On the right panel, we illustrate the training data for three nominal values of the diffusion coefficient, highlighted by green shades. Green arrows indicate the direction of time.

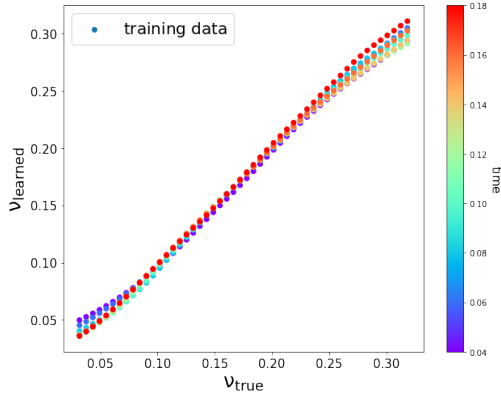
5. Conclusion

We introduced BiPDE networks, a natural architecture to infer hidden parameters in partial differential equations given a limited number of observations. We showed that this approach is versatile as it can be easily applied to arbitrary static or nonlinear time-dependent inverse-PDE problems. We showed the performance of this design on multiple inverse Poisson problems in one and two spatial dimensions as well as on the non-linear time-dependent Burgers' equation in one spatial dimension. Moreover, our results indicate BiPDEs are robust to noise and can be adapted for data collected on unstructured grids by resorting to traditional mesh-free numerical methods for solving partial differential equations. We also showed the applicability of this framework to the discovery of inverse transforms for different inverse-PDE problems.

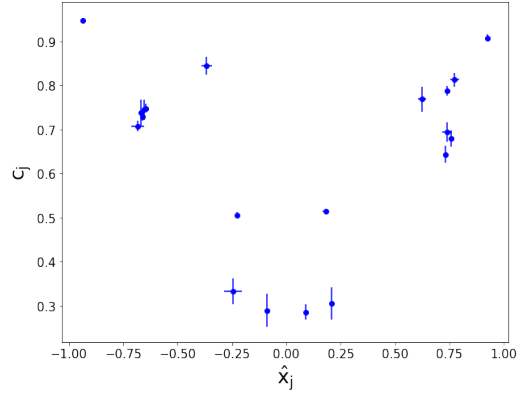
There are many areas of research that could be further investigated, such as considering diffusion maps with discontinuities across subdomains, using more sophisticated neural network architectures for more complex problems, using higher-order numerical solvers and finally tackle more complicated governing PDE problems with a larger number of unknown fields or in higher dimensions.

Acknowledgment

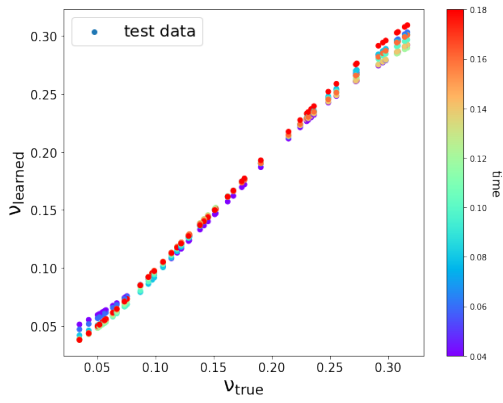
This research was supported by ARO W911NF-16-1-0136 and ONR N00014-17-1-2676.



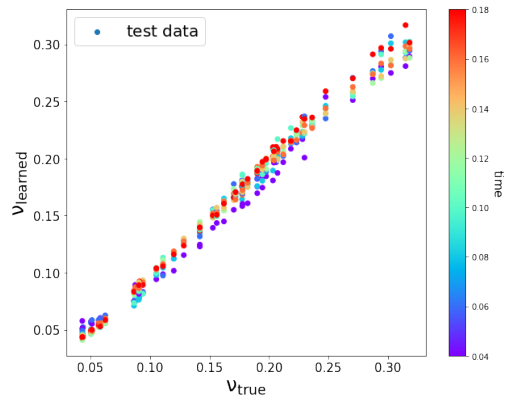
(a) Performance of encoder on training data set.



(b) Distribution of interpolation points and shape parameters discovered by the network.



(c) Performance of the encoder on unseen data.



(d) Performance of the encoder on unseen data with Gaussian noise with standard deviation 0.01.

Figure 18: Semantic autoencoder learns how to discover hidden variables from pairs of solutions. These results are obtained after 500 epochs on 50 data points along the ν -axis.

References

References

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*, 2016.
- [2] M. A. Aragon-Calvo. Self-supervised learning with physics-aware neural networks i: Galaxy model fitting. *arXiv preprint arXiv:1907.03957*, 2019.
- [3] R. R. Bailey and M. Srinath. Orthogonal moment features for use with parametric and non-parametric classifiers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 18(4):389–399, 1996.
- [4] L. Bar and N. Sochen. Unsupervised deep learning algorithm for PDE-based forward and inverse problems. *arXiv preprint arXiv:1904.05417*, 2019.
- [5] S. O. Belkasim, M. Shridhar, and M. Ahmadi. Pattern recognition with moment invariants: a comparative study and new results. *Pattern recognition*, 24(12):1117–1138, 1991.
- [6] J. Berg and K. Nyström. Neural network augmented inverse problems for PDEs. *arXiv preprint arXiv:1712.09685*, 2017.
- [7] J. Berg and K. Nyström. Data-driven discovery of PDEs in complex datasets. *Journal of Computational Physics*, 384:239–252, 2019.
- [8] D. S. Broomhead and D. Lowe. Radial basis functions, multi-variable functional interpolation and adaptive networks. Technical report, Royal Signals and Radar Establishment Malvern (United Kingdom), 1988.
- [9] J. M. Burgers. A mathematical model illustrating the theory of turbulence. In *Advances in applied mechanics*, volume 1, pages 171–199. Elsevier, 1948.
- [10] A. Chandrasekaran, D. Kamal, R. Batra, C. Kim, L. Chen, and R. Ramprasad. Solving the electronic structure problem with machine learning. *npj Computational Materials*, 5(1):22, 2019.
- [11] R. Chartrand. Numerical differentiation of noisy, nonsmooth data. *ISRN Applied Mathematics*, 2011, 2011.
- [12] F. Chollet et al. Keras, 2015.
- [13] B. C. Csáji et al. Approximation with artificial neural networks. *Faculty of Sciences, Eötvös Lornd University, Hungary*, 24(48):7, 2001.
- [14] J. Cullum. Numerical differentiation and regularization. *SIAM Journal on numerical analysis*, 8(2):254–265, 1971.
- [15] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.
- [16] N. Dal Santo, S. Deparis, and L. Pegolotti. Data driven approximation of parametrized PDEs by reduced basis and neural networks. *Journal of Computational Physics*, page 109550, 2020.
- [17] J. Darbon, G. P. Langlois, and T. Meng. Overcoming the curse of dimensionality for some Hamilton–Jacobi partial differential equations via neural network architectures. *Research in the Mathematical Sciences*, 7(3):1–50, 2020.
- [18] L. Debnath. *Nonlinear partial differential equations for scientists and engineers*. Springer Science & Business Media, 2011.

- [19] S. Dong, J. Kettenbach, I. Hinterleitner, H. Bergmann, and W. Birkfellner. The Zernike expansion—an example of a merit function for 2D/3D registration based on orthogonal functions. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 964–971. Springer, 2008.
- [20] C. L. Epstein. *Introduction to the mathematics of medical imaging*. SIAM, 2007.
- [21] R. Franke. Scattered data interpolation: tests of some methods. *Mathematics of computation*, 38(157):181–200, 1982.
- [22] D. V. Gaitonde and M. R. Visbal. High-order schemes for Navier-Stokes equations: algorithm and implementation into FDL3DI. Technical report, Air Force Research Lab Wright-Patterson AFB OH Air Vehicles Directorate, 1998.
- [23] Z. Geng, D. Johnson, and R. Fedkiw. Coercing machine learning to output physically accurate results. *Journal of Computational Physics*, page 109099, 2019.
- [24] F. Gibou, R. Fedkiw, and S. Osher. A review of level-set methods and some recent applications. *Journal of Computational Physics*, 353:82–109, 2018.
- [25] F. Gibou, D. Hyde, and R. Fedkiw. Sharp interface approaches and deep learning techniques for multiphase flows. *Journal of Computational Physics*, 380:442 – 463, 2019.
- [26] R. H. Hahnloser, R. Sarpeshkar, M. A. Mahowald, R. J. Douglas, and H. S. Seung. Digital selection and analogue amplification coexist in a cortex-inspired silicon circuit. *Nature*, 405(6789):947, 2000.
- [27] J. Han, A. Jentzen, and E. Weinan. Solving high-dimensional partial differential equations using deep learning. *Proceedings of the National Academy of Sciences*, 115(34):8505–8510, 2018.
- [28] R. L. Hardy. Multiquadric equations of topography and other irregular surfaces. *Journal of geophysical research*, 76(8):1905–1915, 1971.
- [29] S. He, Y. Li, Y. Feng, S. Ho, S. Ravanbakhsh, W. Chen, and B. Póczos. Learning to predict the cosmological structure formation. *Proceedings of the National Academy of Sciences*, page 201821458, 2019.
- [30] M. Hedayatpour, M. Mehrandezh, and F. Janabi-Sharifi. A unified approach to configuration-based dynamic analysis of quadcopters for optimal stability. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5116–5121. IEEE, 2017.
- [31] M. Hedayatpour, M. Mehrandezh, and F. Janabi-Sharifi. Precision modeling and optimally-safe design of quadcopters for controlled crash landing in case of rotor failure. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5206–5211. IEEE, 2019.
- [32] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012.
- [33] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [34] Y.-C. Hon and X. Mao. An efficient numerical scheme for Burgers’ equation. *Applied Mathematics and Computation*, 95(1):37–50, 1998.
- [35] G. P. J.-P. Fouque, J. Garnier and K. Solna. *Wave Propagation and Time Reversal in Randomly Layered Media*. Springer, 2007.

- [36] K. H. Jin, M. T. McCann, E. Froustey, and M. Unser. Deep convolutional neural network for inverse problems in imaging. *IEEE Transactions on Image Processing*, 26(9):4509–4522, 2017.
- [37] E. J. Kansa. Multiquadrics—a scattered data approximation scheme with applications to computational fluid-dynamics—i surface approximations and partial derivative estimates. *Computers & Mathematics with applications*, 19(8-9):127–145, 1990.
- [38] E. J. Kansa. Multiquadrics—a scattered data approximation scheme with applications to computational fluid-dynamics—ii solutions to parabolic, hyperbolic and elliptic partial differential equations. *Computers & mathematics with applications*, 19(8-9):147–161, 1990.
- [39] E. A. Kaye, Y. Hertzberg, M. Marx, B. Werner, G. Navon, M. Levoy, and K. B. Pauly. Application of zernike polynomials towards accelerated adaptive focusing of transcranial high intensity focused ultrasound. *Medical physics*, 39(10):6254–6263, 2012.
- [40] A. Khotanzad and Y. H. Hong. Invariant image recognition by zernike moments. *IEEE Transactions on pattern analysis and machine intelligence*, 12(5):489–497, 1990.
- [41] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [42] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [43] I. E. Lagaris, A. Likas, and D. I. Fotiadis. Artificial neural networks for solving ordinary and partial differential equations. *IEEE transactions on neural networks*, 9(5):987–1000, 1998.
- [44] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [45] J. Ling, R. Jones, and J. Templeton. Machine learning strategies for systems with invariance properties. *Journal of Computational Physics*, 318:22 – 35, 2016.
- [46] Z. Long, Y. Lu, and B. Dong. PDE-Net 2.0: Learning PDEs from data with a numeric-symbolic hybrid deep network. *Journal of Computational Physics*, 399:108925, 2019.
- [47] Z. Long, Y. Lu, X. Ma, and B. Dong. PDE-Net: Learning PDEs from data. In *International Conference on Machine Learning*, pages 3208–3216, 2018.
- [48] P. Y. Lu, S. Kim, and M. Soljačić. Extracting interpretable physical parameters from spatiotemporal systems using unsupervised learning. *Physical Review X*, 10(3):031056, 2020.
- [49] A. Maas, Q. V. Le, T. M. O’neil, O. Vinyals, P. Nguyen, and A. Y. Ng. Recurrent neural networks for noise reduction in robust ASR. 2012.
- [50] P. Markelj, D. Tomažević, B. Likar, and F. Pernuš. A review of 3d/2d registration methods for image-guided interventions. *Medical image analysis*, 16(3):642–661, 2012.
- [51] R. J. Mathar. Zernike basis to Cartesian transformations. *arXiv preprint arXiv:0809.2368*, 2008.
- [52] P. Mistani, A. Guittet, D. Bochkov, J. Schneider, D. Margetis, C. Ratsch, and F. Gibou. The island dynamics model on parallel quadtree grids. *Journal of Computational Physics*, 361:150–166, 2018.
- [53] P. Mistani, A. Guittet, C. Poignard, and F. Gibou. A parallel Voronoi-based approach for mesoscale simulations of cell aggregate electropermeabilization. *Journal of Computational Physics*, 380:48–64, 2019.
- [54] T. Nagatani. Density waves in traffic flow. *Physical Review E*, 61(4):3564, 2000.

- [55] F. Natterer. *The mathematics of computerized tomography*. SIAM, 2001.
- [56] T. Oliphant. *NumPy: A guide to NumPy*. USA: Trelgol Publishing, 2006–.
- [57] H. Owhadi. Bayesian numerical homogenization. *Multiscale Modeling & Simulation*, 13(3):812–828, 2015.
- [58] H. Owhadi. Calculus for the optimal quantification of uncertainties. Presented at the 2015 SIAM conference in Computational Science and Engineering, Salt Lake City, UT, 2015.
- [59] H. Owhadi. Multigrid with rough coefficients and multiresolution operator decomposition from hierarchical information games. *SIAM Review*, 59(1):99–149, 2017.
- [60] H. Owhadi and C. Scovel. *Operator-Adapted Wavelets, Fast Solvers, and Numerical Homogenization: From a Game Theoretic Approach to Numerical Approximation and Algorithm Design*, volume 35. Cambridge University Press, 2019.
- [61] H. Owhadi, C. Scovel, and F. Schäfer. Statistical numerical approximation. 2019.
- [62] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in PyTorch. 2017.
- [63] T. Poggio and F. Girosi. Networks for approximation and learning. *Proceedings of the IEEE*, 78(9):1481–1497, 1990.
- [64] R. J. Prokop and A. P. Reeves. A survey of moment-based techniques for unoccluded object representation and recognition. *CVGIP: Graphical Models and Image Processing*, 54(5):438–460, 1992.
- [65] R. Ragazzoni, E. Marchetti, and G. Valente. Adaptive-optics corrections available for the whole sky. *Nature*, 403(6765):54, 2000.
- [66] M. Raissi and G. E. Karniadakis. Hidden physics models: Machine learning of nonlinear partial differential equations. *Journal of Computational Physics*, 357:125 – 141, 2018.
- [67] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Machine learning of linear differential equations using gaussian processes. *Journal of Computational Physics*, 348:683 – 693, 2017.
- [68] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics informed deep learning (part i): Data-driven solutions of nonlinear partial differential equations. *arXiv preprint arXiv:1711.10561*, 2017.
- [69] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics informed deep learning (part ii): Data-driven discovery of nonlinear partial differential equations. *ArXiv*, abs/1711.10566, 2017.
- [70] K. Samsami, S. A. Mirbagheri, F. Meshkati, and H. C. Fu. Stability of soft magnetic helical microrobots. *Fluids*, 5(1):19, 2020.
- [71] M. Sari and G. Gürarslan. A sixth-order compact finite difference scheme to the numerical solutions of Burgers’ equation. *Applied Mathematics and Computation*, 208(2):475–483, 2009.
- [72] H. Schaeffer. Learning partial differential equations via data discovery and sparse optimization. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 473(2197):20160446, 2017.
- [73] H. Shen, D. George, E. Huerta, and Z. Zhao. Denoising gravitational waves using deep learning with recurrent denoising autoencoders. *arXiv preprint arXiv:1711.09919*, 2017.
- [74] C.-W. Shu and S. Osher. Efficient implementation of essentially non-oscillatory shock capturing schemes, 2. 1988.

- [75] A. V. Sinitskiy and V. S. Pande. Deep neural network computes electron densities and energies of a large set of organic molecules faster than density functional theory (DFT). *arXiv preprint arXiv:1809.02723*, 2018.
- [76] J. Sirignano and K. Spiliopoulos. DGM: A deep learning algorithm for solving partial differential equations. *Journal of Computational Physics*, 375:1339 – 1364, 2018.
- [77] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- [78] J. J. Stickel. Data smoothing and numerical differentiation by a regularization method. *Computers & chemical engineering*, 34(4):467–475, 2010.
- [79] P. Stinis, T. Hagge, A. M. Tartakovsky, and E. Yeung. Enforcing constraints for interpolation and extrapolation in generative adversarial networks. *Journal of Computational Physics*, 397:108844, 2019.
- [80] V. Tikhomirov. On the representation of continuous functions of several variables as superpositions of continuous functions of one variable and addition. In *Selected Works of AN Kolmogorov*, pages 383–387. Springer, 1991.
- [81] N. Trask, R. G. Patel, B. J. Gross, and P. J. Atzberger. GMLS-Nets: A framework for learning from unstructured data. *arXiv preprint arXiv:1909.05371*, 2019.
- [82] J. Valentin, C. Keskin, P. Pidlypenskyi, A. Makadia, A. Sud, and S. Bouaziz. Tensorflow graphics. Available at: <https://github.com/tensorflow/graphics>, 2019.
- [83] R. J. van Sloun, R. Cohen, and Y. C. Eldar. Deep learning in ultrasound imaging. *Proceedings of the IEEE*, 2019.
- [84] Z. von F. Beugungstheorie des schneidenverfahrens und seiner verbesserten form, der phasenkontrastmethode. *physica*, 1(7-12):689–704, 1934.
- [85] E. W. Weisstein. Zernike polynomial. 2002.
- [86] J. C. Wyant and K. Creath. Basic wavefront aberration theory for optical metrology. *Applied optics and optical engineering*, 11(part 2):28–39, 1992.
- [87] H. Xie and D. Li. A meshless method for Burgers’ equation using MQ-RBF and high-order temporal approximation. *Applied Mathematical Modelling*, 37(22):9215–9222, 2013.
- [88] K. Xu and E. Darve. The neural network approach to inverse problems in differential equations. *arXiv preprint arXiv:1901.07758*, 2019.
- [89] J. Zamudio-Fernandez, A. Okan, F. Villaescusa-Navarro, S. Bilaloglu, A. D. Cengiz, S. He, L. P. Levasseur, and S. Ho. HIGAN: Cosmic neutral hydrogen with generative adversarial networks. *arXiv preprint arXiv:1904.12846*, 2019.
- [90] X. Zhang, Y. Wang, W. Zhang, Y. Sun, S. He, G. Contardo, F. Villaescusa-Navarro, and S. Ho. From dark matter to galaxies with convolutional networks. *arXiv preprint arXiv:1902.05965*, 2019.
- [91] A. Zupanic, B. Kos, and D. Miklavcic. Treatment planning of electroporation-based medical interventions: electrochemotherapy, gene electrotransfer and irreversible electroporation. *Physics in Medicine & Biology*, 57(17):5425, 2012.