

A Mathematical Guide to Operator Learning

Nicolas Boullé

*Department of Applied Mathematics and Theoretical Physics
University of Cambridge
Cambridge, CB3 0WA, UK*

NB690@CAM.AC.UK

Alex Townsend

*Department of Mathematics
Cornell University
Ithaca, NY 14853, USA*

TOWNSEND@CORNELL.EDU

Abstract

Operator learning aims to discover properties of an underlying dynamical system or partial differential equation (PDE) from data. Here, we present a step-by-step guide to operator learning. We explain the types of problems and PDEs amenable to operator learning, discuss various neural network architectures, and explain how to employ numerical PDE solvers effectively. We also give advice on how to create and manage training data and conduct optimization. We offer intuition behind the various neural network architectures employed in operator learning by motivating them from the point-of-view of numerical linear algebra.

Keywords: Scientific machine learning, deep learning, operator learning, partial differential equations

Contents

1	Introduction	2
1.1	What is a neural operator?	4
1.2	Where is operator learning relevant?	4
1.3	Organization of the paper	6
2	From numerical linear algebra to operator learning	6
2.1	Low rank matrix recovery	7
2.2	Circulant matrix recovery	8
2.3	Banded matrix recovery	9
2.4	Hierarchical low rank matrix recovery	10
3	Neural operator architectures	10
3.1	Deep operator networks	12
3.2	Fourier neural operators	14
3.3	Deep Green networks	17
3.4	Graph neural operators	18
3.5	Multipole graph neural operators	19

4	Learning neural operators	22
4.1	Data acquisition	22
4.1.1	Distribution of source terms	22
4.1.2	Numerical PDE solvers	25
4.1.3	Amount of training data	26
4.2	Optimization	28
4.2.1	Loss functions	28
4.2.2	Optimization algorithms and implementation	29
4.2.3	Measuring convergence and super-resolution	29
5	Conclusions and future challenges	32
	References	34

1 Introduction

The recent successes of deep learning (LeCun et al., 2015) in computer vision (Krizhevsky et al., 2012), language model (Brown et al., 2020), and biology (Jumper et al., 2021) have caused a surge of interest in applying these techniques to scientific problems. The field of **scientific machine learning (SciML)** (Karniadakis et al., 2021), which combines the approximation power of machine learning (ML) methodologies and observational data with traditional modeling techniques based on partial differential equations (PDEs), sets out to use ML tools for accelerating scientific discovery.

SciML techniques can roughly be categorized into three main areas: (1) PDE solvers, (2) PDE discovery, and (3) operator learning (see Fig. 1). First, PDE solvers, such as physics-informed neural networks (PINNs) (Raissi et al., 2019; Lu et al., 2021b; Cuomo et al., 2022; Wang et al., 2023), the deep Galerkin method (Sirignano and Spiliopoulos, 2018), and the deep Ritz method (E and Yu, 2018), consist of approximating the solution a known PDE by a neural network by minimizing the solution’s residual. At the same time, PDE discovery aims to identify the coefficients of a PDE from data, such as the SINDy approach (Brunton et al., 2016; Champion et al., 2019), which relies on sparsity-promoting algorithms to determine coefficients of dynamical systems. There are also symbolic regression techniques, such as AI Feynman introduced by (Udrescu and Tegmark, 2020; Udrescu et al., 2020) and genetic algorithms (Schmidt and Lipson, 2009; Searson et al., 2010), that discover physics equations from experimental data.

Here, we focus on the third main area of SciML, called operator learning (Lu et al., 2021a; Kovachki et al., 2023). **Operator learning aims to discover or approximate an unknown operator \mathcal{A} , which often takes the form of the solution operator associated with a differential equation.** In mathematical terms, the problem can be defined as follows. Given pairs of data (f, u) , where $f \in \mathcal{U}$ and $u \in \mathcal{V}$ are from function spaces on a d -dimensional spatial domain $\Omega \subset \mathbb{R}^d$, and a (potentially nonlinear) operator $\mathcal{A} : \mathcal{U} \rightarrow \mathcal{V}$ such that $\mathcal{A}(f) = u$, the objective is to find an approximation of \mathcal{A} , denoted as $\hat{\mathcal{A}}$, such that for any new data $f' \in \mathcal{U}$, we have $\hat{\mathcal{A}}(f') \approx \mathcal{A}(f')$. In other words, the approximation should be accurate for both the training and unseen data, thus demonstrating good generalization.

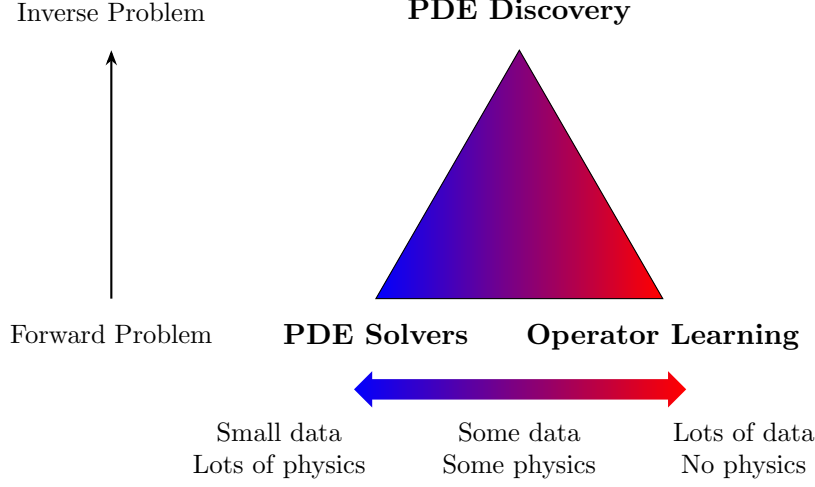


Figure 1: Illustrating the role of operator learning in SciML. Operator learning aims to discover or approximate an unknown operator \mathcal{A} , which often corresponds to the solution operator of an unknown PDE. In contrast, PDE discovery aims to discover coefficients of the PDE itself, while PDE solvers aim to solve a known PDE using ML techniques.

This problem is typically approached by representing $\hat{\mathcal{A}}$ as a neural operator, which is a generalization of neural networks as the inputs and outputs are functions, not vectors. After discretizing the functions at sensor points $x_1, \dots, x_m \in \Omega$, one then parametrizes the neural operator with a set of parameters $\theta \in \mathbb{R}^N$, which could represent the weights and biases of the underlying neural network. Then, one typically formulates an optimization problem to find the best parameters:

$$\min_{\theta \in \mathbb{R}^N} \sum_{(f,u) \in \text{data}} L(\hat{\mathcal{A}}(f; \theta), u), \quad (1)$$

where L is a loss function that measures the discrepancy between $\hat{\mathcal{A}}(f; \theta)$ and u , and the sum is over all available training data pairs (f, u) . The challenges of operator learning often arise from selecting an appropriate neural operator architecture for $\hat{\mathcal{A}}$, the computational complexities of solving the optimization problem, and the ability to generalize to new data.

A typical application of operator learning arises when learning the solution operator associated with a PDE, which maps a forcing function f to a solution u . One can informally think of it as the (right) inverse of a differential operator. One of the simplest examples is the solution operator associated with Poisson's equation with zero Dirichlet conditions:

$$-\nabla^2 u = f, \quad x \in \Omega \subset \mathbb{R}^d, \quad u|_{\partial\Omega} = 0, \quad (2)$$

where $\partial\Omega$ means the boundary of Ω . In this case, the solution operator, \mathcal{A} , can be expressed as an integral operator:

$$\mathcal{A}(f) = \int_{\Omega} G(\cdot, y) f(y) \, dy = u,$$

where G is the Green’s function associated with Eq. (2) (Evans, 2010, Chapt. 2). A neural operator is then trained to approximate the action of \mathcal{A} using training data pairs $(f_1, u_1), \dots, (f_M, u_M)$.

In general, recovering the solution operator is challenging, as it is often nonlinear and high-dimensional, and the available data may be scarce or noisy. Nevertheless, unlike inverse problems, which aim to recover source terms from solutions, the forward problem is usually well-posed. As we shall see, learning solution operators lead to new insights or applications that can complement inverse problem techniques, as described in two surveys (Stuart, 2010) and (Arridge et al., 2019).

1.1 What is a neural operator?

Neural operators (Kovachki et al., 2023; Lu et al., 2021a) are analogues of neural networks with infinite-dimensional inputs. Neural operators were introduced to generalize standard deep learning techniques to learn mappings between function spaces instead of between discrete vector spaces \mathbb{R}^{d_1} to \mathbb{R}^{d_L} , where d_1 is the input dimension of a neural network and d_L is the output dimension. In its most traditional formulation, a fully connected neural network can be written as a succession of affine transformations and nonlinear activation functions as

$$\mathcal{N}(x) = \sigma(A_L(\cdots \sigma(A_1 x + b_1) \cdots) + b_L),$$

where $L \geq 1$ is the number of layers, A_i are the weight matrices, b_i are the bias vectors, and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is the activation function, often chosen to be the ReLU function $\sigma(x) = \max\{x, 0\}$. Neural operators generalize this architecture, where the input and output of the neural network are functions instead of vectors. Hence, the input of a neural operator is a function $f : \Omega \rightarrow \mathbb{R}^{d_1}$, where $\Omega \subset \mathbb{R}^d$ is the domain of the function, and the output is a function $u : \Omega \rightarrow \mathbb{R}^{d_L}$. The neural operator is then defined as a composition of integral operators and nonlinear functions, which results in the following recursive definition at layer i :

$$u_{i+1}(x) = \sigma \left(\int_{\Omega_i} K^{(i)}(x, y) u_i(y) dy + b_i(x) \right), \quad x \in \Omega_{i+1}, \quad (3)$$

where $\Omega_i \subset \mathbb{R}^{d_i}$ is a compact domain, b_i is a bias function, and $K^{(i)}$ is the kernel. The kernels and biases are then parameterized and trained similarly to standard neural networks. However, approximating the kernels or evaluating the integral operators could be computationally expensive. Hence, several neural operator architectures have been proposed to overcome these challenges, such as DeepONets (Lu et al., 2021a) and Fourier neural operators (Li et al., 2021a).

1.2 Where is operator learning relevant?

Operator learning has been successfully applied to many PDEs from different fields, including fluid dynamics with simulations of fluid flow turbulence in the Navier–Stokes equations at high Reynolds number (Li et al., 2023b; Peng et al., 2022), continuum mechanics (You et al., 2022), astrophysics (Mao et al., 2023), quantum mechanics with the Schrödinger equation (Li et al., 2021a), and weather forecasting (Kurth et al., 2023; Lam et al., 2023). The following four types of applications might directly benefit from operator learning.

Speeding up numerical PDE solvers. First, one can use operator learning to build reduced-order models of complex systems that are computationally challenging to simulate with traditional numerical PDE solvers. For example, this situation arises in fluid dynamics applications such as modeling turbulent flows, which require a very fine discretization or the simulation of high dimensional PDEs. Moreover, specific problems in engineering require the evaluation of the solution operator many times, such as in the design of aircraft or wind turbines. In these cases, a fast but less accurate solver provided by operator learning may be used for forecasting or optimization. This is one of the main motivations behind Fourier neural operators in (Li et al., 2021a). There are also applications of operator learning (Zheng et al., 2023) to speed up the sampling process in diffusion models or score-based generative models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2021), which require solving complex differential equations. However, one must be careful when comparing performance against classical numerical PDE solvers, mainly due to the significant training time required by operator learning.

Parameter optimization. In our experience, the computational efficiency of operator learning is mainly seen in downstream applications such as parameter optimization. Once the solution operator has been approximated, it can be exploited in an inverse problem framework to recover unknown parameters of the PDE, which may be computationally challenging to perform with existing numerical PDE solvers. Additionally, neural operators do not rely on a fixed discretization as they are mesh-free and parameterized by a neural network that can be evaluated at any point. This property makes them suitable for solving PDEs on irregular domains or transferring the model to other spatial resolutions (Kovachki et al., 2023).

Benchmarking new techniques. Operator learning may also be used to benchmark and develop new deep learning models. As an example, one can design specific neural network architectures to preserve quantities of interest in PDEs, such as symmetries (Olver, 1993a), conservation laws (Evans, 2010, Sec. 3.4), and discretization independence. This could lead to efficient architectures that are more interpretable and generalize better to unseen data, and exploit geometric priors within datasets (Bronstein et al., 2021). Moreover, the vast literature on PDEs and numerical solvers can be leveraged to create datasets and assess the performance of these models in various settings without requiring significant computational resources for training.

Discovering unknown physics. Last but not least, operator learning is helpful for the discovery of new physics (Lu et al., 2021a). Indeed, the solution operator of a PDE is often unknown, or one may only have access to a few data points without any prior knowledge of the underlying PDE. In this case, operator learning can be used to discover the PDE or a mathematical model to perform predictions. This can lead to new insights into the system’s behavior, such as finding conservation laws, symmetries, shock locations, and singularities (Boullé et al., 2022a). However, the complex nature of neural networks makes them challenging to interpret or explain, and there are many future directions for making SciML more interpretative.

1.3 Organization of the paper

This paper is organized as follows. We begin in Section 2 by exploring the connections between numerical linear algebra and operator learning. Then, we review the main neural network architectures used to approximate operators in Section 3. In Section 4, we focus on the data acquisition process, a crucial step in operator learning. We discuss the choice of the distribution of source terms used to probe the system, the numerical PDE solver, and the number of training data. Along the way, we analyze the optimization pipeline, including the possible choices of loss functions, optimization algorithms, and assessment of the results. Finally, in Section 5, we conclude with a discussion of the remaining challenges in the field that include the development of open-source software and datasets, the theoretical understanding of the optimization procedure, and the discovery of physical properties in operator learning, such as symmetries and conservation laws.

2 From numerical linear algebra to operator learning

There is a strong connection between operator learning and the recovery of structured matrices from matrix-vector products. Suppose one aims to approximate the solution operator associated with a linear PDE using a single layer neural operator in the form of Eq. (3) without the nonlinear activation function. Then, after discretizing the integral operator using a quadrature rule, it can be written as a matrix-vector product, where the integral kernel $K : \Omega \times \Omega \rightarrow \mathbb{R}$ is approximated by a matrix $A \in \mathbb{R}^{N \times N}$. Moreover, the structure of the matrix is inherited from the properties of the Green’s function (see Table 1). The matrix’s underlying structure—whether it is low-rank, circulant, banded, or hierarchical low-rank—plays a crucial role in determining the efficiency and approach of the recovery process. This section describes the matrix recovery problem as a helpful way to gain intuition about operator learning and the design of neural operator architectures (see Section 3). Another motivation for recovering structured solution operators is to ensure that the neural operators are fast to evaluate, which is essential in applications involving parameter optimization and benchmarking (Kovachki et al., 2023).

Property of the PDE	Solution operator’s kernel	Matrix structure
—	Globally smooth	Low rank
Periodic BCs & const. coeffs	Convolution kernel	Circulant
Localized behavior	Off-diagonal decay	Banded
Elliptic / Parabolic	Off-diagonal low rank	Hierarchical

Table 1: Solution operators associated with linear PDEs can often be represented as integral operators with a kernel called a Green’s function. The properties of a linear PDE induce different structures on the Green’s function, such as translation-invariant or off-diagonal low-rank. When these integral operators are discretized, one forms a matrix-vector product, and hence the matrix recovery problem can be viewed as a discrete analogue of operator learning. The dash in the first column and row means no PDE has a solution operator with a globally smooth kernel.

Consider an unknown matrix $A \in \mathbb{R}^{N \times N}$ with a known structure such as rank- k or banded. We assume that the matrix A is a black box and cannot be seen, but that one can probe A for information via matrix-vector products, i.e., the maps $x \mapsto Ax$ and $x \mapsto A^\top x$, with A^\top representing the transpose of A . The matrix recovery problem is the task of approximating the matrix A using as few queries to $x \mapsto Ax$ and $x \mapsto A^\top x$ as possible. Every matrix with N columns can be deduced in a maximum of N matrix-vector product queries, as Ae_j for $1 \leq j \leq N$ returns the j th column, where e_j denotes the j th standard basis vector. However, if the matrix A has a specific structure such as low-rank, circulant, banded, or hierarchical low-rank, it is often possible to recover A using far fewer queries. This section describes how to recover structured matrices efficiently using matrix-vector products. We prefer doing matrix recovery with Gaussian random vectors because the infinite-dimensional analogue of these vectors are random functions drawn from a Gaussian process, which is a widespread choice of training input data in operator learning.

2.1 Low rank matrix recovery

Let $A \in \mathbb{R}^{N \times N}$ be a rank- k matrix, then it can be expressed for some $C \in \mathbb{R}^{N \times k}$ and $R \in \mathbb{R}^{k \times N}$ as

$$A = CR.$$

Halikias and Townsend 2023 showed that at least $2k$ queries are required to capture the k -dimensional row and column spaces and deduce A . Halko et al. 2011; Martinsson and Tropp 2020 introduced the randomized singular value decomposition (SVD) as a method to recover a rank- k matrix with probability one in $2k$ matrix-vector products with random Gaussian vectors. The randomized SVD can be expressed as a recovery algorithm in Algorithm 1.

Algorithm 1 Randomized singular value decomposition.

- 1: Draw a random matrix $X \in \mathbb{R}^{N \times k}$ with i.i.d. standard Gaussian entries.
 - 2: Perform k queries with A : $Y = AX$.
 - 3: Compute the QR factorization of $Y = QR$.
 - 4: Perform k queries with A^\top : $Z = A^\top Q$.
 - 5: Return $A = QZ^\top$.
-

A randomized algorithm is crucial for low-rank matrix recovery to prevent input vectors from lying within the $N - k$ dimensional nullspace of A . Hence, recovering any low-rank matrix with a deterministic algorithm using fixed input vectors is impossible. Then, for the rank- k matrix recovery problem, Algorithm 1 recovers A with probability one. A small oversampling parameter $p \geq 1$ is used for numerical stability, such as $p = 5$. This means that $X \in \mathbb{R}^{N \times (k+p)}$, preventing the chance that a random Gaussian vector might be highly aligned with the nullspace of A .

A convenient feature of Algorithm 1 is that it also works for matrices with numerical rank¹ k , provided that one uses a random matrix X with $k + p$ columns. In particular, a simplified statement of (Halko et al., 2011, Thm. 10.7) shows that the randomized SVD

1. For a fixed $0 < \epsilon < 1$, we say that a matrix A has numerical rank k if $\sigma_{k+1}(A) < \epsilon \sigma_1(A)$ and $\sigma_k(A) \geq \epsilon \sigma_1(A)$, where $\sigma_1(A) \geq \sigma_2(A) \geq \dots \geq \sigma_N(A) \geq 0$ are the singular values of A .

recovers a near-best low-rank matrix in the sense that

$$\mathbb{P} \left[\|A - QZ^\top\|_F \leq \left(1 + 15\sqrt{k+5}\right) \min_{\text{rank}(A_k) \leq k} \|A - A_k\|_F \right] \geq 0.999,$$

where $\|\cdot\|_F$ is the Frobenius norm of A . While other random embeddings can be used to probe A , Gaussian random vectors give the cleanest probability bounds (Martinsson and Tropp, 2020). Moreover, ensuring that the entries in each column of X have some correlation and come from a multivariable Gaussian distribution allows for the infinite-dimensional extension of the randomized SVD and its application to recover Hilbert–Schmidt operators (Boullé and Townsend, 2022, 2023). This analysis allows one to adapt Algorithm 1 to recover solution operators with low-rank kernels.

Low-rank matrix recovery is one of the most straightforward settings to motivate DeepONet (see Section 3.1). One of the core features of DeepONet is to use the trunk net to represent the action of a solution operator on a set of basis functions generated by the so-called branch net. Whereas in low-rank matrix recovery, we often randomly draw the columns of X as input vectors, DeepONet is trained with these functions. However, like DeepONet, low-rank matrix recovery is constructing an accurate approximant whose action is on these vectors. Many operators between function spaces can often be represented to high accuracy with DeepONet in the same way that the kernels of solution operators associated with linear PDEs often have algebraically fast decaying singular values.

2.2 Circulant matrix recovery

The FNO (Li et al., 2021a) structure is closely related to circulant matrix recovery. Consider an $N \times N$ circulant matrix C_c , which is parameterized a vector $c \in \mathbb{R}^N$ as follows:

$$C_c = \begin{bmatrix} c_0 & c_{N-1} & \cdots & c_2 & c_1 \\ c_1 & c_0 & c_{N-1} & \ddots & c_2 \\ \vdots & c_1 & c_0 & \ddots & \vdots \\ c_{N-2} & \ddots & \ddots & \ddots & c_{N-1} \\ c_{N-1} & c_{N-2} & \cdots & c_1 & c_0 \end{bmatrix}.$$

To recover C_c with a random Gaussian vector g , we recall that C_c can be interpreted as a multiplication operator in the Fourier basis. By associativity of multiplication, we have

$$C_c g = C_g c.$$

If we perform the matrix-vector product query $y = C_c g$, we can find the vector c by solving the linear system $C_g c = y$. Since c completely defines C_c , we have recovered the circulant matrix. Moreover, the linear system $C_g c = y$ can be solved efficiently using the fast Fourier transform (FFT) in $\mathcal{O}(N \log N)$ operations. A convenient feature of circulant matrices is that given a new vector $x \in \mathbb{R}^N$, one can compute $C_c x$ in $\mathcal{O}(N \log N)$ operations using the FFT.

Circulant matrix recovery motivates Fourier neural operators. Hence, FNOs leverage the fast Fourier transform to efficiently parameterize the kernel of a solution operator,

essentially capturing the operator in a spectral sense. Similarly, circulant matrices are diagonalized by the discrete Fourier transform matrix. The infinite-dimensional analogue of a circulant matrix is a solution operator with a periodic and translation invariant kernel, and this is the class of solution operators for which the FNO assumptions are fully justified. FNOs are extremely fast to evaluate because of their structure, making them popular for parameter optimization and favorable for benchmarking against reduced-order models.

2.3 Banded matrix recovery

We now consider a banded matrix $A \in \mathbb{R}^{N \times N}$ with a fixed bandwidth w , i.e.,

$$A_{ij} = 0, \quad \text{if } |i - j| > w.$$

The matrix A can be recovered with $w + 2$ matrix-vector products, but not fewer, using the $w + 2$ columns of the following matrix as input vectors:

$$[I_{w+2} \quad \cdots \quad I_{w+2}]^\top,$$

where I_{w+2} is the $(w + 2) \times (w + 2)$ identity matrix. Since every w th column has disjoint support, these input vectors recover the columns of A . Of course, this also means that a A can be recovered with $w + 2$ Gaussian random vectors.

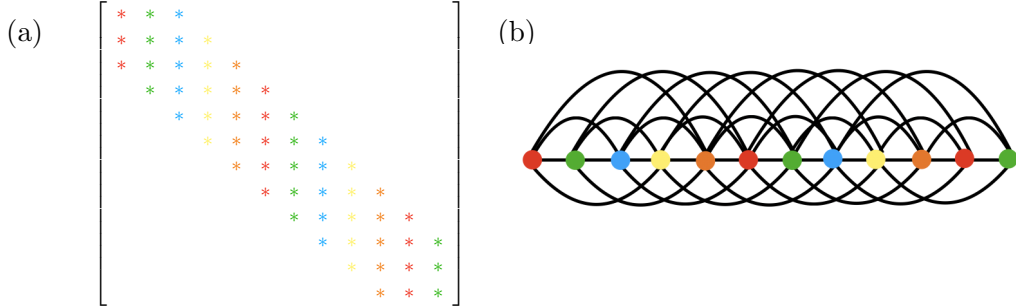


Figure 2: (a) A generic 12×12 banded matrix with bandwidth 2, with a maximum of 5 diagonals, and the corresponding graph (b). Here, each vertex is a column of the banded matrix, and two vertices are connected if their corresponding columns do not have disjoint support. The coloring number of 5 determines the minimum number of matrix-vector products needed to recover the structure. Generally, an $N \times N$ banded matrix with bandwidth w can be recovered in $2w + 1$ matrix-vector products.

There is a way to understand how many queries one needs as a graph-coloring problem. Consider the graph of size N , corresponding to an $N \times N$ banded matrix with bandwidth w , where two vertices are connected if their corresponding columns do not have disjoint support (see Fig. 2). Then, the minimum number of matrix-vector product queries needed to recover A is the graph coloring number of this graph.² One can see why this is the case

2. Recall that the coloring number of a graph is the minimum number of colors required to color the vertices so that no two vertices connected by an edge are identically colored.

because all the columns with the same color can be deduced simultaneously with a single matrix-vector product as they must have disjoint support.

Banded matrix recovery motivates Graph neural operators (GNOs), which we will describe later in Section 3.4, as both techniques exploit localized structures within data. GNOs use the idea that relationships in nature are local and can be represented as graphs with no faraway connections. By only allowing local connections, GNOs can efficiently represent solution operators corresponding to local solution operators, mirroring the way banded matrices are concentrated on the diagonal. Likewise, with a strong locality, GNOs are relatively fast to evaluate, making them useful for parameter optimization and benchmarking. However, they may underperform if the bandwidth increases or the solution operator is not local.

2.4 Hierarchical low rank matrix recovery

An $N \times N$ rank- k hierarchical off-diagonal low rank (HODLR) matrix, denoted as $H_{N,k}$, is a structure that frequently appears in the context of discretized solution operators associated with elliptic and parabolic PDEs (Hackbusch et al., 2004). To understand its recursive structure, we assume N to be a power of 2 and illustrate the structure in Fig. 3(a).

Since a rank- k matrix requires $2k$ matrix-vector product queries to be recovered (see Section 2.1), a naive approach to deducing $H_{N,k}$ is to use $2k$ independent queries on each submatrix. However, one can show that some of the submatrices of $H_{N,k}$ can be recovered concurrently using the same queries. We use a graph coloring approach (Levitt and Martinsson, 2022b) to determine which submatrices can be recovered concurrently. This time, we consider the graph where each vertex is a low-rank submatrix of $H_{N,k}$ and connect two vertices if their corresponding low-rank submatrices occupy the same column as in Fig. 3(b). The low-rank submatrices that are identically colored at each level can be recovered concurrently in only $2k$ queries. Hence, it can be shown that an $N \times N$ hierarchical rank- k matrix can be recovered in fewer than $10k \lceil \log_2(N) \rceil$ matrix-vector products (Halikias and Townsend, 2023). The precise coloring of the graph in Fig. 3(b) can also be used to derive a particular algorithm for hierarchical matrix recovery known as peeling (Levitt and Martinsson, 2022a,b; Lin et al., 2011; Martinsson, 2011). These peeling algorithms have been recently generalized to the infinite-dimensional setting by (Boullé et al., 2023).

HODLR recovery can be seen as the simplest version of a multipole graph neural operator (MGNO) (see Section 3.5) as both emphasize the importance of capturing operators at multiple scales. MGNOs are based on a hierarchical graph with interactions at different scales or levels (see Section 3.4). By incorporating local (near-field) and global (far-field) interactions, MGNOs can effectively learn complex patterns. MGNOs are often great at representing solution operators due to their multiscale nature. The price to pay is that the final neural operator can be computationally expensive to evaluate, and it is a complicated structure to implement.

3 Neural operator architectures

In this section, we review the main neural operator architectures used in the literature, namely DeepONets (Lu et al., 2021a), Fourier neural operators (Li et al., 2021a), and Deep Green networks (Gin et al., 2021; Boullé et al., 2022a). We also refer to the recent survey

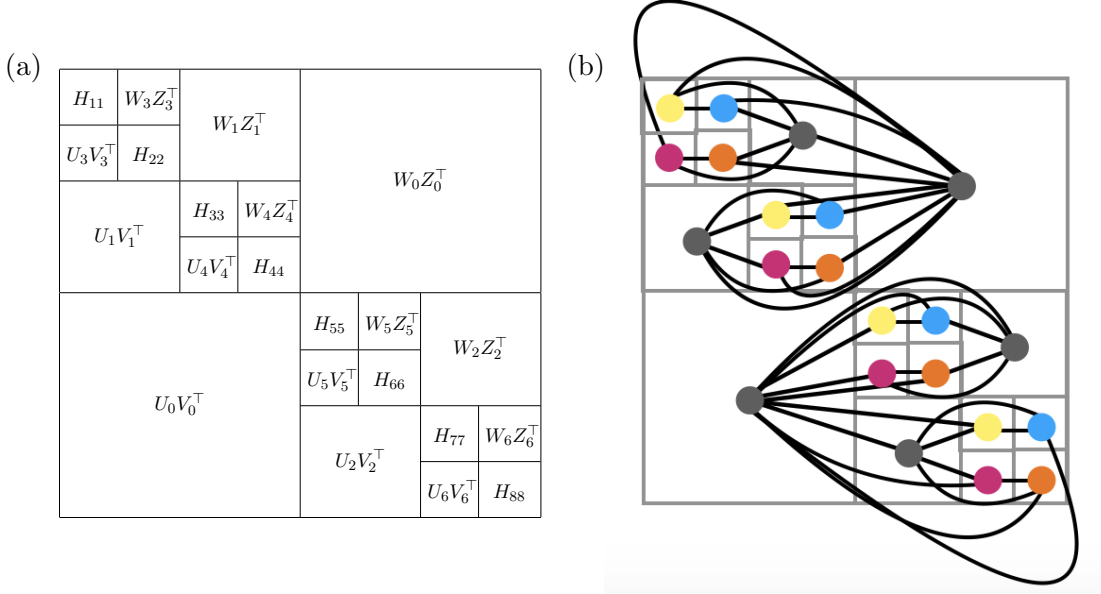


Figure 3: (a) A HODLR matrix $H_{N,k}$ after three levels of partitioning. Since $H_{N,k}$ is a rank- k HODLR matrix, U_i , V_i , W_i , and Z_i have at most k columns. The matrices A_{ii} are themselves rank- k HODLR matrices of size $N/8 \times N/8$ and can be further partitioned. (b) Graph corresponding to a hierarchical low-rank matrix with three levels. Here, each vertex is a low-rank block of the matrix, where two vertices are connected if their low-rank blocks occupy the same row. At each level, the number of required matrix-vector input probes to recover that level is proportional to the coloring number of the graph when restricted to submatrices of the same size. In this case, the submatrices that are identically colored can be recovered simultaneously.

by (Goswami et al., 2023) for a review of the different neural operator architectures and their applications. Each of these architectures employs different discretization and approximation techniques to make the neural operator more efficient and scalable by enforcing certain structures on the kernel such as low-rank, periodicity, translation invariance, or hierarchical low-rank structure (see Table 2).

Most neural operator architectures also come with theoretical guarantees on their approximation power. These theoretical results essentially consist of universal approximation properties for neural operators (Chen and Chen, 1995; Kovachki et al., 2023; Lu et al., 2021a), in a similar manner as neural networks (DeVore, 1998), and quantitative error bounds based on approximation theory to estimate the size, i.e., the number of trainable parameters, of a neural operator needed to approximate a given operator between Banach spaces to within a prescribed accuracy (Lanthaler et al., 2022; Yarotsky, 2017).

Neural operators	Property of the operator	Kernel parameterization
DeepONet	Low-rank	Branch and trunk networks
FNO	Translation-invariant	Fourier coefficients
GreenLearning	Linear	Rational neural network
DeepGreen	Semi-linear	Kernel matrix
Graph neural operator	Diagonally dominant	Message passing network
Multipole GNO	Off-diagonal low rank	Neural network

Table 2: Summary table of neural operator architectures, describing the property assumption on the operator along with the discretization of the integral kernels.

3.1 Deep operator networks

Deep Operator Networks (DeepONets) are a promising model for learning nonlinear operators and capturing the inherent relationships between input and output functions (Lu et al., 2021a). They extend the capabilities of traditional deep learning techniques by leveraging the expressive power of neural networks to approximate operators in differential equations, integral equations, or more broadly, any functional maps from one function space to another. A key theoretical motivation for DeepONet is the universal operator approximation theorem (Chen and Chen, 1995; Lu et al., 2021a). This result can be seen as an infinite dimensional analogue of the universal approximation operator for neural networks (Cybenko, 1989; Hornik, 1991), which guarantee that a sufficiently wide neural network can approximate any continuous function to any accuracy. Since the introduction of DeepONets by (Lu et al., 2021a), several research works focused on deriving error bounds for the approximation of nonlinear operators by DeepONets in various settings, such as learning the solution operator associated with Burger’s equation or the advection-diffusion equation (Deng et al., 2022), and the approximation of nonlinear parabolic PDEs (De Ryck and Mishra, 2022; Lanthaler et al., 2022).

A DeepONet is a two-part deep learning network consisting of a branch network and a trunk network. The branch net encodes the operator’s input functions f into compact, fixed-size latent vectors $b_1(f(x_1), \dots, f(x_m)), \dots, b_p(f(x_1), \dots, f(x_m))$, where $\{x_i\}_{i=1}^m$ are the sensor points at which the input functions are evaluated. The trunk net decodes these latent vectors to produce the final output function at the location $y \in \Omega$ as

$$\mathcal{N}(f)(y) = \sum_{k=1}^p b_k(f(x_1), \dots, f(x_m)) t_k(y).$$

A schematic of a deep operator network is given in Fig. 4. The defining feature of DeepONets is their ability to handle functional input and output, thus enabling them to learn a wide array of mathematical operators effectively. It’s worth mentioning that the branch network and the trunk network can have distinct neural network architectures tailored for different purposes, such as performing a feature expansion on the input of the trunk network as $y \rightarrow (y \ \cos(\pi y) \ \sin(\pi y) \ \dots)$ to take into account any potential oscillatory patterns in the data (Di Leoni et al., 2023). Moreover, while the interplay of the branch and trunk networks is crucial, the output of a DeepONet does not necessarily depend on the specific

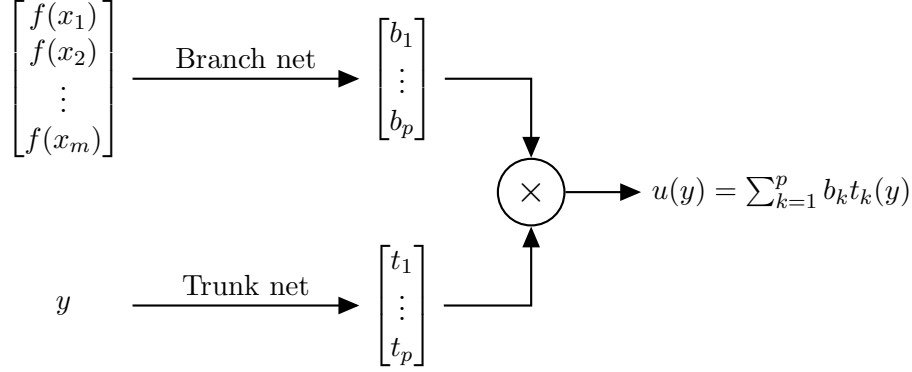


Figure 4: Schematic diagram of a deep operator network (DeepONet). A DeepONet parametrizes a neural operator using a branch network and a truncation (trunk) network. The branch network encodes the input function f as a vector of p features, which is then multiplied by the trunk network to yield a rank- p representation of the solution u .

input points but rather on the global property of the entire input function, which makes it suitable for learning operator maps.

One reason behind the performance of DeepONet might be its connection with the low-rank approximation of operators and the SVD (see Section 2.1). Hence, one can view the trunk network as learning a basis of functions $\{t_k\}_{k=1}^p$ that are used to approximate the operator, while the branch network expresses the output function in this basis by learning the coefficients $\{b_k\}_{k=1}^p$. Moreover, the branch network can be seen as a feature extractor, which encodes the input function into a compact representation, thus reducing the problem’s dimensionality to p , where p is the number of branch networks. Additionally, several architectures, namely the POD-DeepONet (Lu et al., 2022) and SVD-DeepONet (Venturi and Casey, 2023), have been proposed to strengthen the connections between DeepONet and the SVD of the operator and increase its interpretability.

A desirable property for a neural operator architecture is to be discretization invariant in the sense that the model can act on any discretization of the source term and be evaluated at any point of the domain (Kovachki et al., 2023). This property is crucial for the generalization of the model to unseen data and the transferability of the model to other spatial resolutions. While DeepONets can be evaluated at any location of the output domain, DeepONets are not discretization invariant in their original formulation by Lu et al. 2021a as the branch network is evaluated at specific points of the input domain (see Fig. 4). However, this can be resolved using a low-rank neural operator (Kovachki et al., 2023), sampling the input functions at local spatial averages (Lanthaler et al., 2022), or employing a principal component analysis (PCA) alternative of the branch network (de Hoop et al., 2022).

The training of DeepONets is performed using a supervised learning process. It involves minimizing the mean-squared error between the predicted output $\mathcal{N}(f)(y)$ and the actual

output of u the operator on the training functions at random locations $\{y_j\}_{j=1}^n$, i.e.,

$$\min_{\theta \in \mathbb{R}^N} \frac{1}{|\text{data}|} \sum_{(f,u) \in \text{data}} \frac{1}{n} \sum_{j=1}^n |\mathcal{N}(f)(y_j) - u(y_j)|^2. \quad (4)$$

The term inside the first sum approximates the integral of the mean-squared error, $|\mathcal{N}(f) - u|^2$, over the domain Ω using Monte-Carlo integration. The optimization is typically done via backpropagation and gradient descent algorithms, which are the same as in traditional neural networks. Importantly, DeepONets allow for different choices of loss functions, depending on the problem. For example, mean squared error is commonly used for regression tasks, but other loss functions might be defined to act as a regularizer and incorporate prior physical knowledge of the problem (Goswami et al., 2022; Wang et al., 2021b). The selection of an appropriate loss function is a crucial step in defining the learning process of these networks and has a substantial impact on their performance (see Section 4.2.1).

DeepONet has been successfully applied and adapted to a wide range of problems, including predicting cracks in fracture mechanics using a variational formulation of the governing equations (Goswami et al., 2022), simulating the New York-New England power grid behavior with a probabilistic and Bayesian framework to quantify the uncertainty of the trajectories (Moya et al., 2023), as well as predicting linear instabilities in high-speed compressible flows with boundary layers (Di Leoni et al., 2023).

3.2 Fourier neural operators

Fourier neural operators (FNOs) (Li et al., 2021a; Kovachki et al., 2023) are a class of neural operators motivated by Fourier spectral methods. FNOs have found their niche in dealing with high-dimensional PDEs, which are notoriously difficult to solve using traditional numerical methods due to the curse of dimensionality. They’ve demonstrated significant success in learning and predicting solutions to various PDEs, particularly those with periodic boundary conditions or those that can be transformed into the spectral domain via Fourier transform. This capability renders FNOs an invaluable tool in areas where PDEs play a central role, such as fluid dynamics, quantum mechanics, and electromagnetism.

The main idea behind FNOs is to choose the kernels $K^{(i)}$ in Eq. (3) as translation-invariant kernels satisfying $K^{(i)}(x, y) = k^{(i)}(x - y)$ (provided the input and output domains are torus) such that the integration of the kernel can be performed efficiently as a convolution using the Fast Fourier Transform (FFT) (Cooley and Tukey, 1965), i.e., multiplication in the feature space of Fourier coefficients. Hence, the integral operation in Eq. (3) can be performed as

$$\int_{\Omega_i} k^{(i)}(x - y) u_i \, dy = \mathcal{F}^{-1}(\mathcal{F}(k^{(i)})\mathcal{F}(u_i))(x) = \mathcal{F}^{-1}(\mathcal{R} \cdot \mathcal{F}(u_i))(x), \quad x \in \Omega_i,$$

where \mathcal{F} denotes the Fourier transform and \mathcal{F}^{-1} its inverse. The kernel $K^{(i)}$ is parametrized by a periodic function $k^{(i)}$, which is discretized by a (trainable) weight vector of Fourier coefficients \mathcal{R} , and truncated to a finite number of Fourier modes. Then, if the input domain is discretized uniformly with m sensor points, and the vector \mathcal{R} contains at most $k_{\max} \leq m$ modes, the convolution can be performed in quasi-linear complexity in $\mathcal{O}(m \log m)$ operations via the FFT. This is a significant improvement over the $\mathcal{O}(m^2)$ operations required

to evaluate the integral in Eq. (3) using a quadrature rule. In practice, one can restrict the number of Fourier modes to $k_{\max} \ll m$ without significantly affecting the accuracy of the approximation whenever the input and output functions are smooth so that their representation in the Fourier basis enjoy rapid decay of the coefficients, thus further reducing the computational and training complexity of the neural operator.

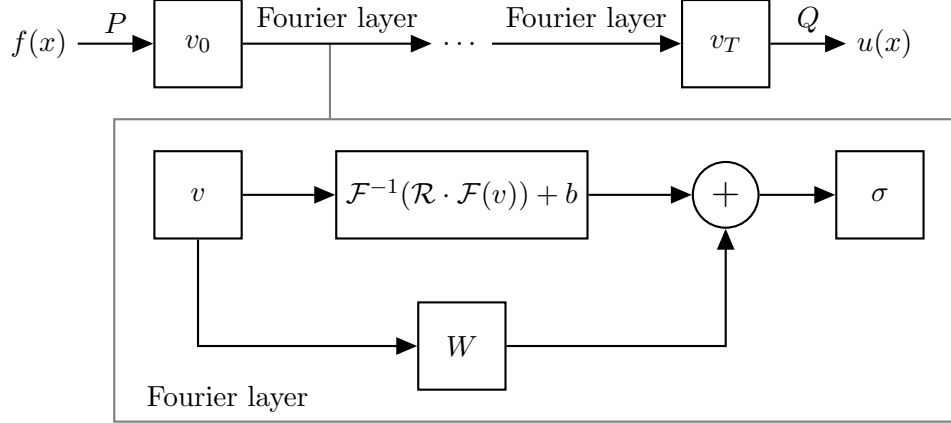


Figure 5: Schematic diagram of a Fourier neural operator (FNO). The networks P and Q , respectively, lift the input function f to a higher dimensional space and project the output of the last Fourier layer to the output dimension. An FNO mainly consists of a succession of Fourier layers, which perform the integral operations in neural operators as a convolution in the Fourier domain and component-wise composition with an activation function σ .

We display a diagram of the architecture of an FNO in Fig. 5. The input function f is first lifted to a higher dimensional space by a neural network P . Then, a succession of Fourier layers is applied to the lifted function, which is parametrized by a vector of Fourier coefficients \mathcal{R}_i , a bias vector b_i , and a weight matrix W . Then, the output of the FNO at the i th layer is given by

$$v_i = \sigma(W_i v_{i-1} + \mathcal{F}^{-1}(\mathcal{R}_i \cdot \mathcal{F}(u_{i-1})) + b_i),$$

where $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is the activation function whose action is defined component-wise, often chosen to be the ReLU function. The weight matrix W_i and bias vector b_i perform a linear transformation of the input function v_i . After the last Fourier layer, the output of the FNO is obtained by applying a final neural network Q on the output of the last Fourier layer to project it to the output dimension.

The training of FNOs, like DeepONets, is carried out via a supervised learning process. It typically involves minimizing a loss function that measures the discrepancy between the predicted and the true output of the operator on the input functions with respect to the trainable parameters of the neural network as in Eq. (4). Here, one needs to perform backpropagation through the Fourier layers, which is enabled by the implementation of fast GPU differentiable FFTs (Mathieu et al., 2014) in deep learning frameworks such as PyTorch (Paszke et al., 2019) and TensorFlow (Abadi et al., 2016).

While FNOs have been proposed initially to alleviate the computational expense of performing integral operations in neural operators by leveraging the FFT, they have a distinctive advantage in learning operators where computations in the spectral domain are more efficient or desirable. This arises naturally when the target operator, along with the input and output functions, are smooth so that their representation as Fourier coefficients decay exponentially fast, yielding an efficient truncation. Hence, by selecting the architecture of the FNO appropriately, such as the number of Fourier modes k_{\max} , or the initialization of the Fourier coefficients, one can obtain a neural operator that preserves specific smoothness properties. However, when the input or output training data is not smooth, FNO might suffer from Runge’s phenomenon near discontinuities (de Hoop et al., 2022).

One main limitation of the FNO architectures is that the FFT should be performed on a uniform grid and rectangular domains, which is not always the case in practice. This can be overcome by applying embedding techniques to transform the input functions to a uniform grid and extend them to simple geometry, using a Fourier analytic continuation technique (Bruno et al., 2007). Recently, several works have been proposed to extend the FNO architecture to more general domains, such as using a zero padding, linear interpolation (Lu et al., 2022), or encoding the geometry to a regular latent space with a neural network (Li et al., 2022, 2023a). However, this might lead to a loss of accuracy and additional computational cost. Moreover, the FFT is only efficient for approximating translation invariant kernels, which do not occur when learning solution operators of PDEs with non-constant coefficients.

Other related architectures aim to approximate neural operators directly in the feature space, such as spectral neural operators (SNO) (Fanaskov and Oseledets, 2022), which are based on spectral methods and employ a simple feedforward neural network to map the input function, represented as a vector of Fourier or Chebyshev coefficients to an output vector of coefficients. Finally, (Raonic et al., 2023) introduces convolutional neural operators (CNOs) to alleviate the aliasing phenomenon of convolutional neural networks (CNNs) by learning the mapping between bandlimited functions. Contrary to FNOs, CNOs parameterize the integral kernel on a $k \times k$ grid and perform the convolution in the physical space as

$$\int_{\Omega} k(x - y)f(y) dy = \sum_{i,j=1}^k k_{ij}f(x - z_{ij}), \quad x \in \Omega,$$

where z_{ij} are the grid points.

Similarly to DeepONets, Fourier neural operators are universal approximators, in the sense that they are dense in the space of continuous operators (Bhattacharya et al., 2021; Kovachki et al., 2023, 2021). However, even while being universal approximators, FNOs could, in theory, require a huge number of parameters to approximate a given operator to a prescribed accuracy $\epsilon > 0$. As an example, (Kovachki et al., 2021) showed that the size of the FNO must grow exponentially fast as ϵ decreases to approximate any operator between rough functions whose Fourier coefficients decay only at a logarithmic rate. Fortunately, these pessimistic lower bounds are not observed in practice when learning solution operators associated with PDEs. Indeed, in this context, one can exploit PDE regularity theory and Sobolev embeddings to derive quantitative bounds on the size of FNOs for approximating solution operators that only grow sublinearly with the error. Here, we refer to the analysis of Darcy flow and the two-dimensional Navier–Stokes equations by (Kovachki et al., 2021).

3.3 Deep Green networks

Deep Green networks (DGN) employ a different approach compared to DeepONets and FNOs to approximate solution operators of PDEs (Gin et al., 2021; Boullé et al., 2022a). Instead of enforcing certain properties on the integral kernel in Eq. (3), such as being low-rank (DeepONet) or translation-invariant (FNO), DGN learns the kernel directly in the physical space. Hence, assume that the underlying differential operator is a linear boundary value problem of the form

$$\mathcal{L}u = f \text{ in } \Omega, \quad \text{and} \quad u = 0 \text{ on } \partial\Omega. \quad (5)$$

Under suitable regularity assumptions on operator \mathcal{L} (e.g., uniform ellipticity or parabolicity), the solution operator \mathcal{A} can be expressed as an integral operator with a Green kernel $G : \Omega \times \Omega \rightarrow \mathbb{R} \cup \{\infty\}$ as (Evans, 2010; Boullé et al., 2022b; Boullé and Townsend, 2023)

$$\mathcal{A}(f)(x) = u(x) = \int_{\Omega} G(x, y) f(y) \, dy, \quad x \in \Omega.$$

(Boullé et al., 2022a) introduced GreenLearning networks (GL) to learn the Green kernel G directly from data. The main idea behind GL is to parameterize the kernel G as a neural network \mathcal{N} and minimize the following relative mean-squared loss function to recover an approximant to G :

$$\min_{\theta \in \mathbb{R}^N} \frac{1}{|\text{data}|} \sum_{(f, u) \in \text{data}} \frac{1}{\|u\|_{L^2(\Omega)}^2} \int_{\Omega} \left(u(x) - \int_{\Omega} \mathcal{N}(x, y) f(y) \, dy \right)^2 \, dx. \quad (6)$$

Once trained, the network \mathcal{N} can be evaluated at any point in the domain, similarly to FNO and DON. A key advantage of this method is that it provides a more interpretable model, as the kernel can be visualized and analyzed to recover properties of the underlying differential operators (Boullé et al., 2022a). However, this comes at the cost of higher computational complexity, as the integral operation in Eq. (6) must be computed accurately using a quadrature rule and typically requires $\mathcal{O}(m^2)$ operations, as opposed to the $\mathcal{O}(m \log m)$ operations required by FNOs, where m is the spatial discretization of the domain Ω .

As Green’s functions may be unbounded or singular, (Boullé et al., 2022a) propose to use a rational neural network (Boullé et al., 2020) to approximate the Green kernel. A rational neural network is a neural network whose activation functions are rational functions, defined as the ratio of two polynomials whose coefficients are learned during the training phase of the network. This choice of architecture is motivated by the fact that rational networks have higher approximation power than the standard ReLU networks (Boullé et al., 2020), in the sense that they require exponentially fewer layers to approximate continuous functions within a given accuracy and may take arbitrary large values, which is a desirable property for approximating Green kernels. A schematic diagram of a GL architecture is available in Fig. 6.

When the underlying differential operator is nonlinear, the solution operator \mathcal{A} cannot be written as an integral operator with a Green’s function. In this case, (Gin et al., 2021) propose to learn the solution operator \mathcal{A} using a dual auto-encoder architecture Deep Green network (DGN), which is a neural operator architecture that learns an invertible coordinate

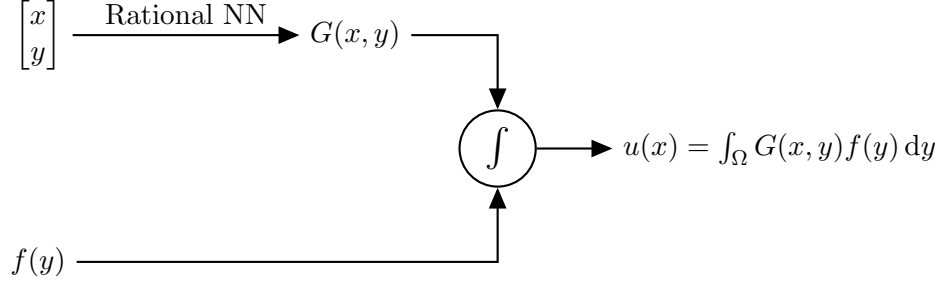


Figure 6: Schematic diagram of a GreenLearning network (GL), which approximates the integral kernel (Green’s function) associated with linear PDEs using rational neural networks.

transform map that linearizes the nonlinear boundary value problem. The resulting linear operator is approximated by a matrix, which represents a discretized Green’s function but could also be represented by a neural network if combined with the GL technique. This approach has been successfully applied to learn the solution operator of the nonlinear cubic Helmholtz equation non-Sturm–Liouville equation and discover an underlying Green’s function (Gin et al., 2021).

Other deep learning-based approaches (Lin et al., 2023; Peng et al., 2023; Sun et al., 2023) have since been introduced to recover Green’s functions using deep learning, but they rely on a PINN technique in the sense that they require the knowledge of the underlying PDE operator. Finally, (Stepaniants, 2023) proposes to learn the Green kernel associated with linear partial differential operators using a reproducible kernel Hilbert space (RKHS) framework, which leads to a convex loss function.

3.4 Graph neural operators

As described in Section 3.3, solution operators associated with linear, elliptic, or parabolic PDEs of the form $\mathcal{L}u = f$ can be written as an integral operator with a Green kernel G (Evans, 2010, Sec. 2.2.4). For simplicity, we consider Green kernels associated with uniformly elliptic operators in divergence form defined on a bounded domain Ω in spatial dimension $d \geq 3$:

$$\mathcal{L}u = -\operatorname{div}(A(x)\nabla u) = f, \quad \text{on } \Omega \subset \mathbb{R}^d, \quad (7)$$

where $A(x)$ is a bounded coefficient matrix satisfying the uniform ellipticity condition $A(x)\xi \cdot \xi \geq \lambda|\xi|^2$ for all $x \in \Omega$ and $\xi \in \mathbb{R}^d$, for some $\lambda > 0$. In this section, we present a neural operator architecture that takes advantage of the local structure of the Green kernel associated with Eq. (7), inferred by PDE regularity theory.

This architecture is called graph neural operator (GNO) (Li et al., 2020a) and is inspired by graph neural network (GNN) models (Scarselli et al., 2008; Wu et al., 2020; Zhou et al., 2020). It focuses on capturing the Green kernel’s short-range interactions to reduce the integral operation’s computational complexity in Eq. (3). The main idea behind GNO is to perform the integral operation in Eq. (5) locally on a small ball of radius r , $B(x, r)$, around

x for each $x \in \Omega$ as follows:

$$\mathcal{A}(f)(x) = u(x) \approx \int_{B(x,r)} G(x,y) f(y) dy, \quad x \in \Omega. \quad (8)$$

Here, (Li et al., 2020a) propose to discretize the domain Ω using a graph, whose nodes represent discretized spatial locations, and use a message passing network architecture (Gilmer et al., 2017) to perform an average aggregation of the nodes as in Eq. (8). The approach introduced by (Li et al., 2020a) aims to approximate the restriction G_r to the Green's function G on a band of radius r along the diagonal of the domain $\Omega \times \Omega$, defined as

$$G_r(x,y) = \begin{cases} G(x,y), & \text{if } |x-y| \leq r, \\ 0, & \text{otherwise,} \end{cases}$$

where $|\cdot|$ is the Euclidean distance in \mathbb{R}^d .

This neural architecture is justified by the following pointwise bound satisfied by the Green's function and proven by Grüter and Widman (1982, Thm. 1.1):

$$|G(x,y)| \leq C(d,A)|x-y|^{2-d}, \quad x,y \in \Omega, \quad (9)$$

where Ω is a compact domain in \mathbb{R}^d for $d \geq 3$, and C is a constant depending only on d and the coefficient matrix $A(x)$. Similar bounds have been derived in spatial dimension $d = 2$ by (Cho et al., 2012; Dong and Kim, 2009) and for Green's functions associated with time-dependent, parabolic, PDEs (Hofmann and Kim, 2004; Cho et al., 2008). Then, integrating Eq. (9) over the domain $\Gamma_r := \{(x,y) \in \Omega \times \Omega: |x-y| > r\}$ yields a bound on the approximation error between G and G_r that decays algebraically fast as r increases:

$$\begin{aligned} \|G - G_r\|_{L^2(\Omega \times \Omega)} &= \left(\int_{\Gamma_r} |G(x,y)|^2 dx dy \right)^{1/2} \leq C(d,A) \left(\int_{\Gamma_r} r^{4-2d} dx dy \right)^{1/2} \\ &\leq |\Omega| C(d,A) r^{2-d}. \end{aligned}$$

This implies that the Green's function can be well approximated by a bandlimited kernel G_r and that the approximation error bound improves in high dimensions. To illustrate this, we plot in Fig. 7 the Green's function associated with the one-dimensional Poisson equation on $\Omega = [0,1]$ with homogeneous Dirichlet boundary conditions, along with the error between the Green's function G and its truncation G_r along a bandwidth of radius r along the diagonal of the domain.

3.5 Multipole graph neural operators

Multipole graph neural operator (MGNO) has been introduced by (Li et al., 2020b) and is a class of multi-scale networks that extends the graph neural operator architecture described in Section 3.4 to capture the long-range interactions of the Green kernel. The main idea behind MGNO is to decompose the Green kernel G into a sum of low-rank kernels as $G = K_1 + \dots + K_L$, which approximates the short and wide-range interactions in the PDEs. This architecture is motivated by the same reasons that led to the development of hierarchical low-rank matrices (see Section 2.4), such as the fast multipole method (Greengard and

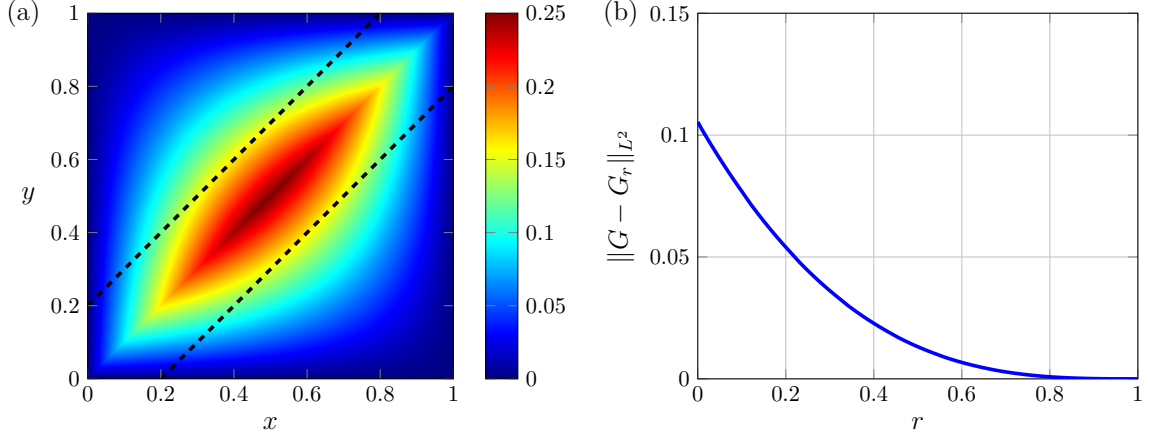


Figure 7: (a) Green's function associated with the one-dimensional Poisson equation on $\Omega = [0, 1]$ with homogeneous Dirichlet boundary conditions. The dashed lines highlight a band of radius $r = \sqrt{2}/10$ around the diagonal. (b) L^2 -norm of the error between the Green's function G and its truncation G_r along a bandwidth of radius r along the diagonal of the domain.

Rokhlin, 1997; Ying et al., 2004). It allows for the evaluation of the integral operation in Eq. (3) in linear complexity.

MGNO is based on low-rank approximations of kernels, similar to DeepONets or low-rank neural operators (see Section 3.1 and Li et al. 2020a), but is more flexible than vanilla DeepONets since it does not require the underlying kernels to be low-rank. Hence, if we consider a Green's function G associated with a uniformly elliptic PDE in the form of Eq. (7), then Weyl's law (Weyl, 1911; Canzani, 2013; Minakshisundaram and Pleijel, 1949) states that the eigenvalues of the solution operator associated with Eq. (7) decay at an algebraic rate of $\lambda_n \sim cn^{-2/d}$ for a constant $c > 0$. This implies that the approximation error between the solution operator and its best rank- k approximant decays only algebraically with k . Moreover, the decay rate deteriorates in high dimensions. In particular, the length p of the feature vector in DeepONets must be significantly large to approximate the solution operator to a prescribed accuracy.

However, Bebendorf and Hackbusch (2003, Thm. 2.8) showed that the Green's function G associated with Eq. (7) can be well approximated by a low-rank kernel when restricted to separated subdomains $D_X \times D_Y$ of $\Omega \times \Omega$, satisfying the strong admissibility condition: $\text{dist}(D_X, D_Y) < \text{diam}(D_Y)$. Here, the distance and diameter in \mathbb{R}^d are defined as

$$\text{dist}(D_X, D_Y) = \inf_{x \in D_X, y \in D_Y} |x - y|, \quad \text{diam}(D_Y) = \sup_{y_1, y_2 \in D_Y} |y_1 - y_2|.$$

Then, for any $\epsilon \in (0, 1)$, there exists a separable approximation of the form $G_k(x, y) = \sum_{i=1}^k u_i(x)v_i(y)$, with $k = \mathcal{O}(\log(1/\epsilon)^{d+1})$, such that

$$\|G - G_k\|_{L^2(D_X \times D_Y)} \leq \epsilon \|G\|_{L^2(D_X \times \hat{D}_Y)},$$

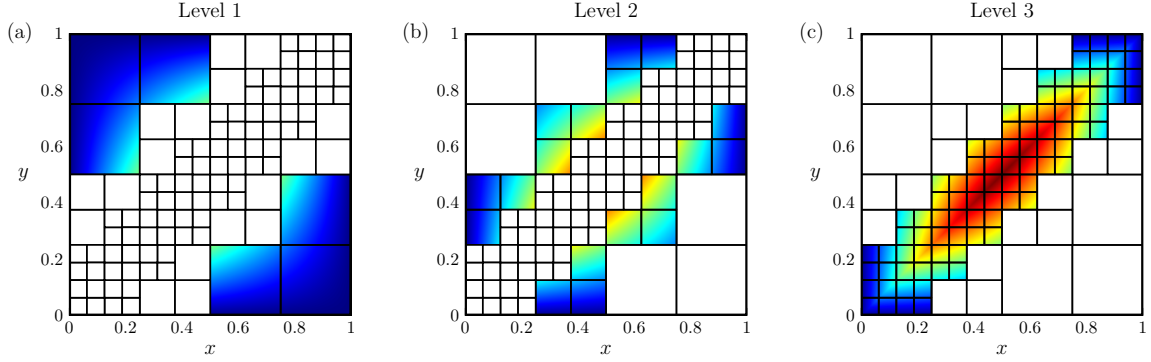


Figure 8: Decomposition of the Green’s function associated with the 1D Poisson equation, displayed in Fig. 7(a), into a hierarchy of kernels capturing different ranges of interactions: from long-range (a) to short-range (c) interactions.

where \hat{D}_Y is a domain slightly larger than D_Y (Bebendorf and Hackbusch, 2003, Thm. 2.8). This property has been exploited by (Boullé and Townsend, 2023; Boullé et al., 2023; Boullé et al., 2022b) to derive sample complexity bounds for learning Green’s functions associated with elliptic and parabolic PDEs. It motivates the decomposition of the Green kernel into a sum of low-rank kernels $G = K_1 + \dots + K_L$ in MGNO architectures. Indeed, one can exploit the low-rank structure of Green’s functions on well-separated domains to perform a hierarchical decomposition of the domain $\Omega \times \Omega$ into a tree of subdomains satisfying the admissibility condition. In Fig. 8, we illustrate the decomposition of the Green’s function associated with the 1D Poisson equation on $\Omega = [0, 1]$ with homogeneous Dirichlet boundary conditions into a hierarchy of $L = 3$ levels of different range of interactions. The first level captures the long-range interactions, while the last level captures the short-range interactions. Then, the integral operation in Eq. (3) can be performed by aggregating the contributions of the subdomains in the tree, starting from the leaves and moving up to the root. This allows for the evaluation of the integral operation in Eq. (3) in linear complexity in the number of subdomains. The key advantage is that the approximation error on each subdomain decays exponentially fast as the rank of the approximating kernel increases.

One alternative approach to MGNO is to encode the different scales of the solution operators using a wavelet basis. This class of operator learning techniques (Feliu-Faba et al., 2020; Gupta et al., 2021; Tripura and Chakraborty, 2022) is based on the wavelet transform and aims to learn the solution operator kernel at multiple scale resolutions. One advantage over MGNO is that it does not require building a hierarchy of meshes, which could be computationally challenging in high dimensions or for complex domain geometries.

Finally, motivated by the success of the self-attention mechanism in transformers architectures for natural language processing (Vaswani et al., 2017) and image recognition (Dosovitskiy et al., 2020), several architectures have been proposed to learn global correlations in solution operators of PDEs. In particular, (Cao, 2021) introduced an architecture based on the self-attention mechanism for operator learning and observed higher performance on benchmark problems when compared against the Fourier Neural Operator. More recently,

(Kissas et al., 2022) propose a Kernel-Coupled Attention mechanism to learn correlations between the entries of a vector feature representation of the output functions. In contrast, (Hao et al., 2023a) designed a general neural operator transformer (GNOT) that allows for multiple input functions and complex meshes.

4 Learning neural operators

In this section, we discuss various technical aspects involved in training neural operators, such as the data acquisition of forcing terms and solutions, the amount of training data required in practice, and the optimization algorithms and loss functions used to train neural operators.

4.1 Data acquisition

This section focuses on the data acquisition process for learning neural operators. In real-world applications, one may not have control over the distribution of source terms and solutions or locations of the sensors to measure the solutions at specific points in the domain. Therefore, we consider an idealized setting where one is interested in generating synthetic data using numerical PDE solvers to develop neural operator architectures. In this case, one has complete control over the distribution of source terms and solutions, as well as the locations of the sensors.

4.1.1 DISTRIBUTION OF SOURCE TERMS

The source terms $\{f_j\}_{j=1}^N$ used to generate pairs of training data to train neural operators are usually chosen to be random functions, sampled from a Gaussian random field (Lu et al., 2021a). Let $\Omega \subset \mathbb{R}^d$ be a domain, then a stochastic process $\{X_x, x \in \Omega\}$ indexed by Ω , is Gaussian if, for every finite set of indices $x_1, \dots, x_n \in \Omega$, the vector of random variables $(X_{x_1}, \dots, X_{x_n})$ follows a multivariate Gaussian distribution. The Gaussian process (GP) distribution is completely determined by the following mean and covariance functions (Adler, 2010, Sec. 1.6):

$$\mu(x) = \mathbb{E}\{X_x\}, \quad K(x, y) = \mathbb{E}\{[X_x - \mu(x)]^\top [X_y - \mu(y)]\}, \quad x, y \in \Omega.$$

In the rest of the paper, we will denote a Gaussian process with mean μ and covariance kernel K by $\mathcal{GP}(\mu, K)$. The mean function μ is usually chosen to be zero, while K is symmetric and positive-definite.

When K is continuous Mercer’s theorem (Mercer, 1909) states that there exists an orthonormal basis of eigenfunctions $\{\psi_j\}_{j=1}^\infty$ of $L^2(\Omega)$, and nonnegative eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots > 0$ such that

$$K(x, y) = \sum_{j=1}^{\infty} \lambda_j \psi_j(x) \psi_j(y), \quad x, y \in \Omega,$$

where the sum is absolutely and uniformly convergent (Hsing and Eubank, 2015, Thm. 4.6.5). Here, the eigenvalues and eigenfunctions of the kernel are defined as solutions to the asso-

ciated Fredholm integral equation:

$$\int_{\Omega} K(x, y) \psi_j(y) dy = \lambda_j \psi_j(x), \quad x \in \Omega.$$

Then, the Karhunen–Loève theorem (Karhunen, 1946; Loève, 1946) ensures that a zero mean square-integrable Gaussian process X_x with continuous covariance function K admits the following representation:

$$X_x = \sum_{j=1}^{\infty} \sqrt{\lambda_j} c_j \psi_j(x), \quad c_j \sim \mathcal{N}(0, 1), \quad x \in \Omega, \quad (10)$$

where c_j are independent and identically distributed (i.i.d.) random variables, and the convergence is uniform in $x \in \Omega$. Suppose the eigenvalue decomposition of the covariance function is known. In that case, one can sample a random function from the associated GP, $\mathcal{GP}(0, K)$, by sampling the coefficients c_j in Eq. (10) from a standard Gaussian distribution and truncated the series up to the desired resolution. Under suitable conditions, one can relate the covariance function K ’s smoothness to the random functions sampled from the associated GP (Adler, 2010, Sec. 3). Moreover, the decay rate of the eigenvalues provides information about the smoothness of the underlying kernel (Ritter et al., 1995; Zhu et al., 1998). In practice, the number of eigenvalues greater than machine precision dictates the dimension of the finite-dimensional vector space spanned by the random functions sampled from $\mathcal{GP}(0, K)$.

One of the most common choices of covariance functions for neural operator learning include the squared-exponential kernel (Lu et al., 2021a; Boullé et al., 2022a), which is defined as

$$K(x, y) = \exp(-|x - y|^2 / (2\ell^2)), \quad x, y \in \Omega,$$

where $\ell > 0$ is the length-scale parameter, which roughly characterizes the distance at which two point values of a sampled random function become uncorrelated (Rasmussen and Williams, 2006, Chapt. 5). Moreover, eigenvalues of the squared-exponential kernel decay exponentially fast at a rate that depends on the choice of ℓ (Zhu et al., 1998; Boullé and Townsend, 2022). After a random function f has been sampled from the GP, one typically discretizes it by performing a piecewise linear interpolation at sensor points $x_1, \dots, x_m \in \Omega$ by evaluating f at these points. The interpolant can then solve the underlying PDE or train a neural operator. The number of sensors is chosen to resolve the underlying random functions and depends on their smoothness. Following the analysis by Lu et al. (2021a, Suppl. Inf. S4), in one dimension, the error between f and its piecewise linear interpolant is of order $\mathcal{O}(1/(m^2\ell^2))$, and one should choose $m \geq 1/\ell$. A typical value of ℓ lies in the range $\ell \in [0.01, 0.1]$ with $m = 100$ sensors (Lu et al., 2021a; Boullé et al., 2022a). We illustrate the eigenvalues of the squared-exponential kernel on $\Omega = [0, 1]$ with length-scale parameters $\ell \in \{0.1, 0.05, 0.01\}$, along with the corresponding random functions sampled from the associated GP in Fig. 9. As the length-scale parameter ℓ decreases, the eigenvalues decay faster, and the sampled random functions become smoother.

Another possible choice of covariance functions for neural operator learning (Benitez et al., 2023; Zhu et al., 2023) comes from the Matérn class of covariance functions (Ras-

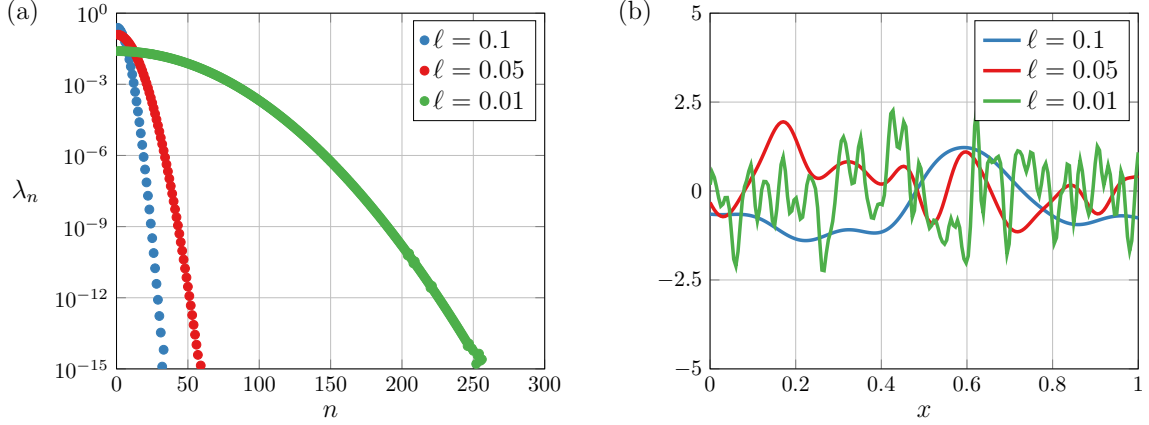


Figure 9: (a) Eigenvalues of the squared-exponential kernel on $\Omega = [0, 1]$ with length-scale parameters $\ell \in \{0.1, 0.05, 0.01\}$. (b) Random functions sampled from the associated Gaussian process $\mathcal{GP}(0, K)$, where K is the squared-exponential kernel with length-scale parameters $\ell \in \{0.1, 0.05, 0.01\}$.

mussen and Williams, 2006; Stein, 1999):

$$K(x, y) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}|x-y|}{\ell} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu}|x-y|}{\ell} \right), \quad x, y \in \Omega,$$

where Γ is the Gamma function, K_ν is a modified Bessel function, and ν, ℓ are positive parameters that enable the control of the smoothness of the sampled random functions. Hence, the resulting Gaussian process is $\lceil \nu \rceil - 1$ times mean-squared differentiable (Rasmussen and Williams, 2006, Sec. 4.2). Moreover, the Matérn kernel converges to the squared-exponential covariance function as $\nu \rightarrow \infty$. We refer to the book by Rasmussen and Williams (2006) for a detailed analysis of other standard covariance functions in Gaussian processes.

Finally, (Li et al., 2021a; Kovachki et al., 2023) propose to use a Green kernel associated with a differential operator, which is a power of the Helmholtz equation, as a covariance function:

$$K = A(-\nabla^2 + cI)^{-\nu}.$$

Here, $A, \nu > 0$, and $c \geq 0$ are parameters that respectively govern the scaling of the Gaussian process, the algebraic decay rate of the spectrum, and the frequency of the first eigenfunctions of the covariance function. One motivation for this choice of distribution is that it allows the enforcement of prior information about the underlying model, such as the order of the differential operator, directly into the source terms. A similar behavior has been observed in a randomized linear algebra context when selecting the distribution of random vectors for approximating matrices from matrix-vector products using the randomized SVD (Boullé and Townsend, 2022). For example, the eigenvalues associated with this covariance kernel decay at an algebraic rate, implying that random functions sampled from this GP would be more oscillatory. This could lead to higher performance of the

neural operator on high-frequency source terms. However, one downside of this approach is that a poor choice of the parameters can affect the training and approximation error of the neural operator. Hence, the studies (Li et al., 2021a; Kovachki et al., 2023) employ different parameter choices in each of the reported numerical experiments, suggesting that the covariance hyperparameters have been heavily optimized.

4.1.2 NUMERICAL PDE SOLVERS

After choosing the covariance function and generating the random source terms from the Gaussian process, one must solve the underlying PDE to generate the corresponding solutions. In general, the PDE is unknown, and one only has access to an oracle (such as physical experiments or black-box numerical solver) that outputs solutions u to the PDE from input source terms f as $\mathcal{L}u = f$. However, generating synthetic data from known mathematical models to design or evaluate neural operator architectures is often convenient. One can use numerical PDE solvers to generate the corresponding solutions in this case. This section briefly describes the different numerical methods that can be used to solve the underlying PDEs, along with their key attributes summarized in Table 3. We want to emphasize that these methods have many variations, and we refer to the most standard ones.

Property	Finite differences	Finite elements	Spectral methods
Domain geometry	Simple	Complex	Simple
Approximation	Local	Local	Global
Linear system	Large sparse	Large sparse	Small dense
Convergence rate	Algebraic	Algebraic	Spectral

Table 3: Summary of the different properties of standard finite difference, finite element, and spectral methods for solving PDEs.

When the PDE does not depend on time, the most common techniques for discretizing and solving it are finite difference methods (FDM), finite element methods (FEM), and spectral methods. The finite difference method consists of discretizing the computational domain Ω into a grid and approximating spatial derivatives of the solution u from linear combinations of the values of u at the grid points using finite difference operators (Iserles, 2009, Chap. 8). This approach is based on a local Taylor expansion of the solution and is very easy to implement on rectangular geometries. However, it usually requires a uniform grid approximation of the domain, which might not be appropriate for complex geometries and boundary conditions. The finite element method (Süli and Mayers, 2003, Chap. 14) employs a different approach than FDM and considers the approximation of the solution u on a finite-dimensional vector space spanned by basis functions with local support on Ω . The spatial discretization of the domain Ω is performed via a mesh representation. The basis functions are often chosen as piecewise polynomials supported on a set of elements, which are adjacent cells in the mesh. This approach is more flexible than FDM and can be used to solve PDEs on complex geometries and boundary conditions. However, it is more challenging to implement and requires the construction of a mesh of the domain Ω , which can be computationally expensive. We highlight that the finite difference and

finite element methods lead to large, sparse, and highly structured linear algebra systems, which can be solved efficiently using iterative methods. Two commonly used finite element software for generating training data for neural operators are FEniCS (Alnæs et al., 2015) and Firedrake (Rathgeber et al., 2016; Ham et al., 2023). These open-source software exploit the Unified Form Language (UFL) developed by (Alnæs et al., 2014) to define the weak form of the PDE in a similar manner as in mathematics and automatically generate the corresponding finite element assembly code before exploiting fast linear and nonlinear solvers through a Python interface with the high-performance PETSc library (Balay et al., 2023).

Finally, spectral methods (Iserles, 2009; Gottlieb and Orszag, 1977; Trefethen, 2000) are based on the approximation of the solution u on a finite-dimensional vector space spanned by basis functions with global support on Ω , usually Chebyshev polynomials or trigonometric functions. Spectral methods are motivated by the fact that the solution u of a PDE defined on a 1D interval is often smooth if the source term is smooth and so can be well-approximated by a Fourier series if it is periodic or a Chebyshev series if it is not periodic. Hence, spectral methods lead to exponential convergence, also called spectral accuracy, to analytic solutions with respect to the number of basis functions, unlike FDM and FEM, which only converge at an algebraic rate. The fast convergence rate of spectral methods implies that a small number of basis functions is usually required to achieve a given accuracy. Therefore, the matrices associated with the resulting linear algebra systems are much smaller than for FEM but are dense. In summary, spectral methods have a competitive advantage over FDM and FEM on simple geometries, such as tensor-product domains, and when the solution is smooth. At the same time, FEM might be difficult to implement but is more flexible. A convenient software for solving simple PDEs using spectral methods is the Chebfun software system (Driscoll et al., 2014), an open-source package written in MATLAB.

For time-dependent PDEs, one typically starts by performing a time-discretization using a time-stepping scheme, such as backward differentiation schemes (e.g. backward Euler) and Runge–Kutta methods (Iserles, 2009; Süli and Mayers, 2003), and then employ a spatial discretization method, such as the techniques described before in this section, to solve the resulting stationary PDE at each time-step.

4.1.3 AMOUNT OF TRAINING DATA

Current neural operator approaches typically require a relatively small amount of training data, in the order of a thousand input-output pairs, to approximate solution operators associated with PDEs (Lu et al., 2021a; Goswami et al., 2023; Kovachki et al., 2023; Boullé et al., 2023). This contrasts with the vast amount of data used to train neural networks for standard supervised learning tasks, such as image classification, which could require hundreds of millions of labeled samples (LeCun et al., 2015). This difference can be explained by the fact that solution operators are often highly structured, which can be exploited to design data-efficient neural operator architectures (see Section 3).

Recent numerical experiments have shown that the rate of convergence of neural operators with respect to the number of training samples evolves in two regimes (Lu et al., 2021a, Fig. S10). In the first regime, we observe a fast decay of the testing error at an exponential

rate (Boullé et al., 2023). Then, the testing error decays at a slower algebraic rate in the second regime for a larger amount of samples and saturates due to discretization error and optimization issues, such as convergence to a suboptimal local minima.

On the theoretical side, several works derived sample complexity bounds that characterize the amount of training data required to learn solution operators associated with certain classes of linear PDEs to within a prescribed accuracy $0 < \epsilon < 1$. In particular, (Boullé et al., 2023; Schäfer and Owhadi, 2021) focus on approximating Green’s functions associated with uniformly elliptic PDEs in divergence form:

$$-\operatorname{div}(A(x)\nabla u) = f, \quad x \in \Omega \subset \mathbb{R}^d, \quad (11)$$

where $A(x)$ is a symmetric bounded coefficient matrix (see Eq. (7)). These studies construct data-efficient algorithms that converge exponentially fast with respect to the number of training pairs. Hence, they can recover the Green’s function associated with Eq. (11) to within ϵ using only $\mathcal{O}(\text{polylog}(1/\epsilon))$ sample pairs. The method employed by (Boullé et al., 2023) consists of recovering the hierarchical low-rank structure satisfied by Green’s function on well-separated subdomains (Bebendorf and Hackbusch, 2003; Lin et al., 2011; Levitt and Martinsson, 2022b) using a generalization of the rSVD to Hilbert–Schmidt operators (Boullé and Townsend, 2022, 2023; Halko et al., 2011; Martinsson and Tropp, 2020). Interestingly, the approach by (Schäfer and Owhadi, 2021) is not based on low-rank techniques but relies on the sparse Cholesky factorization of elliptic solution operators (Schäfer et al., 2021).

Some of the low-rank recovery techniques employed by (Boullé et al., 2023) extend naturally to time-dependent parabolic PDEs (Boullé et al., 2022b) in spatial dimension $d \geq 1$ of the form:

$$\frac{\partial u}{\partial t} - \operatorname{div}(A(x, t)\nabla u) = f, \quad x \in \Omega \subset \mathbb{R}^d, \quad t \in (0, T], \quad (12)$$

where the coefficient matrix $A(x, t) \in \mathbb{R}^{d \times d}$ is symmetric positive definite with bounded coefficient functions in $L^\infty(\Omega \times [0, T])$, for some $0 < T < \infty$, and satisfies the uniform parabolicity condition (Evans, 2010, Sec. 7.1.1). Parabolic systems in the form of Eq. (12) model various time-dependent phenomena, including heat conduction and particle diffusion. Boullé et al. (2022b, Thm. 9) showed that the Green’s function associated with Eq. (12) admits a hierarchical low-rank structure on well-separated subdomains, similarly to the elliptic case (Bebendorf and Hackbusch, 2003). Combining this with the pointwise bounds satisfied by the Green’s function (Cho et al., 2012), one can construct an algorithm that recovers the Green’s function to within ϵ using $\mathcal{O}(\text{poly}(1/\epsilon))$ sample pairs (Boullé et al., 2022b, Thm. 10).

Finally, other approaches (de Hoop et al., 2023; Jin et al., 2023; Stepaniants, 2023) derived convergence rates for a broader class of operators between infinite-dimensional Hilbert spaces, which are not necessarily associated with solution operators of PDEs. In particular, (de Hoop et al., 2023) consider the problem of estimating the eigenvalues of an unknown, and possibly unbounded, self-adjoint operator assuming that the operator is diagonalizable in a known basis of eigenfunctions, and highlight the impact of varying the smoothness of training and test data on the convergence rates. Then, (Stepaniants, 2023; Jin et al., 2023) derive upper and lower bounds on the sample complexity of Hilbert–Schmidt operators between two reproducing kernel Hilbert spaces (RKHS) that depend on the smoothness of the input and output functions.

4.2 Optimization

Once the neural operator architecture has been selected and the training dataset is constituted, the next task is to train the neural operator by solving an optimization problem in the form of Eq. (1). The aim is to identify the best parameters of the underlying neural network so that the output $\hat{\mathcal{A}}(f; \theta)$ of the neural operator evaluated at a forcing term f in the training dataset fits the corresponding ground truth solution u . This section describes the most common choices of loss functions and optimization algorithms employed in current operator learning approaches. Later in Section 4.2.3, we briefly discuss how to measure the convergence and performance of a trained neural operator.

4.2.1 LOSS FUNCTIONS

The choice of loss function in operator learning is a critical step, as it directs the optimization process and ultimately affects the model’s performance. Different types of loss functions can be utilized depending on the task’s nature, the operator’s structure, and the function space’s properties. A common choice of loss function in ML is the mean squared error (MSE), which is defined as

$$L_{\text{MSE}} = \frac{1}{N} \sum_{i=1}^N \frac{1}{m} \sum_{j=1}^m |\hat{\mathcal{A}}(f_i)(x_j) - u_i(x_j)|^2 \approx \frac{1}{N} \sum_{i=1}^N \|\hat{\mathcal{A}}(f_i) - u_i\|_{L^2(\Omega)}^2, \quad (13)$$

and is employed in the original DeepONet study (Lu et al., 2021a). Here, N is the number of training samples, m is the number of sensors, f_i is the i -th forcing term, u_i is the corresponding ground truth solution, $\hat{\mathcal{A}}(f_i)$ is the output of the neural operator evaluated at f_i , and x_j is the j -th sensor location. This loss function discretizes the squared L^2 error between the output of the neural operator and the ground truth solution at the sensor locations. When the sensor grid is regular, one can employ a higher-order quadrature rule to discretize the L^2 norm. Moreover, in most cases, it may be beneficial to use a relative error, especially when the magnitudes of the output function can vary widely. Then, (Boullé et al., 2022a) use the following relative squared L^2 loss function

$$L = \frac{1}{N} \sum_{i=1}^N \frac{\|\hat{\mathcal{A}}(f_i) - u_i\|_{L^2(\Omega)}^2}{\|u_i\|_{L^2(\Omega)}^2}, \quad (14)$$

which is discretized using a trapezoidal rule. The most common loss function in operator learning is the relative L^2 error employed in Fourier neural operator techniques (Li et al., 2021a):

$$L_2 = \frac{1}{N} \sum_{i=1}^N \frac{\|\hat{\mathcal{A}}(f_i) - u_i\|_{L^2(\Omega)}}{\|u_i\|_{L^2(\Omega)}}. \quad (15)$$

(Kovachki et al., 2023) observed a better normalization of the model when using a relative loss function, and that the choice of Eq. (15) decreases the testing error by a factor of two compared to Eq. (14).

For tasks that require robustness to outliers or when it is important to measure the absolute deviation, the L^1 loss can be employed. It is defined as

$$L_1 = \frac{1}{N} \sum_{i=1}^N \frac{\|\hat{A}(f_i) - u_i\|_{L^1(\Omega)}}{\|u_i\|_{L^1(\Omega)}}. \quad (16)$$

This loss function tends to be less sensitive to large deviations than the L^2 loss (Alpak et al., 2023; Lyu et al., 2023; Zhao et al., 2024). Furthermore, Sobolev norms can also be used as a loss function when the unknown operator \mathcal{A} involves functions in Sobolev spaces (Evans, 2010, Chapt. 5), particularly when the derivatives of the input and output functions play a role (Son et al., 2021; Yu et al., 2023; O’Leary-Roseberry et al., 2024). For example, one could perform training with a relative H^1 loss to enforce the smoothness of the neural operator output. Finally, when the underlying PDE is known, one can enforce it as a weak constraint when training the neural operator by adding a PDE residual term to the loss function (Li et al., 2021b; Wang et al., 2021b), similarly to physics-informed neural networks (Raissi et al., 2019).

4.2.2 OPTIMIZATION ALGORITHMS AND IMPLEMENTATION

The training procedure of neural operators is typically performed using Adam optimization algorithm (Kingma and Ba, 2015; Kovachki et al., 2023; Lu et al., 2021a; Li et al., 2021a; Goswami et al., 2023) or one of its variants such as AdamW (Loshchilov and Hutter, 2019; Hao et al., 2023a). Hence, the work introducing DeepONets by (Lu et al., 2021a) employed Adam algorithm to train the neural network architecture with a default learning rate of 0.001. In contrast, (Kovachki et al., 2023) incorporate learning rate decay throughout the optimization of Fourier neural operators. One can also employ a two-step training approach by minimizing the loss function using Adam algorithm for a fixed number of iterations and then fine-tuning the neural operator using the L-BFGS algorithm (Byrd et al., 1995; Cuomo et al., 2022; He et al., 2020; Mao et al., 2020; Boull   et al., 2022a). This approach has been shown to improve the convergence rate of the optimization when little data is available in PINNs applications (He et al., 2020). Popular libraries for implementing and training neural operators include PyTorch (Paszke et al., 2019) and TensorFlow (Abadi et al., 2016).

Thus far, there has been limited focus on the theoretical understanding of convergence and optimization of neural operators. Since neural operators are a generalization of neural networks in infinite dimensions, existing convergence results of physics-informed neural networks (Wang et al., 2021a, 2022b) based on the neural tangent kernel (NTK) framework (Jacot et al., 2018; Du et al., 2019; Allen-Zhu et al., 2019) should naturally extend to neural operators. One notable exception is the study by (Wang et al., 2022a), which analyzes the training of physics-informed DeepONets (Wang et al., 2021b) and derives a weighting scheme guided by NTK theory to balance the data and the PDE residual terms in the loss function.

4.2.3 MEASURING CONVERGENCE AND SUPER-RESOLUTION

After training a neural operator, one typically measures its performance by evaluating the testing error, such as the relative L^2 -error, on a set of unseen data generated using the

procedure described in Section 4. In general, state-of-the-art neural operator architectures report a relative testing error of 1% – 10% depending on the problems considered (Kovachki et al., 2023; Lu et al., 2021a; Li et al., 2021a).

However, it is essential to note that the testing error may not be a good measure of the performance of a neural operator, as it does not provide any information about the generalization properties of the model. Hence, the testing forcing terms are usually sampled from the same distribution as the training forcing terms, so they lie on the same finite-dimensional function space, determined by the spectral decay of the GP covariance kernel eigenvalues (see Section 4.1). Moreover, in real applications, the testing source terms could have different distributions than the ones used for training ones, and extrapolation of the neural operator may be required. (Zhu et al., 2023) investigate the extrapolation properties of DeepONet with respect to the length-scale parameter of the underlying source term GP. They observe that the testing error increases when the length-scale parameter corresponding to the test data decreases. At the same time, the neural operator can extrapolate to unseen data with a larger length scale than the training dataset, i.e., smoother functions.

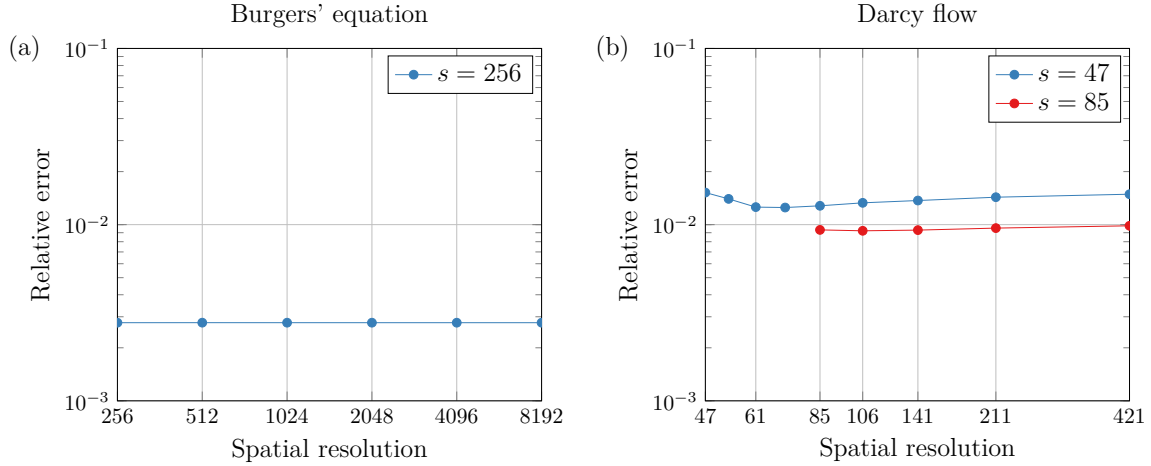


Figure 10: (a) Relative test error at different spatial resolutions of the Fourier neural operator trained to approximate the solution operator of the 1D Burgers' equation (17) with trained resolution of $s = 256$. (b) FNO trained on 2D Darcy flow a resolution of $s = 47$ (blue) and $s = 85$ (red) and evaluated at higher spatial resolutions.

One attractive property of neural operators is their resolution invariance to perform predictions at finer spatial resolutions than the training dataset on which they have been trained. This is usually called zero-shot super-resolution (Kovachki et al., 2023, Sec. 7.2.3). To investigate this property, we reproduce the numerical examples of (Kovachki et al., 2023, Sec. 7.2) and train a Fourier neural operator to approximate the solution operator of Burgers' equation and Darcy flow at a low-resolution data and evaluate the operator at

higher resolutions. We consider the one-dimensional Burgers' equation:

$$\frac{\partial}{\partial t}u(x,t) + \frac{1}{2}\frac{\partial}{\partial x}(u(x,t)^2) = \nu \frac{\partial^2}{\partial x^2}u(x,t), \quad x \in (0, 2\pi), \quad t \in [0, 1], \quad (17)$$

with periodic boundary conditions and viscosity $\nu = 0.1$. We are interested in learning the solution operator $\mathcal{A} : L^2_{\text{per}}((0, 2\pi)) \rightarrow H^1_{\text{per}}((0, 2\pi))$, which maps initial conditions $u_0 \in L^2_{\text{per}}((0, 2\pi))$ to corresponding solutions $u(\cdot, 1) \in H^1_{\text{per}}((0, 2\pi))$ to Eq. (17) at time $t = 1$. We then discretize the source and solution training data on a uniform grid with spatial resolution $s = 256$ and evaluate the trained neural operator at finer spatial resolutions $s \in \{512, 1024, 2048\}$. We observe in Fig. 10(a) that the relative testing error of the neural operator is independent of the spatial resolutions, as reported by Kovachki et al. (2023, Sec. 7.2).

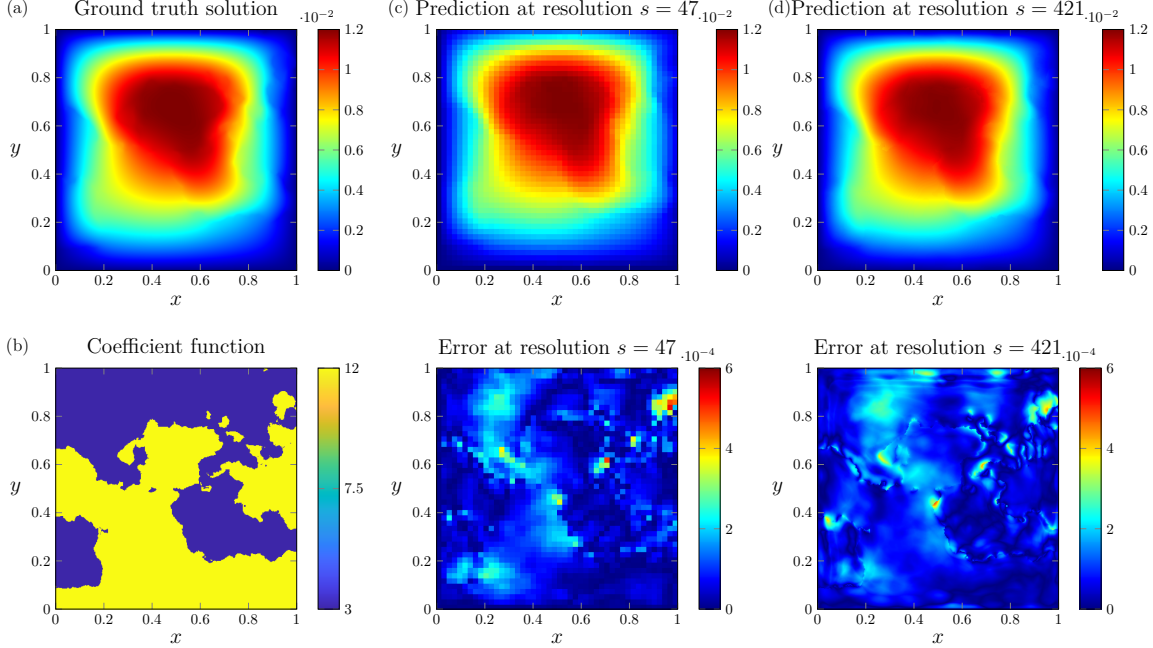


Figure 11: (a) Ground truth solution u to the 2D Darcy flow equation (18) corresponding to the coefficients function plotted in (b). (c)-(d) Predicted solution and approximation error at $s = 47$ and $s = 421$ by a Fourier neural operator trained on a Darcy flow dataset with spatial resolution of $s = 47$.

Next, we consider the two-dimensional Darcy flow equation (18) with constant source term $f = 1$ and homogeneous Dirichlet boundary conditions on a unit square domain $\Omega = [0, 1]^2$:

$$-\operatorname{div}(a(x)\nabla u) = 1, \quad x \in [0, 1]^2. \quad (18)$$

We train a Fourier neural operator to approximate the solution operator, mapping the coefficient function a to the associated solution u to Eq. (18). We reproduce the numerical

experiment in (Kovachki et al., 2023, Sec. 6.2), where the random coefficient functions a are piecewise constant. The random functions a are generated as $a \sim T \circ f$, where $f \sim \mathcal{GP}(0, C)$, with $C = (-\Delta + 9I)^{-2}$ and $T : \mathbb{R} \rightarrow \mathbb{R}^+$ is defined as

$$T(x) = \begin{cases} 12, & \text{if } x \geq 0, \\ 3, & \text{if } x < 0. \end{cases}$$

We discretize the coefficient and solution training data on a $s \times s$ uniform grid with spatial resolution $s = 47$ and evaluate the trained neural operator at higher spatial resolutions in Fig. 10. The relative testing error does not increase as the spatial resolution increases. Moreover, training the neural operator on a higher spatial resolution dataset can decrease the testing error. We also plot the ground truth solution u to Eq. (18) in Fig. 11(a) corresponding to the coefficient function plotted in panel (b), along with the predicted solutions and approximation errors at $s = 47$ and $s = 421$ by the Fourier neural operator in panels (c) and (d). We want to point the reader interested in the discretization properties of neural operators to the recent perspective on representation equivalent neural operators (ReNO) by (Bartolucci et al., 2023).

5 Conclusions and future challenges

In this paper, we provided a comprehensive overview of the recent developments in neural operator learning, a new paradigm at the intersection of scientific computing and ML for learning solution operators of PDEs. Given the recent surge of interest in this field, a key question concerns the choice of neural architectures for different PDEs. Most theoretical studies in the field analyze and compare neural operators through the prism of approximation theory. We proposed a framework based on numerical linear algebra and matrix recovery problems for interpreting the type of neural operator architectures that can be used to learn solution operators of PDEs. Hence, solution operators associated with linear PDEs can often be written as integral operators with a Green’s function and recover by a one-layer neural operator, which after discretization is equivalent to a matrix recovery problem.

Moreover, the choice of architectures, such as FNO or DeepONet, enforces or preserves different properties of the PDE solution operator, such as being translation invariant, low-rank, or off-diagonal low-rank (see Table 2). We then focused on the data acquisition process. We highlighted the importance of the distribution of source terms, usually sampled from a Gaussian process with a tailored covariance kernel, on the resulting performance of the neural operator. Following recent works on elliptic and parabolic PDEs and numerical experiments, we also discussed the relatively small amount of training data needed for operator learning. Finally, we studied the different choices of optimization algorithms and loss functions and highlighted the super-resolution properties of neural operators, i.e., their ability to be evaluated at higher resolution than the training dataset with a minor impact on the performance. There are, however, several remaining challenges in the field.

Distribution of probes. Most applications of neural operators employ source terms that are globally supported on the domain, sampled from a Gaussian process, and whose distribution is fixed before training. However, this might not always apply to real-world engineering

or biological systems, where source terms could be localized in space and time. A significant problem is to study the impact of the distribution of locally supported source terms on the performance of neural operators, both from a practical and theoretical viewpoint. Hence, recent sample complexity works on elliptic and parabolic PDEs exploit structured source terms (Boullé et al., 2023; Boullé et al., 2022b; Schäfer and Owhadi, 2021). Another area of future research is to employ adaptive source terms to fine-tune neural operators for specific applications. This could lead to higher performance by selecting source terms that maximize the training error or allow efficient transfer learning between different applications without retraining a large neural operator.

Software and datasets. An essential step towards democratizing operator learning involves the development of open-source software and datasets for training and comparing neural operators, similar to the role played by the MNIST (LeCun et al., 1998) and ImageNet (Deng et al., 2009) databases in the improvement of computer vision techniques. However, due to the fast emerging methods in operator learning, there have been limited attempts beyond (Lu et al., 2022) to standardize the datasets and software used in the field. Establishing a list of standard PDE problems across different scientific fields, such as fluid dynamics, quantum mechanics, and epidemiology, with other properties (e.g. linear/nonlinear, steady/time-dependent, low/high dimensional, smooth/rough solutions, simple/complex geometry) would allow researchers to compare and identify the neural operator architectures that are the most appropriate for a particular task. A recent benchmark has been proposed to evaluate the performance of physics-informed neural networks for solving PDEs (Hao et al., 2023b).

Real-world applications. Neural operators have been successfully applied to perform weather forecasting and achieve spectacular performance in terms of accuracy and computational time to solutions compared to traditional numerical weather prediction techniques while being trained on historical weather data (Kurth et al., 2023; Lam et al., 2023). An exciting development in the field of operator learning would be to expand the scope of applications to other scientific fields and train the models on real datasets, where the underlying PDE governing the data is unknown to discover new physics.

Theoretical understanding. Following the recent works on the approximation theory of neural operators and sample complexity bounds for different classes of PDEs, there is a growing need for a theoretical understanding of convergence and optimization. In particular, an exciting area of research would be to extend the convergence results of physics-informed neural networks and the neural tangent kernel framework to neural operators. This would enable the derivation of rigorous convergence rates for different types of neural operator architectures and loss functions and new schemes for initializing the weight distributions in the underlying neural networks.

Physical properties. Most neural operator architectures are motivated by obtaining a good approximation of the solution operator of a PDE. However, the resulting neural operator is often highly nonlinear, difficult to interpret mathematically, and might not satisfy the physical properties of the underlying PDE, such as conservation laws or symmetries (Olver, 1993b). There are several promising research directions in operator learning related to symmetries and conservation laws (Otto et al., 2023). One approach would be to enforce

known physical properties when training neural operators, either strongly through structure preserving architectures (Richter-Powell et al., 2022), or weakly by adding a residual term in the loss function (Li et al., 2021b; Wang et al., 2021b). Another direction is to discover new physical properties of the underlying PDEs from the trained neural operator. While (Boullé et al., 2022a) showed that symmetries of linear PDEs can be recovered from the learned Green’s function, this approach has not been extended to nonlinear PDEs. Finally, one could also consider using reinforcement learning techniques for enforcing physical constraints after the optimization procedure, similar to recent applications in large language models (Ouyang et al., 2022).

Acknowledgments and Disclosure of Funding

The work of both authors was supported by the Office of Naval Research (ONR), under grant N00014-23-1-2729. N.B. was supported by an INI-Simons Postdoctoral Research Fellowship. A.T. was supported by National Science Foundation grants DMS-2045646 and a Weiss Junior Fellowship Award.

References

- M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th USENIX Symposium on Operating Systems Design and Implementation*, pages 265–283, 2016.
- R. J. Adler. *The geometry of random fields*. SIAM, 2010.
- Z. Allen-Zhu, Y. Li, and Z. Song. A convergence theory for deep learning via over-parameterization. In *International Conference on Machine Learning*, pages 242–252, 2019.
- M. Alnæs, J. Blechta, J. Hake, A. Johansson, B. Kehlet, A. Logg, C. Richardson, J. Ring, M. E. Rognes, and G. N. Wells. The FEniCS project version 1.5. *Arch. Numer. Softw.*, 3(100), 2015.
- M. S. Alnæs, A. Logg, K. B. Ølgaard, M. E. Rognes, and G. N. Wells. Unified form language: A domain-specific language for weak formulations of partial differential equations. *ACM Trans. Math. Softw.*, 40(2):1–37, 2014.
- F. O. Alpak, J. Vamaraju, J. W. Jennings, S. Pawar, P. Devarakota, and D. Hohl. Augmenting Deep Residual Surrogates with Fourier Neural Operators for Rapid Two-Phase Flow and Transport Simulations. *SPE J.*, pages 1–22, 2023.
- S. Arridge, P. Maass, O. Öktem, and C.-B. Schönlieb. Solving inverse problems using data-driven models. *Acta Numer.*, 28:1–174, 2019.

- S. Balay, S. Abhyankar, M. F. Adams, J. Brown, P. Brune, K. Buschelman, L. Dalcin, V. Eijkhout, D. Kaushik, M. G. Knepley, L. C. McInnes, W. D. Gropp, K. Rupp, B. F. Smith, S. Zampini, H. Zhang, and H. Z. *PETSc Users Manual*. Argonne National Laboratory, 2023.
- F. Bartolucci, E. de Bézenac, B. Raonić, R. Molinaro, S. Mishra, and R. Alaifari. Are Neural Operators Really Neural Operators? Frame Theory Meets Operator Learning. *arXiv preprint arXiv:2305.19913*, 2023.
- M. Bebendorf and W. Hackbusch. Existence of \mathcal{H} -matrix approximants to the inverse FE-matrix of elliptic operators with L^∞ -coefficients. *Numer. Math.*, 95(1):1–28, 2003.
- J. A. L. Benitez, T. Furuya, F. Faucher, A. Kratsios, X. Tricoche, and M. V. de Hoop. Out-of-distributional risk bounds for neural operators with applications to the Helmholtz equation. *arXiv preprint arXiv:2301.11509*, 2023.
- K. Bhattacharya, B. Hosseini, N. B. Kovachki, and A. M. Stuart. Model reduction and neural networks for parametric PDEs. *SMAI J. Comput. Math.*, 7:121–157, 2021.
- N. Boullé and A. Townsend. A generalization of the randomized singular value decomposition. In *International Conference on Learning Representations*, 2022.
- N. Boullé and A. Townsend. Learning elliptic partial differential equations with randomized linear algebra. *Found. Comput. Math.*, 23(2):709–739, 2023.
- N. Boullé, Y. Nakatsukasa, and A. Townsend. Rational neural networks. In *Advances in Neural Information Processing Systems*, volume 33, pages 14243–14253, 2020.
- N. Boullé, C. J. Earls, and A. Townsend. Data-driven discovery of Green’s functions with human-understandable deep learning. *Sci. Rep.*, 12(1):4824, 2022a.
- N. Boullé, S. Kim, T. Shi, and A. Townsend. Learning Green’s functions associated with time-dependent partial differential equations. *J. Mach. Learn. Res.*, 23(218):1–34, 2022b.
- N. Boullé, D. Halikias, and A. Townsend. Elliptic PDE learning is provably data-efficient. *Proc. Natl. Acad. Sci. USA*, 120(39):e2303904120, 2023.
- M. M. Bronstein, J. Bruna, T. Cohen, and P. Velićković. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges. *arXiv preprint arXiv:2104.13478*, 2021.
- T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.
- O. P. Bruno, Y. Han, and M. M. Pohlman. Accurate, high-order representation of complex three-dimensional surfaces via Fourier continuation analysis. *J. Comput. Phys.*, 227(2):1094–1125, 2007.
- S. L. Brunton, J. L. Proctor, and J. N. Kutz. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proc. Natl. Acad. Sci. U.S.A.*, 113(15):3932–3937, 2016.

- R. H. Byrd, P. Lu, J. Nocedal, and C. Zhu. A limited memory algorithm for bound constrained optimization. *SIAM J. Sci. Comput.*, 16(5):1190–1208, 1995.
- Y. Canzani. Analysis on manifolds via the Laplacian. Harvard University, 2013.
- S. Cao. Choose a transformer: Fourier or Galerkin. In *Advances in Neural Information Processing Systems*, volume 34, pages 24924–24940, 2021.
- K. Champion, B. Lusch, J. N. Kutz, and S. L. Brunton. Data-driven discovery of coordinates and governing equations. *Proc. Natl. Acad. Sci. U.S.A.*, 116(45):22445–22451, 2019.
- T. Chen and H. Chen. Universal approximation to nonlinear operators by neural networks with arbitrary activation functions and its application to dynamical systems. *IEEE Trans. Neur. Netw.*, 6(4):911–917, 1995.
- S. Cho, H. Dong, and S. Kim. On the Green’s matrices of strongly parabolic systems of second order. *Indiana Univ. Math. J.*, 57(4):1633–1677, 2008.
- S. Cho, H. Dong, and S. Kim. Global estimates for Green’s matrix of second order parabolic systems with application to elliptic systems in two dimensional domains. *Potential Anal.*, 36(2):339–372, 2012.
- J. W. Cooley and J. W. Tukey. An algorithm for the machine calculation of complex Fourier series. *Math. Comput.*, 19(90):297–301, 1965.
- S. Cuomo, V. S. Di Cola, F. Giampaolo, G. Rozza, M. Raissi, and F. Piccialli. Scientific machine learning through physics-informed neural networks: Where we are and what’s next. *J. Sci. Comput.*, 92(3):88, 2022.
- G. Cybenko. Approximation by superpositions of a sigmoidal function. *Math. Control Signals Syst.*, 2(4):303–314, 1989.
- M. V. de Hoop, D. Z. Huang, E. Qian, and A. M. Stuart. The cost-accuracy trade-off in operator learning with neural networks. *arXiv preprint arXiv:2203.13181*, 2022.
- M. V. de Hoop, N. B. Kovachki, N. H. Nelsen, and A. M. Stuart. Convergence rates for learning linear operators from noisy data. *SIAM-ASA J. Uncertain. Quantif.*, 11(2):480–513, 2023.
- T. De Ryck and S. Mishra. Generic bounds on the approximation error for physics-informed (and) operator learning. In *Advances in Neural Information Processing Systems*, volume 35, pages 10945–10958, 2022.
- B. Deng, Y. Shin, L. Lu, Z. Zhang, and G. E. Karniadakis. Approximation rates of DeepONets for learning operators arising from advection–diffusion equations. *Neur. Netw.*, 153:411–426, 2022.
- J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Conference on Computer Vision and Pattern Recognition*, pages 248–255. IEEE, 2009.

- R. A. DeVore. Nonlinear approximation. *Acta Numer.*, 7:51–150, 1998.
- P. C. Di Leoni, L. Lu, C. Meneveau, G. E. Karniadakis, and T. A. Zaki. Neural operator prediction of linear instability waves in high-speed boundary layers. *J. Comput. Phys.*, 474:111793, 2023.
- H. Dong and S. Kim. Green’s matrices of second order elliptic systems with measurable coefficients in two dimensional domains. *Trans. Amer. Math. Soc.*, 361(6):3303–3323, 2009.
- A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- T. A. Driscoll, N. Hale, and L. N. Trefethen. *Chebfun Guide*. Pafnuty Publications, 2014. URL <http://www.chebfun.org/docs/guide/>.
- S. Du, J. Lee, H. Li, L. Wang, and X. Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, pages 1675–1685, 2019.
- W. E and B. Yu. The deep Ritz method: a deep learning-based numerical algorithm for solving variational problems. *Commun. Math. Stat.*, 6(1):1–12, 2018.
- L. C. Evans. *Partial differential equations*. American Mathematical Society, 2nd edition, 2010.
- V. Fanaskov and I. Oseledets. Spectral neural operators. *arXiv preprint arXiv:2205.10573*, 2022.
- J. Feliu-Faba, Y. Fan, and L. Ying. Meta-learning pseudo-differential operators with deep neural networks. *J. Comput. Phys.*, 408:109309, 2020.
- J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing for quantum chemistry. In *International Conference on Machine Learning*, pages 1263–1272, 2017.
- C. R. Gin, D. E. Shea, S. L. Brunton, and J. N. Kutz. DeepGreen: deep learning of Green’s functions for nonlinear boundary value problems. *Sci. Rep.*, 11(1):1–14, 2021.
- S. Goswami, M. Yin, Y. Yu, and G. E. Karniadakis. A physics-informed variational deeponet for predicting crack path in quasi-brittle materials. *Comput. Methods Appl. Mech. Eng.*, 391:114587, 2022.
- S. Goswami, A. Bora, Y. Yu, and G. E. Karniadakis. Physics-informed deep neural operator networks. In *Machine Learning in Modeling and Simulation: Methods and Applications*, pages 219–254. Springer, 2023.
- D. Gottlieb and S. A. Orszag. *Numerical analysis of spectral methods: theory and applications*. SIAM, 1977.

- L. Greengard and V. Rokhlin. A new version of the Fast Multipole Method for the Laplace equation in three dimensions. *Acta Numer.*, 6:229–269, 1997.
- M. Grüter and K.-O. Widman. The Green function for uniformly elliptic equations. *Manuscripta Math.*, 37(3):303–342, 1982.
- G. Gupta, X. Xiao, and P. Bogdan. Multiwavelet-based Operator Learning for Differential Equations. In *Advances in Neural Information Processing Systems*, volume 34, pages 24048–24062, 2021.
- W. Hackbusch, B. N. Khoromskij, and R. Kriemann. Hierarchical matrices based on a weak admissibility criterion. *Computing*, 73(3):207–243, 2004.
- D. Halikias and A. Townsend. Structured matrix recovery from matrix-vector products. *Numer. Linear Algebra Appl.*, page e2531, 2023.
- N. Halko, P.-G. Martinsson, and J. A. Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Rev.*, 53(2):217–288, 2011.
- D. A. Ham, P. H. J. Kelly, L. Mitchell, C. J. Cotter, R. C. Kirby, K. Sagiya, N. Bouziani, S. Vorderwuelbecke, T. J. Gregory, J. Betteridge, D. R. Shapero, R. W. Nixon-Hill, C. J. Ward, P. E. Farrell, P. D. Brubeck, I. Marsden, T. H. Gibson, M. Homolya, T. Sun, A. T. T. McRae, F. Luporini, A. Gregory, M. Lange, S. W. Funke, F. Rathgeber, G.-T. Bercea, and G. R. Markall. *Firedrake User Manual*. Imperial College London and University of Oxford and Baylor University and University of Washington, 1st edition, 2023.
- Z. Hao, Z. Wang, H. Su, C. Ying, Y. Dong, S. Liu, Z. Cheng, J. Song, and J. Zhu. GNOT: A general neural operator transformer for operator learning. In *International Conference on Machine Learning*, pages 12556–12569, 2023a.
- Z. Hao, J. Yao, C. Su, H. Su, Z. Wang, F. Lu, Z. Xia, Y. Zhang, S. Liu, L. Lu, et al. PIN-Nacle: A Comprehensive Benchmark of Physics-Informed Neural Networks for Solving PDEs. *arXiv preprint arXiv:2306.08827*, 2023b.
- Q. He, D. Barajas-Solano, G. Tartakovsky, and A. M. Tartakovsky. Physics-informed neural networks for multiphysics data assimilation with application to subsurface transport. *Adv. Water Resour.*, 141:103610, 2020.
- J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851, 2020.
- S. Hofmann and S. Kim. Gaussian estimates for fundamental solutions to certain parabolic systems. *Publ. Mat.*, pages 481–496, 2004.
- K. Hornik. Approximation capabilities of multilayer feedforward networks. *Neur. Netw.*, 4(2):251–257, 1991.

- T. Hsing and R. Eubank. *Theoretical Foundations of Functional Data Analysis, with an Introduction to Linear Operators*. John Wiley & Sons, 2015.
- A. Iserles. *A first course in the numerical analysis of differential equations*. Cambridge University Press, 2009.
- A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, volume 31, 2018.
- J. Jin, Y. Lu, J. Blanchet, and L. Ying. Minimax optimal kernel operator learning via multilevel training. In *International Conference on Learning Representations*, 2023.
- J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Židek, A. Potapenko, et al. Highly accurate protein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, 2021.
- K. Karhunen. Über lineare methoden in der wahrscheinlichkeitsrechnung. *Ann. Acad. Science Fenn., Ser. A. I.*, 37:3–79, 1946.
- G. E. Karniadakis, I. G. Kevrekidis, L. Lu, P. Perdikaris, S. Wang, and L. Yang. Physics-informed machine learning. *Nat. Rev. Phys.*, 3(6):422–440, 2021.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *Proc. 3rd International Conference on Learning Representation*, 2015.
- G. Kissas, J. H. Seidman, L. F. Guilhoto, V. M. Preciado, G. J. Pappas, and P. Perdikaris. Learning operators with coupled attention. *J. Mach. Learn. Res.*, 23(1):9636–9698, 2022.
- N. Kovachki, S. Lanthaler, and S. Mishra. On universal approximation and error bounds for Fourier Neural Operators. *J. Mach. Learn. Res.*, 22:1–76, 2021.
- N. Kovachki, Z. Li, B. Liu, K. Azizzadenesheli, K. Bhattacharya, A. Stuart, and A. Anandkumar. Neural operator: Learning maps between function spaces with applications to PDEs. *J. Mach. Learn. Res.*, 24(89):1–97, 2023.
- A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, volume 25, 2012.
- T. Kurth, S. Subramanian, P. Harrington, J. Pathak, M. Mardani, D. Hall, A. Miele, K. Kashinath, and A. Anandkumar. Fourcastnet: Accelerating global high-resolution weather forecasting using adaptive fourier neural operators. In *Proceedings of the Platform for Advanced Scientific Computing Conference*, pages 1–11, 2023.
- R. Lam, A. Sanchez-Gonzalez, M. Willson, P. Wirnsberger, M. Fortunato, F. Alet, S. Ravuri, T. Ewalds, Z. Eaton-Rosen, W. Hu, et al. Learning skillful medium-range global weather forecasting. *Science*, page eadi2336, 2023.
- S. Lanthaler, S. Mishra, and G. E. Karniadakis. Error estimates for DeepONets: A deep learning framework in infinite dimensions. *Trans. Math. Appl.*, 6(1), 2022.

- Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proc. IEEE*, 86(11):2278–2324, 1998.
- Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- J. Levitt and P.-G. Martinsson. Linear-complexity black-box randomized compression of hierarchically block separable matrices. *arXiv preprint arXiv:2205.02990*, 2022a.
- J. Levitt and P.-G. Martinsson. Randomized compression of rank-structured matrices accelerated with graph coloring. *arXiv preprint arXiv:2205.03406*, 2022b.
- Z. Li, N. Kovachki, K. Azizzadenesheli, B. Liu, K. Bhattacharya, A. Stuart, and A. Anandkumar. Neural operator: Graph kernel network for partial differential equations. *arXiv preprint arXiv:2003.03485*, 2020a.
- Z. Li, N. Kovachki, K. Azizzadenesheli, B. Liu, A. Stuart, K. Bhattacharya, and A. Anandkumar. Multipole graph neural operator for parametric partial differential equations. In *Advances in Neural Information Processing Systems*, volume 33, pages 6755–6766, 2020b.
- Z. Li, N. Kovachki, K. Azizzadenesheli, B. Liu, K. Bhattacharya, A. Stuart, and A. Anandkumar. Fourier neural operator for parametric partial differential equations. In *International Conference on Learning Representations*, 2021a.
- Z. Li, H. Zheng, N. Kovachki, D. Jin, H. Chen, B. Liu, K. Azizzadenesheli, and A. Anandkumar. Physics-informed neural operator for learning partial differential equations. *arXiv preprint arXiv:2111.03794*, 2021b.
- Z. Li, D. Z. Huang, B. Liu, and A. Anandkumar. Fourier neural operator with learned deformations for PDEs on general geometries. *arXiv preprint arXiv:2207.05209*, 2022.
- Z. Li, N. B. Kovachki, C. Choy, B. Li, J. Kossaifi, S. P. Otta, M. A. Nabian, M. Stadler, C. Hundt, K. Azizzadenesheli, et al. Geometry-Informed Neural Operator for Large-Scale 3D PDEs. *arXiv preprint arXiv:2309.00583*, 2023a.
- Z. Li, W. Peng, Z. Yuan, and J. Wang. Long-term predictions of turbulence by implicit U-Net enhanced Fourier neural operator. *Phys. Fluids*, 35(7), 2023b.
- G. Lin, F. Chen, P. Hu, X. Chen, J. Chen, J. Wang, and Z. Shi. BI-GreenNet: learning Green’s functions by boundary integral network. *Commun. Math. Stat.*, 11(1):103–129, 2023.
- L. Lin, J. Lu, and L. Ying. Fast construction of hierarchical matrix representation from matrix–vector multiplication. *J. Comput. Phys.*, 230(10):4071–4087, 2011.
- M. Loève. Fonctions aleatoire de second ordre. *Rev. Sci.*, 84:195–206, 1946.
- I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.

- L. Lu, P. Jin, G. Pang, Z. Zhang, and G. E. Karniadakis. Learning nonlinear operators via DeepONet based on the universal approximation theorem of operators. *Nat. Mach. Intell.*, 3(3):218–229, 2021a.
- L. Lu, X. Meng, Z. Mao, and G. E. Karniadakis. DeepXDE: A deep learning library for solving differential equations. *SIAM Rev.*, 63(1):208–228, 2021b.
- L. Lu, X. Meng, S. Cai, Z. Mao, S. Goswami, Z. Zhang, and G. E. Karniadakis. A comprehensive and fair comparison of two neural operators (with practical extensions) based on fair data. *Comput. Methods Appl. Mech. Eng.*, 393:114778, 2022.
- Y. Lyu, X. Zhao, Z. Gong, X. Kang, and W. Yao. Multi-fidelity prediction of fluid flow based on transfer learning using Fourier neural operator. *Phys. Fluids*, 35(7), 2023.
- S. Mao, R. Dong, L. Lu, K. M. Yi, S. Wang, and P. Perdikaris. PPDONet: Deep Operator Networks for Fast Prediction of Steady-state Solutions in Disk–Planet Systems. *Astrophys. J. Lett.*, 950(2):L12, 2023.
- Z. Mao, A. D. Jagtap, and G. E. Karniadakis. Physics-informed neural networks for high-speed flows. *Comput. Methods Appl. Mech. Eng.*, 360:112789, 2020.
- P.-G. Martinsson. A fast randomized algorithm for computing a hierarchically semiseparable representation of a matrix. *SIAM J. Matrix Anal. Appl.*, 32(4):1251–1274, 2011.
- P.-G. Martinsson and J. A. Tropp. Randomized numerical linear algebra: Foundations and algorithms. *Acta Numer.*, 29:403–572, 2020.
- M. Mathieu, M. Henaff, and Y. LeCun. Fast training of convolutional networks through FFTs. In *International Conference on Learning Representations*, 2014.
- J. Mercer. Functions of positive and negative type, and their connection with the theory of integral equations. *Philos. Trans. Royal Soc. A*, 209:415–446, 1909.
- S. Minakshisundaram and Å. Pleijel. Some properties of the eigenfunctions of the Laplace-operator on Riemannian manifolds. *Canad. J. Math.*, 1(3):242–256, 1949.
- C. Moya, S. Zhang, G. Lin, and M. Yue. Deeponet-grid-uq: A trustworthy deep operator framework for predicting the power grid’s post-fault trajectories. *Neurocomputing*, 535:166–182, 2023.
- T. O’Leary-Roseberry, P. Chen, U. Villa, and O. Ghattas. Derivative-informed neural operator: an efficient framework for high-dimensional parametric derivative learning. *J. Comput. Phys.*, 496:112555, 2024.
- P. J. Olver. *Applications of Lie groups to differential equations*. Springer Science & Business Media, 1993a.
- P. J. Olver. *Applications of Lie groups to differential equations*. Springer-Verlag, 2nd edition, 1993b.

- S. E. Otto, N. Zolman, J. N. Kutz, and S. L. Brunton. A unified framework to enforce, discover, and promote symmetry in machine learning. *arXiv preprint arXiv:2311.00212*, 2023.
- L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744, 2022.
- A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems*, 32, 2019.
- R. Peng, J. Dong, J. Malof, W. J. Padilla, and V. Tarokh. Deep Generalized Green’s Functions. *arXiv preprint arXiv:2306.02925*, 2023.
- W. Peng, Z. Yuan, and J. Wang. Attention-enhanced neural network models for turbulence simulation. *Phys. Fluids*, 34(2), 2022.
- M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *J. Comput. Phys.*, 378:686–707, 2019.
- B. Raonic, R. Molinaro, T. Rohner, S. Mishra, and E. de Bezenac. Convolutional neural operators. In *ICLR 2023 Workshop on Physics for Machine Learning*, 2023.
- C. E. Rasmussen and C. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- F. Rathgeber, D. A. Ham, L. Mitchell, M. Lange, F. Luporini, A. T. McRae, G.-T. Bercea, G. R. Markall, and P. H. Kelly. Firedrake: automating the finite element method by composing abstractions. *ACM Trans. Math. Softw.*, 43(3):1–27, 2016.
- J. Richter-Powell, Y. Lipman, and R. T. Chen. Neural conservation laws: A divergence-free perspective. In *Advances in Neural Information Processing Systems*, volume 35, pages 38075–38088, 2022.
- K. Ritter, G. W. Wasilkowski, and H. Woźniakowski. Multivariate integration and approximation for random fields satisfying sacks-ylvisaker conditions. *Ann. Appl. Probab.*, pages 518–540, 1995.
- F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini. The graph neural network model. *IEEE Trans. Neural Netw.*, 20(1):61–80, 2008.
- F. Schäfer and H. Owhadi. Sparse recovery of elliptic solvers from matrix-vector products. *arXiv preprint arXiv:2110.05351*, 2021.
- F. Schäfer, T. J. Sullivan, and H. Owhadi. Compression, inversion, and approximate PCA of dense kernel matrices at near-linear computational complexity. *Multiscale Model. Simul.*, 19(2):688–730, 2021.

- M. Schmidt and H. Lipson. Distilling free-form natural laws from experimental data. *Science*, 324(5923):81–85, 2009.
- D. P. Searson, D. E. Leahy, and M. J. Willis. GPTIPS: an open source genetic programming toolbox for multigene symbolic regression. In *Proceedings of the International multiconference of engineers and computer scientists*, volume 1, pages 77–80. Citeseer, 2010.
- J. Sirignano and K. Spiliopoulos. DGM: A deep learning algorithm for solving partial differential equations. *J. Comput. Phys.*, 375:1339–1364, 2018.
- J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265, 2015.
- H. Son, J. W. Jang, W. J. Han, and H. J. Hwang. Sobolev training for physics informed neural networks. *arXiv preprint arXiv:2101.08932*, 2021.
- Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- M. L. Stein. *Interpolation of spatial data: some theory for kriging*. Springer Science & Business Media, 1999.
- G. Stepaniants. Learning partial differential equations in reproducing kernel Hilbert spaces. *J. Mach. Learn. Res.*, 24(86):1–72, 2023.
- A. M. Stuart. Inverse problems: a Bayesian perspective. *Acta Numer.*, 19:451–559, 2010.
- E. Süli and D. F. Mayers. *An Introduction to Numerical Analysis*. Cambridge University Press, 2003.
- J. Sun, Y. Liu, Y. Wang, Z. Yao, and X. Zheng. BINN: A deep learning approach for computational mechanics problems based on boundary integral equations. *Comput. Methods Appl. Mech. Eng.*, 410:116012, 2023.
- L. N. Trefethen. *Spectral methods in MATLAB*. SIAM, 2000.
- T. Tripura and S. Chakraborty. Wavelet neural operator: a neural operator for parametric partial differential equations. *arXiv preprint arXiv:2205.02191*, 2022.
- S.-M. Udrescu and M. Tegmark. AI Feynman: A physics-inspired method for symbolic regression. *Sci. Adv.*, 6(16):eaay2631, 2020.
- S.-M. Udrescu, A. Tan, J. Feng, O. Neto, T. Wu, and M. Tegmark. AI Feynman 2.0: Pareto-optimal symbolic regression exploiting graph modularity. In *Advances in Neural Information Processing Systems*, volume 33, pages 4860–4871, 2020.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30, 2017.

- S. Venturi and T. Casey. Svd perspectives for augmenting deeponet flexibility and interpretability. *Comput. Methods Appl. Mech. Eng.*, 403:115718, 2023.
- S. Wang, H. Wang, and P. Perdikaris. On the eigenvector bias of Fourier feature networks: From regression to solving multi-scale PDEs with physics-informed neural networks. *Comput. Methods Appl. Mech. Eng.*, 384:113938, 2021a.
- S. Wang, H. Wang, and P. Perdikaris. Learning the solution operator of parametric partial differential equations with physics-informed DeepONets. *Sci. Adv.*, 7(40):eabi8605, 2021b.
- S. Wang, H. Wang, and P. Perdikaris. Improved architectures and training algorithms for deep operator networks. *J. Sci. Comput.*, 92(2):35, 2022a.
- S. Wang, X. Yu, and P. Perdikaris. When and why PINNs fail to train: A neural tangent kernel perspective. *J. Comput. Phys.*, 449:110768, 2022b.
- S. Wang, S. Sankaran, H. Wang, and P. Perdikaris. An expert’s guide to training physics-informed neural networks. *arXiv preprint arXiv:2308.08468*, 2023.
- H. Weyl. Über die asymptotische verteilung der eigenwerte. *Nachr. Ges. Wiss. Göttingen, Math.-Phys. Kl.*, 1911:110–117, 1911.
- Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and S. Y. Philip. A comprehensive survey on graph neural networks. *IEEE Trans. Neural Netw. Learn. Syst.*, 32(1):4–24, 2020.
- D. Yarotsky. Error bounds for approximations with deep ReLU networks. *Neural Netw.*, 94:103–114, 2017.
- L. Ying, G. Biros, and D. Zorin. A kernel-independent adaptive fast multipole algorithm in two and three dimensions. *J. Comput. Phys.*, 196(2):591–626, 2004.
- H. You, Q. Zhang, C. J. Ross, C.-H. Lee, and Y. Yu. Learning deep implicit Fourier neural operators (IFNOs) with applications to heterogeneous material modeling. *Comput. Methods Appl. Mech. Eng.*, 398:115296, 2022.
- A. Yu, Y. Yang, and A. Townsend. Tuning frequency bias in neural network training with nonuniform data. In *International Conference on Learning Representations*, 2023.
- X. Zhao, X. Chen, Z. Gong, W. Zhou, W. Yao, and Y. Zhang. RecFNO: a resolution-invariant flow and heat field reconstruction method from sparse observations via Fourier neural operator. *Int. J. Therm. Sci.*, 195:108619, 2024.
- H. Zheng, W. Nie, A. Vahdat, K. Azizzadenesheli, and A. Anandkumar. Fast sampling of diffusion models via operator learning. In *International Conference on Machine Learning*, pages 42390–42402, 2023.
- J. Zhou, G. Cui, S. Hu, Z. Zhang, C. Yang, Z. Liu, L. Wang, C. Li, and M. Sun. Graph neural networks: A review of methods and applications. *AI Open*, 1:57–81, 2020.

- H. Zhu, C. K. Williams, R. Rohwer, and M. Morciniec. Gaussian regression and optimal finite dimensional linear models. In *Neural Networks and Machine Learning*. Springer-Verlag, 1998.
- M. Zhu, H. Zhang, A. Jiao, G. E. Karniadakis, and L. Lu. Reliable extrapolation of deep neural operators informed by physics or sparse observations. *Comput. Methods Appl. Mech. Eng.*, 412:116064, 2023.