



Kernel methods are competitive for operator learning

Pau Batlle^a, Matthieu Darcy^{a,*}, Bamdad Hosseini^b, Houman Owhadi^a

^a Computing and Mathematical Sciences, Caltech, Pasadena, CA, United States of America

^b Department of Applied Mathematics, University of Washington, Seattle, WA, United States of America

ARTICLE INFO

Keywords:

Operator learning
Optimal recovery
Kernel methods
Gaussian processes
Functional regression
Partial differential equations

ABSTRACT

We present a general kernel-based framework for learning operators between Banach spaces along with a priori error analysis and comprehensive numerical comparisons with popular neural net (NN) approaches such as Deep Operator Networks (DeepONet) [46] and Fourier Neural Operator (FNO) [45]. We consider the setting where the input/output spaces of target operator $\mathcal{G}^\dagger : \mathcal{U} \rightarrow \mathcal{V}$ are reproducing kernel Hilbert spaces (RKHS), the data comes in the form of partial observations $\phi(u_i), \varphi(v_i)$ of input/output functions $v_i = \mathcal{G}^\dagger(u_i)$ ($i = 1, \dots, N$), and the measurement operators $\phi : \mathcal{U} \rightarrow \mathbb{R}^n$ and $\varphi : \mathcal{V} \rightarrow \mathbb{R}^m$ are linear. Writing $\psi : \mathbb{R}^n \rightarrow \mathcal{U}$ and $\chi : \mathbb{R}^m \rightarrow \mathcal{V}$ for the optimal recovery maps associated with ϕ and φ , we approximate \mathcal{G}^\dagger with $\tilde{\mathcal{G}} = \chi \circ \tilde{f} \circ \phi$ where \tilde{f} is an optimal recovery approximation of $f^\dagger := \varphi \circ \mathcal{G}^\dagger \circ \psi : \mathbb{R}^n \rightarrow \mathbb{R}^m$. We show that, even when using vanilla kernels (e.g., linear or Matérn), our approach is competitive in terms of cost-accuracy trade-off and either matches or beats the performance of NN methods on a majority of benchmarks. Additionally, our framework offers several advantages inherited from kernel methods: simplicity, interpretability, convergence guarantees, a priori error estimates, and Bayesian uncertainty quantification. As such, it can serve as a natural benchmark for operator learning.

1. Introduction

Operator learning is a well-established field going back at least to the 1970s with the articles [1,56] who introduced the reduced basis method as a way speeding up expensive model evaluations. In the most broad sense operator learning arises in the solution of stochastic PDEs [28], emulation of computer codes [37], reduced order modeling (ROM) [48], and numerical homogenization [61]. In recent years, and with the rise of machine learning, operator learning has become the focus of extensive research with the development of neural net (NN) methods such as Deep Operator Nets [46] and Fourier Neural Nets [45] among many others. While these NN methods are often benchmarked against each other [47], they are rarely compared with the aforementioned classical approaches. Furthermore, the theoretical analysis of NN methods is often limited to density/universal approximation results; showing the existence of a network of a requisite size achieving a certain error rate, without guarantees whether this network is computable in practice (see for example [20,41]).

In order to alleviate the aforementioned shortcomings we present a mathematical framework for approximation of mappings between Banach spaces using the theory of operator valued reproducing Kernel Hilbert spaces (RKHS) and Gaussian Processes (GPs). Our abstract framework is: (1) mathematically simple and interpretable, (2) convenient to implement, (3) encompasses some of the

* Corresponding author.

E-mail addresses: pbattlef@caltech.edu (P. Batlle), mdarcy@caltech.edu (M. Darcy), bamdadh@uw.edu (B. Hosseini), owhadi@caltech.edu (H. Owhadi).

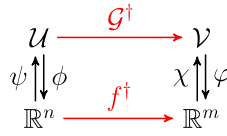


Fig. 1. Commutative diagram of our operator learning setup.

classical approaches such as linear methods; and (4) comes with a priori error analysis and convergence theory. We further present extensive benchmarking of our kernel method with the DeepONet and FNO approaches and show that the kernel approach either matches or outperforms NN methods in most benchmark examples.

In the remainder of this section we give a summary of our methodology and results: We pose the operator learning problem in Subsection 1.1 before presenting a running example in Subsection 1.2 which is used to outline our proposed framework and main theoretical results in Subsections 1.3 and 1.4 as well as brief numerical results in Subsection 1.5. Our main contributions are summarized in Subsection 1.6 followed by a literature review in Subsection 1.7.

1.1. The operator learning problem

Let \mathcal{U} and \mathcal{V} be two (possibly infinite-dimensional) separable Banach spaces and suppose that

$$\mathcal{G}^\dagger : \mathcal{U} \rightarrow \mathcal{V} \quad (1.1)$$

is an arbitrary (possibly nonlinear) operator. Then, broadly speaking, the goal of operator learning is to approximate \mathcal{G}^\dagger from a finite number N of input/output data on \mathcal{G}^\dagger . For our framework, we consider the setting where the input/output data are only partially observed through a finite collection of linear measurements which we formalize as follows:

Problem 1. Let $\{u_i, v_i\}_{i=1}^N$ be N elements of $\mathcal{U} \times \mathcal{V}$ such that

$$\mathcal{G}^\dagger(u_i) = v_i, \quad \text{for } i = 1, \dots, N. \quad (1.2)$$

Let $\phi : \mathcal{U} \rightarrow \mathbb{R}^n$ and $\varphi : \mathcal{V} \rightarrow \mathbb{R}^m$ be bounded linear operators. Given the data $\{\phi(u_i), \varphi(v_i)\}_{i=1}^N$ approximate \mathcal{G}^\dagger .

1.2. Running example

To give context to the above problem and our solution method we briefly outline a running example to which the reader can refer to throughout the rest of this section. Consider the following elliptic PDE, which is of broad interest in geosciences and material science:

$$\begin{cases} -\operatorname{div} e^u \nabla v = w, & \text{in } \Omega, \\ v = 0, & \text{on } \partial\Omega, \end{cases} \quad (1.3)$$

where $\Omega = (0, 1)^2$, $u \in H^3(\Omega)$, $w \in H^1(\Omega)$ and $v \in H^3(\Omega) \cap H_0^1(\Omega)$. For a fixed forcing term w , we wish to approximate the nonlinear operator mapping the diffusion coefficient u to the solution v , i.e., $\mathcal{G}^\dagger : u \mapsto v$. In this case we may take $\mathcal{U} \equiv H^3(\Omega)$ and $\mathcal{V} \equiv H^3(\Omega) \cap H_0^1(\Omega)$. We further assume that a training data set is available in the form of limited observations of input-out pairs. As a canonical example, consider the evaluation bounded and linear operators

$$\phi : u \mapsto (u(X_1), u(X_2), \dots, u(X_n))^T \quad \text{and} \quad \varphi : v \mapsto (v(Y_1), v(Y_2), \dots, v(Y_m))^T, \quad (1.4)$$

where the $\{X_j\}_{j=1}^n$ and $\{Y_j\}_{j=1}^m$ are distinct collocation points in the domain Ω as well as pairs $\{u_i, v_i\}_{i=1}^N$ that satisfy the PDE (1.3). Then our goal is to approximate \mathcal{G}^\dagger from the training data set $\{\phi(u_i), \varphi(v_i)\}_{i=1}^N$.¹

1.3. The proposed solution

Our setup naturally gives rise to a commutative diagram depicted in Fig. 1. Here the map $f^\dagger : \mathbb{R}^n \rightarrow \mathbb{R}^m$ explicitly defined as

$$f^\dagger := \varphi \circ \mathcal{G}^\dagger \circ \psi \quad (1.5)$$

is a mapping between finite-dimensional Euclidean spaces, and is therefore amenable to numerical approximation. However, in order to approximate \mathcal{G}^\dagger we also need the reconstruction maps $\psi : \mathbb{R}^n \rightarrow \mathcal{U}$ and $\chi : \mathbb{R}^m \rightarrow \mathcal{V}$.

Our proposed solution is to endow \mathcal{U} and \mathcal{V} with an RKHS structure and use kernel/GP regression to identify the maps ψ and χ . As a prototypical example we consider the situation where \mathcal{U} is an RKHS of functions $u : \Omega \rightarrow \mathbb{R}$ defined by a kernel $Q : \Omega \times \Omega \rightarrow \mathbb{R}$

¹ Choosing ϕ, φ as pointwise evaluation functionals is common to many applications, although our abstract framework readily accommodates other choices such as integral operators and basis projections.

and \mathcal{V} is an RKHS of functions $u : D \rightarrow \mathbb{R}$ defined by a kernel $K : D \times D \rightarrow \mathbb{R}$. For our running example, we have $D = \Omega$, and we can take Q and K to be Matérn like kernels, e.g., the Green's function of elliptic PDEs (possibly on Ω or restricted to Ω) with appropriate regularity. One can also choose Q, K to be smoother kernels such that their RKHSs are embedded in \mathcal{U} and \mathcal{V} .

We then define ψ and χ as the following optimal recovery maps²:

$$\begin{aligned}\psi(U) &:= \arg \min_{w \in \mathcal{U}} \|w\|_Q \quad \text{s.t.} \quad \phi(w) = U, \\ \chi(V) &:= \arg \min_{w \in \mathcal{V}} \|w\|_K \quad \text{s.t.} \quad \varphi(w) = V,\end{aligned}\tag{1.6}$$

where $\|\cdot\|_Q$ and $\|\cdot\|_K$ are the RKHS norms arising from their pertinent kernels.

In the case where ϕ and φ are pointwise evaluation maps ($\phi(u) = (u(X_1), \dots, u(X_n))$ and $\varphi(v) = (v(Y_1), \dots, v(Y_m))$ where the X_i and Y_j are pairwise distinct collocation points in Ω and D), our optimal recovery maps can be expressed in closed form using standard representer theorems for kernel interpolation [71]:

$$\psi(U)(x) = Q(x, X)Q(X, X)^{-1}U, \quad \chi(V)(y) = K(y, Y)K(Y, Y)^{-1}V,\tag{1.7}$$

where $Q(X, X)$ and $K(Y, Y)$ are kernel matrices with entries $Q(X, X)_{ij} = Q(X_i, X_j)$ and $K(Y, Y)_{ij} = K(Y_i, Y_j)$ respectively, while $Q(x, X)$ and $K(y, Y)$ denote row-vector fields with entries $Q(x, X)_i = Q(x, X_i)$ and $K(y, Y)_i = K(y, Y_i)$.

We further propose to approximate f^\dagger by optimal recovery in a vector-valued RKHS. Let $\Gamma : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathcal{L}(\mathbb{R}^m)$ be a matrix valued kernel [3]; here $\mathcal{L}(\mathbb{R}^m)$ is the space of $m \times m$ matrices) with RKHS \mathcal{H}_Γ equipped with the norm $\|\cdot\|_\Gamma$ ³ and proceed to approximate f^\dagger by the map \tilde{f} defined as

$$\tilde{f} := \arg \min_{f \in \mathcal{H}_\Gamma} \|f\|_\Gamma \quad \text{s.t.} \quad f(\phi(u_i)) = \varphi(v_i) \quad \text{for } i = 1, \dots, N.$$

A simple and practical choice for Γ is the diagonal kernel

$$\Gamma(U, U') = S(U, U')I\tag{1.8}$$

where $S : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ is an arbitrary scalar-valued kernel, such as RBF, Laplace, or Matérn, and I is the $m \times m$ identity matrix. More complicated choices, such as sums of kernels or replacing the identity matrix for a fixed positive definite matrix, implying correlations between various input or output correlations, are also possible. However, these may lead to greater computational cost and we observe empirically that the simple choice of the identity matrix already provides good performance. Then we can approximate the components of \tilde{f} via the independent optimal recovery problems

$$\tilde{f}_j := \arg \min_{g \in \mathcal{H}_S} \|g\|_S \quad \text{s.t.} \quad g(\phi(u_i)) = \varphi_j(v_i), \quad \text{for } i = 1, \dots, N\tag{1.9}$$

for $j = 1, \dots, m$. Here we wrote $\varphi_j(v_i)$ for the entry j of the vector $\varphi(v_i)$ and, as our notation suggests, \mathcal{H}_S is the RKHS of S equipped with the norm $\|\cdot\|_S$. Since (1.9) is a standard optimal recovery problem, each \tilde{f}_j can be identified by the usual representer formula:

$$\tilde{f}_j(U) = S(U, \mathbf{U})S(\mathbf{U}, \mathbf{U})^{-1}\mathbf{V}_{\cdot,j},\tag{1.10}$$

where $\mathbf{U} := (\phi(u_1), \dots, \phi(u_N))$ and $\mathbf{V}_{\cdot,j} := (\varphi_j(v_1), \dots, \varphi_j(v_N))^T$ and $S(U, \mathbf{U})$ is a block-vector and $S(\mathbf{U}, \mathbf{U})$ is a block-matrix defined in an analogous manner to those in (1.7). By combining equations (1.7) and (1.10) we obtain the operator

$$\tilde{\mathcal{G}} := \chi \circ \tilde{f} \circ \phi\tag{1.11}$$

as an approximation to \mathcal{G}^\dagger . We provide further details and generalize the proposed framework in Section 2 to the setting where ϕ and φ are obtained from arbitrary linear measurements (e.g., integral operators as in tomography) and \mathcal{U} and \mathcal{V} may not be spaces of continuous functions.

1.4. Convergence guarantee

Under suitable regularity assumptions on \mathcal{G}^\dagger , our method comes with worst-case convergence guarantees as the number of data points N , i.e., input-output pairs and the number of collocations points n and m go to infinity. We present here a condensed version of this result and defer the proof to Section 3. Below we write $B_R(\mathcal{H})$ for the ball of radius $R > 0$ in a normed space \mathcal{H} .

Theorem 1.1 (Condensed version of Theorem 3.4). *Suppose it holds that:*

(1.1.1) (Regularity of the domains Ω and D) Ω and D are compact sets of finite dimensions d_Ω and d_D and with Lipschitz boundary.

² It is possible to define the optimal recovery maps ψ, χ in the setting where ϕ and ψ are nonlinear, following the general framework of [14,59,60]. However, in this setting the closed form formulae (1.7) no longer hold.

³ See Appendix A for a review of operator-valued kernels or the reference [36].

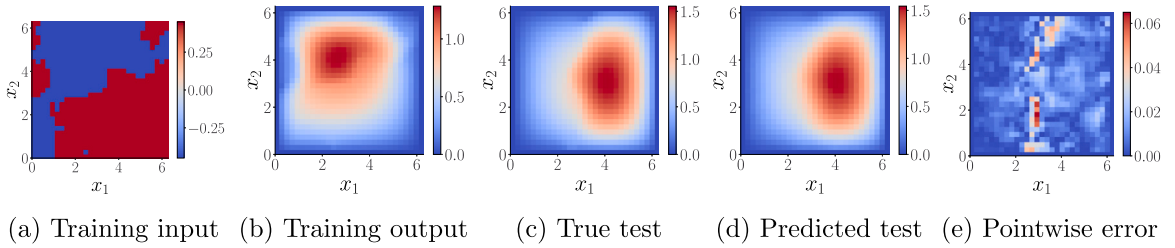


Fig. 2. Example of training data and test prediction and pointwise errors for the Darcy flow problem (1.3). (For interpretation of the colors in the figure(s), the reader is referred to the web version of this article.)

- (1.1.2) *Regularity of the kernels Q and K .* Assume that $\mathcal{H}_Q \subset H^s(\Omega)$ and $\mathcal{H}_K \subset H^t(D)$ for some $s > d_\Omega/2$ and some $t > d_D/2$ with inclusions indicating continuous embeddings.
- (1.1.3) *(Space filling property of collocation points)* The fill distance between the collocation points $\{X_i\}_{i=1}^n \subset \Omega$ and the $\{Y_j\}_{j=1}^m \subset D$ goes to zero as $n \rightarrow \infty$ and $m \rightarrow \infty$.
- (1.1.4) *(Regularity of the operator \mathcal{G}^\dagger)* The operator \mathcal{G}^\dagger is continuous from $H^{s'}(\Omega)$ to \mathcal{H}_K for some $s' \in (0, s)$ as well as from \mathcal{U} to \mathcal{V} and all its Fréchet derivatives are bounded on $B_R(\mathcal{H}_Q)$ for any $R > 0$.
- (1.1.5) *(Regularity of the kernels S^n)* Assume that for any $n \geq 1$ and any compact subset Y of \mathbb{R}^n , the RKHS of S^n restricted to Y is contained in $H^r(Y)$ for some $r > n/2$ and contains $H^{r'}(Y)$ for some $r' > 0$ that may depend on n .
- (1.1.6) *(Resolution and space-filling property of the data)* Assume that for n sufficiently large, the data points $(u_i)_{i=1}^N \subset B_R(\mathcal{H}_Q)$ belong to the range of ψ^n and are space filling in the sense that they become dense in $\phi^n(B_R(\mathcal{H}_Q))$ as $N \rightarrow \infty$.

Then, for all $t' \in (0, t)$,

$$\lim_{n,m \rightarrow \infty} \lim_{N \rightarrow \infty} \sup_{u \in B_R(\mathcal{H}_Q)} \|\mathcal{G}^\dagger(u) - \chi^m \circ \bar{f}_N^{m,n} \circ \phi^n(u)\|_{H^{t'}(D)} \rightarrow 0, \quad (1.12)$$

where our notation makes the dependence of ψ, ϕ, χ, S and \bar{f} on n, m and N explicit.

We note that Assumptions (1.1.1)–(1.1.3) are standard, and concern the accuracy of the optimal recovery maps ϕ^n and χ^m as $n, m \rightarrow \infty$. Assumptions (1.1.4)–(1.1.5) are less standard and amount to regularity assumptions on the map \mathcal{G}^\dagger while Assumption (1.1.6) concerns the acquisition and regularity of the training data set.

In Section 3 we also present Theorem 3.3 as the quantitative analogue of the above result which characterizes how the speed of convergence depends on the regularity of the operator \mathcal{G}^\dagger and the choice of ϕ and φ in the setting of pointwise measurement operators. We also comment on how this analysis could be extended to other linear measurements.

1.5. Numerical framework

Returning to our running example, we implement the proposed framework for learning the non-linear operator mapping u to v in (1.3). We consider 1000 inputs and outputs of u and v . The data is taken from [47] and the experimental setup is discussed further in Subsection 4.3.2. We take φ to be of the form (1.4) with $m = 841$ while we define ϕ through a PCA pre-processing step. More precisely, let $\phi_{\text{pointwise}}$ be of the form (1.4) with $n = 841$. Choose $n_{\text{PCA}} = 202$ (this value captures 95% of the empirical variance of our training data) and define

$$\phi(u) = \Pi_{\text{PCA}} \circ \phi_{\text{pointwise}}(u) \in \mathbb{R}^{202}. \quad (1.13)$$

In other words, we take our ϕ map to be the linear map that computes the first 202 PCA coefficients of the input functions u given on a uniform grid; observe that we do not use PCA pre-processing on the output data here although we do this for some of our other examples in Section 4 for better performance.

With ϕ and φ identified (recall Fig. 1) we proceed to implement our kernel method using the simple choice of a diagonal kernel $S(U, U')I$ where S is a rational quadratic (RQ) kernel (see Appendix C). This choice transforms the problem into 841 independent kernel regression problems, each corresponding to one component of f^\dagger (i.e., the f_j^\dagger 's).

We used the PCA and kernel regression modules of the `scikit-learn` Python library [65] to implement our algorithm. This implementation automatically selects the best kernel parameters by maximizing the marginal likelihood function [67] jointly for all problems. Our proposed method can therefore be implemented conveniently using off-the-shelf software. Fig. 2 illustrates examples of the inputs and outputs of our operator learning problem. Despite the simple implementation of our method, we are able to obtain competitive accuracy as shown in Table 1 where the relative testing L^2 loss of our method is compared to other popular algorithms. Moreover, our approach is amenable to well-known numerical analysis techniques, such as sparse or low-rank approximation of kernel matrices, to reduce its complexity. For the present example (and those in Section 4) we only consider “vanilla” kernel methods which compute (1.10) by computing the full Cholesky factors of the matrix $S(U, U)$.

Table 1

The L^2 relative test error of the Darcy flow problem in our running example. The kernel approach is compared with variations of DeepONet and FNO. Results of our kernel method are presented below the dashed line with the pertinent choice of the kernel S .

Method	Accuracy
DeepONet	2.91%
FNO	2.41%
POD-DeepONet	2.32%
Linear	6.74%
Rational quadratic	2.87%

1.6. Summary of contributions

The main results of the article concern the properties, performance, and error analysis of the map $\tilde{\mathcal{G}}$ defined in (1.11). Our contributions can be summarized under four categories:

1. **An abstract kernel framework for operator learning:** In Section 2, we propose a framework for operator learning using kernel methods with several desirable properties. A family of methods of increasing complexity is proposed that includes linear models and diagonal kernels as well as non-diagonal kernels which capture output correlations. These properties make our approach ideal for benchmarking purposes. Furthermore, the methodology is: (i) applicable to any choice of the linear functionals φ and ψ ; (ii) minimax optimal with respect to an implicitly defined operator-valued kernel; and (iii) is mesh-invariant. We emphasize in Remark 2.1 that our optimal recovery maps can be applied to *any* operator learning after training to obtain a mesh-invariant pipeline.
2. **Error analysis and convergence rates for $\tilde{\mathcal{G}}$:** In Section 3, we develop rigorous worst-case a priori error bounds and convergence guarantees for our method: Theorem 3.3 provides quantitative error bounds while Theorem 3.4 (the detailed version of Theorem 3.3) shows the convergence of $\tilde{\mathcal{G}} \rightarrow \mathcal{G}$ under appropriate conditions.
3. **A simple to use vanilla kernel method:** While our abstract kernel method is quite general, our numerical implementation in Section 4 focuses on a simple, easy-to-implement version using diagonal kernels of the form (1.8). Off-the-shelf software, such as the kernel regression modules of `scikit-learn`, can be employed for this task. We empirically observe low training times and robust choice of hyperparameters. These properties further suggest that kernel methods are a good baseline for benchmarking of more complex methods.
4. **Competitive performance.** In Section 4 we present a series of numerical experiments on benchmark PDE problems from the literature and observe that our simple implementation of the kernel approach is competitive in terms of complexity-accuracy tradeoffs in comparison to several NN-based methods. Since kernel methods can be interpreted as an infinite-width, one-layer NN, the results raise the question of how much of a role the depth of a deep NN plays in the performance of algorithms for the purposes of operator learning.

1.7. Review of relevant literature

In the most broad sense, operator learning is the problem of approximating a mapping between two infinite-dimensional function spaces [9,19]. In recent years, this problem has become an area of intense research in the scientific machine learning community with a particular focus on parametric or stochastic PDEs. However, the approximation of such parameter to solution maps has been an area of intense research in the computational mathematics and engineering communities, going back at least to the reduced basis method introduced in the 1970s [1,56] as a way of speeding up the solution of families of parametric PDEs in applications that require many PDE solves such as design [21,52,8,10], uncertainty quantification (UQ) [72,51,35], and multi-scale modeling [76,27,26,42]. In what follows we give a brief summary of the various areas and methodologies that overlap with operator learning; we cannot provide an exhaustive list of references due to space, but refer the reader to key contributions and surveys where further references can be found.

Deep learning techniques The use of NNs for operator learning goes back at least to the 90s and the seminal works of Chen and Chen [13,12] who proved a universal approximation theorem for NN approximations to operators. The use and design of NNs for operator learning has become popular in the last five years as a consequence of growing interest in NNs for scientific computing starting with the article [81] which used autoencoders to build surrogates for UQ of subsurface flow models. Since then many different approaches have been proposed some of which use specific architectures or target particular families of PDEs [33,39,45,38,46,29,11,44,40]. The most relevant of among these methods to our proposed framework are the DeepONet family [46,74,47,75], FNO [45], and PCA-Net [33,9] where the main novelty appears to be the use of novel, flexible, and expressive NN architectures that allow the algorithm to learn and adapt the bases that are selected for the input and outputs of the solution map as well as possible nonlinear dependencies between the basis coefficients. Although not part of our comparisons, we note that [22–24] obtained competitive accuracy by using deep neural networks with architectures inspired by conventional fast solvers.

Classical numerical approximation methods Operator learning has been the subject of intense research in the computational mathematics literature in the context of stochastic Galerkin methods [28,80], polynomial chaos [78,79], reduced basis methods [56,49] and numerical homogenization [63,57,61,2]. In the setting of stochastic and parametric PDEs, the goal is often to approximate the solution of a PDE as a function of a random or uncertain parameter. The well-established approach to such problems is to pick or construct appropriate bases for the input parameter and the solution of the PDE and then construct a parametric, high-dimensional map, that transforms the input basis coefficients to the output coefficients. Well-established methods such as polynomial chaos, stochastic finite element methods, reduced basis methods [28,78,18,34,48] fall within this category. A vast amount of literature in applied mathematics exists on this subject, and the theoretical analysis of these methods is extensive; see for example [7,16,17,55,54,30] and references therein.

Operator compression For solving PDEs, the objectives of operator learning are also similar to those of operator compression [25,44] as formulated in numerical homogenization [61,2] and reduced order modeling (ROM) [4,48], i.e., the approximation of the solution operator from pairs of solutions and source/forcing terms. While both ROM and numerical homogenization seek operator compression through the identification of reduced basis functions that are as accurate as possible (this translates into low-rank approximations with SVD and its variants [11]), numerical homogenization also requires those functions to be as localized as possible [50] and in turn leverages both low rank and sparse approximations. These localized reduced basis functions are known as Wannier functions in the physics literature [53], and can be interpreted as linear combinations of eigenfunctions that are localized in both frequency space and the physical domain, akin to wavelets. The hierarchical generalization of numerical homogenization [58] (gamblots) has led to the current state-of-the-art for operator compression of linear elliptic [68,70] and parabolic/hyperbolic PDEs [64]. In particular, for arbitrary (and possibly unknown) elliptic PDEs [69] shows that the solution operator (i.e., the Green's function) can be approximated in near-linear complexity to accuracy ϵ from only $\mathcal{O}(\log^{d+1}(\frac{1}{\epsilon}))$ solutions of the PDE.

GP emulators In the case where the range of the operator of interest is finite dimensional, then operator learning coincides with surrogate modeling techniques that were developed in the UQ literature, such as GP surrogate modeling/emulation [37,6]. When the kernels of the underlying GPs are also learned from data [62,15], GP surrogate modeling has been shown to offer a simple, low-cost, and accurate solution to learning dynamical systems [32], geophysical forecasting [31], and radiative transfer emulation [73], and the inference of the structure of convective storms from passive microwave observations [66]. Indeed, our proposed kernel framework for operator learning can be interpreted as an extension of these well-established GP surrogates to the setting where the range of the operator is a function space.

1.8. Outline of the article

The remainder of the article is organized as follows: We present our operator learning framework in Section 2 for the generalized setting where ϕ, φ can be any collection of bounded and linear operators along with an interpretation of our method from the GP perspective. Our convergence analysis and quantitative error bounds are presented in Section 3 where we present the full version of Theorem 3.4. Our numerical experiments, implementation details, and benchmarks against FNO and DeepONet are collected in Section 4. We discuss future directions and open problems in Section 5. The appendix collects a review of operator valued kernels and GPs along with other auxiliary details.

2. The RKHS/GP framework for operator learning

We now present our general kernel framework for operator learning, i.e., the proposed solution to Problem 1. We emphasize that here we do not require the spaces \mathcal{U} and \mathcal{V} to be spaces of continuous functions and in particular, we do not require the maps ϕ and φ to be obtained from pointwise measurements. To describe this, we will introduce the dual spaces of \mathcal{U} and \mathcal{V} to define optimal recovery with respect to kernel operators rather than just kernel functions.

Write \mathcal{U}^* and \mathcal{V}^* for the duals of \mathcal{U} and \mathcal{V} , and write $[\cdot, \cdot]$ for the pertinent duality pairings. Assume that \mathcal{U} is endowed with a quadratic norm $\|\cdot\|_Q$, i.e., there exists a linear bijection $Q : \mathcal{U}^* \rightarrow \mathcal{U}$ that is symmetric ($[\phi_a, Q\phi_b] = [\phi_b, Q\phi_a]$), positive ($[\phi_a, Q\phi_a] > 0$ for $\phi_a \neq 0$), and such that $\|u\|_Q^2 = [Q^{-1}u, u]$, $\forall u \in \mathcal{U}$.

As in [61, Ch. 11], although \mathcal{U} , and \mathcal{U}^* are also Hilbert spaces under $\|\cdot\|_Q$ and its dual norm $\|\cdot\|_Q^*$ (with inner products $\langle u, v \rangle_Q = [Q^{-1}u, v]$ and $\langle \phi_a, \phi_b \rangle_Q^* = [\phi_a, Q\phi_b]$), we will keep using the Banach space terminology to emphasize the fact that our dual pairings will not be based on the inner product through the Riesz representation theorem, but on a different realization of the dual space, as this setting is more practical.

If \mathcal{U} is a space of continuous functions on a subset $\Omega \subset \mathbb{R}^{d_\Omega}$ then \mathcal{U}^* contains delta Dirac functions and, to simplify notations, we also write $Q(x, y) := [\delta_x, Q\delta_y]$ for $x, y \in \mathbb{R}^{d_\Omega}$ to denote the kernel induced by the operator Q . Note that in that case, \mathcal{U} is a RKHS with norm $\|\cdot\|_Q$ induced by the kernel Q . Since ϕ is bounded and linear, its entries ϕ_i (write $\phi := (\phi_1, \dots, \phi_n)$) must be elements of \mathcal{U}^* . We assume those elements to be linearly independent. Write $\psi : \mathbb{R}^n \rightarrow \mathcal{U}$ for the linear operator defined by

$$\psi(Y) := (Q\phi)Q(\phi, \phi)^{-1}Y \text{ for } Y \in \mathbb{R}^n, \quad (2.1)$$

where we write $Q(\phi, \phi)$ for the $n \times n$ symmetric positive definite (SPD) matrix with entries $Q(\phi_i, \phi_j) := [\phi_i, Q\phi_j]^4$ and $Q\phi$ for $(Q\phi_1, \dots, Q\phi_n) \in \mathcal{U}^n$. As described in [61, Chap. 11], for $u \in \mathcal{U}$, given $\phi(u) = Y$, $\psi(Y)$ is the minmax optimal recovery of u when using the relative error in $\|\cdot\|_Q$ -norm as a loss.

Similarly, assume that \mathcal{V} is endowed with a quadratic norm $\|\cdot\|_K$, defined by the symmetric positive linear bijection $K : \mathcal{V}^* \rightarrow \mathcal{V}$. Write $\varphi := (\varphi_1, \dots, \varphi_m)$ and assume the entries of φ to be linearly independent elements of \mathcal{V}^* . Using the same notations as in (2.1) write $\chi : \mathbb{R}^m \rightarrow \mathcal{V}$ for the linear operator defined by

$$\chi(Z) := (K\varphi)K(\varphi, \varphi)^{-1}Z \text{ for } Z \in \mathbb{R}^m. \quad (2.2)$$

Then, as above, for $v \in \mathcal{V}$, given $\varphi(v) = Z$, $\chi(Z)$ is the minmax optimal recovery of v when using the relative error in $\|\cdot\|_K$ -norm as a loss.

Write $\mathcal{L}(\mathbb{R}^m)$ for the space of bounded linear operators mapping \mathbb{R}^m to itself, i.e., $m \times m$ matrices. Let $\Gamma : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathcal{L}(\mathbb{R}^m)$ be a matrix-valued kernel [3] defining an RKHS \mathcal{H}_Γ of continuous functions $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ equipped with an RKHS norm $\|\cdot\|_\Gamma$. For $i \in \{1, \dots, N\}$, write $U_i := \phi(u_i)$ and $V_i := \varphi(v_i)$. Write \mathbf{U} and \mathbf{V} for the block-vectors with entries U_i and V_i . Write $\Gamma(\mathbf{U}, \mathbf{U})$ for the $N \times N$ block-matrix with entries $\Gamma(U_i, U_j)$ and assume $\Gamma(\mathbf{U}, \mathbf{U})$ to be invertible (which is satisfied if Γ is non-degenerate and $U_i \neq U_j$ for $i \neq j$). Let f^\dagger be an element of \mathcal{H}_Γ and write $f^\dagger(\mathbf{U})$ for the block vector with entries $f^\dagger(U_i)$. Then given $f^\dagger(\mathbf{U}) = \mathbf{V}$ it follows that

$$\bar{f}(\mathbf{U}) := \Gamma(\mathbf{U}, \mathbf{U})\Gamma(\mathbf{U}, \mathbf{U})^{-1}\mathbf{V}, \quad (2.3)$$

is the minimax optimal recovery of f^\dagger , where $\Gamma(\cdot, \mathbf{U})$ is the block-vector with entries $\Gamma(\cdot, U_i)$.

To this end, we propose to approximate the ground truth operator \bar{G}^\dagger with

$$\bar{G} := \chi \circ \bar{f} \circ \phi, \quad (2.4)$$

also recall Fig. 1. Combining (2.2) and (2.3) we further infer that \bar{G} admits the following explicit representer formula

$$\bar{G}(u) = (K\varphi)K(\varphi, \varphi)^{-1}\Gamma(\phi(u), \mathbf{U})\Gamma(\mathbf{U}, \mathbf{U})^{-1}\mathbf{V}. \quad (2.5)$$

In the remainder of this section we will provide more details and observations regarding our approximate operator \bar{G} that is useful later in Section 3 and of independent interest.

2.1. The kernel and RKHS associated with \bar{G}

The explicit formula (2.5) suggests that the operator \bar{G} is an element of an RKHS defined by an operator-valued kernel, which we now characterize. For $u_1, u_2 \in \mathcal{U}$ and $v \in \mathcal{V}$ write

$$G(u_1, u_2)v := (K\varphi)(K(\varphi, \varphi))^{-1}\Gamma(\phi(u_1), \phi(u_2))(K(\varphi, \varphi))^{-1}\varphi(v). \quad (2.6)$$

It turns out that $G : \mathcal{U} \times \mathcal{U} \rightarrow \mathcal{L}(\mathcal{V})$ is a well-defined operator-valued kernel whose RKHS contains operators of the form \bar{G} .

Proposition 2.1. *The kernel G in (2.6) is an operator-valued kernel. Write \mathcal{H}_G for its RKHS and $\|\cdot\|_G$ for the associated norm. Then it holds that $\bar{G} \in \mathcal{H}_G$ if and only if $\bar{G} = \chi \circ f \circ \phi$ for $f = \varphi \circ \mathcal{G} \circ \psi \in \mathcal{H}_\Gamma$ and $\|\bar{G}\|_G = \|f\|_\Gamma$.*

Proof. Since G is Hermitian and positive, we deduce that G is an operator-valued kernel. Indeed for $\tilde{u}_1, \dots, \tilde{u}_m \in \mathcal{U}$ and $\tilde{v}_1, \dots, \tilde{v}_m \in \mathcal{V}$, using $\langle \tilde{v}_i, K\varphi_s \rangle_K = \varphi_s(\tilde{v}_i)$ and the fact that Γ is a matrix-valued kernel we have

$$\begin{aligned} \langle \tilde{v}_i, G(\tilde{u}_i, \tilde{u}_j)\tilde{v}_j \rangle_K &= \varphi(\tilde{v}_i)^T (K\varphi)(K(\varphi, \varphi))^{-1}\Gamma(\phi(\tilde{u}_i), \phi(\tilde{u}_j))(K(\varphi, \varphi))^{-1}\varphi(\tilde{v}_j) \\ &= \langle G(\tilde{u}_j, \tilde{u}_i)\tilde{v}_i, \tilde{v}_j \rangle_K, \end{aligned} \quad (2.7)$$

where we used $\langle \tilde{v}_i, K\varphi_s \rangle_K = \varphi_s(\tilde{v}_i)$ and the fact that Γ is a matrix-valued kernel. Furthermore, summing (2.7), we deduce that $\sum_{i,j=1}^m \langle \tilde{v}_i, G(\tilde{u}_i, \tilde{u}_j)\tilde{v}_j \rangle_K \geq 0$. From (2.6) we infer

$$\sum_{j=1}^m G(u, \tilde{u}_j)\tilde{v}_j = \chi \circ f \circ \phi(u) \quad (2.8)$$

with the function

$$f(\mathbf{U}) = \sum_{j=1}^m \Gamma(\mathbf{U}, \phi(\tilde{u}_j))(K(\varphi, \varphi))^{-1}\varphi(\tilde{v}_j). \quad (2.9)$$

⁴ For linear measurements involving derivatives the computation of these kernel matrices requires the computation of derivatives of the kernels; see [14] for practical examples and considerations.

Furthermore using the reproducing property of G and (2.7) we have $\left\| \sum_{j=1}^m G(u, \tilde{u}_j) \tilde{v}_j \right\|_G^2 = \|f\|_\Gamma^2$. Therefore the closure of the space of operators of the form (2.8) with respect to the RKHS norm induced by G is the space of functions of the form $\chi \circ f \circ \phi$ where f lives in the closure of functions of the form (2.9) with respect to the RKHS norm induced by Γ . We deduce that $\mathcal{H}_G = \{\chi \circ f \circ \phi \mid f \in \mathcal{H}_\Gamma\}$. The uniqueness of f in the representation $\mathcal{G} = \chi \circ f \circ \phi$ for $f \in \mathcal{H}_G$ follows from $f = \varphi \circ \mathcal{G} \circ \psi$ following the identities $\varphi \circ \chi = I_d$ and $\phi \circ \psi = I_d$. ■

Using the above result we can further characterize $\bar{\mathcal{G}}$ and \bar{f} via optimal recovery problems in \mathcal{H}_G and \mathcal{H}_Γ respectively. In what follows we will write \mathbf{u} for the N vector whose entries are the u_i , and $\mathcal{G}(\mathbf{u})$ for the N vector whose entries are $\mathcal{G}^\dagger(u_i)$.

Proposition 2.2. *The operator $\bar{\mathcal{G}}$ is the minimizer of*

$$\begin{cases} \text{Minimize} & \|\mathcal{G}\|_G^2 \\ \text{Over} & \mathcal{G} \in \mathcal{H}_G \text{ such that } \varphi \circ \mathcal{G}(\mathbf{u}) = \varphi \circ \mathcal{G}^\dagger(\mathbf{u}), \end{cases} \quad (2.10)$$

while the map \bar{f} is the minimizer of

$$\begin{cases} \text{Minimize} & \|f\|_\Gamma^2 \\ \text{Over} & f \in \mathcal{H}_\Gamma \text{ such that } f \circ \phi(\mathbf{u}) = \varphi \circ \mathcal{G}^\dagger(\mathbf{u}). \end{cases} \quad (2.11)$$

Proof. By Proposition 2.1 $\bar{\mathcal{G}}$ is completely identified by \bar{f} and $\|\bar{\mathcal{G}}\|_G = \|\bar{f}\|_\Gamma$. Then solving (2.10) is equivalent to solving (2.11). The statement regarding \bar{f} follows directly from representer formulae for optimal recovery with matrix-valued kernels. ■

2.2. Regularizing $\bar{\mathcal{G}}$ by operator regression

As is often the case with optimal recovery/kernel regression the estimator for \bar{f} in (2.3) is susceptible to numerical error due to ill-conditioning of the kernel matrix $\Gamma(\mathbf{U}, \mathbf{U})$. To overcome this issue we regularize our estimator by adding a small diagonal perturbation to this matrix. More precisely, let $\gamma > 0$ and write I for the identity matrix. We then define the regularized map

$$\bar{f}_\gamma(\mathbf{U}) := \Gamma(\mathbf{U}, \mathbf{U})(\Gamma(\mathbf{U}, \mathbf{U}) + \gamma I)^{-1} \mathbf{V}. \quad (2.12)$$

This regularized map gives rise to the regularized approximate operator

$$\bar{\mathcal{G}}_\gamma := \chi \circ \bar{f}_\gamma \circ \phi,$$

which admits the following representer formula

$$\bar{\mathcal{G}}_\gamma(u) = (K\varphi)K(\varphi, \varphi)^{-1}\Gamma(\phi(u), \mathbf{U})(\Gamma(\mathbf{U}, \mathbf{U}) + \gamma I)^{-1} \mathbf{V}. \quad (2.13)$$

We can further characterize this operator as the solution to an operator regression problem.

Proposition 2.3. *$\bar{\mathcal{G}}_\gamma$ is the solution to*

$$\text{Minimize}_{\mathcal{G} \in \mathcal{H}_G} \|\mathcal{G}\|_G^2 + \gamma^{-1} \|\varphi \circ \mathcal{G}(\mathbf{u}) - \varphi \circ \mathcal{G}^\dagger(\mathbf{u})\|^2. \quad (2.14)$$

Proof. By Proposition 2.1, $\mathcal{G} = \chi \circ f \circ \phi$ solves (2.14) if and if f solves

$$\text{Minimize}_{f \in \mathcal{H}_\Gamma} \|f\|_\Gamma^2 + \gamma^{-1} \|f(\mathbf{U}) - \mathbf{V}\|^2. \quad (2.15)$$

It then follows, by standard representer theorems for matrix-valued kernel regression (see Section A.5) that \bar{f}_γ is the minimizer of (2.15). ■

2.3. Interpretation as conditioned operator valued GPs

Our kernel approach to operator learning has a natural GP regression interpretation that is compatible with Bayesian inference and UQ pipelines. We present some facts and observations in this direction.

Write $\xi \sim \mathcal{N}(0, G)$ for the centered operator-valued GP with covariance kernel G^5 and $\zeta \sim \mathcal{N}(0, \Gamma)$ for a centered vector valued GP with covariance kernel Γ . Then it is straightforward to show that the law of ξ is equivalent to that of $\chi \circ \zeta \circ \phi$. Let $Z = (Z_1, \dots, Z_N)$ be a random block-vector, independent from ξ , with i.i.d. entries $Z_j \sim \mathcal{N}(0, \gamma I_m)$ for $j = 1, \dots, N$; here $\gamma \geq 0$ and I_m is the $m \times m$ identity matrix.

⁵ See Appendix A.6 for a review of operator valued GPs.

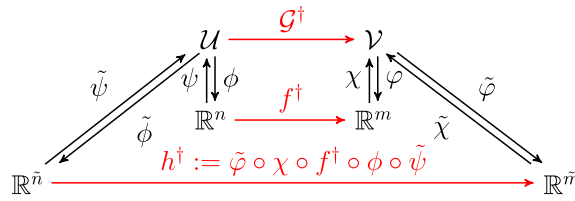


Fig. 3. Generalization of Fig. 1 to the mesh invariant setting where the measurement functionals are different at test time.

Then ξ conditioned on $\varphi \circ \xi(\mathbf{u}) = \varphi(\mathbf{v}) + Z$ is an operator-valued GP with mean $\bar{\mathcal{G}}_\gamma$, as in (2.13), and conditional covariance kernel

$$G^\perp(u, u')v = (K\varphi)(K(\varphi, \varphi))^{-1}\Gamma(\phi(u_1), \phi(u_2)) \\ (\Gamma(\phi(u), \phi(u')) - \Gamma(\phi(u), \mathbf{U})\Gamma(\mathbf{U}, \mathbf{U}) + \gamma I)^{-1}\Gamma(\mathbf{U}, \phi(u')))(K(\varphi, \varphi))^{-1}\varphi(v)$$

Furthermore, the law of ξ conditioned on $\varphi \circ \xi(\mathbf{u}) = \varphi(\mathbf{v}) + Z$ is equivalent to that of $\chi \circ \xi^\perp \circ \phi$ where $\xi^\perp \sim \mathcal{N}(\bar{f}_\gamma, \Gamma^\perp)$ is the GP ξ conditioned on $\zeta(\mathbf{U}) = \mathbf{V} + Z'$, whose mean is \bar{f}_γ as in (2.12) and conditional covariance kernel is

$$\Gamma^\perp(U, U') = \Gamma(U, U') - \Gamma(U, \mathbf{U})\Gamma(\mathbf{U}, \mathbf{U}) + \gamma I)^{-1}\Gamma(\mathbf{U}, U').$$

We also use the GP approach to derive an alternative regularization of (2.14) in Appendix B.

2.4. Measurement and mesh invariance

As argued in [45], mesh invariance is a key property for operator learning methods, i.e., the learned operator should be generalizable at test time beyond the specific discretization that was used during training. In our framework, this translates to being able to predict the output of a test input function \tilde{u} given only a linear measurement $\tilde{\phi}(\tilde{u})$, where $\tilde{\phi}$ was unknown at training time. For example $\tilde{\phi}$ could be of the same form as ϕ (say (1.4)) but on a finer or coarser grid. Similarly, we may choose to output with an operator $\tilde{\varphi}$ which is a coarse/fine version of φ . Our proposed framework can easily provide mesh invariance using additional optimal recovery and measurement operators at the input and outputs of the operator $\tilde{\mathcal{G}}$ as depicted in Fig. 3. In fact, we can not only accommodate modification of the grid but completely different measurement operators at testing time. For example, while ϕ, φ may be of the form (1.4) we may take $\tilde{\phi}$ and $\tilde{\varphi}$ to be integral operators such as Fourier or Radon transforms.

Let us describe our approach to mesh invariance in detail. Given bounded and linear operators $\tilde{\phi} : \mathcal{U} \rightarrow \mathbb{R}^{\tilde{n}}$ and $\tilde{\varphi} : \mathcal{V} \rightarrow \mathbb{R}^{\tilde{m}}$ we can approximate $\tilde{\varphi}(\tilde{\mathcal{G}}^*(\tilde{u}))$ using the map \tilde{f} obtained from (2.3) defined in terms of our training. To achieve mesh invariance we simply need a consistent approach to interpolate/extend the testing measurement operators to those used for training and we achieve this using the optimal recovery map $\tilde{\psi}$ that is defined from $\tilde{\phi}$ analogously to ψ in (2.1). This setup gives rise to a natural approximation of $\tilde{\mathcal{G}}^*$ in terms of the function $h^\dagger : \mathbb{R}^{\tilde{n}} \rightarrow \mathbb{R}^{\tilde{m}}$ depicted in Fig. 3 which in turn can be approximated with $\tilde{h} := \tilde{\varphi} \circ \chi \circ \tilde{f} \circ \phi \circ \tilde{\psi} \equiv \tilde{\varphi} \circ \tilde{\mathcal{G}} \circ \tilde{\psi}$. This expression further gives rise to another approximation to $\tilde{\mathcal{G}}^*$ given by the operator $\tilde{\mathcal{G}} = \tilde{\chi} \circ \tilde{h} \circ \tilde{\phi}$.

Remark 2.1. Observe that the definition of \tilde{h} (and consequently $\tilde{\mathcal{G}}$) is independent of the fact that \tilde{f} is constructed using the kernel approach. Thus, the optimal recovery maps χ and $\tilde{\psi}$ can be used to retrofit any fixed-mesh operator learning algorithm, to become mesh-invariant and able to use arbitrary linear measurements of the function \tilde{u} at test time.

3. Convergence and error analysis

In this section, we present convergence guarantees and rigorous a priori error bounds for our proposed kernel method for operator learning and give a detailed statement and proof of Theorem 3.3. We assume that \mathcal{H}_Q is a space of continuous functions from $\Omega \subset \mathbb{R}^{d_\Omega}$ and that \mathcal{H}_K is a space of continuous functions from $D \subset \mathbb{R}^{d_D}$. Abusing notations we write $Q : \Omega \times \Omega \rightarrow \mathbb{R}^{d_\Omega}$ and $K : D \times D \rightarrow \mathbb{R}^{d_D}$ for the kernels induced by the operators Q and K . Let $X = (X_1, \dots, X_n) \subset \Omega$ and $Y = (Y_1, \dots, Y_m) \subset D$ be distinct collections of points and define their fill-distances

$$h_X := \max_{x' \in \Omega} \min_{x \in X} |x - x'|, \quad h_Y := \max_{y' \in D} \min_{y \in Y} |y - y'|.$$

This section focuses on operators ϕ and φ that are linear combinations of pointwise measurements in X and Y . The presented results can be extended by using analogs of the sampling inequalities for other linear measurements, see [61, Theorem 4.11, Lemma 14.34] for a general framework that allows one to obtain such inequalities.

Let L_Q and L_K be invertible $n \times n$ and $m \times m$ matrices. For $u \in \mathcal{H}_Q$ write $u(X)$ for the n -vector with entries $u(X_i)$ and let $\phi : \mathcal{H}_Q \rightarrow \mathbb{R}^n$ be the bounded linear map defined by

$$\phi(u) = L_Q u(X). \quad (3.1)$$

For $v \in \mathcal{H}_K$ write $v(Y)$ for the m -vector with entries $v(Y_j)$ and let $\varphi : \mathcal{H}_K \rightarrow \mathbb{R}^m$ be the bounded linear map defined by

$$\varphi(v) = L_K v(Y). \quad (3.2)$$

Write $\|\phi\| := \sup_{u \in \mathcal{H}_Q} |\phi(u)|/\|u\|_Q$ and $\|\psi\| := \sup_{U' \in \mathbb{R}^n} \|\psi(U')\|_Q/|U'|$, and similarly $\|\varphi\| := \sup_{v \in \mathcal{H}_K} |\varphi(v)|/\|v\|_K$ and $\|\chi\| := \sup_{V' \in \mathbb{R}^m} \|\chi(V')\|_K/|V'|$. We will also assume the following regularity conditions on the domains Ω, D , the kernels Q, K , and the operator \mathcal{G}^\dagger .

Condition 3.1. Assume that the following conditions hold.)

(3.1.1) Ω and D are compact sets with Lipschitz boundary.

(3.1.2) There exist indices $s > d_\Omega/2$ and $t > d_D/2$ so that $\mathcal{H}_Q \subset H^s(\Omega)$ and $\mathcal{H}_K \subset H^t(D)$, with inclusions indicating continuous embeddings.

(3.1.3) \mathcal{G}^\dagger is a (possibly) nonlinear operator from $H^{s'}(\Omega)$ to \mathcal{H}_K with $s' < s$ that satisfies,

$$\|\mathcal{G}^\dagger(u) - \mathcal{G}^\dagger(v)\|_K \leq \omega\left(\|u - v\|_{H^{s'}(\Omega)}\right), \quad (3.3)$$

where $\omega : \mathbb{R} \rightarrow \mathbb{R}_+$ is the modulus of continuity of \mathcal{G}^\dagger .

Note that conditions (3.1.2) and 3.3 imply

$$\|\mathcal{G}^\dagger\|_{B_R(\mathcal{H}_Q) \rightarrow \mathcal{H}_K} := \sup_{u \in B_R(\mathcal{H}_Q)} \|\mathcal{G}^\dagger(u)\|_K < +\infty.$$

Proposition 3.1. Suppose that Condition 3.1 holds. Let $0 < t' < t$. Then there exist constants $h_\Omega, h_D, C_\Omega, C_D > 0$ such that if $h_X < h_\Omega$ and $h_Y < h_D$, then

$$\|\mathcal{G}^\dagger(u) - \chi \circ f^\dagger \circ \phi(u)\|_{H^{t'}(D)} \leq C_D \omega\left(C_\Omega h_X^{s-s'} R\right) + C_D h_Y^{t-t'} (\|\mathcal{G}^\dagger(0)\|_K + \omega(C_\Omega R)),$$

for any $u \in B_R(\mathcal{H}_Q)$, where f^\dagger is defined as in (1.5).

Proof. By the definition of f^\dagger and the triangle inequality we have

$$\begin{aligned} \|\mathcal{G}^\dagger(u) - \chi \circ \varphi \circ \mathcal{G}^\dagger \circ \psi^\dagger \circ \phi(u)\|_{H^{t'}(\Gamma)} &\leq \|\mathcal{G}^\dagger(u) - \mathcal{G}^\dagger \circ \psi \circ \phi(u)\|_{H^{t'}(\Gamma)} \\ &\quad + \|\mathcal{G}^\dagger \circ \psi \circ \phi(u) - \chi \circ \varphi \circ \mathcal{G}^\dagger \circ \psi \circ \phi(u)\|_{H^{t'}(\Gamma)} \\ &=: T_1 + T_2. \end{aligned}$$

Let us first bound T_1 : By conditions (3.1.2) and 3.3, we have

$$T_1 \leq C_D \|\mathcal{G}^\dagger(u) - \mathcal{G}^\dagger \circ \psi \circ \phi(u)\|_K \leq C_D \omega\left(\|u - \psi \circ \phi(u)\|_{H^{s'}(\Omega)}\right).$$

At the same time, since $(u - \psi \circ \phi(u))(X) = 0$, condition (3.1.1) and the sampling inequality for interpolation in Sobolev spaces [5, Thm. 4.1], and condition (3.1.2) imply that there exists a constant $h_\Omega > 0$ so that if $h_X < h_\Omega$ then

$$\|u - \psi \circ \phi(u)\|_{H^{s'}(\Omega)} \leq C'_\Omega h_X^{s-s'} \|u - \psi \circ \phi(u)\|_{H^s(\Omega)} \leq C_\Omega h_X^{s-s'} \|u - \psi \circ \phi(u)\|_Q, \quad (3.4)$$

where $C'_\Omega, C_\Omega > 0$ are constants that are independent of u . Using $\|u - \psi \circ \phi(u)\|_Q \leq \|u\|_Q$ [61, Thm. 12.3] we deduce the desired bound

$$T_1 \leq C_D \omega\left(C_\Omega h_\Omega^{s-s'} \|u\|_Q\right). \quad (3.5)$$

Let us now bound T_2 : Once again, by the continuous embedding of condition (3.1.2) and the sampling inequality for interpolation in Sobolev spaces, we have that, there exists $h_D > 0$ so that if $h_Y < h_D$, then for any $v \in H^t(D)$ it holds that

$$\|v - \chi \circ \varphi(v)\|_{H^{t'}(D)} \leq C'_D h_Y^{t-t'} \|v - \chi \circ \varphi(v)\|_{H^t(D)} \leq C_D h_Y^{t-t'} \|v - \chi \circ \varphi(v)\|_K \leq C_D h_Y^{t-t'} \|v\|_K.$$

Taking $v \equiv \mathcal{G}^\dagger \circ \psi \circ \phi(u)$, we deduce that

$$\begin{aligned} T_2 &\leq C_D h_Y^{t-t'} \|\mathcal{G}^\dagger \circ \psi \circ \phi(u)\|_K, \\ &\leq C_D h_Y^{t-t'} (\|\mathcal{G}^\dagger(0)\|_K + \omega(\|\psi \circ \phi(u)\|_{H^{s'}(\Omega)})). \end{aligned}$$

Using $\|\psi \circ \phi(u)\|_{H^{s'}(\Omega)} \leq C_\Omega \|\psi \circ \phi(u)\|_Q \leq C_\Omega \|u\|_Q$ concludes the proof. ■

While Proposition 3.1 gives an error bound for the distance between the maps \mathcal{G}^\dagger and $\varphi \circ f^\dagger \circ \phi$, we can never compute this map when $N < \infty$ and so we have to approximate this map as well. Given the kernel Γ , our optimal recovery approximant for the map f^\dagger is \bar{f} as in (2.3), which we recall is the minimizer of (2.11).

To proceed, we need to consider another intermediary problem that defines an approximation \hat{f} to the map f^\dagger :

$$\hat{f} := \begin{cases} \text{Minimize} & \|f\|_\Gamma^2 \\ \text{Over} & f \in \mathcal{H}_\Gamma \text{ such that } f \circ \phi(\mathbf{u}) = f^\dagger \circ \phi(\mathbf{u}). \end{cases} \quad (3.6)$$

We emphasize that the difference between the problems (2.11) and (3.6) is simply in the training data that is injected in the equality constraints, and this difference is quite subtle:

In practical applications, observations may be taken from $\mathcal{G}^\dagger(u_i)$, which is different from $f^\dagger \circ \phi(u_i) \equiv \varphi \circ \mathcal{G}^\dagger \circ \psi \circ \phi(u_i)$. To make our analysis simple, henceforth we assume the following condition on our input data.

Condition 3.2. The input data points u_i satisfy

$$u_i = \psi \circ \phi(u_i) \text{ for } i = 1, \dots, N,$$

We observe that this condition implies $\mathcal{G}^\dagger(u_i) = f^\dagger \circ \phi(u_i)$ and $\bar{f} = \hat{f}$. Removing this assumption requires bounding some norm of the error $f^\dagger - \bar{f}$, and we postpone that analysis to a sequel paper as this step can become very technical.

The next step in our convergence analysis is then to control the error between the maps \hat{f} and f^\dagger which we will achieve using similar arguments as in the proof of Proposition 3.1. For our analysis, we take Γ to be a diagonal, matrix-valued kernel, of the form (1.8) which we recall for reference

$$\Gamma(U, U') = S(U, U')I \quad (3.7)$$

where I is the $m \times m$ identity matrix and $S : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ is a real valued kernel.

Proposition 3.2. Suppose that Condition 3.2 holds. Let $Y \subset \mathbb{R}^n$ be a compact set with Lipschitz boundary and consider $U = (U_1, \dots, U_N) \subset Y$ with fill distance

$$h_Y := \max_{U' \in Y} \min_{1 \leq i \leq N} |U_i - U'|.$$

Let Γ be of the form (3.7), with S restricted to the set Y , and suppose $H_S \subset H^r(Y)$ for $r > n/2$ and that $f_j^\dagger \in H_S$ for $j = 1, \dots, m$. Then there exist constants $h_Y', C_Y > 0$ so that whenever $h_Y < h_Y'$ then for any $r' < r$ it holds that

$$\|f_j^\dagger - \hat{f}_j\|_{H^{r'}(Y)} \leq C_Y h_Y^{r-r'} \|f_j^\dagger\|_S.$$

Proof. The proof is a direct consequence of the fact that the components of \hat{f} are given by the optimal recovery problems (3.6) and the sampling inequality for interpolation in Sobolev spaces [5, Thm. 4.1] following the same arguments used in the proof of Theorem 3.1. ■

We can now combine the above results to obtain the following theorem.

Theorem 3.3. Suppose that Conditions 3.1 and 3.2 hold in addition to those of Proposition 3.2 with a set of inputs $(u_i)_{i=1}^N \subset B_R(\mathcal{H}_Q)$, the set $Y = \phi(B_R(\mathcal{H}_Q))$, and index $n/2 < r' < r$. Then for any $u \in B_R(\mathcal{H}_Q)$, it holds that

$$\begin{aligned} \|\mathcal{G}^\dagger(u) - \chi \circ \bar{f} \circ \phi(u)\|_{H^{r'}(D)} &\leq C_D \omega\left(C_\Omega h_X^{s-s'} R\right) + C_D h_Y^{t-t'} (\|\mathcal{G}^\dagger(0)\|_K + \omega(C_\Omega R)) \\ &\quad + \sqrt{m} C_D C_Y \|\chi\| h_Y^{(r-r')} \max_{1 \leq j \leq m} \|f_j^\dagger\|_S \end{aligned} \quad (3.8)$$

Proof. An application of the triangle inequality yields

$$\begin{aligned} \|\mathcal{G}^\dagger(u) - \chi \circ \bar{f} \circ \phi(u)\|_{H^{r'}(D)} &\leq \|\mathcal{G}^\dagger(u) - \chi \circ f^\dagger \circ \phi(u)\|_{H^{r'}(D)} \\ &\quad + \|\chi \circ f^\dagger \circ \phi(u) - \chi \circ \hat{f} \circ \phi(u)\|_{H^{r'}(D)} \\ &\quad + \|\chi \circ \hat{f} \circ \phi(u) - \chi \circ \bar{f} \circ \phi(u)\|_{H^{r'}(D)} =: I_1 + I_2 + I_3. \end{aligned}$$

We can bound I_1 immediately using Proposition 3.1. Furthermore, by Condition 3.2 we have that $I_3 = 0$. So it remains for us to bound I_2 : By the continuous embedding of \mathcal{H}_K into $H^{r'}(D)$ we can write

$$\begin{aligned} I_2 &\leq C_D \|\chi \circ f^\dagger \circ \phi(u) - \chi \circ \hat{f} \circ \phi(u)\|_K \leq C_D \|\chi\| \|f^\dagger \circ \phi(u) - \hat{f} \circ \phi(u)\| \\ &\leq C_D \|\chi\| \sqrt{\sum_{j=1}^m \|f_j^\dagger - \hat{f}_j\|_{H^{r'}(Y)}^2}, \end{aligned}$$

where the last line follows from the Sobolev embedding theorem and the assumption that $r' > n/2$. Then an application of Proposition 3.2 yields,

$$I_2 \leq \sqrt{m} C_D C_Y \|\chi\| h_Y^{(r-r')} \max_{1 \leq j \leq m} \|f_j^\dagger\|_S. \quad \blacksquare$$

3.1. Convergence theorem

Our next step will be to consider the limits $N, n, m \rightarrow \infty$ and show the convergence of $\bar{\mathcal{G}}$ to \mathcal{G}^\dagger . To obtain this result we first need to make assumptions on the regularity of the true operator \mathcal{G}^\dagger .

For $k \geq 1$ write $D^k \mathcal{G}^\dagger$ for the functional derivative of \mathcal{G}^\dagger of order k . Recall that for $u \in \mathcal{H}_Q$, $D^k \mathcal{G}^\dagger(u)$ is a multilinear operator mapping $\otimes_{i=1}^k \mathcal{H}_Q$ to \mathcal{H}_K . For $w_1, \dots, w_k \in \mathcal{H}_Q$ write $[D^k \mathcal{G}^\dagger(u), \otimes_{i=1}^k w_i]$ for the (multilinear) action of $D^k \mathcal{G}^\dagger(u)$ on $\otimes_{i=1}^k w_i$ and write $\|D^k \mathcal{G}^\dagger(u)\|$ for the smallest constant such that for $w_1, \dots, w_k \in \mathcal{H}_Q$,

$$\left\| [D^k \mathcal{G}^\dagger(u), \otimes_{i=1}^k w_i] \right\|_{\mathcal{H}_K} \leq \|D^k \mathcal{G}^\dagger(u)\| \prod_{i=1}^k \|w_i\|_{\mathcal{H}_Q} \quad (3.9)$$

Similarly, for $k \geq 1$ write $D^k f^\dagger$ for the derivation tensor of f^\dagger of order k (the gradient for $k = 1$ and the Hessian for $k = 2$, etc). Recall that for $U \in \mathbb{R}^n$, $D^k f^\dagger(U)$ is a multilinear operator mapping $\otimes_{i=1}^k \mathbb{R}^n$ to \mathbb{R}^m . For $W_1, \dots, W_k \in \mathbb{R}^n$ write $[D^k f^\dagger(U), \otimes_{i=1}^k W_i]$ for the (multilinear) action of $D^k f^\dagger(U)$ on $\otimes_{i=1}^k W_i$ and write $\|D^k f^\dagger(U)\|$ for the smallest constant such that for $W_1, \dots, W_k \in \mathbb{R}^n$,

$$\left| [D^k f^\dagger(U), \otimes_{i=1}^k W_i] \right| \leq \|D^k f^\dagger(U)\| \prod_{i=1}^k |W_i| \quad (3.10)$$

where $|\cdot|$ is the Euclidean norm.

Lemma 3.3. *It holds true that $\|D^k f^\dagger(U)\| \leq \|\varphi\| \|\psi\|^k \|D^k \mathcal{G}^\dagger \circ \psi(U)\|$, $\forall U \in \mathbb{R}^n$.*

Proof. The chain rule and the linearity of φ and ψ imply that

$$[D^k f^\dagger(U), \otimes_{i=1}^k W_i] = \varphi[D^k \mathcal{G}^\dagger \circ \psi(U), \otimes_{i=1}^k \psi(W_i)].$$

We then conclude the proof by writing

$$\begin{aligned} \left| [D^k f^\dagger(U), \otimes_{i=1}^k W_i] \right| &\leq \|\varphi\| \|D^k \mathcal{G}^\dagger \circ \psi(U)\| \prod_{i=1}^k \|\psi(W_i)\|_{\mathcal{H}_Q} \\ &\leq \|\varphi\| \|\psi\|^k \|D^k \mathcal{G}^\dagger \circ \psi(U)\| \prod_{i=1}^k |W_i|. \quad \blacksquare \end{aligned}$$

Let us now consider an infinite and dense sequence of points X_1, X_2, X_3, \dots of Ω , such that the closure of $\cup_{i=1}^\infty \{X_i\}$ is the closure of Ω . Write X^n for the n -vector formed by the first n points, i.e.,

$$X^n := (X_1, \dots, X_n) \quad (3.11)$$

and let L_Q^n be an arbitrary invertible $n \times n$ matrix. Further let $\phi^n : \mathcal{H}_Q \rightarrow \mathbb{R}^n$ be defined by

$$\phi^n(u) = L_Q^n u(X^n). \quad (3.12)$$

Write ψ^n for the corresponding optimal recovery ψ -map. Similarly, we assume that we are given an infinite and dense sequence of points Y_1, Y_2, Y_3, \dots of D , such that the closure of $\cup_{i=1}^\infty \{Y_i\}$ is the closure of D . Write Y^m for the m -vector formed by the first m points, i.e.,

$$Y^m := (Y_1, \dots, Y_m). \quad (3.13)$$

Let L_K^m be an arbitrary invertible $m \times m$ matrix and let $\varphi^m : \mathcal{H}_K \rightarrow \mathbb{R}^m$ be defined by

$$\varphi^m(v) = L_K^m v(Y^m). \quad (3.14)$$

Write χ^m for the corresponding optimal recovery χ -map. We also assume that we are given a sequence of diagonal matrix-valued kernels $\Gamma^{m,n} : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathcal{L}(\mathbb{R}^m)$ with scalar-valued kernels $S^n : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ as diagonal entries. Write $\bar{f}_N^{m,n}$ for the corresponding minimizer of (2.11) (also identified by the formula (2.3)) for the above setup.

Theorem 3.4. *Let m, n be the dimensionality of the input and output observations $\phi : \mathcal{U} \rightarrow \mathbb{R}^n$ and $\varphi : \mathcal{V} \rightarrow \mathbb{R}^m$. Suppose that the closure of $\lim_{n \uparrow \infty} \cup_{i=1}^n \{X_i\}$ is equal to the closure of Ω and that the closure of $\lim_{m \uparrow \infty} \cup_{i=1}^m \{Y_i\}$ is equal to the closure of D . Suppose Condition 3.1 is satisfied and that*

$$\sup_{u \in B_R(\mathcal{H}_Q)} \|D^k \mathcal{G}^\dagger(u)\| < \infty \text{ for all } k \geq 1, \quad (3.15)$$

for an arbitrary $R > 0$. Assume that for any $n \geq 1$ and any compact set Y of \mathbb{R}^n , the RKHS of S^n restricted to Y (which we write $\mathcal{H}_{S^n}(Y)$) is contained in $H^r(Y)$ for some $r > n/2$ and contains $H^{r'}(Y)$ for some $r' > 0$ that may depend on n . Let $(u_i)_{i=1}^N$ be a sequence of inputs in $B_R(\mathcal{H}_Q)$. Assume that there exists an integer n_0 such that for $n \geq n_0$, the data points $(u_i)_{i=1}^N$ satisfy Condition 3.2, i.e., they satisfy $u_i = \psi^n \circ \phi^n(u_i)$ for all $i \geq 1$. Further assume that the $(\phi^n(u_i))_{1 \leq i \leq N}$ are space filling in the sense that for any $n \geq n_0$ we have

$$\lim_{N \rightarrow \infty} \sup_{u \in B_R(\mathcal{H}_Q)} \min_{1 \leq i \leq N} |u_i(X^n) - u(X^n)| = 0. \quad (3.16)$$

Then for any $t' \in (0, t)$, it holds that

$$\lim_{n, m \rightarrow \infty} \lim_{N \rightarrow \infty} \sup_{u \in B_R(\mathcal{H}_Q)} \|\mathcal{G}^\dagger(u) - \chi^m \circ \bar{f}_N^{m, n} \circ \phi^n(u)\|_{H^{t'}(D)} = 0 \quad (3.17)$$

Proof. Following [61, Chap. 12.1] define the projection $P_n^U = \psi^n \circ \phi^n$ onto the range of ψ^n . Since the points X_i and Y_j are dense in Ω and D we have $h_{Y^n} \downarrow 0$ as $n \rightarrow \infty$ and $h_{Y^m} \downarrow 0$ as $m \rightarrow \infty$. Given n , take $Y = \phi^n(B_R(\mathcal{H}_Q))$. Then Lemma 3.3 and (3.15) imply that $\bar{f}_j^{m, n} \in H^{r'}(Y)$ for all $r' \geq 0$. Therefore $\bar{f}_j^{m, n} \in \mathcal{H}_{S^n}(Y)$. Now (3.16) implies that for any n , the fill distance, in $\phi^n(B_R(\mathcal{H}_Q))$, between the points $(\phi^n(u_i))_{1 \leq i \leq N}$ goes to zero as $N \rightarrow \infty$. Since the conditions of Proposition 3.2 are satisfied, we conclude by taking the limit $N \rightarrow \infty$ in (3.8) before taking the limit $m, n \rightarrow \infty$. ■

3.2. The effect of the L_Q and L_K preconditioners

We conclude this section and our discussion of convergence results, by highlighting the importance of the choice of the matrices L_Q^n and L_K^m in (3.12) and (3.14). It is clear from the bounds (3.8) and (3.10) that our error estimates depend on the norms of the linear operators φ^m, ψ^n and χ^m . To ensure that those norms do not blow up as $n, m \rightarrow \infty$ we can select the matrices L_Q^n and L_K^m to be the Cholesky factors of the precision matrices obtained from pointwise measurements of the kernels Q and K , i.e.,

$$L_Q^n (L_Q^n)^T = Q(X^n, X^n)^{-1} \quad \text{and} \quad L_K^m (L_K^m)^T = K(Y^m, Y^m)^{-1}. \quad (3.18)$$

We now obtain the following proposition.

Proposition 3.4. *If ϕ^n is as in (3.12) and L_Q^n as in (3.18), then $\|\phi^n\| = 1$ and $\|\psi^n\| = 1$. If φ^m is as in (3.2) and L_K^m as in (3.18), then $\|\varphi^m\| = 1$ and $\|\chi^m\| = 1$.*

Proof. For $u \in \mathcal{H}_Q$, $|\phi^n(u)|^2 = u(X^n)^T Q(X^n, X^n)^{-1} u(X^n) = \|\psi^n \circ \phi^n(u)\|_Q^2$. Since $\psi^n \circ \phi^n$ is a projection [61, Chap. 12.1] we deduce that $\|\phi^n\| = 1$. Using $\psi^n(U') = Q(\cdot, X^n) L_Q^n U'$ leads to $\|\psi^n(U')\|_Q^2 = |U'|^2$ and $\|\psi^n\| = 1$. The proof of $\|\varphi^m\| = 1$ and $\|\chi^m\| = 1$ is similar. ■

We note that although useful for obtaining tighter approximation errors, this particular choice for the matrices L_Q^n and L_K^m is not required for convergence if one first takes the limit $N \rightarrow \infty$ as in Theorem 3.4, which does not put any requirements on the matrices L_Q^n and L_K^m beyond invertibility.

4. Numerics

In this section, we present numerical experiments and benchmarks that compare a straightforward implementation of our kernel operator learning framework to state-of-the-art NN-based techniques. We discuss some implementation details of our method in Subsection 4.1 followed by the setup of experiments and test problems in Subsections 4.2 and 4.3. A detailed discussion of our findings is presented in Subsection 4.4.

4.1. Implementation considerations

Below we summarize some of the key details in the implementation of our kernel approach for operator learning for benchmark examples. Our code to reproduce the experiments can be found in a public repository.⁶

4.1.1. Choice of the kernel Γ

Following our theoretical discussions in Sections 2 and 3, we primarily take Γ to be a diagonal kernel of the form (3.7). This implies that our estimation of \bar{f} can be split into independent problems for each of its components \bar{f}_j in the RKHS of the scalar kernel S . In our experiments, we investigate different choices of S belonging to the families of the linear kernel, rational quadratic, and Matérn; see Appendix C for detailed expressions of these kernels. The rational quadratic kernel has two parameters, the lengthscale

⁶ <https://github.com/MatthieuDarcy/KernelsOperatorLearning/>.

Table 2

Summary of datasets used for benchmarking. The first three examples were considered in [47], and the last four were taken from [19].

Equation	Input	Output	Input Distribution μ
Burger's	Initial condition	Solution at time T	Gaussian field (GF)
Darcy problem	Coefficient	Solution	Binary function of GF
Advection I	Initial condition	Solution at time T	Random square waves
Advection II	Initial condition	Solution at time T	Binary function of GF
Helmholtz	Coefficient	Solution	Function of Gaussian field
Structural mechanics	Initial force	Stress field	Gaussian field
Navier Stokes	Forcing term	Solution at time T	Gaussian field

l and the exponent α . We tuned these parameters using standard cross validation or log marginal likelihood maximization over the training data (see [67, p.112] for a detailed description). The Matérn kernel is parameterized by two positive parameters: a smoothness parameter ν and the length scale l . The smoothness parameter ν controls the regularity of the RKHS and we considered $\nu \in \{\frac{1}{2}, \frac{3}{2}, \frac{5}{2}, \frac{7}{2}, \infty\}$. In practice we found that $\nu = \frac{5}{2}$ almost always had the best performance. For a fixed choice of ν we tuned the length scale l similarly to the rational quadratic kernel. We implemented the kernel regressions of the \tilde{f}_j and parameter tuning algorithms in scikit-learn for low-dimensional examples and manually in JAX for high-dimensional examples.

4.1.2. Preconditioning and dimensionality reduction

Following (3.1) and (3.2) and the discussion in Subsection 3.2, we consider two preconditioning strategies for our pointwise measurements, i.e., choices of the matrices L_Q and L_K : (1) we consider the Cholesky factors of the underlying covariance matrices as in (3.18); (2) we use PCA projection matrices of the input and output functions computed from the training data. We truncated the PCA expansions to preserve (0.90, 0.95, 0.99) of the variance. The use of PCA in learning mappings between infinite dimensional spaces was proposed in [43] and recently revisited in [9,33].

4.2. Experimental setup

We compare the test performance of our method with different choices of the kernel S of increasing complexity using the examples in [19] and [47] and their reported test relative L^2 loss (see (4.2) below). We use the data provided by these papers for the training set and the test set.⁷ Both articles provide performance comparisons between different variants of Neural Operators (most notably FNO and DeepONet) on a variety of PDE operator learning tasks, where the data is sampled independently from a distribution $(Id, \mathcal{G}^\dagger)^\# \mu$ supported on $\mathcal{U} \times \mathcal{V}$, where μ is a specified (input) distribution on \mathcal{U} . The example problems are outlined in detail in Subsection 4.3; a summary of the specific PDEs, problem type, and distribution μ for each test is given in Table 2. In some instances the train-test split of the data was not clear from the available online repositories in which case we re-sampled them from the assumed distribution μ . The datasets from [47] contain 1000 training data-points per problem (which we will refer to as the “low-data regime”), whereas the datasets from [19] contain 20000 training data-points (which we will refer to as the “high-data” regime). We make this distinction because the complexity of kernel methods, unlike that of neural networks, may depend on the number of data-points.

Following the suggestion of [19] we not only compare test errors and training complexity but also the complexity of operator learning at the inference/evaluation stage in Subsection 4.2.2. For the examples in [19], we investigate the accuracy-complexity trade-off of our method against the reported values of that article.

4.2.1. Measures of accuracy

As our first performance metric we measured the accuracy of models by a relative loss on the output space \mathcal{V} :

$$\mathcal{R}(\mathcal{G}) = \mathbb{E}_{u \sim \mu} \left[\frac{\|\mathcal{G}^\dagger(u) - \mathcal{G}(u)\|_{\mathcal{V}}}{\|\mathcal{G}^\dagger(u)\|_{\mathcal{V}}} \right] \quad (4.1)$$

where \mathcal{G}^\dagger is true operator and \mathcal{G} is a candidate operator. Following previous works, we often took $\|u\|_{\mathcal{V}} = \|u\|_{L^2} := (\int u(x)^2 dx)^{\frac{1}{2}}$, which in turn is discretized using the trapezoidal rule. In practice, we do not have the access to the underlying probability measure μ and we compute the empirical loss on a withheld test set:

$$\mathcal{R}_N(\mathcal{G}) = \frac{1}{N} \sum_{n=1}^N \left[\frac{\|\mathcal{G}^\dagger(u^n) - \mathcal{G}(u^n)\|_{\mathcal{V}}}{\|\mathcal{G}^\dagger(u^n)\|_{\mathcal{V}}} \right], \quad u^i \sim \mu. \quad (4.2)$$

4.2.2. Measures of complexity

For our second performance metric we considered the complexity of operator learning algorithms at the inference stage (i.e., evaluating the learned operator). Complexity at inference time is the main metric used in [19] to compare numerical methods for operator learning. The motivation is that training of the methods can be performed in an offline fashion, and therefore the cost

⁷ See <https://github.com/Zhengyu-Huang/Operator-Learning> and <https://github.com/lu-group/deepnet-fno>, respectively, for the data.

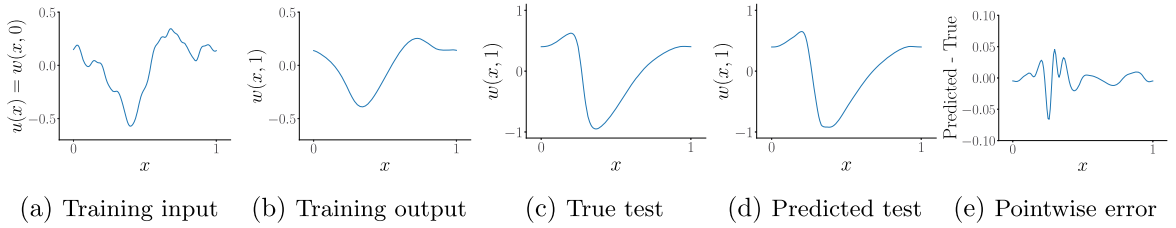


Fig. 4. Example of training data and test prediction and pointwise errors for the Burger's equation (4.3).

per test example dominates in the limit of many test queries. In particular, they compare the online evaluation costs of the neural networks by computing the requisite floating point operations (FLOPs) per test example. We adopt this metric as well for the methods not based on neural networks that we develop in this work, and we compare, when available, the cost-accuracy tradeoff with the numbers reported in [19]. We computed the FLOPs with the same assumptions as in the original work: a matrix-vector product where the input vector is in \mathbb{R}^n and the output vector is in \mathbb{R}^m amounts to $m(2-1)$ flops, and non-linear functions with n -dimensional inputs (activation functions for neural networks, kernel computations for kernel methods) are assumed to have cost $\mathcal{O}(n)$.

Remark 4.1 (Training complexity). While the inference complexity of a model eventually dominates the cost of training during applications, the training cost cannot be ignored since the allocated computational resources during this stage may still be limited and the resulting errors will have a profound impact on the quality and performance of the learned operators. Therefore numerical methods in which the offline data assimilation step is cheaper, faster, and more robust will always be preferred. Computing the exact number of FLOPs at training time is difficult to estimate for NN methods, as it depends on the optimization algorithms used, the hyperparameters and the optimization over such hyperparameters, among many other factors. Therefore in this work we limit the training complexity evaluation to the qualitative observation that kernel methods provided in this work are significantly simpler at training time, as they have no NN weights, they do not require the use of stochastic gradient descent, and have few or no hyperparameters which can be tuned using standard methods such as grid search or gradient descent in a low-dimensional space.

4.3. Test problems and qualitative results

Below we outline the setup of each of our benchmark problems. In all cases, \mathcal{U} and \mathcal{V} are spaces of real-valued functions with input domains $\Omega, D \subset \mathbb{R}^d$ for $k = 1$ or 2 . Whenever $\Omega = D$, we simply write D for both.

4.3.1. Burger's equation

Consider the one-dimensional Burger's equation:

$$\begin{aligned} \frac{\partial w}{\partial t} + w \frac{\partial w}{\partial x} &= \nu \frac{\partial^2 w}{\partial x^2}, \quad (x, t) \in (0, 1) \times (0, 1], \\ w(x, 0) &= u(x), \quad x \in (0, 1) \end{aligned} \quad (4.3)$$

with $D = (0, 1)$, and periodic boundary conditions. The viscosity parameter ν is set to 0.1. We learn the operator mapping the initial condition u to $v = w(\cdot, 1)$, the solution at time $t = 1$, i.e., $\mathcal{G}^\dagger : w(\cdot, 0) \mapsto w(\cdot, 1)$.

The training data is generated by sampling the initial condition u from a GP with a Riesz kernel, denoted by $\mu = \mathcal{GP}(0, 625(-\Delta + 25I)^{-2})$. As in [47], we used a spatial resolution with 128 grid points to represent the input and output functions, and used 1000 instances for training and 200 instances for testing. Fig. 4 shows an example of training input and output pairs as well as a test example along with its pointwise error.

4.3.2. Darcy flow

Consider the two-dimensional Darcy flow problem (1.3). Recall that in this example, we are interested in learning the mapping from the permeability field u to the solution v and the source term w is assumed to be fixed, hence $D \equiv \Omega = (0, 1)^2$ and $\mathcal{G}^\dagger : u \mapsto v$. The coefficient u is sampled by setting $u \sim \log \circ h_\mu \mu$ where $\mu = \mathcal{GP}(0, (-\Delta + 9I)^{-2})$ is a GP and h is binary function mapping positive inputs to 12 and negative inputs to 3. The resulting permeability/diffusion coefficient e^u is therefore piecewise constant. As in [47], we use a discretized grid of resolution 29×29 , with the data generated by the MATLAB PDE Toolbox. We use 1000 points for training and 200 points for testing. Fig. 2 shows an example of training input and output of the map \mathcal{G}^\dagger , and an example of predictions along with pointwise error at the test stage.

4.3.3. Advection equations (I and II)

Consider the one-dimensional advection equation:

$$\begin{aligned} \frac{\partial w}{\partial t} + \frac{\partial w}{\partial x} &= 0 \quad x \in (0, 1), t \in (0, 1] \\ w(x, 0) &= u(x) \quad x \in (0, 1) \end{aligned} \quad (4.4)$$

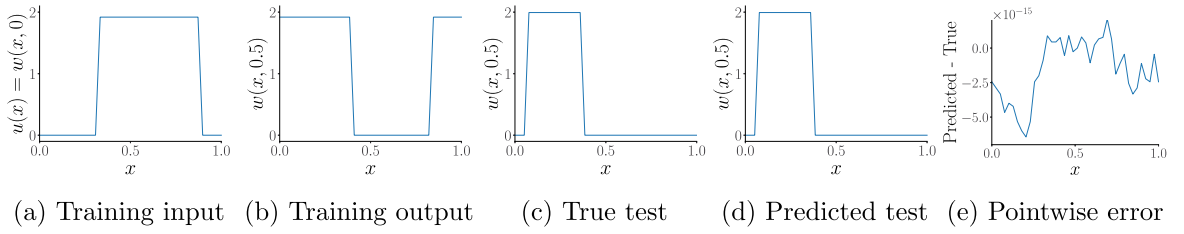


Fig. 5. Example of training data and test prediction and pointwise errors for the Advection problem (4.4)-I.

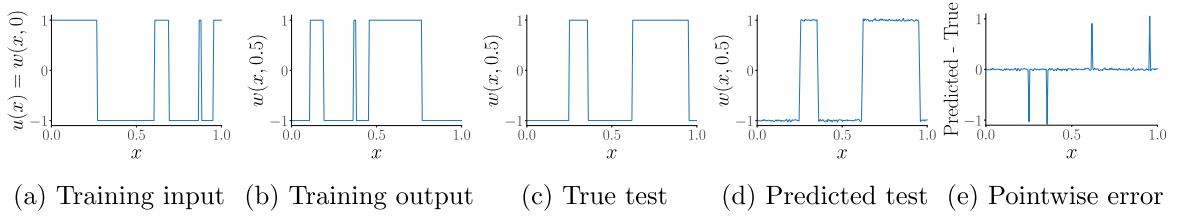


Fig. 6. Example of training data and test prediction and pointwise errors for the Advection problem (4.4)-II.

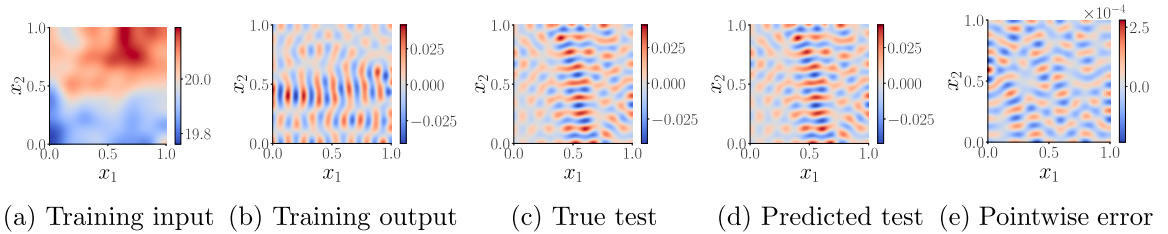


Fig. 7. Example of training data and test prediction and pointwise errors for the Helmholtz problem (4.7).

with $D = (0, 1)$ and periodic boundary conditions. Similar to the example for Burgers' equation, we learn the mapping from the initial condition u to $v = w(\cdot, 0.5)$, the solution at $t = 0.5$, i.e., $\mathcal{G}^t : w(\cdot, 0) \mapsto w(\cdot, 0.5)$.

This problem was considered in [47,19] with different distributions μ for the initial condition. We will show in the following section how these different distributions lead to different performances. In [47], henceforth referred to as Advection I, the initial condition is a square wave centered at $x = c$ of width b and height h :

$$u(x) = h \mathbf{1}_{[c-\frac{b}{2}, c+\frac{b}{2}]}, \quad (4.5)$$

where the parameters $(c, b, h) \sim \mathcal{U}([0.3, 0.7] \times [0.3, 0.6] \times [1, 2])$. In [19], henceforth referred to as Advection II, the initial condition is

$$u = -1 + 2\mathbf{1}_{\{\tilde{u}_0 \geq 0\}} \quad (4.6)$$

where $\tilde{u}_0 \sim \mathcal{GP}(0, (-\Delta + 3^2 I)^{-2})$.

For Advection I, the spatial grid was of resolution 40, and we used 1000 instances for training and 200 instances for testing. For Advection II, the resolution was of 200 and we used 20000 training and test instances, following [19].

Figs. 5 and 6 show an example of training input and output for Advection the I and II problems, respectively. Observe that the functional samples from the distribution in Advection I will have exactly two discontinuities almost surely, but the samples for Advection II can have many more jumps. We observe that prediction is challenging around discontinuities, hence Advection II is a significantly harder problem (across all benchmarked methods) than Advection I. Figs. 5 and 6 also show an instance of a test sample, along with a prediction and the pointwise errors.

4.3.4. Helmholtz's equation

For a given frequency ω and wavespeed field $u : D \rightarrow \mathbb{R}$, with $D = (0, 1)^2$, the excitation field $v : D \rightarrow \mathbb{R}$ solves

$$\left(-\Delta - \frac{\omega^2}{u^2(x)} \right) v = 0, \quad x \in (0, 1)^2 \quad (4.7)$$

$$\frac{\partial v}{\partial n} = 0, \quad x \in \{0, 1\} \times [0, 1] \cup [0, 1] \times \{0\} \quad \text{and} \quad \frac{\partial v}{\partial n} = v_N, \quad x \in [0, 1] \times \{1\}$$

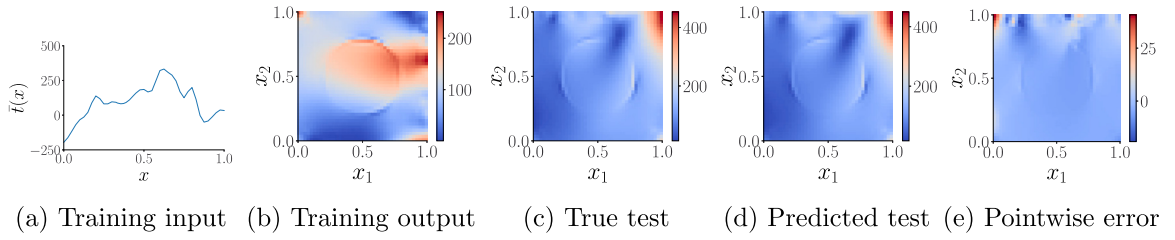


Fig. 8. Example of training data and test prediction and pointwise errors for the Structural Mechanics problem (4.8).

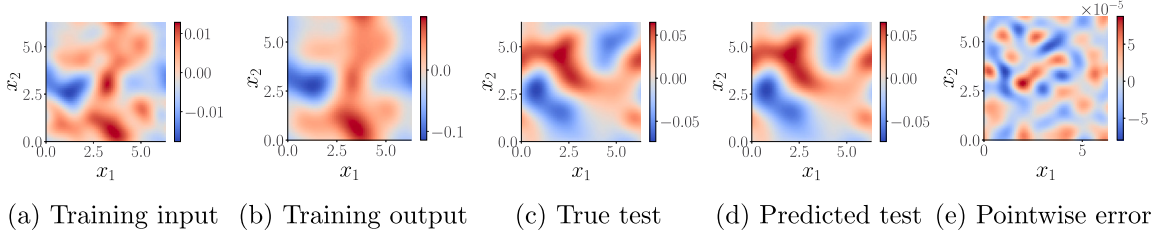


Fig. 9. Example of training data and test prediction and pointwise errors for the Navier-Stokes problem (4.9).

In the results that follow, we take $\omega = 10^3$, $v_N = \mathbf{1}_{\{0.35 \leq x \leq 0.65\}}$, and we aim to learn the map $\mathcal{G} : u \mapsto v$, i.e., the mapping from the wavespeed field to the excitation field. The distribution μ is specified as the law of $u(x) = 20 + \tanh(\tilde{u}(x))$, where \tilde{u} is drawn from the GP, $\mathcal{GP}(0, (-\Delta + 3^2 I)^{-2})$. The training and test data were generated by solving (4.7) with a Finite Element Method on a discretization of size 100×100 of the unit square. Fig. 7 shows an example of training input and output, a test prediction, and pointwise errors.

4.3.5. Structural mechanics

We let $\Omega = [0, 1] \times [0, 1]$, $D = [0, 1]^2$, the equation that governs the displacement vector w in an elastic solid undergoing infinitesimal deformations is

$$\nabla \cdot \sigma = 0 \quad \text{in } (0, 1)^2, \quad w = \bar{w}, \quad \text{on } \Gamma_w, \quad \nabla \cdot n = u \quad \text{on } \Gamma_u, \quad (4.8)$$

where the boundary ∂D is split in $[0, 1] \times \{0, 1\} = \Gamma_l$ (the part of the boundary subject to stress) and its complement Γ_u .

The goal is to learn the operator that maps the one-dimensional load u on Γ_u to the two-dimensional von Mises stress field v on Ω , i.e., $\mathcal{G} : u \mapsto v$. Here the distribution μ is $\mathcal{GP}(100, 400^2(-\Delta + 3^2 I)^{-1})$, with Δ being the Laplacian subject to homogeneous Neumann boundary conditions on the space of zero-mean functions. The function v was obtained by a finite element code, see [19] for implementation details and the constitutive model used. Fig. 8 shows an example of training input and outputs, a test prediction, and pointwise errors.

4.3.6. Navier-Stokes equations

Consider the vorticity-stream (ω, ψ) formulation of the incompressible Navier-Stokes equations:

$$\frac{\partial \omega}{\partial t} + (c \cdot \nabla) \omega - \nu \Delta \omega = u, \quad \omega = -\Delta \psi, \quad \int_D \psi = 0, \quad c = \left(\frac{\partial \psi}{\partial x_2}, -\frac{\partial \psi}{\partial x_1} \right) \quad (4.9)$$

where $D = [0, 2\pi]^2$, periodic boundary conditions are considered and the initial condition $w(\cdot, 0)$ is fixed. Here we are interested in the mapping from the forcing term u to $v = \omega(\cdot, T)$, the vorticity field at a given time $t = T$, i.e., $\mathcal{G}^\dagger : u \mapsto \omega(\cdot, T)$.

The distribution μ is $\mathcal{GP}(0, (-\Delta + 3^2 I)^{-4})$. The viscosity ν is fixed and equal to 0.025, and the equation is solved on a 64×64 grid with a pseudo-spectral method and Crank-Nicholson time integration; see [19] for further implementation details. Fig. 9 shows an example of input and output in the test set, along with an example of test prediction and pointwise errors.

4.4. Results and discussion

Below we discuss our main findings in benchmarking our kernel method against state-of-the-art NN based techniques

4.4.1. Performance against NNs

Table 3 summarizes the L^2 relative test error of our vanilla implementation of the kernel method along with those of DeepONet, FNO, PCA-Net, and PARA-Net. We observed that our vanilla kernel method was reliable in terms of accuracy across all examples. In particular, observe that between the Matérn or rational quadratic kernel, we always managed to get close to the other methods, see for example the results for the Burgers' equation or Darcy problem, and even outperform them in several examples such as Navier-Stokes and Helmholtz. Overall we observed that the performance of the kernel method is stable across all examples suggesting

Table 3

Summary of numerical results: we report the L^2 relative test error of our numerical experiments and compare the kernel approach with variations of DeepONet, FNO, PCA-Net and PARA-Net. We considered two choices of the kernel S , the rational quadratic and the Matérn, but we observed little difference between the two.

	Low-data regime			High-data regime			
	Burger's	Darcy problem	Advection I	Advection II	Hemholtz	Structural Mechanics	Navier Stokes
DeepONet	2.15%	2.91%	0.66%	15.24%	5.88%	5.20%	3.63%
POD-DeepONet	1.94%	2.32%	0.04%	n/a	n/a	n/a	n/a
FNO	1.93%	2.41%	0.22%	13.49%	1.86%	4.76%	0.26%
PCA-Net	n/a	n/a	n/a	12.53%	2.13%	4.67%	2.65%
PARA-Net	n/a	n/a	n/a	16.64%	12.54%	4.55%	4.09%
Linear	36.24%	6.74%	$2.15 \times 10^{-13}\%$	11.28%	10.59%	27.11%	5.41%
Best of Matérn/RQ	2.15%	2.75%	$2.75 \times 10^{-3}\%$	11.44%	1.00%	5.18%	0.12%

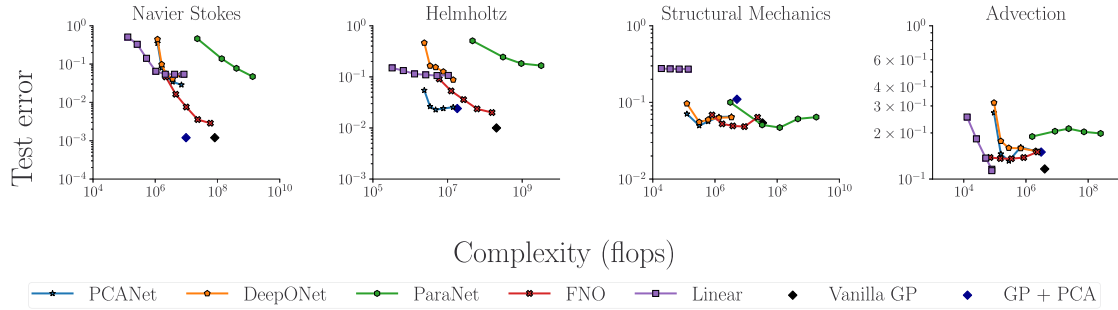


Fig. 10. Accuracy-complexity tradeoff achieved in the problems in [19]. Data for NNs was obtained from the aforementioned article. Linear model refers to the linear kernel, vanilla GP is our implementation with the nonlinear kernels and minimal preprocessing, GP+PCA corresponds to preprocessing through PCA both the input and the output to reduce complexity.

Table 4

Comparison between Cholesky preconditioning and PCA dimensionality reduction on three examples for our vanilla kernel implementation with the Matérn kernel.

	Advection II	Burger's	Darcy problem
No preprocessing	14.37%	3.04%	4.47%
PCA	14.50%	2.41%	2.89%
Cholesky	11.44%	2.15%	2.75%

that our method is reliable and provides a good baseline for a large class of problems. Moreover, we did not observe a significant difference in performance in terms of the choice of the particular kernel family once the hyper-parameters were tuned. This indicates that a large class of kernels are effective for these problems. Furthermore, we found the hyper-parameter tuning to be robust, i.e., results were consistent in a reasonable range of parameters such as length scales.

In the high data regime, we found the vanilla kernel method to be the most accurate, although this comes with a greater cost, as seen in Fig. 10. However, the kernel method appears to provide the highest accuracy for its level of complexity as the accuracy of NNs typically stagnates or even decreases after a certain level of complexity; see the Navier-Stokes and Helmholtz panels of Fig. 10 where most of the NN methods seem to plateau after a certain complexity level.

We also observed that the linear model did not provide the best accuracy as it quickly saturated in performance. Nonetheless, it provided surprisingly good accuracy at low levels of complexity: for example, in the case of Navier-Stokes, the linear kernel provided the best accuracy below 10^6 FLOPS of complexity. This indicates that while simple, the linear model can be a valuable low-complexity model. Another notable example is the Advection equation (both I and II), where the operator \mathcal{G}^\dagger is linear. In this case, the linear kernel had the best accuracy and the best complexity-accuracy tradeoff. We note, however, that while the linear model was close to machine precision on Advection I (error on the order $10^{-13}\%$), its performance was significantly worse on Advection II (error on the order of 10%). Moreover, the gap between the linear kernel and all other models was significantly smaller for Advection I; we conjecture this difference in performance is likely due to the setup of these problems.

Finally, we note that the most challenging problem for our kernel method was the Structural Mechanics example. In this case, the vanilla kernel method has higher complexity but did not beat the NNs. In fact, the NNs seem to be able to reduce complexity without loss of accuracy compared to our method.

4.4.2. Effect of preconditioners

Table 4 compares the performance of our method with the Matérn kernel family using various preconditioning steps. Overall we observed that both PCA and Cholesky preconditioning improved the performance of our vanilla kernel method.

The Cholesky preconditioning generally offers the greatest improvement. However, we observed that getting the best results from the Cholesky approach required careful tuning of the parameters of the kernels K and Q which we did using cross-validation. While tuning the parameters does not increase the inference complexity, it does increase the training complexity.

On the other hand, the PCA approach was more robust to changes in hyperparameters, i.e., the number of PCA components following Subsection 4.1.2. We observed that applying PCA on the input and output reduces complexity and has varying levels of effectiveness in providing a better cost-accuracy tradeoff. For example, for Navier-Stokes, it greatly reduced the complexity without affecting accuracy. But for the Helmholtz and Advection equations, PCA reduced the accuracy while remaining competitive with NN models. For structural mechanics, however, PCA significantly reduced accuracy and was worse than other models. We hypothesize that the loss in accuracy can be related to the decay of the eigenvalues of the PCA matrix in that example.

5. Conclusions

In this work we presented a kernel/GP framework for the learning of operators between function spaces. We presented an abstract formulation of our kernel framework along with convergence proofs and error bounds in certain asymptotic limits. Numerical experiments and benchmarking against popular NN based algorithms revealed that our vanilla implementation of the kernel approach is competitive and either matches the performance of NN methods or beats them in several benchmarks. Due to simplicity of implementation, flexibility, and the empirical results, we suggest that the proposed kernel methods are a good benchmark for future, perhaps more sophisticated, algorithms. Furthermore, these methods can be used to guide practitioners in the design of new and challenging benchmarks (e.g., identify problems where vanilla kernel methods do not perform well). Numerous directions of future research exist. In the theoretical direction it is interesting to remove the stringent Condition 3.2 and we anticipate this to require a particular selection of the kernel employed to obtain the map \tilde{f} . Moreover, obtaining error bounds for more general measurement functionals beyond pointwise evaluations would be interesting. One could also adapt our framework to non-vanilla kernel methods such as random features or inducing point methods to provide a low-complexity alternative to NNs in the large-data regime. Finally, since the proposed approach is essentially a generalization of GP Regression to the infinite-dimensional setting, we anticipate that some of the hierarchical techniques of [58,68,70] could be extended to this setting and provide a better cost-accuracy trade-off than current methods.

CRedit authorship contribution statement

Pau Batlle: Conceptualization, Investigation, Methodology, Software, Visualization, Writing – original draft. **Matthieu Darcy:** Conceptualization, Investigation, Methodology, Software, Visualization, Writing – original draft. **Bamdad Hosseini:** Conceptualization, Formal analysis, Methodology, Writing – original draft. **Houman Owahdi:** Conceptualization, Formal analysis, Funding acquisition, Methodology, Writing – original draft, Writing – review & editing.

Declaration of competing interest

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Matthieu Darcy reports financial support was provided by Air Force Office of Scientific Research under MURI award number FA9550-20-1-0358 (Machine Learning and Physics-Based Modeling and Simulation). Pau Batlle reports financial support was provided by Air Force Office of Scientific Research under MURI award number FA9550-20-1-0358 (Machine Learning and Physics-Based Modeling and Simulation). Houman Owahdi reports financial support was provided by Air Force Office of Scientific Research under MURI award number FA9550-20-1-0358 (Machine Learning and Physics-Based Modeling and Simulation). Bamdad Hosseini reports financial support was provided by National Science Foundation under grant number NSF-DMS-2208535 (Machine Learning for Bayesian Inverse Problems). Houman Owahdi reports financial support was provided by Department of Energy under award number DE-SC0023163 (SEA-CROGS: Scalable, Efficient and Accelerated Causal Reasoning Operators, Graphs and Spikes for Earth and Embedded Systems).

Data availability

The data is available publicly on GitHub, with a link provided in the article.

Acknowledgements

MD, PB, and HO acknowledge support by the Air Force Office of Scientific Research under MURI award number FA9550-20-1-0358 (Machine Learning and Physics-Based Modeling and Simulation). BH acknowledges support by the National Science Foundation grant number NSF-DMS-2208535 (Machine Learning for Bayesian Inverse Problems). HO also acknowledges support by the Department of Energy under award number DE-SC0023163 (SEA-CROGS: Scalable, Efficient and Accelerated Causal Reasoning Operators, Graphs and Spikes for Earth and Embedded Systems). We thank F. Schäfer for comments and references.

Appendix A. Review of operator valued kernels and GPs

We review the theory of operator valued kernels and GPs [60] as these are utilized throughout the article. Operator-valued kernels were introduced in [36] as a generalization of vector-valued kernels [3].

A.1. Operator valued kernels

Let \mathcal{U} and \mathcal{V} be separable Hilbert spaces endowed with the inner products $\langle \cdot, \cdot \rangle_{\mathcal{U}}$ and $\langle \cdot, \cdot \rangle_{\mathcal{V}}$. Write $\mathcal{L}(\mathcal{V})$ for the set of bounded linear operators mapping \mathcal{V} to \mathcal{V} .

Definition A.1. We call $G : \mathcal{U} \times \mathcal{U} \rightarrow \mathcal{L}(\mathcal{V})$ an “operator-valued kernel” if

1. G is Hermitian, i.e. $G(u, u') = G(u', u)^T$ for all $u, u' \in \mathcal{U}$, writing A^T for the adjoint of the operator A with respect to $\langle \cdot, \cdot \rangle_{\mathcal{V}}$.
2. G is non-negative, i.e., for all $m \in \mathbb{N}$ and any set of points $(u_i, v_i)_{i=1}^m \subset \mathcal{U} \times \mathcal{V}$ it holds that $\sum_{i,j=1}^m \langle v_i, G(u_i, u_j) v_j \rangle_{\mathcal{V}} \geq 0$.

We call G non-degenerate if $\sum_{i,j=1}^m \langle v_i, G(u_i, u_j) v_j \rangle_{\mathcal{V}} = 0$ implies $v_i = 0$ for all i whenever $u_i \neq u_j$ for $i \neq j$.

A.2. RKHSs

Each non-degenerate, locally bounded and separately continuous operator-valued kernel G is in one to one correspondence with an RKHS \mathcal{H} of continuous operators $\mathcal{G} : \mathcal{U} \rightarrow \mathcal{V}$ obtained as the closure of the linear span of the maps $z \mapsto G(z, u)v$ with respect to the inner product identified by the reproducing property

$$\langle g, G(\cdot, u)v \rangle_{\mathcal{H}} = \langle g(u), v \rangle_{\mathcal{V}} \quad (\text{A.1})$$

A.3. Feature maps

Let \mathcal{F} be a separable Hilbert space (with inner product $\langle \cdot, \cdot \rangle_{\mathcal{F}}$ and norm $\|\cdot\|_{\mathcal{F}}$) and let $\psi : \mathcal{U} \rightarrow \mathcal{L}(\mathcal{V}, \mathcal{F})$ be a continuous function mapping \mathcal{U} to the space of bounded linear operators from \mathcal{V} to \mathcal{F} .

Definition A.2. We say that \mathcal{F} and $\psi : \mathcal{U} \rightarrow \mathcal{L}(\mathcal{V}, \mathcal{F})$ are a *feature space* and a *feature map* for the kernel G if, for all $(u, u', v, v') \in \mathcal{U}^2 \times \mathcal{V}^2$,

$$\langle v, G(u, u') v' \rangle = \langle \psi(u)v, \psi(u')v' \rangle_{\mathcal{F}}.$$

Write $\psi^T(u)$, for the adjoint of $\psi(u)$ defined as the linear function mapping \mathcal{F} to \mathcal{V} satisfying

$$\langle \psi(u)v, \alpha \rangle_{\mathcal{F}} = \langle v, \psi^T(u)\alpha \rangle_{\mathcal{V}}$$

for $u, v, \alpha \in \mathcal{U} \times \mathcal{V} \times \mathcal{F}$. Note that $\psi^T : \mathcal{U} \rightarrow \mathcal{L}(\mathcal{F}, \mathcal{V})$ is therefore a function mapping \mathcal{U} to the space of bounded linear functions from \mathcal{F} to \mathcal{V} . Writing $\alpha^T \alpha' := \langle \alpha, \alpha' \rangle_{\mathcal{F}}$ for the inner product in \mathcal{F} we can ease our notations by writing

$$G(u, u') = \psi^T(u)\psi(u') \quad (\text{A.2})$$

which is consistent with the finite-dimensional setting and $v^T G(u, u') v' = (\psi(u)v)^T (\psi(u')v')$ (writing $v^T v'$ for the inner product in \mathcal{V}). For $\alpha \in \mathcal{F}$ write $\psi^T \alpha$ for the function $\mathcal{U} \rightarrow \mathcal{V}$ mapping $u \in \mathcal{U}$ to the element $v \in \mathcal{V}$ such that

$$\langle v', v \rangle_{\mathcal{V}} = \langle v', \psi^T(u)\alpha \rangle_{\mathcal{V}} = \langle \psi(u)v', \alpha \rangle_{\mathcal{F}} \text{ for all } v' \in \mathcal{V}.$$

We can, without loss of generality, restrict \mathcal{F} to be the range of $(u, v) \rightarrow \psi(u)v$ so that the RKHS \mathcal{H} defined by G is the closure of the pre-Hilbert space spanned by $\psi^T \alpha$ for $\alpha \in \mathcal{F}$. Note that the reproducing property (A.1) implies that for $\alpha \in \mathcal{F}$

$$\langle \psi^T(\cdot)\alpha, \psi^T(\cdot)\psi(u)v \rangle_{\mathcal{H}} = \langle \psi^T(u)\alpha, v \rangle_{\mathcal{V}} = \langle \alpha, \psi(u)v \rangle_{\mathcal{F}}$$

for all $u, v \in \mathcal{U} \times \mathcal{V}$, which leads to the following theorem.

Theorem A.3. The RKHS \mathcal{H} defined by the kernel (A.2) is the linear span of $\psi^T \alpha$ over $\alpha \in \mathcal{F}$ such that $\|\alpha\|_{\mathcal{F}} < \infty$. Furthermore, $\langle \psi^T(\cdot)\alpha, \psi^T(\cdot)\alpha' \rangle_{\mathcal{H}} = \langle \alpha, \alpha' \rangle_{\mathcal{F}}$ and

$$\|\psi^T(\cdot)\alpha\|_{\mathcal{H}}^2 = \|\alpha\|_{\mathcal{F}}^2 \text{ for } \alpha, \alpha' \in \mathcal{F}.$$

A.4. Interpolation

Let us consider the interpolation problem in operator valued RKHSs.

Problem 2. Let \mathcal{G}^\dagger be an unknown continuous operator mapping \mathcal{U} to \mathcal{V} . Given the information⁸ $\mathcal{G}^\dagger(\mathbf{u}) = \mathbf{v}$ with the data $(\mathbf{u}, \mathbf{v}) \in \mathcal{U}^N \times \mathcal{V}^N$, approximate \mathcal{G}^\dagger .

Using the relative error in $\|\cdot\|_{\mathcal{H}}$ -norm as a loss, the minimax optimal recovery solution of Problem 2 is, by [61, Thm. 12.4,12.5], given by

$$\begin{cases} \text{Minimize} & \|\mathcal{G}\|_{\mathcal{H}}^2 \\ \text{subject to} & \mathcal{G}(\mathbf{u}) = \mathbf{v} \end{cases} \quad (\text{A.3})$$

The minimizer is then of the form $\mathcal{G}(\cdot) = \sum_{j=1}^N G(\cdot, u_j) w_j$, where the coefficients $w_j \in \mathcal{V}$ are identified by solving the system of linear equations $\sum_{j=1}^N G(u_i, u_j) w_j = v_i$ for all $i \in \{1, \dots, N\}$. Using our compressed notation we can rewrite this equation as $G(\mathbf{u}, \mathbf{u}) \mathbf{w} = \mathbf{v}$ where $\mathbf{w} = (w_1, \dots, w_N)$, $\mathbf{v} = (v_1, \dots, v_N) \in \mathcal{V}^N$ and $G(\mathbf{u}, \mathbf{u})$ is the $N \times N$ block-operator matrix⁹ with entries $G(u_i, u_j)$. Therefore, writing $G(\cdot, \mathbf{u})$ for the vector $(G(\cdot, u_1), \dots, G(\cdot, u_N)) \in \mathcal{H}^N$, the optimal recovery interpolant is given by

$$\bar{\mathcal{G}}(\cdot) = G(\cdot, \mathbf{u}) G(\mathbf{u}, \mathbf{u})^{-1} \mathbf{v}, \quad (\text{A.4})$$

which implies that the value of (A.3) at the minimum is

$$\|\bar{\mathcal{G}}\|_{\mathcal{H}}^2 = \mathbf{v}^T G(\mathbf{u}, \mathbf{u})^{-1} \mathbf{v}, \quad (\text{A.5})$$

where $G(\mathbf{u}, \mathbf{u})^{-1}$ is the inverse of $G(\mathbf{u}, \mathbf{u})$, whose existence is implied by the non-degeneracy of G combined with $u_i \neq u_j$ for $i \neq j$.

A.5. Ridge regression

Let $\gamma > 0$. A ridge regression (approximate) solution to Problem 2 can be found as the minimizer of

$$\inf_{\mathcal{G} \in \mathcal{H}} \lambda \|\mathcal{G}\|_{\mathcal{H}}^2 + \gamma^{-1} \sum_{i=1}^N \|v_i - \mathcal{G}(u_i)\|_{\mathcal{V}}^2. \quad (\text{A.6})$$

This minimizer is given by the formula

$$\bar{\mathcal{G}}(u) = G(u, \mathbf{u}) (G(\mathbf{u}, \mathbf{u}) + \gamma I)^{-1} \mathbf{v}, \quad (\text{A.7})$$

writing I for the identity matrix. We can further compute directly

$$\|\bar{\mathcal{G}}\|_{\mathcal{H}}^2 = \mathbf{v}^T (G(\mathbf{u}, \mathbf{u}) + \gamma I)^{-1} \mathbf{v}.$$

A.6. Operator-valued GPs

The following definition of operator-valued Gaussian processes is a natural extension of scalar-valued Gaussian fields [61].

Definition A.4. [60, Def. 5.1] Let $G : \mathcal{U} \times \mathcal{U} \rightarrow \mathcal{L}(\mathcal{V})$ be an operator-valued kernel. Let m be a function mapping \mathcal{U} to \mathcal{V} . We call $\xi : \mathcal{U} \rightarrow \mathcal{L}(\mathcal{V}, \mathbf{H})$ an operator-valued GP if ξ is a function mapping $u \in \mathcal{U}$ to $\xi(u) \in \mathcal{L}(\mathcal{V}, \mathbf{H})$ where \mathbf{H} is a Gaussian space and $\mathcal{L}(\mathcal{V}, \mathbf{H})$ is the space of bounded linear operators from \mathcal{V} to \mathbf{H} . Abusing notations we write $\langle \xi(u), v \rangle_{\mathcal{V}}$ for $\xi(u)v$. We say that ξ has mean m and covariance kernel G and write $\xi \sim \mathcal{N}(m, G)$ if $\langle \xi(u), v \rangle_{\mathcal{V}} \sim \mathcal{N}(m(u), v^T G(u, u) v)$ and

$$\text{Cov}(\langle \xi(u), v \rangle_{\mathcal{V}}, \langle \xi(u'), v' \rangle_{\mathcal{V}}) = v^T G(u, u') v'. \quad (\text{A.8})$$

We say that ξ is centered if it is of zero mean.

If $G(u, u)$ is trace class ($\text{Tr}[G(u, u)] < \infty$) then $\xi(u)$ defines a measure on \mathcal{V} , i.e. a \mathcal{V} -valued random variable.¹⁰

Theorem A.5. [60, Thm. 5.2] The law of an operator-valued GP is uniquely determined by its mean m and covariance kernel G . Conversely given m and G there exists an operator-valued GP having mean m and covariance kernel G . In particular if G has feature space F and map ψ , the e_i form an orthonormal basis of F , and the Z_i are i.i.d. $\mathcal{N}(0, 1)$ random variables, then $\xi = m + \sum_i Z_i \psi^T e_i$ is an operator-valued GP with mean m and covariance kernel G .

⁸ For a N -vector $\mathbf{u} = (u_1, \dots, u_N) \in \mathcal{U}^N$ and a function $\mathcal{G} : \mathcal{U} \rightarrow \mathcal{V}$, write $\mathcal{G}(\mathbf{u})$ for the N vector with entries $(\mathcal{G}(u_1), \dots, \mathcal{G}(u_N))$.

⁹ For $N \geq 1$ let \mathcal{V}^N be the N -fold product space endowed with the inner-product $\langle \mathbf{v}, \mathbf{w} \rangle_{\mathcal{V}^N} := \sum_{i,j=1}^N \langle v_i, w_j \rangle_{\mathcal{V}}$ for $\mathbf{v} = (v_1, \dots, v_N)$, $\mathbf{w} = (w_1, \dots, w_N) \in \mathcal{V}^N$. $\mathbf{A} \in \mathcal{L}(\mathcal{V}^N)$

given by $\mathbf{A} = \begin{pmatrix} A_{1,1} & \cdots & A_{1,N} \\ \vdots & & \vdots \\ A_{N,1} & \cdots & A_{N,N} \end{pmatrix}$ where $A_{i,j} \in \mathcal{L}(\mathcal{V})$, is called a block-operator matrix. Its adjoint \mathbf{A}^T with respect to $\langle \cdot, \cdot \rangle_{\mathcal{V}^N}$ is the block-operator matrix with entries $(\mathbf{A}^T)_{i,j} = (A_{j,i})^T$.

¹⁰ Otherwise it only defines a (weak) cylinder-measure in the sense of Gaussian fields.

Theorem A.6. [60, Thm. 5.3] Let ξ be a centered operator-valued GP with covariance kernel $G : \mathcal{U} \times \mathcal{U} \rightarrow \mathcal{L}(\mathcal{V})$. Let $\mathbf{u}, \mathbf{v} \in \mathcal{U}^N \times \mathcal{V}^N$. Let $Z = (Z_1, \dots, Z_N)$ be a random Gaussian vector, independent from ξ , with i.i.d. $\mathcal{N}(0, \gamma I_{\mathcal{V}})$ entries ($\gamma \geq 0$ and $I_{\mathcal{V}}$ is the identity map on \mathcal{V}). Then ξ conditioned on $\xi(\mathbf{u}) + Z$ is an operator-valued GP with mean

$$\mathbb{E}[\xi(u) | \xi(\mathbf{u}) + Z = \mathbf{v}] = G(u, \mathbf{u}) (G(\mathbf{u}, \mathbf{u}) + \gamma I_{\mathcal{V}})^{-1} \mathbf{v} = (\text{A.7}) \quad (\text{A.9})$$

and conditional covariance operator

$$G^{\perp}(u, u') := G(u, u') - G(u, \mathbf{u}) (G(\mathbf{u}, \mathbf{u}) + \gamma I_{\mathcal{V}})^{-1} G(\mathbf{u}, u'). \quad (\text{A.10})$$

In particular, if G is trace class, then

$$\sigma^2(u) := \mathbb{E} \left[\left\| \xi(u) - \mathbb{E}[\xi(u) | \xi(\mathbf{u}) + Z = \mathbf{v}] \right\|_{\mathcal{V}}^2 \middle| \xi(\mathbf{u}) + Z = \mathbf{v} \right] = \text{Tr} [G^{\perp}(u, u)]. \quad (\text{A.11})$$

A.7. Deterministic error estimates for operator-valued regression

The following theorem shows that the standard deviation (A.11) provides deterministic a priori error bounds on the accuracy of the ridge regressor (A.9) to \mathcal{G}^{\dagger} in Problem 2. Local error estimates such as (A.12) below are classical in the Kriging literature [77] where $\sigma^2(u)$ is known as the power function/kriging variance; see also [57][Thm. 5.1] for applications to PDEs.

Theorem A.7. [60, Thm. 5.4] Let \mathcal{G}^{\dagger} be the unknown function of Problem 2 and let $\mathcal{G}(u) = (\text{A.9}) = (\text{A.7})$ be its ridge regressor. Let \mathcal{H} be the RKHS associated with G and let \mathcal{H}_{γ} be the RKHS associated with the kernel $G_{\gamma} := G + \gamma I_{\mathcal{V}}$. It holds true that

$$\left\| \mathcal{G}^{\dagger}(u) - \mathcal{G}(u) \right\|_{\mathcal{V}} \leq \sigma(u) \|\mathcal{G}^{\dagger}\|_{\mathcal{H}} \quad (\text{A.12})$$

and

$$\left\| \mathcal{G}^{\dagger}(u) - \mathcal{G}(u) \right\|_{\mathcal{V}} \leq \sqrt{\sigma^2(u) + \gamma \dim(\mathcal{V})} \|\mathcal{G}^{\dagger}\|_{\mathcal{H}_{\gamma}}, \quad (\text{A.13})$$

where $\sigma(u)$ is the standard deviation (A.11).

Appendix B. An alternative regularization of operator regression

For $\gamma > 0$, the regularization implied by (2.14) is equivalent to adding noise on the $\varphi(\mathbf{v})$ measurements. If one could observe \mathbf{v} (and not just $\varphi(\mathbf{v})$), then an alternative approach to regularizing the problem is to add noise to $\xi(\mathbf{u})$. To describe this let $Z' = (Z'_1, \dots, Z'_N)$ be a random block-vector, independent from ξ , with i.i.d. entries $Z'_j \sim \mathcal{N}(0, \gamma I_{\mathcal{V}})$ for $j = 1, \dots, N$ (where $I_{\mathcal{V}}$ denotes the identity map on \mathcal{V}). Then the GP ξ conditioned on $\xi(\mathbf{u}) = \mathbf{v} + Z'$ is a GP with conditional covariance kernel (A.10) and conditional mean $\tilde{\mathcal{G}}_{\gamma} = (\text{A.7})$ that is also the minimizer of (A.6). Observing¹¹ that $\varphi(Z'_i) \sim \mathcal{N}(0, \gamma K(\varphi, \varphi))$, we deduce that $\tilde{\mathcal{G}}_{\gamma} = \chi \circ \tilde{f}_{\gamma} \circ \phi$ where \tilde{f}_{γ} minimizes

$$\begin{cases} \text{Minimize} & \|f\|_{\Gamma}^2 + \gamma^{-1} \sum_{i=1}^N (f(U_i) - V_i)^T K(\varphi, \varphi)^{-1} (f(U_i) - V_i). \\ \text{Over} & f \in \mathcal{H}_{\Gamma}. \end{cases} \quad (\text{B.1})$$

Furthermore, the distribution of ξ conditioned on $\xi(\mathbf{u}) = \mathbf{v} + Z'$ is that of $\chi \circ \tilde{\xi}^{\perp} \circ \phi$ where $\tilde{\xi}^{\perp} \sim \mathcal{N}(\tilde{f}_{\gamma}, \tilde{\Gamma}^{\perp})$ is the GP ξ conditioned on $\xi(\mathbf{U}) = \mathbf{V} + \varphi(Z')$, whose mean is \tilde{f}_{γ} and conditional covariance kernel is $\tilde{\Gamma}^{\perp}(U, U') = \Gamma(U, U') - \Gamma(U, \mathbf{U})(\Gamma(\mathbf{U}, \mathbf{U}) + \gamma A)^{-1} \Gamma(\mathbf{U}, U')$ where A is a $N \times N$ block diagonal matrix with $K(\varphi, \varphi)$ as diagonal entries.

Appendix C. Expressions for the kernels used in experiments

Below we collect the expressions for the kernels that were referred to in the article or utilized for our numerical experiments. These can be found in many standard textbooks on GPs such as [67].

C.1. The linear kernel

The linear kernel has the simple expression $K_{\text{linear}}(x, x') = \langle x, x' \rangle$ and may be defined on any inner product space. It has no hyper-parameters.

¹¹ This follows from $\varphi(Z'_i) \sim \mathcal{N}(0, \gamma \varphi \varphi^T)$ where φ^T is the adjoint of φ identified as the linear map from \mathbb{R}^m to \mathcal{V} satisfying $\langle W, \varphi(w) \rangle_{\mathbb{R}^m} = \langle \varphi^{\perp} W, w \rangle_{\mathcal{V}}$ for $w \in \mathcal{V}$ and $W \in \mathbb{R}^m$ (i.e., $\varphi^{\perp}(W) = (K(\varphi, \varphi))^{-1} \varphi^T(W)$).

C.2. The rational quadratic kernel

The rational quadratic kernel has the expression $K(x, x') = k_{\text{RQ}}(\|x - x'\|)$ where

$$k_{\text{RQ}}(r) = \left(1 + \frac{r^2}{2l^2}\right)^{-\alpha}. \quad (\text{C.1})$$

It has hyper-parameters $\alpha > 0$ and l .

C.3. The Matérn parametric family

The Matérn kernel family is of the form $K(x, x') = k(\|x - x'\|)$ where

$$k_\nu(r) = \exp\left(-\frac{\sqrt{2\nu r}}{l}\right) \frac{\Gamma(p+1)}{\Gamma(2p+1)} \sum_{i=0}^p \frac{(p+1)!}{i!(p-i)!} \left(\frac{\sqrt{8\nu r}}{l}\right)^{p-i}, \quad (\text{C.2})$$

for $\nu = p + \frac{1}{2}$. This kernel has hyper-parameters $p \in \mathbb{Z}_+$ and $l > 0$. In the limiting case where $\nu \rightarrow \infty$, the Matérn kernel, we obtain the Gaussian or squared exponential kernel:

$$k_\infty(r) = \exp\left(-\frac{r^2}{2l^2}\right), \quad (\text{C.3})$$

with hyper-parameter $l > 0$.

References

- [1] B.O. Almroth, P. Stern, F.A. Brogan, Automatic choice of global shape functions in structural analysis, *AIAA J.* 16 (1978) 525–528.
- [2] R. Altmann, P. Henning, D. Peterseim, Numerical homogenization beyond scale separation, *Acta Numer.* 30 (2021) 1–86.
- [3] M.A. Alvarez, L. Rosasco, N.D. Lawrence, et al., Kernels for vector-valued functions: a review, *Found. Trends Mach. Learn.* 4 (2012) 195–266.
- [4] D. Amsallem, C. Farhat, Interpolation method for adapting reduced-order models and application to aeroelasticity, *AIAA J.* 46 (2008) 1803–1813.
- [5] R. Arcangéli, M.C. López de Silanes, J.J. Torrens, An extension of a bound for functions in Sobolev spaces, with applications to (m, s)-spline interpolation and smoothing, *Numer. Math.* 107 (2007) 181–211.
- [6] L.S. Bastos, A. O'hagan, Diagnostics for Gaussian process emulators, *Technometrics* 51 (2009) 425–438.
- [7] J. Beck, R. Tempone, F. Nobile, L. Tamellini, On the optimal polynomial approximation of stochastic PDEs by Galerkin and collocation methods, *Math. Models Methods Appl. Sci.* 22 (2012) 1250023.
- [8] M.P. Bendsoe, O. Sigmund, *Topology Optimization: Theory, Methods, and Applications*, Springer Science & Business Media, 2003.
- [9] K. Bhattacharya, B. Hosseini, N.B. Kovachki, A.M. Stuart, Model reduction and neural networks for parametric PDEs, *SMAI J. Comput. Math.* 7 (2021) 121–157.
- [10] G. Boncoraglio, C. Farhat, Active manifold and model-order reduction to accelerate multidisciplinary analysis and optimization, *AIAA J.* 59 (2021) 4739–4753.
- [11] N. Boullé, A. Townsend, Learning elliptic partial differential equations with randomized linear algebra, *Found. Comput. Math.* (2022) 1–31.
- [12] T. Chen, H. Chen, Approximation capability to functions of several variables, nonlinear functionals, and operators by radial basis function neural networks, *IEEE Trans. Neural Netw.* 6 (1995) 904–910.
- [13] T. Chen, H. Chen, Universal approximation to nonlinear operators by neural networks with arbitrary activation functions and its application to dynamical systems, *IEEE Trans. Neural Netw.* 6 (1995) 911–917.
- [14] Y. Chen, B. Hosseini, H. Owahdi, A.M. Stuart, Solving and Learning Nonlinear PDEs with Gaussian Processes, 2021.
- [15] Y. Chen, H. Owahdi, A. Stuart, Consistency of empirical Bayes and kernel flow for hierarchical parameter estimation, *Math. Comput.* 90 (2021) 2527–2578.
- [16] A. Chkifa, A. Cohen, R. DeVore, C. Schwab, Sparse adaptive Taylor approximation algorithms for parametric and stochastic elliptic PDEs, *ESAIM: Math. Model. Numer. Anal.* 47 (2012) 253–280.
- [17] A. Chkifa, A. Cohen, C. Schwab, High-dimensional adaptive sparse polynomial interpolation and applications to parametric PDEs, *Found. Comput. Math.* 14 (2014) 601–633.
- [18] A. Cohen, R. DeVore, Approximation of high-dimensional parametric PDEs, *Acta Numer.* 24 (2015) 1–159.
- [19] M. De Hoop, D.Z. Huang, E. Qian, A.M. Stuart, The cost-accuracy trade-off in operator learning with neural networks, *arXiv preprint, arXiv:2203.13181*, 2022.
- [20] B. Deng, Y. Shin, L. Lu, Z. Zhang, G.E. Karniadakis, Approximation rates of deepnets for learning operators arising from advection–diffusion equations, *Neural Netw.* 153 (2022) 411–426.
- [21] T.D. Economou, F. Palacios, S.R. Copeland, T.W. Lukaczyk, J.J. Alonso, Su2: an open-source suite for multiphysics simulation and design, *AIAA J.* 54 (2016) 828–846.
- [22] Y. Fan, C.O. Bohorquez, L. Ying, BCR-Net: a neural network based on the nonstandard wavelet form, *J. Comput. Phys.* 384 (2019) 1–15.
- [23] Y. Fan, J. Feliu-Faba, L. Lin, L. Ying, L. Zepeda-Núñez, A multiscale neural network based on hierarchical nested bases, *Res. Math. Sci.* 6 (2019) 1–28.
- [24] Y. Fan, L. Lin, L. Ying, L. Zepeda-Núñez, A multiscale neural network based on hierarchical matrices, *Multiscale Model. Simul.* 17 (2019) 1189–1213.
- [25] M. Feischl, D. Peterseim, Sparse compression of expected solution operators, *SIAM J. Numer. Anal.* 58 (2020) 3144–3164.
- [26] F. Feyel, J.-L. Chaboche, FE^2 multiscale approach for modelling the elastoviscoplastic behaviour of long fibre SiC/Ti composite materials, *Comput. Methods Appl. Mech. Eng.* 183 (2000) 309–330.
- [27] J. Fish, K. Shek, M. Pandheeradi, M.S. Shephard, Computational plasticity for composite structures based on mathematical homogenization: theory and practice, *Comput. Methods Appl. Mech. Eng.* 148 (1997) 53–73.
- [28] R.G. Ghanem, P.D. Spanos, *Stochastic Finite Elements: a Spectral Approach*, Dover Publications, 2003.
- [29] C.R. Gin, D.E. Shea, S.L. Brunton, J.N. Kutz, Deepgreen: deep learning of Green's functions for nonlinear boundary value problems, *Sci. Rep.* 11 (2021) 21614.
- [30] M.D. Gunzburger, C.G. Webster, G. Zhang, Stochastic finite element methods for partial differential equations with random input data, *Acta Numer.* 23 (2014) 521–650.
- [31] B. Hamzi, R. Maulik, H. Owahdi, Simple, low-cost and accurate data-driven geophysical forecasting with learned kernels, *Proc. R. Soc. A, Math. Phys. Eng. Sci.* 477 (2021) 20210326.
- [32] B. Hamzi, H. Owahdi, Learning dynamical systems from data: a simple cross-validation perspective, Part I: parametric kernel flows, *Physica D* 421 (2021) 132817.
- [33] J. Hesthaven, S. Ubbiali, Non-intrusive reduced order modeling of nonlinear problems using neural networks, *J. Comput. Phys.* 363 (2018) 55–78.

- [34] J.S. Hesthaven, G. Rozza, B. Stamm, et al., *Certified Reduced Basis Methods for Parametrized Partial Differential Equations*, vol. 590, Springer, 2016.
- [35] D.Z. Huang, T. Schneider, A.M. Stuart, Iterated Kalman methodology for inverse problems, *J. Comput. Phys.* 463 (2022) 111262.
- [36] H. Kadri, E. Duflos, P. Preux, S. Canu, A. Rakotomamonjy, J. Audiffren, *Operator-Valued Kernels for Learning from Functional Response Data*, 2016.
- [37] M.C. Kennedy, A. O'Hagan, Bayesian calibration of computer models, *J. R. Stat. Soc., Ser. B, Stat. Methodol.* 63 (2001) 425–464.
- [38] Y. Khoo, J. Lu, L. Ying, Solving parametric PDE problems with artificial neural networks, *Eur. J. Appl. Math.* 32 (2021) 421–435.
- [39] Y. Khoo, L. Ying, Switchnet: a neural network model for forward and inverse scattering problems, *SIAM J. Sci. Comput.* 41 (2019) A3182–A3201.
- [40] G. Kassis, J.H. Seidman, L.F. Guilhoto, V.M. Preciado, G.J. Pappas, P. Perdikaris, Learning operators with coupled attention, *J. Mach. Learn. Res.* 23 (2022) 1–63.
- [41] N. Kovachki, S. Lanthaler, S. Mishra, On universal approximation and error bounds for Fourier neural operators, *J. Mach. Learn. Res.* 22 (2021) 13237–13312.
- [42] N. Kovachki, B. Liu, X. Sun, H. Zhou, K. Bhattacharya, M. Ortiz, A. Stuart, Multiscale modeling of materials: computing, data science, uncertainty and goal-oriented optimization, *Mech. Mater.* 165 (2022) 104156.
- [43] K. Krischer, R. Rico-Martínez, I. Kevrekidis, H. Rotermund, G. Ertl, J. Hudson, Model identification of a spatiotemporally varying catalytic reaction, *AIChE J.* 39 (1993) 89–98.
- [44] F. Kröpl, R. Maier, D. Peterseim, Operator compression with deep neural networks, *Adv. Cont. Discr. Mod.* 2022 (2022) 1–23.
- [45] Z. Li, N. Kovachki, K. Azizzadenesheli, B. Liu, K. Bhattacharya, A. Stuart, A. Anandkumar, Fourier Neural Operator for Parametric Partial Differential Equations, 2020.
- [46] L. Lu, P. Jin, G. Pang, Z. Zhang, G.E. Karniadakis, Learning nonlinear operators via DeepONet based on the universal approximation theorem of operators, *Nat. Mach. Intell.* 3 (2021) 218–229.
- [47] L. Lu, X. Meng, S. Cai, Z. Mao, S. Goswami, Z. Zhang, 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 (2022) 114778.
- [48] D.J. Lucia, P.S. Beran, W.A. Silva, Reduced-order modeling: new approaches for computational physics, *Prog. Aerosp. Sci.* 40 (2004) 51–117.
- [49] Y. Maday, A.T. Patera, G. Turinici, A priori convergence theory for reduced-basis approximations of single-parameter elliptic partial differential equations, *J. Sci. Comput.* 17 (2002) 437–446.
- [50] A. Målqvist, D. Peterseim, Localization of elliptic multiscale problems, *Math. Comput.* 83 (2014) 2583–2603.
- [51] J. Martin, L.C. Wilcox, C. Burstedde, O. Ghattas, A stochastic Newton MCMC method for large-scale statistical inverse problems with application to seismic inversion, *SIAM J. Sci. Comput.* 34 (2012) A1460–A1487.
- [52] J.R. Martins, A.B. Lambe, Multidisciplinary design optimization: a survey of architectures, *AIAA J.* 51 (2013) 2049–2075.
- [53] N. Marzari, A.A. Mostofi, J.R. Yates, I. Souza, D. Vanderbilt, Maximally localized Wannier functions: theory and applications, *Rev. Mod. Phys.* 84 (2012) 1419.
- [54] F. Nobile, R. Tempone, C.G. Webster, An anisotropic sparse grid stochastic collocation method for partial differential equations with random input data, *SIAM J. Numer. Anal.* 46 (2008) 2411–2442.
- [55] F. Nobile, R. Tempone, C.G. Webster, A sparse grid stochastic collocation method for partial differential equations with random input data, *SIAM J. Numer. Anal.* 46 (2008) 2309–2345.
- [56] A.K. Noor, J.M. Peters, Reduced basis technique for nonlinear analysis of structures, *AIAA J.* 18 (1980) 455–462.
- [57] H. Owadi, Bayesian numerical homogenization, *Multiscale Model. Simul.* 13 (2015) 812–828.
- [58] H. Owadi, Multigrid with rough coefficients and multiresolution operator decomposition from hierarchical information games, *SIAM Rev.* 59 (2017) 99–149.
- [59] H. Owadi, Computational graph completion, *Res. Math. Sci.* 9 (2022) 27.
- [60] H. Owadi, Do ideas have shape? Idea registration as the continuous limit of artificial neural networks, *Physica D* 444 (2023) 133592.
- [61] H. Owadi, C. Scovel, *Operator-Adapted Wavelets, Fast Solvers, and Numerical Homogenization: From a Game Theoretic Approach to Numerical Approximation and Algorithm Design*, Cambridge Monographs on Applied and Computational Mathematics, Cambridge University Press, 2019.
- [62] H. Owadi, G.R. Yoo, Kernel flows: from learning kernels from data into the abyss, *J. Comput. Phys.* 389 (2019) 22–47.
- [63] H. Owadi, L. Zhang, Metric-based upscaling, *Commun. Pure Appl. Math.* 60 (2007) 675–723.
- [64] H. Owadi, L. Zhang, Gamblets for opening the complexity-bottleneck of implicit schemes for hyperbolic and parabolic odes/PDEs with rough coefficients, *J. Comput. Phys.* 347 (2017) 99–128.
- [65] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay, Scikit-learn: machine learning in Python, *J. Mach. Learn. Res.* 12 (2011) 2825–2830.
- [66] S. Prasanth, Z. Haddad, J. Susiluo, A. Braverman, H. Owadi, B. Hamzi, S. Hristova-Veleva, J. Turk, Kernel flows to infer the structure of convective storms from satellite passive microwave observations, in: *AGU Fall Meeting Abstracts*, vol. 2021, 2021, A55F-1445.
- [67] C.E. Rasmussen, C.K.I. Williams, *Gaussian Processes for Machine Learning*, Adaptive Computation and Machine Learning, MIT Press, 2006.
- [68] F. Schäfer, M. Katzfuss, H. Owadi, Sparse Cholesky factorization by Kullback–Leibler minimization, *SIAM J. Sci. Comput.* 43 (2021) A2019–A2046.
- [69] F. Schäfer, H. Owadi, Sparse recovery of elliptic solvers from matrix-vector products, arXiv preprint, arXiv:2110.05351, 2021.
- [70] F. Schäfer, T.J. Sullivan, H. Owadi, Compression, inversion, and approximate PCA of dense kernel matrices at near-linear computational complexity, *Multiscale Model. Simul.* 19 (2021) 688–730.
- [71] B. Schölkopf, R. Herbrich, A.J. Smola, A generalized representer theorem, in: D. Helmbold, B. Williamson (Eds.), *Computational Learning Theory*, Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 416–426, 2001.
- [72] B. Sudret, S. Marelli, J. Wiart, Surrogate models for uncertainty quantification: an overview, in: *2017 11th European Conference on Antennas and Propagation, EUCAP, IEEE*, 2017, pp. 793–797.
- [73] J. Susiluo, A. Braverman, P. Brodrick, B. Hamzi, M. Johnson, O. Lamminpää, H. Owadi, C. Scovel, J. Teixeira, M. Turmon, Radiative transfer emulation for hyperspectral imaging retrievals with advanced kernel flows-based Gaussian process emulation, in: *AGU Fall Meeting Abstracts*, vol. 2021, 2021, pp. NG25A–0506.
- [74] S. Wang, H. Wang, P. Perdikaris, Learning the solution operator of parametric partial differential equations with physics-informed deepONets, *Sci. Adv.* 7 (2021), eabi8605.
- [75] S. Wang, H. Wang, P. Perdikaris, Improved architectures and training algorithms for deep operator networks, *J. Sci. Comput.* 92 (2022) 35.
- [76] E. Weinan, *Principles of Multiscale Modeling*, Cambridge University Press, 2011.
- [77] Z.-m. Wu, R. Schaback, Local error estimates for radial basis function interpolation of scattered data, *IMA J. Numer. Anal.* 13 (1993) 13–27.
- [78] D. Xiu, *Numerical Methods for Stochastic Computations: A Spectral Method Approach*, Princeton University Press, 2010.
- [79] D. Xiu, G.E. Karniadakis, The Wiener–Askey polynomial chaos for stochastic differential equations, *SIAM J. Sci. Comput.* 24 (2002) 619–644.
- [80] D. Xiu, J. Shen, Efficient stochastic Galerkin methods for random diffusion equations, *J. Comput. Phys.* 228 (2009) 266–281.
- [81] Y. Zhu, N. Zabaras, Bayesian deep convolutional encoder–decoder networks for surrogate modeling and uncertainty quantification, *J. Comput. Phys.* 366 (2018) 415–447.