

# Infinite Neural Operators: Gaussian processes on functions

**Daniel Augusto de Souza\***  
University College London

**Yuchen Zhu**  
University College London

**Harry Jake Cunningham**  
University College London

**Yuri Saporito**  
Fundação Getulio Vargas

**Diego Mesquita**  
Fundação Getulio Vargas

**Marc Peter Deisenroth**  
University College London

## Abstract

A variety of infinitely wide neural architectures (e.g., dense NNs, CNNs, and transformers) induce Gaussian process (GP) priors over their outputs. These relationships provide both an accurate characterization of the prior predictive distribution and enable the use of GP machinery to improve the uncertainty quantification of deep neural networks. In this work, we extend this connection to neural operators (NOs), a class of models designed to learn mappings between function spaces. Specifically, we show conditions for when arbitrary-depth NOs with Gaussian-distributed convolution kernels converge to function-valued GPs. Based on this result, we show how to compute the covariance functions of these NO-GPs for two NO parametrizations, including the popular Fourier neural operator (FNO). With this, we compute the posteriors of these GPs in regression scenarios, including PDE solution operators. This work is an important step towards uncovering the inductive biases of current FNO architectures and opens a path to incorporate novel inductive biases for use in kernel-based operator learning methods.

## 1 Introduction

Neural Operators (NOs, Kovachki et al., 2023) are deep learning architectures designed to learn mappings between function spaces—with direct applications in many areas of science and engineering (Pathak et al., 2022; Li et al., 2024). NOs generalize conventional convolutional neural networks using *kernel integral operators*, which integrate the input function against a learnable kernel at each layer. Importantly, unlike CNNs, NOs can be trained with inputs of mixed, arbitrary resolutions and output predictions in discretizations of arbitrary granularity.

Despite their growing adoption, most works on NOs are primarily empirical, and most of the theoretical properties of NOs are still unexplored. In contrast, the convergence of Bayesian neural networks to Gaussian processes as their width goes to infinity has been amply studied (Neal, 1995; Novak et al., 2019; Yang, 2019). However, due to the infinite dimensionality of function spaces, it is unclear whether GPs are a limiting case for NOs and, if this is the case, how to characterize them.

In this work, we elucidate this question and present a set of assumptions that guarantee the existence of the infinite limit of NOs as Gaussian elements in the space of operators. Additionally, we present how to derive the covariance function for infinite-width NOs in an analogous fashion to the covariance functions of infinitely wide, densely connected NNs. Finally, we characterize the infinite-width limit of Fourier neural operators (FNOs) and propose a novel Bayesian NO architecture based on Matérn GP-distributed integral kernels.

\*Corresponding author: [daniel.souza.21@ucl.ac.uk](mailto:daniel.souza.21@ucl.ac.uk)

Our experiments reinforce our theoretical results, showcasing the agreement between increasingly wide NOs and our derived expressions for the infinite limit at initialization. Additionally, we compare the performance of these models in a regression setting.

## 2 Background

This section provides a brief background on NOs (Section 2.1), along with basic notions of probability in Hilbert spaces (Section 2.2) and Gaussian processes on functions.

### 2.1 Operator learning and neural operators

Kovachki et al. (2023) propose neural operators, a family of parametrized operators. Recall that multilayer perceptrons transform vectors using successive layers of sums of linear transformations followed by element-wise non-linear activation functions. Analogously, Kovachki et al. (2023) define the building layers of neural operators (NOs) as sums of both point-wise linear operations and kernel integral operators, possibly followed by point-wise element-wise non-linear activation functions.

Well-defined dot products in function spaces are central to coherently defining NOs. Thus, we will often assume functions lie in a vector space in which their dot product is finite wrt some measure  $\mu_{\mathcal{X}}$  over their domain  $\mathcal{X}$ . We define the Lebesgue space  $L^2(\mathcal{X}, \mu_{\mathcal{X}}; \mathbb{R}^d)$  as the equivalence classes of functions in this vector space that agree almost everywhere in  $\mathcal{X}$  with respect to  $\mu_{\mathcal{X}}$ . When clear from context, we will simply denote this vector space by  $L^2(\mathcal{X})$ . Whenever needed to evaluate functions point-wise, we further assume the function lies in an appropriate Reproducing Kernel Hilbert Space (RKHSs). In this work, we will be using both the Lebesgue space  $L^2(\mathcal{X})$  and RKHSs, when adequate.

**Point-wise operators.** These operations are carried over from standard neural networks. Thus, given a function  $\mathbf{f}: \mathcal{X} \rightarrow \mathbb{R}^d$ , we consider dense layer-operations, with parameters  $\mathbf{W} \in \mathbb{R}^{b \times d}$ , defined as  $(\mathbf{W}\mathbf{f})(\mathbf{x}): \mathcal{X} \rightarrow \mathbb{R}^b := \mathbf{W}\mathbf{f}(\mathbf{x})$ , and element-wise activations, with a given  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ , to be defined as  $\sigma[\mathbf{f}]_j(\mathbf{x}) := \sigma(f_j(\mathbf{x}))$ . By composing and adding results between layers, we can build neural operators that basically act just on the output of the functions.

**Kernel integral operator.** The majority of interesting behaviors require expanding the receptive field and aggregate results from different function evaluations into one. The *kernel integral operator*  $A_{\mathbf{K}}: (\mathcal{X} \rightarrow \mathbb{R}^d) \rightarrow (\mathcal{Y} \rightarrow \mathbb{R}^b)$ , parametrized by a matrix-valued kernel function  $\mathbf{K}: \mathcal{Y} \times \mathcal{X} \rightarrow \mathbb{R}^{b \times d}$  together with a measure  $\mu_{\mathcal{X}}$  on  $\mathcal{X}$ , is defined as:

$$A_{\mathbf{K}}[\mathbf{f}](\mathbf{y}): \mathcal{Y} \rightarrow \mathbb{R}^b = \int_{\mathcal{X}} \mathbf{K}(\mathbf{y}, \mathbf{x}) \mathbf{f}(\mathbf{x}) d\mu_{\mathcal{X}}(\mathbf{x}). \quad (1)$$

Under this operation, the function evaluated at a single evaluating point  $\mathbf{y}$  linearly aggregates information on all evaluating points in the domain  $\mathcal{X}$  as modulated by the kernel  $\mathbf{K}$  and the measure  $\mu_{\mathcal{X}}$ . Note, that this function may not converge for all values of  $\mathbf{y}$ , but, for any kernel  $\mathbf{K} \in L^2(\mathcal{Y} \times \mathcal{X}, \mu_{\mathcal{Y}} \times \mu_{\mathcal{X}}; \mathbb{R}^{b \times d})$ , the operator  $A_{\mathbf{K}}: L^2(\mathcal{X}, \mu_{\mathcal{X}}) \rightarrow L^2(\mathcal{Y}, \mu_{\mathcal{Y}})$  is well-defined.

**Constructing neural operators.** Given these building blocks, Kovachki et al. (2023) describe a neural operator as a three-part layered model. Firstly, a sequence of point-wise operators are applied to pre-process the function and change the dimension of its output. This is called the *Lift layer*. The second component is a combination of point-wise and kernel integral operators, in the so-called *Neural Operator layer*. Finally, the *Projection layer*, a sequence of point-wise operators is applied to the final result.

Specifically, a neural operator layer combines a matrix-valued kernel  $\mathbf{K}$  and matrix  $\mathbf{W}$  into

$$H[\mathbf{f}](\mathbf{x}): \mathcal{X} \rightarrow \mathbb{R}^h = A_{\mathbf{K}}[\mathbf{f}](\mathbf{x}) + \mathbf{W}\mathbf{f}(\mathbf{x}) = \int_{\mathcal{X}} \mathbf{K}(\mathbf{x}, \mathbf{z}) \mathbf{f}(\mathbf{z}) d\mu_{\mathcal{X}}(\mathbf{z}) + \mathbf{W}\mathbf{f}(\mathbf{x}). \quad (2)$$

Setting the matrix-valued kernel to zero recovers the lift and projection layers. Therefore, a neural operator with depth  $d$  and scalar output can be written succinctly as the composition:

$$Z[\mathbf{f}](\mathbf{x}): \mathcal{X} \rightarrow \mathbb{R} = (\mathbf{w}^{\top} \circ \sigma \circ H_d \circ \cdots \sigma \circ H_1)[\mathbf{f}](\mathbf{x}). \quad (3)$$

### 2.2 Probability in Hilbert spaces

Given a probability space  $(\Omega, \Sigma, \mathbb{P})$ , and a Hilbert space  $\mathcal{H}$ , random elements in  $\mathcal{H}$  are functions  $x: \Omega \rightarrow \mathcal{H}$ , such that the inner product  $\omega \in \Omega \mapsto \langle y, x(\omega) \rangle_{\mathcal{H}}$  is a real-valued random variable, for

any  $y \in \mathcal{H}$ . As usual, we follow the standard notation of denoting the random elements/variables not as functions  $x$  but as elements  $x$ . Likewise, expectation is defined in terms of the random variables  $\langle y, x(\omega) \rangle$ , for each  $y \in \mathcal{H}$ . We say that the expectation of  $x$ , when it exists, is the element of  $\mathcal{H}$ , denoted by  $\mathbb{E}[x]$ , such that  $\mathbb{E}[\langle y, x \rangle] = \langle y, \mathbb{E}[x] \rangle$ , for any  $y \in \mathcal{H}$ .

We denote the space of Hilbert–Schmidt (HS) operators mapping elements from a Hilbert space  $\mathcal{A}$  to  $\mathcal{B}$  by  $\text{HS}(\mathcal{A}; \mathcal{B})$ . This space is the completion of the span of rank-one operators of the form  $a \otimes b: \mathcal{A} \rightarrow \mathcal{B}$ , defined as  $(a \otimes b)(x) = \langle x, a \rangle_{\mathcal{A}} b$  for all  $a \in \mathcal{A}$  and  $b \in \mathcal{B}$ . For  $L^2$  spaces, we have the isomorphism  $\text{HS}(L^2(\mathcal{X}; \mathbb{R}^d), L^2(\mathcal{Y}; \mathbb{R}^b)) \cong L^2(\mathcal{X} \times \mathcal{Y}; \mathbb{R}^{b \times d})$ , under which  $(\mathbf{f} \otimes \mathbf{g})[\mathbf{h}](\cdot) = \int_{\mathcal{X}} \mathbf{g} \mathbf{f}^\top(\cdot, \mathbf{x}) \mathbf{h}(\mathbf{x}) d\mu_{\mathcal{X}}(\mathbf{x})$ , where  $\mathbf{f} \in L^2(\mathcal{X}; \mathbb{R}^d)$ ,  $\mathbf{g} \in L^2(\mathcal{Y}; \mathbb{R}^b)$ , and  $\mathbf{g} \mathbf{f}^\top \in L^2(\mathcal{X} \times \mathcal{Y}; \mathbb{R}^{b \times d})$ .

Moreover, the (cross-)covariance operator of two centered variables  $x$  and  $y$  is defined as the expectation of the tensor product  $\mathbb{E}[x \otimes y]$ . When this expectation exists, it is also a HS operator denoted as  $\text{Cov}(x, y)$ . From these definitions, we have that  $\langle z_2, \text{Cov}(x, y)[z_1] \rangle = \text{cov}(\langle z_2, y \rangle \langle z_1, x \rangle)$ , for any  $z_1, z_2 \in \mathcal{H}$ . In  $L^2$  spaces, we will make use of the isomorphism above and represent the covariance operator by its integration kernel. So, for any random elements  $\mathbf{f} \in L^2(\mathcal{X}; \mathbb{R}^d)$  and  $\mathbf{g} \in L^2(\mathcal{Y}; \mathbb{R}^b)$ , we introduce the function  $\mathbf{C}[\mathbf{f}, \mathbf{g}]: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^{b \times d}$  such that  $\text{Cov}(\mathbf{f}, \mathbf{g})[\mathbf{h}](\cdot) = \int_{\mathcal{X}} \mathbf{C}[\mathbf{f}, \mathbf{g}](\cdot, \mathbf{x}) \mathbf{h}(\mathbf{x}) d\mu_{\mathcal{X}}(\mathbf{x})$ .

In this work, we will make use of an extension of the strong law of large numbers to random elements:

**Theorem 2.1** (Strong law of large numbers (Mourier, 1956)). *Let  $\mathcal{H}$  be a separable Hilbert space and  $\{x_j\}_{j \in \mathbb{N}}$  be a countable sequence of identically distributed random elements. Consider the sample average  $y_N = (1/N) \sum_{j=1}^N x_j$ . If, for any  $j$ , the expected norm  $\mathbb{E}[\|x_j\|]$  exists, then, the sequence  $\{y_N\}_{N \in \mathbb{N}}$  converges almost surely to the constant random element  $y_\infty = \mathbb{E}[x_j]$ .*

### 2.3 Operator valued kernels and Hilbert space valued Gaussian processes

Now, given a set  $\mathcal{X}$  and a separable Hilbert space  $\mathcal{H}$ , an *operator-valued kernel*  $\mathbf{C}: \mathcal{X} \times \mathcal{X} \rightarrow \text{HS}(\mathcal{H}; \mathcal{H})$  is any Hermitian positive-definite function, i.e., for all  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ ,  $\mathbf{C}(\mathbf{x}, \mathbf{x}') = \mathbf{C}(\mathbf{x}', \mathbf{x})^\top$ , and, for any  $n > 0$ ,  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n \subset \mathcal{X} \times \mathcal{H}$  and  $\{\alpha_{ij}\}_{i,j=1}^n \subset \mathbb{R}$ , we have that  $\sum_{i,j=1}^n \alpha_{ij} \langle \mathbf{y}_j, \mathbf{C}(\mathbf{x}_i, \mathbf{x}_j)[\mathbf{y}_i] \rangle > 0$  (Kadri et al., 2016).

Consider an operator-valued kernel  $\mathbf{C}: \mathcal{X} \times \mathcal{X} \rightarrow \text{HS}(\mathcal{H}; \mathcal{H})$  such that  $\mathbf{x} \mapsto \mathbf{C}(\mathbf{x}, \mathbf{x})$  is of trace-class. We say  $\mathbf{f}: \mathcal{X} \times \Omega \rightarrow \mathcal{H}$  is a centered Gaussian process with covariance function  $\mathbf{C}$  if, for any  $n > 0$  and  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n \subset \mathcal{X} \times \mathcal{H}$ , the vector  $(\langle \mathbf{y}_1, \mathbf{f}(\mathbf{x}_1, \cdot) \rangle, \dots, \langle \mathbf{y}_n, \mathbf{f}(\mathbf{x}_n, \cdot) \rangle)$  is a random element distributed as an  $n$ -dimensional Gaussian with covariance

$$\mathbb{E}[\langle \mathbf{y}_2, \mathbf{f}(\mathbf{x}_2, \cdot) \rangle \langle \mathbf{y}_1, \mathbf{f}(\mathbf{x}_1, \cdot) \rangle] = \langle \mathbf{y}_2, \mathbf{C}(\mathbf{x}_1, \mathbf{x}_2)[\mathbf{y}_1] \rangle. \quad (4)$$

We denote this by  $\mathbf{f} \sim \text{GP}(0, \mathbf{C})$ . For simplicity, we also define  $\mathbf{f}(\mathbf{x}) := \mathbf{f}(\mathbf{x}, \cdot)$ .

## 3 Infinite-width neural operators as Gaussian processes

It is well known that infinite-width limits of various Bayesian neural networks are Gaussian processes (Neal, 1995; Matthews et al., 2018). We generalize this connection and show that infinite-width neural operators are function-valued Gaussian processes.

Analogous to Novak et al. (2019), who place Gaussian priors on the convolution kernels of a CNN, the natural step towards function-valued GPs is to put independent GP priors on the component operators. Similarly, we require the weights and kernel for any component operator to be i.i.d. and with covariance shrinking with width. Theorem 3.1 states the main result of this work.

**Theorem 3.1** (Infinite-width neural operators are Gaussian processes). *Let  $\mathcal{X} \subseteq \mathbb{R}^{d_x}$  be a measurable space and let  $\mathcal{H}(\mathcal{X}; \mathbb{R}^J) \subset L^2(\mathcal{X}; \mathbb{R}^J)$  be an RKHS for any  $J \in \mathbb{N}^+$ . Then, for a given depth  $D \in \mathbb{N}^+$ , consider a vector of positive integers  $\mathbf{J} = [J_0, J_1, \dots, J_{D-1}, 1]^\top \in \mathbb{N}^{D+1}$  and a  $\mathbf{J}$ -indexed neural operators  $Z_{\mathbf{J}}^{(D)}$  of depth  $D$ :*

$$Z_{\mathbf{J}}^{(D)} := H^{(D)} \circ \sigma \circ Z_{\mathbf{J}}^{(D-1)} \in (\mathcal{X} \rightarrow \mathbb{R}^{J_0}) \rightarrow (\mathcal{X} \rightarrow \mathbb{R}), \quad (5)$$

where,

$$Z_{\mathbf{J}}^{(1)} := H^{(1)} \in L^2(\mathcal{X}; \mathbb{R}^{J_0}) \rightarrow \mathcal{H}(\mathcal{X}; \mathbb{R}^{J_1}), \quad (6)$$

$$H^{(\ell)} := (A_{\mathbf{K}^{(\ell)}} + \mathbf{W}^{(\ell)}) \in L^2(\mathcal{X}; \mathbb{R}^{J_{\ell-1}}) \rightarrow \mathcal{H}(\mathcal{X}; \mathbb{R}^{J_{\ell}}), \quad (7)$$

$$\mathbf{W}^{(\ell)} \in \mathbb{R}^{J_{\ell} \times J_{\ell-1}}, \quad (8)$$

$$\mathbf{K}^{(\ell)} \in \mathcal{H}(\mathcal{X} \times \mathcal{X}; \mathbb{R}^{J_{\ell} \times J_{\ell-1}}), \quad (9)$$

and  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  such that  $(\sigma \circ f) \in L^2(\mathcal{X})$  for any  $f \in L^2(\mathcal{X})$ .

When all parameters are independently distributed a priori according to

$$\mathbf{W}^{(\ell)} \sim \mathcal{N}(\mathbf{0}, \sigma_{\ell}^2 / J_{\ell-1} \mathbf{I}), \text{ and } \mathbf{K}^{(\ell)} \sim \text{GP}(\mathbf{0}, c_{\mathbf{K}^{(\ell)}} / J_{\ell-1} \mathbf{I}), \quad \text{for } \ell \in \{1, \dots, d\}, \quad (10)$$

then, the iterated limit  $\lim_{J_{D-1} \rightarrow \infty} \dots \lim_{J_1 \rightarrow \infty} Z_J^{(D)}$ , in the sense of Definition B.1, is equal to a function-valued GP  $Z_{\infty}^{(D)} \sim \text{GP}(0, c_{\infty})$ , where  $c_{\infty}[\mathbf{f}, \mathbf{g}]$  is available in closed-form.

An outline of the proof is presented in Section 3.2, where we present the explicit formula for  $c_{\infty}$ , which depends on the conditional covariance function between layers. Before delving into these details, we introduce the compositionality property of covariance functions in Section 3.1. This property enables the closed-form computation of the conditional covariances, thereby fully characterizing the limiting covariance function  $c_{\infty}$ .

### 3.1 Operator-valued covariance functions

We realize the following crucial points: i) The covariance function only depends on the inner product of the values of the input functions, and ii) Using the strong law of large numbers, the covariance of the composition of operators can be described by composing its covariance functions. This is presented in the next lemma, with proof postponed to Appendix B.2.

**Lemma 3.2** (Compositionality of covariance functions). *Let  $B_1: L^2(\mathcal{X}; \mathbb{R}^d) \rightarrow L^2(\mathcal{X}; \mathbb{R}^J)$  be a random operator and  $B_2: L^2(\mathcal{X}; \mathbb{R}^J) \rightarrow L^2(\mathcal{X})$  be a centered function-valued Gaussian process. If the following assumptions hold:*

- For all  $\mathbf{f} \in L^2(\mathcal{X}; \mathbb{R}^d)$  and  $\mathbf{x} \in \mathcal{X}$ , each component of  $B_1[\mathbf{f}](\mathbf{x}) \in \mathbb{R}^J$  is independent and identically distributed such that the covariance function  $\mathbf{C}_{B_1}[\mathbf{f}, \mathbf{g}] = c_{B_1}[\mathbf{f}, \mathbf{g}] \mathbf{I}_J$ ;
- The covariance function of  $B_2$  can be expressed, for all  $\mathbf{f}, \mathbf{g} \in L^2(\mathcal{X}; \mathbb{R}^J)$  as  $c_{B_2}[\mathbf{f}, \mathbf{g}] = c_{B_2}[\frac{1}{J} \mathbf{g}^T \mathbf{f}]$  and the function  $h \mapsto c_{B_2}[h]$  is a continuous map from  $L^2(\mathcal{X} \times \mathcal{X})$  to itself.

Then,  $B_2 \circ B_1$  converges in distribution to a function-valued Gaussian process as  $J \rightarrow \infty$ , and

$$c_{B_2 \circ B_1}[\mathbf{f}_1, \mathbf{f}_2] = c_{B_2}[c_{B_1}[\mathbf{f}_1, \mathbf{f}_2]]. \quad (11)$$

For each operator discussed in Section 2.1, below we state the conditions under which they are function-valued Gaussian processes, and derive their covariance functions.

**Point-wise linear operator.** Given a vector  $\mathbf{w} \in \mathbb{R}^d$  and a function  $\mathbf{f}: \mathcal{X} \rightarrow \mathbb{R}^d$ , then, define the linear operator  $(\mathbf{w}^T \mathbf{f}): \mathcal{X} \rightarrow \mathbb{R}$  such that  $(\mathbf{w}^T \mathbf{f})(\mathbf{x}) = \sum_{p=1}^d w_p f_p(\mathbf{x})$ . If the entries of the weight vector follow an i.i.d. Gaussian distribution, i.e.  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ , then, this is a centered Gaussian process taking values from  $L^2(\mathcal{X}; \mathbb{R}^d)$  to  $L^2(\mathcal{X}; \mathbb{R})$  with covariance function:

$$c_{\mathbf{w}}[\mathbf{f}_1, \mathbf{f}_2](\mathbf{x}_1, \mathbf{x}_2) = \sigma^2 \mathbf{f}_2^T(\mathbf{x}_2) \mathbf{f}_1(\mathbf{x}_1), \text{ such that,} \quad (12)$$

$$\text{cov}(\langle \mathbf{h}_2, \mathbf{w}^T \mathbf{f}_2 \rangle \langle \mathbf{h}_1, \mathbf{w}^T \mathbf{f}_1 \rangle) = \iint_{\mathcal{X}} \mathbf{h}_2(\mathbf{x}_2) \mathbf{h}_1(\mathbf{x}_1) c_{\mathbf{w}}[\mathbf{f}_1, \mathbf{f}_2](\mathbf{x}_1, \mathbf{x}_2) d\mu_{\mathcal{X}}(\mathbf{x}_1) d\mu_{\mathcal{X}}(\mathbf{x}_2). \quad (13)$$

Note that  $c_{\mathbf{w}}$  only depends on the function  $\mathbf{f}_2^T \mathbf{f}_1: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , so we abuse notation and write  $c_{\mathbf{w}}[\mathbf{f}_1, \mathbf{f}_2] = c_{\mathbf{w}}[\mathbf{f}_2^T \mathbf{f}_1]$ . Moreover, this function is homogeneous:  $\alpha c_{\mathbf{w}}[\mathbf{f}_2^T \mathbf{f}_1] = c_{\mathbf{w}}[\alpha \mathbf{f}_2^T \mathbf{f}_1]$ , for  $\alpha > 0$ .

**Kernel integral operator.** As defined in Section 2.1, given a function  $\mathbf{k}: \mathcal{Y} \times \mathcal{X} \rightarrow \mathbb{R}^d$  and an input function  $\mathbf{f}: \mathcal{X} \rightarrow \mathbb{R}^d$ , we consider the linear operator  $A_{\mathbf{k}^T}[\mathbf{f}]: \mathcal{Y} \rightarrow \mathbb{R}$ . If  $\mathbf{k}$  follows an i.i.d. GP such that  $\mathbf{k} \in L^2(\mathcal{Y} \times \mathcal{X}) \sim \text{GP}(0, c_{\mathbf{k}})$ , then we have that  $A_{\mathbf{k}^T}$  is a centered function-valued GP mapping from  $L^2(\mathcal{X}; \mathbb{R}^d)$  to  $L^2(\mathcal{Y})$  with covariance function, denoted here by:

$$c_{A_{\mathbf{k}^T}}[\mathbf{f}_1, \mathbf{f}_2](\mathbf{y}_1, \mathbf{y}_2) = \iint_{\mathcal{X}} c_{\mathbf{k}}(\mathbf{y}_1, \mathbf{x}_1, \mathbf{y}_2, \mathbf{x}_2) \mathbf{f}_2^T(\mathbf{x}_2) \mathbf{f}_1(\mathbf{x}_1) d\mu_{\mathcal{X}}(\mathbf{x}_1) d\mu_{\mathcal{X}}(\mathbf{x}_2) \quad (14)$$

$$= A_{c_k}[\mathbf{f}_2^\top \mathbf{f}_1](\mathbf{y}_1, \mathbf{y}_2). \quad (15)$$

Note  $c_{A_{k\tau}}$  also only depends on the inner product of  $\mathbf{f}_2$  and  $\mathbf{f}_1$  and is homogeneous. Thus, we denote  $c_{A_{k\tau}}[\mathbf{f}_1, \mathbf{f}_2] = c_{A_{k\tau}}[\mathbf{f}_2^\top \mathbf{f}_1]$ .

**Point-wise element-wise activation.** Given a non-linear function  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ , we abuse the notation and define the non-linear operator  $\sigma[\cdot]: L^2(\mathcal{X}) \rightarrow L^2(\mathcal{X})$  as

$$\sigma[f](\mathbf{x}) = \sigma(f(\mathbf{x})). \quad (16)$$

Note that some restrictions on  $\sigma$  need to be placed for this to be a well-defined operator in  $L^2(\mathcal{X})$ . As an example of such condition, for their theoretical analysis, Kovachki et al. (2023) restricts activations to measurable linearly bounded functions, noting that the popular ReLU, ELU, tanh, and sigmoid activations satisfy this condition. In Appendix B.1, we provide a proof that this condition is sufficient for finite measure domains.

Consider a centered Gaussian operator  $B: L^2(\mathcal{X}) \rightarrow \mathcal{H}(\mathcal{X})$  with covariance function  $c_B$  such that  $\mathcal{H}(\mathcal{X}) \subset L^2(\mathcal{X})$  is an RKHS with reproducing kernel  $k_{\mathcal{H}}$ . When  $\sigma[\cdot]$  is a well-defined operator, the operator  $(\sigma \circ B)$  is a random operator in  $L^2(\mathcal{X}) \rightarrow L^2(\mathcal{X})$  with covariance function:

$$c_{(\sigma \circ B)}(\mathbf{x}_1, \mathbf{x}_2) = \text{cov}((\sigma \circ B)[f_1](\mathbf{x}_1), (\sigma \circ B)[f_2](\mathbf{x}_2)). \quad (17)$$

Now, since  $B[f_1]$  and  $B[f_2]$  are Gaussian processes with outputs in an RKHS  $\mathcal{H}$ , we can consider the following bivariate Gaussian r.v.  $\mathbf{b}_{[f_1, f_2]} = [B[f_1](\mathbf{x}_1), B[f_2](\mathbf{x}_2)]^\top$ :

$$\mathbf{b}_{[f_1, f_2]} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} c_B[f_1, f_1](\mathbf{x}_1, \mathbf{x}_1) & c_B[f_1, f_2](\mathbf{x}_1, \mathbf{x}_2) \\ c_B[f_2, f_1](\mathbf{x}_2, \mathbf{x}_1) & c_B[f_2, f_2](\mathbf{x}_2, \mathbf{x}_2) \end{bmatrix}\right). \quad (18)$$

This random variable is well-defined due to the reproducing property,  $B[f](\mathbf{x}) = \langle k_{\mathcal{H}}(\cdot, \mathbf{x}), B[f] \rangle$ .

Thus, we can continue to conclude

$$c_{(\sigma \circ B)}(\mathbf{x}_1, \mathbf{x}_2) = \int_{\mathbb{R}^2} \sigma(b_{f_1}) \sigma(b_{f_2}) \mathcal{N}(\mathbf{b}_{[f_1, f_2]} \mid \mathbf{0}, \mathbf{L}^\top \mathbf{L}) d\mathbf{b}_{[f_1, f_2]} \quad (19)$$

$$= \int_{\mathbb{R}^2} \sigma(\mathbf{l}_1^\top \boldsymbol{\xi}) \sigma(\mathbf{l}_2^\top \boldsymbol{\xi}) \mathcal{N}(\boldsymbol{\xi} \mid \mathbf{0}, \mathbf{I}) d\boldsymbol{\xi} \quad (20)$$

$$=: c_\sigma[c_B[f_1, f_2]](\mathbf{x}_1, \mathbf{x}_2), \quad (21)$$

where  $\mathbf{L}$  is a square-root of the covariance matrix of  $\mathbf{b}_{[f_1, f_2]}$  and  $\mathbf{l}_i$  is the  $i$ -th row of this matrix.

The expected value  $c_\sigma$  as a function of  $\mathbf{l}_1$  and  $\mathbf{l}_2$  in Eq. (20) is known as the *dual kernel* of  $\sigma$ . The dual kernels for many activation functions have closed-form solutions (e.g., sigmoid (Williams, 1996, Eq. 10) and ReLU (Cho and Saul, 2009, Eq. 1)) or can be efficiently approximated (Han et al., 2022). Any of these solutions can be directly used in our context by computing the covariance matrix of  $\mathbf{b}_{[f_1, f_2]}$  and applying the rows of its square-root as arguments.

In conclusion, we construct an covariance function  $c_\sigma: \mathcal{H}(\mathcal{X} \times \mathcal{X}) \rightarrow L^2(\mathcal{X} \times \mathcal{X})$  such that, for a given covariance function  $c_B: L^2(\mathcal{X} \times \mathcal{X}) \rightarrow \mathcal{H}(\mathcal{X} \times \mathcal{X})$ :

$$\langle h_2 h_1, c_\sigma[c_B[f_1, f_2]] \rangle = \text{cov}(\langle h_1, (\sigma \circ B)[f_1] \rangle, \langle h_2, (\sigma \circ B)[f_2] \rangle), \quad (22)$$

for all  $f_1, f_2, h_1, h_2 \in L^2(\mathcal{X})$

### 3.2 Outline of the proof for Theorem 3.1

We now describe a sketch for the proof, we refer the readers to Appendix B.3 for the complete proof.

**Step 1.** We start by showing that, under the conditions of Theorem 3.1, each linear layer in a neural operator layer is a function-valued Gaussian process when conditioned on its inputs. Moreover, as discussed in Section 3.1, the conditional covariance function of each node on each layer only depends on the *empirical covariance function* of its inputs  $\bar{c}[\mathbf{f}, \mathbf{g}](\mathbf{x}', \mathbf{x}) = (1/J) \sum_{j=1}^J g_j(\mathbf{x}) f_j(\mathbf{x}')$ . We denote this dependency by writing the conditional covariance function as  $c^{(\ell|\ell-1)}[\cdot](\mathbf{x}, \mathbf{x}')$ .

**Step 2.** Due to the chosen prior distribution of each layer, we know that each node in  $H_\ell[\cdot] \in \mathbb{R}^{J_\ell}$  is i.i.d. and, therefore, we can apply Lemma 3.2 to conclude that, as  $J_{\ell-1} \rightarrow \infty$ , the covariance  $c^{(\ell|\ell-1)}[H_{\ell-1}[\mathbf{g}]^\top H_{\ell-1}[\mathbf{f}]/J_{\ell-1}]$  converges almost surely to  $c^{(\ell|\ell-1)}[c_{H_{\ell-1}}[\mathbf{f}, \mathbf{g}]]$ .

**Step 3.** Combining both steps, we show, by induction on  $\ell$  up until  $\ell = d$ , that, as  $J \rightarrow \infty$ , the covariance function of  $Z_J[\mathbf{f}]$  is simply the composition of all the previous covariances as denoted in Step 1. So, we have that the covariance function of  $Z_\infty$  is:

$$c^{(d)}[\mathbf{f}, \mathbf{g}] = c^{(d|d-1)}[c^{(d-1|d-2)}[\dots c^{(2|1)}[c_{H_1}[\mathbf{f}, \mathbf{g}]] \dots]]. \quad (23)$$

Finally, denote  $c^{(d)}$  as  $c_\infty$ .

## 4 Parametrizations and computations

To apply the results of Theorem 3.1, we must specify a covariance function for the integral kernel operators  $A_{\mathbf{K}}$ . This choice corresponds to selecting a particular neural operator parameterization, following the approach of Kovachki et al. (2023).

In this section, we derive the operator-valued covariance functions for  $A_{\mathbf{K}}$  under two parametrizations of the integral operator. The first is based on the band-limited Fourier Neural Operator (Section 4.1); the second models the kernel as a non-stationary process, with a prior distribution derived from the classical Matérn family of covariance functions (Section 4.2).

A common assumption for both cases is that the input domain is compact. This ensures that samples of the kernel components  $\mathbf{k}_j$  reside in a  $L^2$  space. By further choosing the domain to be the  $d_x$ -dimensional flat torus  $\mathbb{T}^{d_x} = \mathbb{R}^{d_x} / 2\pi\mathbb{Z}^{d_x}$ , we are able to exploit Fourier analysis tools. In particular, by assuming that the input functions are band-limited enables tractable computations through the connection of Fourier series with discrete Fourier transforms for evaluations in regular grids.

### 4.1 Fourier neural operator

Out of the parametrizations proposed by Kovachki et al. (2023), the Fourier neural operator is the most popular due to its computational benefits. By imposing three assumptions into the convolutions kernel – periodicity, shift-invariance, and band-limitedness – we can use the convolution theorem to compute the integrals using sums up to the chosen band-limit of the kernel in the Fourier space.

Concretely, assuming periodicity is equivalent to choosing the domain to be some  $d_x$ -dimensional flat torus  $\mathcal{X} = \mathbb{T}^{d_x}$ , and shift-invariance means kernels satisfy  $\mathbf{k}_j(\mathbf{w}, \mathbf{x}) = \mathbf{k}_j(\mathbf{w} - \mathbf{x})$ , where we abuse notation and represent the kernel as a univariate function of the same name  $\mathbf{k}_j : \mathbb{T}^{d_x} \rightarrow \mathbb{R}^d$ . Under these conditions, any  $\mathbf{k}_j$  admits a Fourier series representation:

$$\mathbf{k}_j(\mathbf{w} - \mathbf{x}) = \sum_{\mathbf{s} \in \mathbb{Z}^{d_x}} \text{FS}_{\mathbf{s}}[\mathbf{k}_j] \psi_{-\mathbf{s}}(\mathbf{w} - \mathbf{x}), \quad (24)$$

where  $\text{FS}_{\mathbf{s}}$  is the  $(s_1, \dots, s_{d_x})$ -th coefficient of the Fourier series and  $\psi_{\mathbf{s}}(\mathbf{x}) = \exp[-i \cdot \mathbf{s}^\top \mathbf{x}]$ , with  $i = \sqrt{-1}$  being the imaginary unit. Moreover, to have a band-limited kernel implies that only finitely many Fourier coefficients are non-zero, i.e. there is some  $B_j \in \mathbb{N}$ ,  $1 \leq j \leq d_x$ , such that  $\text{FS}_{\mathbf{s}}[\mathbf{k}_j] = 0$ , if  $|s_j| > B_j$ , for all  $1 \leq j \leq d_x$ .

Under these conditions, despite all input functions  $\mathbf{f}$  being represented with a (potentially infinite) Fourier series, by the convolution theorem, the NO layer  $H_j[\mathbf{f}]$  is band-limited and its Fourier series coefficients can be computed directly from the product of Fourier coefficients of the kernel function  $\mathbf{k}$  and the input function  $\mathbf{f}$ . Thus, we have that:

$$\text{FS}_{\mathbf{s}}[H_j[\mathbf{f}]] = \text{FS}_{\mathbf{s}}[\mathbf{k}_j]^\top \text{FS}_{\mathbf{s}}[\mathbf{f}] + \mathbf{w}_j^{(1)\top} \text{FS}_{\mathbf{s}}[\mathbf{f}]. \quad (25)$$

**Parameterization of an FNO.** Following Section 3.1, when  $\mathbf{k}$  is a  $\mathbb{R}^d$ -valued GP, the kernel integral operator  $A_{\mathbf{k}}$  is a function-valued Gaussian process with covariance function of  $A_{\mathbf{k}}$  in terms of the covariance function of  $\mathbf{k}$ ,  $\mathbf{C}_{\mathbf{k}}$ :

$$c_{A_{\mathbf{k}}\tau}[\mathbf{f}_1, \mathbf{f}_2](\mathbf{z}_1, \mathbf{z}_2) = A_{\mathbf{C}_{\mathbf{k}}}[\mathbf{f}_2^\top \mathbf{f}_1](\mathbf{z}_1, \mathbf{z}_2). \quad (26)$$

The most popular choice proposed by Kovachki et al. (2023) is to directly parametrize the Fourier coefficients of the kernel. Thus, we let these  $2B + 1$  Fourier coefficients follow i.i.d. centered complex Gaussian distributions with variance  $\sigma_{\mathbf{k}}^2$  (Appendix A.1), obtaining the covariance function  $\mathbf{C}_{\mathbf{k}}$ :

$$\mathbf{k}(\mathbf{z} - \mathbf{x}) = \sum_{\mathbf{s} \in \{-B, \dots, B\}^{d_x}} \text{FS}_{\mathbf{s}}[\mathbf{k}] \cdot \psi_{-\mathbf{s}}(\mathbf{z} - \mathbf{x}) \sim \text{GP}(0, \mathbf{C}_{\mathbf{k}}), \quad (27)$$

$$\mathbf{C}_k((z - x), (z' - x')) = \sigma_k^2 \mathbf{I}_d \sum_{s \in \{-B, \dots, B\}^{d_x}} \psi_{-s}(z - x) \psi_s(z' - x'), \quad (28)$$

where  $B$  is a hyperparameter of the model controlling the band-limit of the integral kernel.

This allows us to derive a finite-sum representation of the covariance of  $A_k$  parameterized by  $\sigma_k^2$ .

$$c_{A_k^\top}[\mathbf{f}_1, \mathbf{f}_2](z, z') = \sigma_k^2 (2\pi)^{2d_x} \sum_{s \in \{-B, \dots, B\}^{d_x}} \text{FS}_{-s}[\mathbf{f}_2]^\top \text{FS}_s[\mathbf{f}_1] \psi_{-s}(z - z'). \quad (29)$$

## 4.2 Toroidal Matérn operator

In this section, we propose a model in which the kernel does not admit a shift-invariant decomposition. Another popular decomposition used in the Gaussian process literature is the tensor-product factorization, where the covariance function of a GP factorizes over the input dimension. That is,  $f: \mathbb{R}^{d_x} \rightarrow \mathbb{R} \sim \text{GP}(0, c)$ , where  $c(\mathbf{a}, \mathbf{b}) = \prod_j c_j(a_j, b_j)$ ; although the covariance factorizes over the input dimensions, in general, samples from  $f$  do not.

Our proposal will make use of the ubiquitous Matérn family of covariance functions, which are characterized by the smoothness parameter  $\nu$ . Following Borovitskiy et al. (2020), we define the Matérn covariance functions in the  $d_x$ -dimensional flat torus  $\mathbb{T}^{d_x} = \mathbb{T} \otimes \dots \otimes \mathbb{T}$  as:

$$c(\mathbf{x}, \mathbf{x}'; \nu, \ell) = (2\pi)^{-d_x} \sum_{\mathbf{n} \in \mathbb{Z}^{d_x}} \psi_{\mathbf{n}}(\mathbf{x}) \psi_{-\mathbf{n}}(\mathbf{x}') \hat{c}\left(\sum_{j=1}^{d_x} n_j^2; \nu, \ell\right). \quad (30)$$

where  $\ell$  is the lengthscale hyper-parameter and the spectral density  $\hat{c}$  is defined as:

$$\hat{c}(\lambda; \nu, \ell) = \begin{cases} \exp\left[-\frac{\ell^2}{2}\lambda\right], & \text{if } \nu = \infty, \\ \left(\frac{2\nu}{\ell^2} + \lambda\right)^{-\nu - \frac{d}{2}}, & \text{otherwise.} \end{cases} \quad (31)$$

In general, this kernel is not tensor-product factorized, but for the special case of  $\nu = \infty$ , the squared exponential covariance function, the factorization holds (Appendix A.2.1). Thus, in general, we enforce the tensor-product factorization:

$$c(\mathbf{x}, \mathbf{x}'; \nu, \ell) = \prod_{j=1}^{d_x} c(x_j, x'_j; \nu, \ell_j) = (2\pi)^{-d_x} \sum_{\mathbf{n} \in \mathbb{Z}^{d_x}} \psi_{\mathbf{n}}(\mathbf{x}) \psi_{-\mathbf{n}}(\mathbf{x}') \prod_{j=1}^{d_x} \hat{c}(n_j^2; \nu, \ell_j), \quad (32)$$

where  $\ell$  is the automatic relevance determination (ARD) lengthscale hyper-parameter.

**Parameterization of a toroidal Matérn operator.** So, we consider a convolution kernel  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^d$  defined as the product of Matérn covariance functions:

$$\mathbf{C}_k(z, x, z', x') = c(z, z'; \nu_z, \ell_z) c(x, x'; \nu_x, \ell_x) \mathbf{I}_d. \quad (33)$$

where  $c(\cdot, \cdot; \nu, \ell): \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is the Matérn covariance functions with smoothness parameter  $\nu$  and length-scale  $\ell$ .

Again, following Section 3.1, we express the covariance of the operator as:

$$c_{A_k^\top}[\mathbf{f}_1, \mathbf{f}_2](z_1, z_2) = A_{C_k}[\mathbf{f}_2^\top \mathbf{f}_1](z_1, z_2). \quad (34)$$

Thus, by using the identity Eq. (33), we can derive:

$$c_{A_k^\top}[\mathbf{f}_1, \mathbf{f}_2](z, z') = c(z, z'; \nu_z, \ell_z) (2\pi)^{d_x} \sum_{\mathbf{n} \in \mathbb{Z}^{d_x}} \text{FS}_{\mathbf{n}}[\mathbf{f}_1]^\top \text{FS}_{-\mathbf{n}}[\mathbf{f}_2] \prod_{j=1}^{d_x} \hat{c}(n_j^2; \nu_x, \ell_{x,j}). \quad (35)$$

## 5 Experimental validation

In this section, we empirically show i) the agreement between finite width neural operators with increasing width and their corresponding infinite-width neural operator, and ii) evaluate our model against FNO in a regression task.

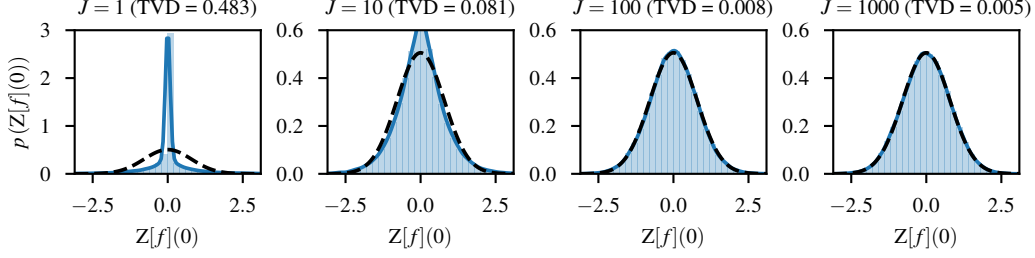


Figure 1: A density estimation of the empirical distribution of the output of increasing channel dimension compared to the infinite width distribution. On top of each plot we show the total variation distance of the empirical distribution against the infinite width distribution.

Section 5.1 explores the distribution of untrained randomly initialized Fourier neural operators of varying width and the distribution of the infinite-width FNO ( $\infty$ -FNO). As expected from the theoretical results, these distributions should eventually match as the hidden dimension increases.

Section 5.2 considers two tasks: a synthetic regression example, where the task is to predict the output of a non-linear operator, and the task of predicting the final evolution of Burgers' equation given the initial state. This situation is not covered in our theory, since it only applies to the distribution of the neural operators at initialization, but our experiments show the behavior of the posteriors of infinite-width neural operators against Adam trained finite-width neural operators of increasing width.

Throughout this section, our starting point is a single hidden-layer neural operator  $Z[f]: \mathbb{T} \rightarrow \mathbb{R} := (\mathbf{w}_2^\top \circ \text{ReLU} \circ (A_{\mathbf{K}} + \mathbf{W}_1))[f]$ . More details for each experiment can be found in Appendix C. All experiments were implemented in Python, mainly based on the GPyTorch (Gardner et al., 2018) library, and run in a desktop computer using a Titan RTX. Code is available at <https://github.com/spectraldani/infinite-neural-operator>.

## 5.1 Empirical demonstration of results

In this experiment, we demonstrate that our analytical computation of the variance for a neural operator layer  $H$  agrees with empirical estimates, and we validate Theorem 3.1 by showing that the output of a neural operator  $Z$  converges to a Gaussian distribution as the number of hidden channels  $J$  increases.

Throughout all experiments, the input function  $f: \mathbb{T} \rightarrow \mathbb{R}$  has band-limit  $B = 3$ , with its output values  $f(x)$  sampled from a uniform distribution  $\mathcal{U}(-1, 1)$ . Both operators are evaluated at  $x = 0$ , so we analyze the empirical distributions of  $H[f](0)$  and  $Z[f](0)$ , respectively.

Following Section 4.1, we parametrize the integral kernel operators using band-limited functions. The band-limit of the kernel is set equal to that of the input  $f$ , and the kernel coefficients are drawn from a Gaussian distribution with unit variance scaled by the inverse of the number of hidden channels.

As shown in Fig. 2, the empirical estimate of the variance converges to the theoretical value as the number of Monte Carlo samples increases, supporting the correctness of our variance computation. Furthermore, Fig. 1 shows that, as the number of hidden channels grows, the total variation distance (TVD) between the empirical distribution and a Gaussian distribution approaches zero, thereby further verifying the validity of Theorem 3.1.

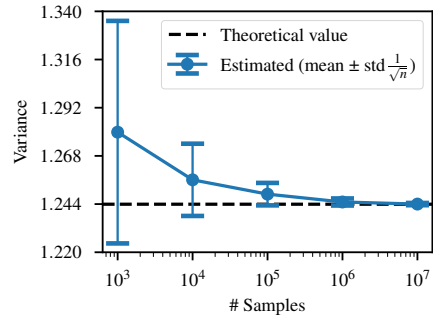


Figure 2: Plot of the MC estimate for the variance of  $H[f](0)$  against our analytical computation (Sec. 3.1).

## 5.2 Regression tasks

In this task, we're given  $n$  pairs of 1D functions  $\{f_i, g_i\}_{i=1}^n$  evaluated in a grid with  $m = 2B_m + 1$  points. We consider FNOs of increasing width  $J \in \mathbb{N}^+$ , as well as  $\infty$ -FNOs, both with increasing kernel band-limits  $B \in \{1, 5, 20\}$ . These models will be trained on two datasets: (a) A operator



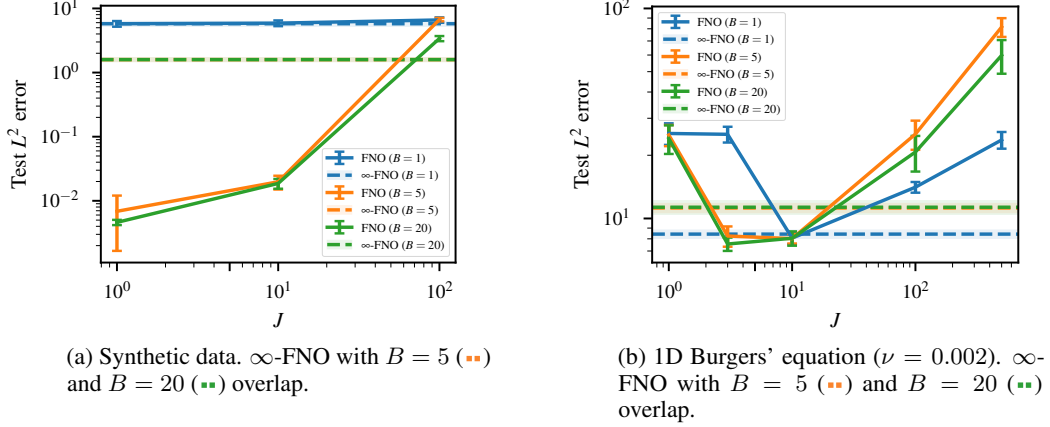


Figure 3: Results for the regression experiments. Mean and std. of test  $L^2$  loss as a function of width  $J$  for different band-limits  $B$ .

generated by a randomly-initialized ground truth FNO  $Z_{\text{true}}$  with band-limit  $B = 5$  and width  $J = 1$ . We sample  $n = 100$  input functions  $f_i: \mathbb{T} \rightarrow \mathbb{R}$  with uniformly-distributed outputs  $\mathcal{U}(-1, 1)$  and band-limit  $B_m = 5$ . (b) 1D Burgers' equation dataset from Takamoto et al. (2022) with  $\nu = 0.002$ . The task is to predict the end state ( $t = 2$ ) given the initial condition ( $t = 0$ ). Due to memory constraints, we subsample the total dataset data to  $n = 100$  functions and a grid size of  $m = 103$ .

The hyperparameters of the  $\infty$ -FNO are estimated using L-BFGS, while the parameters of the FNOs are optimized with Adam using a step size of 0.001. We evaluate all models using 5-fold cross-validation and report the average and standard deviation of the empirical  $L^2$  norm of the prediction error. For  $\infty$ -FNOs, we use the posterior mean as the prediction.

In general, we do not expect close agreement between the predictive performance of  $\infty$ -FNOs and finite-width FNOs, as the former corresponds to a Bayesian estimate while the latter are trained by minimizing an empirical risk, nonetheless, as observed in Figs. 3a and 3b, there is consistency between the gap of hyperparameters in the same model class.

In the synthetic case, as we know the band-limits of the ground truth operator, Fig. 3a shows that the models are only able to accurately predict the output when their band-limits exceed that of the ground truth.

## 6 Related works

**Infinite limits of stochastic NNs.** The study of infinite-width Bayesian neural networks began with the seminal work of Neal (1995) and was later extended to deep architectures (Lee et al., 2018; Yang, 2019; Matthews et al., 2018). Our analysis builds on the ideas developed by Matthews et al. (2018). From the outset, these infinite-width models were considered "disappointing" (Neal, 1995), a view reinforced by findings that neither the Bayesian limit nor the neural tangent kernel limit learns features from data (Aitchison, 2020). However, recent work shows these models still reflect the different inductive biases of their finite-width counterparts (Novak et al., 2019), and that alternative initialization distributions can enable feature learning in the infinite-width setting (Yang and Hu, 2021).

**Bayesian neural operators.** Several works have investigated approximate Bayesian uncertainty quantification in finite-width neural operators using function-valued Gaussian processes. Magnani et al. (2022, 2024) both employ last-layer Laplace approximations to construct GP approximations of the posterior distribution. In addition, Magnani et al. (2022) considers the case where the kernel  $\mathbf{K}$  of the integral operator  $A_{\mathbf{K}}$  follows a Matérn GP prior. However, their analysis is restricted to the finite-width regime on compact subsets of Euclidean space, whereas our work focuses on the flat torus.

**Kernel methods for operator learning.** Batlle et al. (2024) propose the use of kernel methods for operator learning, leveraging operator-valued kernels and the representer theorem in their corresponding RKHS. Their results are promising and highlight the potential of kernel-based approaches in this

domain. Our contribution introduces an additional way to construct operator-valued kernels based on neural operators, enabling new kernel-based models for operator learning.

## 7 Discussion

In this work, we formalized the concept of infinite-width Bayesian neural operators, established their existence (Theorem 3.1), and described how to compute their associated covariance functions (Section 4). We validated these results empirically (Section 5.1) and further assessed the performance of these models in a regression setting (Section 5.2).

Our contributions lay a foundation for future investigations, particularly in bridging the gap between SGD-trained finite-width neural operators and their infinite-width counterparts. Addressing this challenge will require extending the neural tangent kernel framework (Jacot et al., 2018; Lee et al., 2019) to settings involving Hilbert space-valued functions. Moreover, while we focused on the ubiquitous FNO architecture, deriving covariance functions for other architectures, such as the graph neural operator (Kovachki et al., 2023), remains an open direction.

**Limitations.** Our current implementation for computing the required kernel quantities scales with cubically in both the evaluation grid size and the number of training functions. We anticipate that future work can improve computational efficiency by leveraging advances from the Gaussian process literature to improve scalability and efficiency (Borovitskiy et al., 2020; Gilboa et al., 2015).

## Acknowledgments and Disclosure of Funding

YZ acknowledges support by the Engineering and Physical Sciences Research Council with grant number EP/S021566/1. YS was supported by Fundação Carlos Chagas Filho de Amparo à Pesquisa do Estado do Rio de Janeiro (FAPERJ) through the Jovem Cientista do Nosso Estado Program (E-26/201.375/2022 (272760)) and by Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq) through the Productivity in Research Scholarship (306695/2021-9, 305159/2025-9). DM was supported by the Fundação Carlos Chagas Filho de Amparo à Pesquisa do Estado do Rio de Janeiro (FAPERJ) (SEI-260003/000709/2023) and the Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq) (404336/2023-0, 305692/2025-9).

## References

- Laurence Aitchison. Why bigger is not always better: on finite and infinite neural networks. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, 2020. URL <https://proceedings.mlr.press/v119/aitchison20a.html>.
- Pau Batlle, Matthieu Darcy, Bamdad Hosseini, and Houman Owhadi. Kernel methods are competitive for operator learning. *Journal of Computational Physics*, 496, 2024. URL <https://doi.org/10.1016/j.jcp.2023.112549>.
- Viacheslav Borovitskiy, Alexander Terenin, Peter Mostowsky, and Marc Deisenroth. Matérn Gaussian processes on Riemannian manifolds. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. URL <https://papers.nips.cc/paper/2020/hash/92bf5e6240737e0326ea59846a83e076-Abstract.html>.
- Youngmin Cho and Lawrence Saul. Kernel methods for deep learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2009. URL <https://proceedings.neurips.cc/paper/2009/hash/5751ec3e9a4feab575962e78e006250d-Abstract.html>.
- Jacob R. Gardner, Geoff Pleiss, David Bindel, Kilian Q. Weinberger, and Andrew Gordon Wilson. GPYtorch: Blackbox matrix-matrix Gaussian process inference with gpu acceleration. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. URL [https://papers.nips.cc/paper\\_files/paper/2018/hash/27e8e17134dd7083b050476733207ea1-Abstract.html](https://papers.nips.cc/paper_files/paper/2018/hash/27e8e17134dd7083b050476733207ea1-Abstract.html).
- Elad Gilboa, Yunus Saatçi, and John P. Cunningham. Scaling multidimensional inference for structured Gaussian processes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37, 2015. URL <http://dx.doi.org/10.1109/TPAMI.2013.192>.

- Insu Han, Amir Zandieh, Jaehoon Lee, Roman Novak, Lechao Xiao, and Amin Karbasi. Fast neural kernel embeddings for general activations. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/hash/e7be1f4c6212c24919cd743512477c13-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/e7be1f4c6212c24919cd743512477c13-Abstract-Conference.html).
- Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. URL <https://proceedings.neurips.cc/paper/2018/hash/5a4be1fa34e62bb8a6ec6b91d2462f5a-Abstract.html>.
- Hachem Kadri, Emmanuel Duflos, Philippe Preux, Stéphane Canu, Alain Rakotomamonjy, and Julien Audiffren. Operator-valued kernels for learning from functional response data. *Journal of Machine Learning Research (JMLR)*, 17, 2016. URL <http://jmlr.org/papers/v17/11-315.html>.
- Nikola Kovachki, Zongyi Li, Burigede Liu, Kamyar Azizzadenesheli, Kaushik Bhattacharya, Andrew Stuart, and Anima Anandkumar. Neural operator: Learning maps between function spaces with applications to PDEs. *Journal of Machine Learning Research (JMLR)*, 24, 2023. URL <http://jmlr.org/papers/v24/21-1524.html>.
- Jaehoon Lee, Jascha Sohl-Dickstein, Jeffrey Pennington, Roman Novak, Sam Schoenholz, and Yasaman Bahri. Deep neural networks as Gaussian processes. In *International Conference on Learning Representations (ICLR)*, 2018. URL <https://openreview.net/forum?id=B1EA-M-OZ>.
- Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/0d1a9651497a38d8b1c3871c84528bd4-Abstract.html>.
- Zongyi Li, Hongkai Zheng, Nikola Kovachki, David Jin, Haoxuan Chen, Burigede Liu, Kamyar Azizzadenesheli, and Anima Anandkumar. Physics-informed neural operator for learning partial differential equations. *ACM/IMS Journal of Data Science*, 1(3), 2024. URL <https://doi.org/10.1145/3648506>.
- Emilia Magnani, Nicholas Krämer, Runa Eschenhagen, Lorenzo Rosasco, and Philipp Hennig. Approximate Bayesian neural operators: Uncertainty quantification for parametric PDEs. 2022. URL <https://arxiv.org/abs/2208.01565>.
- Emilia Magnani, Marvin Pförtner, Tobias Weber, and Philipp Hennig. Linearization turns neural operators into function-valued Gaussian processes. 2024. URL <https://arxiv.org/abs/2406.05072>.
- Alexander G. de G. Matthews, Mark Rowland, Jiri Hron, Richard E. Turner, and Zoubin Ghahramani. Gaussian process behaviour in wide deep neural networks. 2018. URL <https://arxiv.org/abs/1804.11271>.
- Edith Mourier. L-random elements and l'-random elements in Banach spaces. In *Contributions to Probability Theory*, pages 231–242. University of California Press, December 1956. URL <https://doi.org/10.1525/9780520350670-017>.
- Radford M. Neal. *Bayesian Learning for Neural Networks*. PhD thesis, University of Toronto, 1995.
- Roman Novak, Lechao Xiao, Yasaman Bahri, Jaehoon Lee, Greg Yang, Jiri Hron, Daniel A. Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. Bayesian deep convolutional networks with many channels are Gaussian processes. In *International Conference on Learning Representations (ICLR)*, 2019. URL <https://openreview.net/forum?id=B1g30jOqF7>.
- Jaideep Pathak, Shashank Subramanian, Peter Harrington, Sanjeev Raja, Ashesh Chattopadhyay, Morteza Mardani, Thorsten Kurth, David Hall, Zongyi Li, Kamyar Azizzadenesheli, Pedram Hassanzadeh, Karthik Kashinath, and Animashree Anandkumar. Fourcstnet: A global data-driven high-resolution weather model using adaptive fourier neural operators. 2022. URL <https://arxiv.org/abs/2202.11214>.
- Makoto Takamoto, Timothy Praditia, Raphael Leiteritz, Dan MacKinlay, Francesco Alesiani, Dirk Pflüger, and Mathias Niepert. PDEBench: An extensive benchmark for scientific machine learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. URL [https://papers.neurips.cc/paper\\_files/paper/2022/hash/0a9747136d411fb83f0cf81820d44afb-Abstract-Datasets\\_and\\_Benchmarks.html](https://papers.neurips.cc/paper_files/paper/2022/hash/0a9747136d411fb83f0cf81820d44afb-Abstract-Datasets_and_Benchmarks.html).
- Christopher Williams. Computing with infinite networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 1996. URL <https://proceedings.neurips.cc/paper/1996/hash/ae5e3ce40e0404a45ecacaaf05e5f735-Abstract.html>.

Greg Yang. Wide feedforward or recurrent neural networks of any architecture are Gaussian processes. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. URL [https://papers.neurips.cc/paper\\_files/paper/2019/hash/5e69fda38cda2060819766569fd93aa5-Abstract.html](https://papers.neurips.cc/paper_files/paper/2019/hash/5e69fda38cda2060819766569fd93aa5-Abstract.html).

Greg Yang and Edward J. Hu. Tensor programs IV: Feature learning in infinite-width neural networks. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, 2021. URL <https://proceedings.mlr.press/v139/yang21c.html>.

---

# Infinite Neural Operators: Gaussian processes on functions (Supplemental Materials)

---

## A Covariance Function Computation

In this section, we will work out in detail the computations of Section 4, first for the Fourier neural operator (FNO) case and later for the toroidal Matérn operator.

### A.1 Fourier neural operator

Under the direct parametrization of the integral kernel operator, the coefficients of the kernel's Fourier series (FS) are parametrized and randomly sampled at initialization. Therefore, our first step is to derive what the Gaussian process distribution of a band-limited function with i.i.d. Gaussian FS coefficients is.

**Fourier series.** Given a function on the  $d_x$ -dimensional torus  $\mathbf{f}(\cdot): \mathbb{T}^{d_x} \rightarrow \mathbb{R}^d$ ,  $\mathbf{f} = (f_1, \dots, f_d)$ , it can be represented in terms of a Fourier series:

$$f_p(\mathbf{x}) = \sum_{\mathbf{s} \in \mathbb{Z}^{d_x}} \text{FS}_{\mathbf{s}}[f_p] \psi_{-\mathbf{s}}(\mathbf{x}), \quad (36)$$

for  $p \in \{1, \dots, d\}$ , where,

$$\text{FS}_{\mathbf{s}}[f_p] = (2\pi)^{-d_x} \int_{[-\pi, \pi]^{d_x}} f_p(\mathbf{t}) \psi_{\mathbf{s}}(\mathbf{t}) d\mathbf{t}, \quad (37)$$

$$\psi_{\mathbf{s}}(\mathbf{x}) = \exp[-i \cdot \mathbf{s}^\top \mathbf{x}], \quad (38)$$

and  $i = \sqrt{-1}$  is the imaginary unit.

Note that, as  $f_p$  is a real-valued function, we also have that  $\text{FS}_{\mathbf{s}}[f_p] = \overline{\text{FS}_{-\mathbf{s}}[f_p]}$ , where  $\bar{z}$  is the complex conjugate.

**Gaussian distributed band-limited functions.** Consider the sequence  $\hat{\mathbf{f}}: [-B, \dots, B]^{d_x} \rightarrow \mathbb{C}$  defined as:

$$\Re \hat{f}_{\mathbf{0}} \sim \mathcal{N}(0, \sigma^2), \quad \Im \hat{f}_{\mathbf{0}} = 0, \quad (39)$$

$$\Re \hat{f}_{\mathbf{s}} \sim \mathcal{N}(0, \sigma^2/2), \quad \Im \hat{f}_{\mathbf{s}} \sim \mathcal{N}(0, \sigma^2/2), \quad (40)$$

$$\Re \hat{f}_{-\mathbf{s}} = \Re \hat{f}_{\mathbf{s}}, \quad \Im \hat{f}_{-\mathbf{s}} = -\Im \hat{f}_{\mathbf{s}}, \quad (41)$$

where  $\Re z$  and  $\Im z$  are the real and imaginary parts of the complex number  $z$ , respectively, and all random variables are independent of each other, except the conjugate duals  $\hat{f}_{\mathbf{s}}$  and  $\hat{f}_{-\mathbf{s}}$ . For  $\mathbf{s} \neq \mathbf{0}$ , the equations above can also be expressed as:

$$\begin{bmatrix} \Re \hat{f}_{\mathbf{s}} \\ \Im \hat{f}_{\mathbf{s}} \\ \Re \hat{f}_{-\mathbf{s}} \\ \Im \hat{f}_{-\mathbf{s}} \end{bmatrix} = \mathcal{N} \left( \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \frac{1}{2} \begin{bmatrix} \sigma^2 & 0 & \sigma^2 & 0 \\ 0 & \sigma^2 & 0 & -\sigma^2 \\ \sigma^2 & 0 & \sigma^2 & 0 \\ 0 & -\sigma^2 & 0 & \sigma^2 \end{bmatrix} \right). \quad (42)$$

With this in mind, the expectation of the product of two elements is:

$$\mathbb{E}[\hat{f}_{\mathbf{s}} \cdot \hat{f}_{\mathbf{s}'}] = \mathbb{E}[\Re \hat{f}_{\mathbf{s}} \Re \hat{f}_{\mathbf{s}'}] - \mathbb{E}[\Im \hat{f}_{\mathbf{s}} \Im \hat{f}_{\mathbf{s}'}] + i\mathbb{E}[\Re \hat{f}_{\mathbf{s}'} \Im \hat{f}_{\mathbf{s}}] + i\mathbb{E}[\Re \hat{f}_{\mathbf{s}} \Im \hat{f}_{\mathbf{s}'}] \quad (43)$$

$$= \mathbb{E}[\Re \hat{f}_{\mathbf{s}} \Re \hat{f}_{\mathbf{s}'}] - \mathbb{E}[\Im \hat{f}_{\mathbf{s}} \Im \hat{f}_{\mathbf{s}'}] \quad (44)$$

$$= \begin{cases} \mathbb{E}[\Re \hat{f}_{\mathbf{s}} \Re \hat{f}_{\mathbf{s}}] - \mathbb{E}[\Im \hat{f}_{\mathbf{s}} \Im \hat{f}_{\mathbf{s}}] & \text{if } \mathbf{s}' = \mathbf{s}, \\ \mathbb{E}[\Re \hat{f}_{\mathbf{s}} \Re \hat{f}_{\mathbf{s}}] + \mathbb{E}[\Im \hat{f}_{\mathbf{s}} \Im \hat{f}_{\mathbf{s}}] & \text{if } \mathbf{s}' = -\mathbf{s}, \\ 0 & \text{otherwise;} \end{cases} \quad (45)$$

$$= \begin{cases} \sigma^2 & \text{if } \mathbf{s}' = \mathbf{0} \text{ and } \mathbf{s} = \mathbf{0}, \\ \sigma^2/2 - \sigma^2/2 & \text{if } \mathbf{s}' = \mathbf{s}, \\ \sigma^2/2 + \sigma^2/2 & \text{if } \mathbf{s}' = -\mathbf{s}, \\ 0 & \text{otherwise;} \end{cases} \quad (46)$$

$$= \begin{cases} \sigma^2 & \text{if } \mathbf{s}' = -\mathbf{s}, \\ 0 & \text{otherwise.} \end{cases} \quad (47)$$

Thus, we can define the Gaussian process  $f: \mathbb{T}^{d_x} \rightarrow \mathbb{R}$  through a Fourier series representation:

$$f(\mathbf{x}) = \sum_{\mathbf{s} \in \{-B, \dots, B\}^{d_x}} \hat{f}_{\mathbf{s}} \psi_{-\mathbf{s}}(\mathbf{x}). \quad (48)$$

We compute the covariance function of  $f$  as:

$$c_f(\mathbf{x}) = \mathbb{E}[f(\mathbf{x}) \cdot f(\mathbf{x}')] \quad (49)$$

$$= \mathbb{E} \left[ \sum_{\mathbf{s} \in \{-B, \dots, B\}^{d_x}} \hat{f}_{\mathbf{s}} \psi_{-\mathbf{s}}(\mathbf{x}) \sum_{\mathbf{s}' \in \{-B, \dots, B\}^{d_x}} \hat{f}_{\mathbf{s}'} \psi_{-\mathbf{s}'}(\mathbf{x}') \right] \quad (50)$$

$$= \sum_{\mathbf{s}, \mathbf{s}' \in \{-B, \dots, B\}^{d_x}} \mathbb{E}[\hat{f}_{\mathbf{s}} \psi_{-\mathbf{s}}(\mathbf{x}) \hat{f}_{\mathbf{s}'} \psi_{-\mathbf{s}'}(\mathbf{x}')] \quad (51)$$

$$= \sum_{\mathbf{s}, \mathbf{s}' \in \{-B, \dots, B\}^{d_x}} \mathbb{E}[\hat{f}_{\mathbf{s}} \cdot \hat{f}_{\mathbf{s}'}] \psi_{-\mathbf{s}}(\mathbf{x}) \psi_{-\mathbf{s}'}(\mathbf{x}') \quad (52)$$

$$= \sigma^2 \sum_{\mathbf{s} \in \{-B, \dots, B\}^{d_x}} \psi_{-\mathbf{s}}(\mathbf{x}) \psi_{\mathbf{s}}(\mathbf{x}') \quad (53)$$

$$= \sigma^2 \sum_{\mathbf{s} \in \{-B, \dots, B\}^{d_x}} \psi_{-\mathbf{s}}(\mathbf{x} - \mathbf{x}'). \quad (54)$$

#### A.1.1 Covariance after convolution $c_{A_k \tau}$

Let us place a centered Gaussian distribution on the Fourier series of the band-limited convolution kernel  $\mathbf{k}: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^d$ , so that:

$$\mathbf{C}_{\mathbf{k}}((\mathbf{z} - \mathbf{x}), (\mathbf{z}' - \mathbf{x}')) = \sigma^2 \mathbf{I}_d \sum_{\mathbf{s} \in \{-B, \dots, B\}^{d_x}} \psi_{-\mathbf{s}}((\mathbf{z} - \mathbf{x}) - (\mathbf{z}' - \mathbf{x}')) \quad (55)$$

$$= \sigma^2 \mathbf{I}_d \sum_{\mathbf{s} \in \{-B, \dots, B\}^{d_x}} \psi_{-\mathbf{s}}(\mathbf{z} - \mathbf{x}) \psi_{\mathbf{s}}(\mathbf{z}' - \mathbf{x}') \quad (56)$$

$$= \sigma^2 \mathbf{I}_d \sum_{\mathbf{s} \in \{-B, \dots, B\}^{d_x}} \psi_{-\mathbf{s}}(\mathbf{z} - \mathbf{z}') \psi_{\mathbf{s}}(\mathbf{x} - \mathbf{x}') \quad (57)$$

So, let us consider the quantity  $c_{A_k}[\mathbf{f}_1, \mathbf{f}_2]$  for arbitrary functions  $\mathbf{f}_1$  and  $\mathbf{f}_2$ :

$$A_{c_k}[\mathbf{f}_2^\top \mathbf{f}_1](\mathbf{z}, \mathbf{z}') = \iint_{\mathcal{X}} \mathbf{f}_2^\top(\mathbf{x}') \mathbf{C}_{\mathbf{k}}((\mathbf{z} - \mathbf{x}), (\mathbf{z}' - \mathbf{x}')) \mathbf{f}_1(\mathbf{x}) d\mathbf{x} d\mathbf{x}' \quad (58)$$

$$= \iint_{\mathcal{X}} \mathbf{f}_2^\top(\mathbf{x}') \left( \sigma^2 \mathbf{I}_d \sum_{\mathbf{s} \in \{-B, \dots, B\}^{d_x}} \psi_{-\mathbf{s}}(\mathbf{z} - \mathbf{z}') \psi_{\mathbf{s}}(\mathbf{x} - \mathbf{x}') \right) \mathbf{f}_1(\mathbf{x}) d\mathbf{x} d\mathbf{x}' \quad (59)$$

$$= \sigma^2 \sum_{\mathbf{s} \in \{-B, \dots, B\}^{d_x}} \psi_{-\mathbf{s}}(\mathbf{z} - \mathbf{z}') \iint_{\mathcal{X}} \psi_{\mathbf{s}}(\mathbf{x} - \mathbf{x}') \mathbf{f}_2^T(\mathbf{x}') \mathbf{f}_1(\mathbf{x}) d\mathbf{x} d\mathbf{x}' \quad (60)$$

$$= \sigma^2 \sum_{\mathbf{s} \in \{-B, \dots, B\}^{d_x}} \psi_{-\mathbf{s}}(\mathbf{z} - \mathbf{z}') (2\pi)^{2d_x} \text{FS}_{[\mathbf{s}, -\mathbf{s}]}[\mathbf{f}_2^T \mathbf{f}_1] \quad (61)$$

$$= \sigma^2 (2\pi)^{2d_x} \sum_{\mathbf{s} \in \{-B, \dots, B\}^{d_x}} \psi_{-\mathbf{s}}(\mathbf{z} - \mathbf{z}') \text{FS}_{-\mathbf{s}}[\mathbf{f}_2]^T \text{FS}_{\mathbf{s}}[\mathbf{f}_1]. \quad (62)$$

## A.2 Toroidal Matérn operator

**Definition A.1** (Matérn family of kernels on a closed manifold). The Matérn family of kernels  $c$  with lengthscale  $\ell$  in a  $d$ -dimensional closed manifold  $\mathcal{M}$  are described as:

$$c(a, b; \nu, \ell) = \sum_{k=1}^{\infty} \hat{c}(\lambda_k; \nu, \ell) \cdot \phi_k(a) \cdot \phi_k(b), \quad (63)$$

$$\hat{c}(\lambda_k; \nu, \ell) = \begin{cases} \exp\left[-\frac{\ell^2}{2} \lambda_k\right] & \text{if } \nu = \infty, \\ \left(\frac{2\nu}{\ell^2} + \lambda_k\right)^{-\nu - \frac{d}{2}} & \text{otherwise.} \end{cases}, \quad (64)$$

where,  $\lambda_k$  and  $\phi_k$  are the  $k$ -th eigenvalues and eigenfunctions, respectively, of the Laplace-Beltrami operator of the manifold  $\mathcal{M}$ .

For a 1-dimensional flat torus, an orthonormal eigensystem for the Laplace-Beltrami operator is:

$$\lambda_k = [k/2]^2; \quad \phi_k(x) = \begin{cases} 1/\sqrt{2\pi} & \text{if } k = 1, \\ \cos(\sqrt{\lambda_k}x)/\sqrt{\pi} & \text{if } k = 2n, \\ \sin(\sqrt{\lambda_k}x)/\sqrt{\pi} & \text{if } k = 2n + 1. \end{cases} \quad (65)$$

Additionally, for a  $d_x$ -dimensional flat torus, an orthonormal eigensystem for the Laplace-Beltrami operator is given by the sum and product of the 1-dimensional eigensystem such that, given an index  $\mathbf{k} = [k_1, \dots, k_{d_x}]$ , we have that  $\lambda_{\mathbf{k}} = \sum_{j=1}^{d_x} \lambda_{k_j}$  and  $\phi_{\mathbf{k}}(\mathbf{x}) = \prod_{j=1}^{d_x} \phi_{k_j}(x_j)$ .

**Expression of 1-dimensional toroidal kernel using complex exponentials.** For convenience, we will rewrite the series expansion of this kernel to use exponentials of complex numbers.

Start by noting that:

$$\phi_{2n}(a)\phi_{2n}(b) = \frac{1}{\pi} \cos(na) \cos(nb), \quad \phi_{2n+1}(a)\phi_{2n+1}(b) = \frac{1}{\pi} \sin(na) \sin(nb), \quad (66)$$

Therefore,

$$\phi_{2n}(a)\phi_{2n}(b) + \phi_{2n+1}(a)\phi_{2n+1}(b) = \frac{1}{\pi} \cos(na) \cos(nb) + \frac{1}{\pi} \sin(na) \sin(nb) \quad (67)$$

$$= \frac{1}{\pi} \cos(n(a - b)) \quad (68)$$

$$= \frac{1}{2\pi} (\exp[in(a - b)] + \exp[-in(a - b)]) \quad (69)$$

$$= \frac{1}{2\pi} (\exp[ina] \exp[-inb] + \exp[-ina] \exp[inb]) \quad (70)$$

$$= \frac{1}{2\pi} (\psi_{-n}(a)\psi_n(b) + \psi_n(a)\psi_{-n}(b)) \quad (71)$$

Now, we can rewrite the Matérn kernel expression as:

$$c(a, b; \nu, \ell) = \sum_{k=1}^{\infty} \hat{c}(\lambda_k; \nu, \ell) \cdot \phi_k(a) \cdot \phi_k(b) \quad (72)$$

$$= \frac{1}{2\pi} \hat{c}(\lambda_1; \nu, \ell) + \sum_{n=1}^{\infty} \hat{c}(\lambda_{2n}; \nu, \ell) (\phi_{2n}(a)\phi_{2n}(b) + \phi_{2n+1}(a)\phi_{2n+1}(b)) \quad (73)$$

$$= \frac{1}{2\pi} \hat{c}(0; \nu, \ell) + \sum_{n=1}^{\infty} \frac{1}{2\pi} \hat{c}(n^2; \nu, \ell) (\psi_{-n}(a) \psi_n(b) + \psi_n(a) \psi_{-n}(b)) \quad (74)$$

$$= \frac{1}{2\pi} \hat{c}(0; \nu, \ell) + \sum_{n=1}^{\infty} \frac{1}{2\pi} \hat{c}(n^2; \nu, \ell) \psi_{-n}(a) \psi_n(b) + \sum_{n=1}^{\infty} \frac{1}{2\pi} \hat{c}(n^2; \nu, \ell) \psi_n(a) \psi_{-n}(b) \quad (75)$$

$$= \frac{1}{2\pi} \hat{c}(0; \nu, \ell) + \sum_{n=-1}^{-\infty} \frac{1}{2\pi} \hat{c}(n^2; \nu, \ell) \psi_n(a) \psi_{-n}(b) + \sum_{n=1}^{\infty} \frac{1}{2\pi} \hat{c}(n^2; \nu, \ell) \psi_n(a) \psi_{-n}(b) \quad (76)$$

$$= \frac{1}{2\pi} \hat{c}(0) \psi_0(a) \psi_0(b; \nu, \ell) + \sum_{n \in \mathbb{Z} \setminus \{0\}} \frac{1}{2\pi} \hat{c}(n^2; \nu, \ell) \psi_n(a) \psi_{-n}(b) \quad (77)$$

$$= \frac{1}{2\pi} \sum_{n \in \mathbb{Z}} \hat{c}(n^2; \nu, \ell) \psi_n(a) \psi_{-n}(b). \quad (78)$$

**Product kernel for  $\mathbb{T}^{d_x}$**  In order to have one lengthscale per dimension, we will make a tensor product kernel where the kernel of a  $d_x$ -dimensional torus  $\mathbb{T}^{d_x}$  is the product of the 1-d toroidal kernel for each dimension:

$$c(\mathbf{a}, \mathbf{b}; \nu, \ell_j) = \prod_{j=1}^{d_x} c(a_j, b_j; \nu, \ell_j) \quad (79)$$

$$= \prod_{j=1}^{d_x} \frac{1}{2\pi} \sum_{n \in \mathbb{Z}} \hat{c}(n^2; \nu, \ell_j) \psi_n(a_j) \psi_{-n}(b_j) \quad (80)$$

$$= (2\pi)^{-d_x} \prod_{j=1}^{d_x} \sum_{n \in \mathbb{Z}} \hat{c}(n^2; \nu, \ell_j) \psi_n(a_j) \psi_{-n}(b_j) \quad (81)$$

$$= (2\pi)^{-d_x} \sum_{\mathbf{n} \in \mathbb{Z}^{d_x}} \prod_{j=1}^{d_x} \hat{c}(n_j^2; \nu, \ell_j) \psi_{n_j}(a_j) \psi_{-n_j}(b_j) \quad (82)$$

$$= (2\pi)^{-d_x} \sum_{\mathbf{n} \in \mathbb{Z}^{d_x}} \psi_{\mathbf{n}}(\mathbf{a}) \psi_{-\mathbf{n}}(\mathbf{b}) \prod_{j=1}^{d_x} \hat{c}(n_j^2; \nu, \ell_j). \quad (83)$$

### A.2.1 $\nu = \infty$ lets Matérn kernel be a product kernel

Notice that when  $\nu = \infty$  and  $\ell_j = \ell$ , we have

$$c(\mathbf{a}, \mathbf{b}; \nu, \ell) = (2\pi)^{-d_x} \sum_{\mathbf{n} \in \mathbb{Z}^{d_x}} \psi_{\mathbf{n}}(\mathbf{a}) \psi_{-\mathbf{n}}(\mathbf{b}) \prod_{j=1}^{d_x} \exp[-\ell^2 n_j^2 / 2]. \quad (84)$$

$$= (2\pi)^{-d_x} \sum_{\mathbf{n} \in \mathbb{Z}^{d_x}} \psi_{\mathbf{n}}(\mathbf{a}) \psi_{-\mathbf{n}}(\mathbf{b}) \exp \left[ - \sum_{j=1}^{d_x} \ell^2 n_j^2 / 2 \right]. \quad (85)$$

$$= (2\pi)^{-d_x} \sum_{\mathbf{n} \in \mathbb{Z}^{d_x}} \psi_{\mathbf{n}}(\mathbf{a}) \psi_{-\mathbf{n}}(\mathbf{b}) \exp \left[ - \frac{\ell^2}{2} \sum_{j=1}^{d_x} n_j^2 \right]. \quad (86)$$

$$= (2\pi)^{-d_x} \sum_{\mathbf{n} \in \mathbb{Z}^{d_x}} \psi_{\mathbf{n}}(\mathbf{a}) \psi_{-\mathbf{n}}(\mathbf{b}) \hat{c} \left( \sum_{j=1}^{d_x} n_j^2; \nu, \ell \right). \quad (87)$$

With the proper rearrangement, we can see that this fits the definition of a Matérn kernel in the  $\mathbb{T}^{d_x}$ , as the eigenvalues of its Beltrami-Laplace operator can be expressed as the sum of the eigenvalues for the 1-dimensional flat torus  $\mathbb{T}$ .



### A.2.2 Covariance after convolution $c_{A_k}$

Let us place a centered factored Matérn prior in the convolution kernel  $\mathbf{k}: \mathcal{Z} \times \mathcal{X} \rightarrow \mathbb{R}^d$ , so that:

$$\mathbf{C}_k(\mathbf{z}, \mathbf{x}, \mathbf{z}', \mathbf{x}') = c(\mathbf{z}, \mathbf{z}'; \nu_z, \ell_z) c(\mathbf{x}, \mathbf{x}'; \nu_x, \ell_x) \mathbf{I}_d. \quad (88)$$

When clear from context, we will suppress the dependency on the hyper-parameters of the Matérn kernel.

So, let us consider the quantity  $c_{A_k}[\mathbf{f}_1, \mathbf{f}_2]$  for arbitrary functions  $\mathbf{f}_1$  and  $\mathbf{f}_2$ :

$$\begin{aligned} A_{c_k \tau}[\mathbf{f}_2^\top \mathbf{f}_1](\mathbf{z}, \mathbf{z}') &= \iint_{\mathcal{X}} \mathbf{f}_2^\top(\mathbf{x}') \mathbf{C}_k(\mathbf{z}, \mathbf{x}, \mathbf{z}', \mathbf{x}') \mathbf{f}_1(\mathbf{x}) d\mathbf{x} d\mathbf{x}' \end{aligned} \quad (89)$$

$$= \iint_{\mathcal{X}} \mathbf{f}_2^\top(\mathbf{x}') \left( c(\mathbf{z}, \mathbf{z}'; \nu_z, \ell_z) c(\mathbf{x}, \mathbf{x}'; \nu_x, \ell_x) \mathbf{I}_d \right) \mathbf{f}_1(\mathbf{x}) d\mathbf{x} d\mathbf{x}' \quad (90)$$

$$= c(\mathbf{z}, \mathbf{z}') \iint_{\mathcal{X}} \mathbf{f}_2^\top(\mathbf{x}') \mathbf{f}_1(\mathbf{x}) c(\mathbf{x}, \mathbf{x}') d\mathbf{x} d\mathbf{x}' \quad (91)$$

$$= c(\mathbf{z}, \mathbf{z}') \sum_{\mathbf{n} \in \mathbb{Z}^{d_x}} (2\pi)^{-d_x} \iint_{\mathcal{X}} \mathbf{f}_2^\top(\mathbf{x}') \mathbf{f}_1(\mathbf{x}) \psi_{\mathbf{n}}(\mathbf{a}) \psi_{-\mathbf{n}}(\mathbf{b}) d\mathbf{x} d\mathbf{x}' \prod_{j=1}^{d_x} \hat{c}(n_j^2; \nu, \ell_j) \quad (92)$$

$$= c(\mathbf{z}, \mathbf{z}'; \nu_z, \ell_z) (2\pi)^{d_x} \sum_{\mathbf{n} \in \mathbb{Z}^{d_x}} \text{FS}_{[\mathbf{n}, -\mathbf{n}]}[\mathbf{f}_2^\top \mathbf{f}_1] \prod_{j=1}^{d_x} \hat{c}(n_j^2; \nu, \ell_j) \quad (93)$$

$$= (2\pi)^{d_x} c(\mathbf{z}, \mathbf{z}'; \nu_z, \ell_z) \sum_{\mathbf{n} \in \mathbb{Z}^{d_x}} \text{FS}_{-\mathbf{n}}[\mathbf{f}_2]^\top \text{FS}_{\mathbf{n}}[\mathbf{f}_1] \prod_{j=1}^{d_x} \hat{c}(n_j^2; \nu, \ell_j) \quad (94)$$

## B Proofs

In this section, we include the proofs for Lemma 3.2, Theorem 3.1, and a short lemma on the well-defined-ness of the activation operator.

### B.1 Well-defined-ness of the point-wise element-wise activation operator

**Lemma B.1** *Let  $(\mathcal{X}, \Sigma, \mu_{\mathcal{X}})$  be a finite measure space, i.e.  $\mu_{\mathcal{X}}(\mathcal{X}) \leq \infty$ , and  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  a Borel measurable function such that*

$$\sup_{x \in \mathbb{R}} \frac{|\sigma(x)|}{1 + |x|} < C \quad (95)$$

*for some constant  $C \in \mathbb{R}$ . Then, the operator  $\sigma[\mathbf{f}]: L^2(\mathcal{X}) \rightarrow L^2(\mathcal{X}) = \sigma \circ \mathbf{f}$  is well defined.*

*Proof.* Remember that a function  $\mathbf{f}$  is in  $L^2(\mathcal{X})$  if, and only if,

$$\int_{\mathcal{X}} |\mathbf{f}(\mathbf{x})|^2 d\mu_{\mathcal{X}}(\mathbf{x}) < \infty. \quad (96)$$

Now, from the linear boundedness condition, we know that for any  $\mathbf{f} \in L^2(\mathcal{X})$  and any  $\mathbf{x} \in \mathcal{X}$ :

$$|\sigma(\mathbf{f}(\mathbf{x}))| < C(1 + |\mathbf{f}(\mathbf{x})|), \quad (97)$$

by squaring both sides and taking integrals,

$$\int_{\mathcal{X}} |\sigma(\mathbf{f}(\mathbf{x}))|^2 d\mu_{\mathcal{X}}(\mathbf{x}) < C^2 \int_{\mathcal{X}} (1 + |\mathbf{f}(\mathbf{x})|)^2 d\mu_{\mathcal{X}}(\mathbf{x}) \quad (98)$$

Now, note that the constant function 1 is in  $L^2(\mathcal{X})$  since  $\int_{\mathcal{X}} 1 d\mu_{\mathcal{X}}(\mathbf{x}) = \mu_{\mathcal{X}}(\mathcal{X}) < \infty$  and that  $|\mathbf{f}(\cdot)|$  is in  $L^2 \mathcal{X}$ . Thus, from linearity,  $1 + |\mathbf{f}(\cdot)|$  is also in  $L^2 \mathcal{X}$  and  $\int_{\mathcal{X}} 1 + |\mathbf{f}(\mathbf{x})| d\mu_{\mathcal{X}}(\mathbf{x}) < \infty$ . Therefore,

$$\int_{\mathcal{X}} |\sigma(\mathbf{f}(\mathbf{x}))|^2 d\mu_{\mathcal{X}}(\mathbf{x}) < \infty. \quad (99)$$

□

## B.2 Compositionality of covariance functions

**Lemma 3.2.** Let  $B_1: L^2(\mathcal{X}; \mathbb{R}^d) \rightarrow L^2(\mathcal{X}; \mathbb{R}^J)$  be a random operator and  $B_2: L^2(\mathcal{X}; \mathbb{R}^J) \rightarrow L^2(\mathcal{X})$  be a centered function-valued Gaussian process. If the following assumptions hold:

- For all  $\mathbf{f} \in L^2(\mathcal{X}; \mathbb{R}^d)$  and  $\mathbf{x} \in \mathcal{X}$ , each component of  $B_1[\mathbf{f}](\mathbf{x}) \in \mathbb{R}^J$  is independent and identically distributed such that the covariance function  $\mathbf{C}_{B_1}[\mathbf{f}, \mathbf{g}] = c_{B_1}[\mathbf{f}, \mathbf{g}] \mathbf{I}_J$ ;
- The covariance function of  $B_2$  can be expressed, for all  $\mathbf{f}, \mathbf{g} \in L^2(\mathcal{X}; \mathbb{R}^J)$  as  $c_{B_2}[\mathbf{f}, \mathbf{g}] = c_{B_2}[\frac{1}{J} \mathbf{g}^\top \mathbf{f}]$  and the function  $h \mapsto c_{B_2}[h]$  is a continuous map from  $L^2(\mathcal{X} \times \mathcal{X})$  to itself.

Then,  $B_2 \circ B_1$  converges in distribution to a function-valued Gaussian process as  $J \rightarrow \infty$ , and

$$c_{B_2 \circ B_1}[\mathbf{f}_1, \mathbf{f}_2] = c_{B_2}[c_{B_1}[\mathbf{f}_1, \mathbf{f}_2]]. \quad (100)$$

*Proof.* Consider a set of size  $N \in \mathbb{N}^+$ ,  $\{(\mathbf{f}_n, h_n)\}_{n=1}^N \subset L^2(\mathcal{X}; \mathbb{R}^d) \times L^2(\mathcal{X})$ , then define the  $N$ -dimensional vector:

$$\mathbf{z} := [\langle h_1, (B_2 \circ B_1)[\mathbf{f}_1] \rangle, \dots, \langle h_N, (B_2 \circ B_1)[\mathbf{f}_N] \rangle]^\top \in \mathbb{R}^N. \quad (101)$$

Additionally, define the function:

$$\bar{c}_{B_1}[\mathbf{f}_j, \mathbf{f}_k]: L^2(\mathcal{X}; \mathbb{R}^d) \times L^2(\mathcal{X}; \mathbb{R}^d) \rightarrow L^2(\mathcal{X} \times \mathcal{X}) = \frac{1}{J} B_1[\mathbf{f}_k]^\top B_1[\mathbf{f}_j]. \quad (102)$$

Then, the conditional random variable  $\mathbf{z} \mid \{\bar{c}_{B_1}[\mathbf{f}_i, \mathbf{f}_j]\}_{i,j=1}^N$  is Gaussian distributed with zero mean and covariance:

$$\text{cov}(z_i, z_j \mid \bar{c}_{B_1}[\mathbf{f}_i, \mathbf{f}_j]) = \langle h_j, \langle h_i, c_{B_2}[\bar{c}_{B_1}[\mathbf{f}_i, \mathbf{f}_j]] \rangle \rangle. \quad (103)$$

We want to show that every  $\mathbf{z}$  converges in distribution to a Gaussian distribution when  $J \rightarrow \infty$ , thus, it is useful to remember the following facts:

- **Multivariate Levy's continuity theorem.** A sequence of random variables  $\{\mathbf{x}_j\}_{j=1}^\infty$  converges to another one  $\mathbf{x}_\infty$  if and only if the sequence of characteristic functions  $\phi_{\mathbf{x}_j}(\mathbf{t}) = \mathbb{E}[\exp(i \cdot \mathbf{t}^\top \mathbf{x}_j)]$ , where  $i = \sqrt{-1}$ , converges point-wise to  $\phi_{\mathbf{x}_\infty}$ .
- **Characteristic function of a  $N$ -dimensional Gaussian distribution.** If  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma)$ , then  $\phi_{\mathbf{x}}(\mathbf{t}) = \exp(\mathbf{t}^\top \Sigma \mathbf{t})$  and  $\phi_{\mathbf{x}}(\mathbf{t}) \leq 1$ , for all  $\mathbf{t} \in \mathbb{R}^N$ .
- **Strong law of large numbers.** As  $J \rightarrow \infty$ , the random element  $\mathbf{K}[\mathbf{f}_i, \mathbf{f}_j]$  converges strongly to the constant  $c_{B_1}[\mathbf{f}_i, \mathbf{f}_j]$ , for all  $\mathbf{f}_i$ .
- **Portmanteau theorem.** Given a sequence of random elements in  $\mathcal{H}$  converging in distribution  $\{x_i\}_{i=1}^\infty \rightarrow x_\infty$ , then,  $\lim_{i \rightarrow \infty} \mathbb{E}[f(x_i)] = \mathbb{E}[f(x_\infty)]$ , for all bounded and continuous functions  $f: \mathcal{H} \rightarrow \mathbb{R}$ .

Thus, we begin with the characteristic function of the variable  $\mathbf{z}$ :

$$\phi_{\mathbf{z}}(\mathbf{t}) = \mathbb{E}[\exp(i \cdot \mathbf{t}^\top \mathbf{z})], \quad (104)$$

by the tower rule, we can write

$$= \mathbb{E}[\mathbb{E}[\exp(i \cdot \mathbf{t}^\top \mathbf{z}) \mid \{\bar{c}_{B_1}[\mathbf{f}_j, \mathbf{f}_k]\}_{j,k=1}^N]] = \mathbb{E}[\exp(\mathbf{t}^\top \Sigma \mathbf{t})], \quad (105)$$

where,  $[\Sigma]_{jk} = \text{cov}(z_j, z_k \mid \bar{c}_{B_1}[\mathbf{f}_j, \mathbf{f}_k])$  is a random variable.

Now, because of the continuity of inner products and the assumption that  $c_{B_2}$  is continuous, we know that the mapping  $h \mapsto \text{cov}(z_j, z_k \mid \bar{c}_{B_1}[\mathbf{f}_j, \mathbf{f}_k] = h)$  is continuous. With this we take the limit:

$$\lim_{J \rightarrow \infty} \phi_{\mathbf{z}}(\mathbf{t}) = \lim_{J \rightarrow \infty} \mathbb{E}[\exp(\mathbf{t}^\top \Sigma \mathbf{t})]; \quad (106)$$

using the portmanteau theorem, we get that:

$$\lim_{J \rightarrow \infty} \phi_{\mathbf{z}}(\mathbf{t}) = \mathbb{E} \left[ \exp \left( \mathbf{t}^\top \begin{bmatrix} \langle h_1, \langle h_1, c_{B_2}[\lim_{J \rightarrow \infty} \bar{c}_{B_1}[\mathbf{f}_1, \mathbf{f}_1]] \rangle \rangle & \cdots \\ \vdots & \ddots \end{bmatrix} \mathbf{t} \right) \right], \quad (107)$$

and, finally, by the strong law of large numbers, we write the expectation as:

$$\lim_{J \rightarrow \infty} \phi_{\mathbf{z}}(\mathbf{t}) = \exp \left( \mathbf{t}^\top \begin{bmatrix} \langle \mathbf{h}_1, \langle \mathbf{h}_1, c_{B_2}[c_{B_1}[\mathbf{f}_1, \mathbf{f}_1]] \rangle \rangle & \cdots \\ \vdots & \ddots \end{bmatrix} \mathbf{t} \right). \quad (108)$$

Therefore, we have shown that  $\mathbf{z}$  converges to a centered Gaussian distribution with:

$$\text{cov}(z_j, z_k) = \text{cov}(z_j, z_k \mid \bar{c}_{B_1}[\mathbf{f}_j, \mathbf{f}_k] = c_{B_1}[\mathbf{f}_j, \mathbf{f}_k]) = \langle \mathbf{h}_k, \langle \mathbf{h}_j, c_{B_2}[c_{B_1}[\mathbf{f}_j, \mathbf{f}_k]] \rangle \rangle \quad (109)$$

Since the set  $\{\mathbf{f}_n, \mathbf{h}_n\}_{n=1}^N$  is arbitrary, we have shown that  $B_2 \circ B_1$  is a centered function-valued Gaussian process with covariance function  $c_{B_2}[c_{B_1}[\cdot, \cdot]]$ .  $\square$

### B.3 Infinite-width neural operators are Gaussian processes

**Definition B.1** (Iterated convergence in distribution). Let  $X_{\mathbf{i}}$  be a random variable for each  $\mathbf{i} = [i_1, \dots, i_k] \subset \mathbb{N}^+$ . The iterated limit

$$\lim_{i_k \rightarrow \infty} \left( \lim_{i_{k-1} \rightarrow \infty} \left( \cdots \lim_{i_1 \rightarrow \infty} (X_{\mathbf{i}}) \cdots \right) \right) \quad (110)$$

whenever exists, is defined as the iterated limit in distribution. That is, suppose there are random variables  $X_{[\infty, i_2, \dots, i_k]}, X_{[\infty, \infty, \dots, i_k]}, \dots, X_{[\infty, \infty, \dots, \infty]}$  such that for every  $i_2, \dots, i_k$

$$X_{[i_1, i_2, i_3, \dots, i_{k-1}, i_k]} \xrightarrow{d} X_{[\infty, i_2, i_3, \dots, i_{k-1}, i_k]} \quad \text{as } i_1 \rightarrow \infty. \quad (111)$$

$$X_{[\infty, i_2, i_3, \dots, i_{k-1}, i_k]} \xrightarrow{d} X_{[\infty, \infty, i_3, \dots, i_{k-1}, i_k]} \quad (112)$$

$\vdots$

$$X_{[\infty, \infty, \infty, \dots, \infty, i_k]} \xrightarrow{d} X_{[\infty, \infty, \infty, \dots, \infty, \infty]} \quad (113)$$

Then we define the iterated limit of  $X_{\mathbf{i}}$  as

$$\lim_{i_k \rightarrow \infty} \left( \lim_{i_{k-1} \rightarrow \infty} \left( \cdots \lim_{i_1 \rightarrow \infty} (X_{\mathbf{i}}) \cdots \right) \right) = X_{[\infty, \infty, \dots, \infty]} \quad (114)$$

**Theorem 3.1.** Let  $\mathcal{X} \subseteq \mathbb{R}^{d_x}$  be a measurable space and let  $\mathcal{H}(\mathcal{X}; \mathbb{R}^J) \subset L^2(\mathcal{X}; \mathbb{R}^J)$  be an RKHS for any  $J \in \mathbb{N}^+$ . Then, for a given depth  $D \in \mathbb{N}^+$ , consider a vector positive integers  $\mathbf{J} = [J_0, J_1, \dots, J_{D-1}, 1]^\top \in \mathbb{N}^{D+1}$  and a  $\mathbf{J}$ -indexed neural operators  $Z_{\mathbf{J}}^{(D)}$  of depth  $D$ :

$$Z_{\mathbf{J}}^{(D)} := H^{(D)} \circ \sigma \circ Z_{\mathbf{J}}^{(D-1)} \in (\mathcal{X} \rightarrow \mathbb{R}^{J_0}) \rightarrow (\mathcal{X} \rightarrow \mathbb{R}), \quad (115)$$

where,

$$Z_{\mathbf{J}}^{(1)} := H^{(1)} \in L^2(\mathcal{X}; \mathbb{R}^{J_0}) \rightarrow \mathcal{H}(\mathcal{X}; \mathbb{R}^{J_1}), \text{ and} \quad (116)$$

$$H^{(\ell)} := (A_{\mathbf{K}^{(\ell)}} + \mathbf{W}^{(\ell)}) \in L^2(\mathcal{X}; \mathbb{R}^{J_{\ell-1}}) \rightarrow \mathcal{H}(\mathcal{X}; \mathbb{R}^{J_\ell}), \quad (117)$$

with  $\mathbf{W}^{(\ell)} \in \mathbb{R}^{J_\ell \times J_{\ell-1}}$ , and  $\mathbf{K}^{(\ell)} \in \mathcal{H}(\mathcal{X} \times \mathcal{X}; \mathbb{R}^{J_\ell \times J_{\ell-1}})$ .

When all parameters are independently distributed a priori according to

$$\mathbf{W}^{(\ell)} \sim \mathcal{N}(\mathbf{0}, \sigma_\ell^2 / J_{\ell-1} \mathbf{I}), \text{ and } \mathbf{K}^{(\ell)} \sim \text{GP}(\mathbf{0}, c_{\mathbf{K}^{(\ell)}} / J_{\ell-1} \mathbf{I}), \quad \text{for } \ell \in \{1, \dots, d\}, \quad (118)$$

then, the iterated limit  $\lim_{J_{D-1} \rightarrow \infty} \cdots \lim_{J_1 \rightarrow \infty} Z_{\mathbf{J}}^{(D)}$ , in the sense of Definition B.1, is equal to a function-valued GP  $Z_\infty^{(D)} \sim \text{GP}(0, c_\infty)$ , where  $c_\infty[\mathbf{f}, \mathbf{g}]$  is available in closed-form.

*Proof.* First, we note from Section 3.1 that the covariances  $c_{\mathbf{W}^{(\ell)}}[\mathbf{f}, \mathbf{g}]$  and  $c_{A_{\mathbf{K}^{(\ell)}}}[\mathbf{f}, \mathbf{g}]$  are equal to:

$$\mathbf{C}_{\mathbf{W}^{(\ell)}}[\mathbf{f}, \mathbf{g}] = \sigma_\ell^2 \frac{1}{J_{\ell-1}} \mathbf{g}^\top \mathbf{f} \mathbf{I}_{J_\ell}, \text{ and } \mathbf{C}_{A_{\mathbf{K}^{(\ell)}}}[\mathbf{f}, \mathbf{g}] = A_{c_{\mathbf{K}^{(\ell)}}} \left[ \frac{1}{J_{\ell-1}} \mathbf{g}^\top \mathbf{f} \right] \mathbf{I}_{J_\ell}, \quad (119)$$

such that both depend on the empirical covariance  $\frac{1}{J_{\ell-1}} \mathbf{g}^\top \mathbf{f}$ , for all  $\mathbf{f}, \mathbf{g} \in L^2(\mathcal{X}; \mathbb{R}^{J_{\ell-1}})$  and  $\ell \in \mathbb{N}^+$ . Therefore, since  $H^{(\ell)}$  is the sum of these two independent function-valued Gaussian processes, we have that  $H^{(\ell)} \sim \text{GP}(0, c^{(\ell|\ell-1)} \mathbf{I}_{J_\ell})$  such that:

$$c^{(\ell|\ell-1)}[\mathbf{f}, \mathbf{g}] = c^{(\ell|\ell-1)}[\mathbf{g}^\top \mathbf{f} / J_{\ell-1}] = A_{c_{k(\ell)}}[\mathbf{g}^\top \mathbf{f} / J_{\ell-1}] + \sigma_\ell^2 \mathbf{g}^\top \mathbf{f} / J_{\ell-1} \quad (120)$$

With this in mind, we proceed the proof by induction on the depth  $D$ .

**Base case.** For the base case  $D = 1$ , we consider the operator  $Z_{\mathbf{J}}^{(1)}$ . Therefore, there are no limits to consider in this step. Nonetheless, as discussed in the previous paragraph, this quantity is a function-valued GP with covariance:

$$\mathbf{C}^{(1)}[\mathbf{f}, \mathbf{g}] = c^{(1|0)}[\mathbf{f}, \mathbf{g}] \mathbf{I}_{J_1} = A_{c_{k(1)}}[\mathbf{g}^\top \mathbf{f} / J_0] + \sigma_1^2 \mathbf{g}^\top \mathbf{f} / J_0. \quad (121)$$

Therefore, our claim is proven.

**Inductive step.** Our inductive hypothesis says that, for a specific  $\ell \in \mathbb{N}^+$ , we have that the iterated limit  $\lim_{J_{\ell-1} \rightarrow \infty} \cdots \lim_{J_1 \rightarrow \infty} Z_{\mathbf{J}}^{(\ell)}$  converges in distribution to a  $Z_\infty^{(\ell)} \sim \text{GP}(\mathbf{0}, c^{(\ell)} \mathbf{I}_{J_\ell})$ .

As a first step, we would like to prove that

$$\lim_{J_{\ell-1} \rightarrow \infty} \cdots \lim_{J_1 \rightarrow \infty} H^{(\ell+1)} \circ \sigma \circ Z_{\mathbf{J}}^{(\ell)} \quad (122)$$

converges in distribution to

$$H^{(\ell+1)} \circ \sigma \circ Z_\infty^{(\ell)}. \quad (123)$$

Consider an arbitrary set of size  $N \in \mathbb{N}^+$ ,

$$(\mathcal{F}, \mathcal{H}) = \{(\mathbf{f}_1, \mathbf{h}_1), \dots, (\mathbf{f}_N, \mathbf{h}_N)\} \subset L^2(\mathcal{X}; \mathbb{R}^{J_0}) \times L^2(\mathcal{X}; \mathbb{R}^{J_{\ell+1}}), \quad (124)$$

and define the variables  $\mathbf{z}[\mathcal{F}, \mathcal{H}] \in \mathbb{R}^N$  and  $\mathbf{Z}_{\mathbf{J}}^{(\ell)}[\mathcal{F}] \in L^2(\mathcal{X}; \mathbb{R}^{N \times J_\ell})$  such that:

$$[\mathbf{z}[\mathcal{F}, \mathcal{H}]]_n := \left\langle \mathbf{h}_n, H^{(\ell+1)} \left[ \sigma(Z_{\mathbf{J}}^{(\ell)}[\mathbf{f}_n]) \right] \right\rangle, \text{ and}, \quad (125)$$

$$[\mathbf{Z}_{\mathbf{J}}^{(\ell)}[\mathcal{F}]]_n := Z_{\mathbf{J}}^{(\ell)}[\mathbf{f}_n]. \quad (126)$$

By definition,  $\mathbf{z}[\mathcal{F}, \mathcal{H}]$  conditioned on  $\mathbf{Z}_{\mathbf{J}}^{(\ell)}[\mathcal{F}]$  follows a multivariate centered Gaussian distribution  $\mathcal{N}(\mathbf{0}, \Sigma(\mathbf{Z}_{\mathbf{J}}^{(\ell)}[\mathcal{F}]))$  with covariance matrix:

$$[\Sigma(\mathbf{A})]_{jk} = \left\langle \mathbf{h}_j, \left\langle \mathbf{h}_k, c_{H^{(\ell+1)}}[\sigma([\mathbf{A}]_j)^\top \sigma([\mathbf{A}]_k) / J_\ell] \mathbf{I}_{J_{\ell+1}} \right\rangle \right\rangle. \quad (127)$$

Thus, by the tower rule, the characteristic function of the marginal distribution of  $\mathbf{z}[\mathcal{F}, \mathcal{H}]$  is:

$$\phi_{\mathbf{z}[\mathcal{F}, \mathcal{H}]}(\mathbf{t}) = \mathbb{E} \left[ \mathbb{E}[\exp(i \mathbf{t}^\top \mathbf{z}[\mathcal{F}, \mathcal{H}]) \mid \mathbf{Z}_{\mathbf{J}}^{(\ell)}[\mathcal{F}]] \right] = \mathbb{E} \left[ \exp \left( \mathbf{t}^\top \Sigma(\mathbf{Z}_{\mathbf{J}}^{(\ell)}[\mathcal{F}]) \mathbf{t} \right) \right]. \quad (128)$$

Now, consider the point-wise convergence of the characteristic function:

$$\lim_{J_{\ell-1} \rightarrow \infty} \cdots \lim_{J_1 \rightarrow \infty} \phi_{\mathbf{z}[\mathcal{F}, \mathcal{H}]}(\mathbf{t}) =: \phi_\infty(\mathbf{t}) \quad (129)$$

Using the portmanteau theorem and continuity of  $\Sigma(\cdot)$ , we have that:

$$\phi_\infty(\mathbf{t}) = \mathbb{E} \left[ \exp \left( \mathbf{t}^\top \Sigma \left( \lim_{J_{\ell-1} \rightarrow \infty} \cdots \lim_{J_1 \rightarrow \infty} \mathbf{Z}_{\mathbf{J}}^{(\ell)}[\mathcal{F}] \right) \mathbf{t} \right) \right]. \quad (130)$$

Now, our inductive hypothesis says that  $Z_{\mathbf{J}}^{(\ell)}$  converges in distribution to a function-valued Gaussian process  $Z_\infty^{(\ell)}$  with each output being i.i.d. With this fact, we can conclude that  $\mathbf{Z}_{\mathbf{J}}^{(\ell)}[\mathcal{F}]$  also converges in distribution to the corresponding variable:  $[\mathbf{Z}_\infty^{(\ell)}[\mathcal{F}]]_n := Z_\infty^{(\ell)}[\mathbf{f}_n]$ . This means that:

$$\phi_\infty(\mathbf{t}) = \mathbb{E} \left[ \exp \left( \mathbf{t}^\top \Sigma \left( \mathbf{Z}_\infty^{(\ell)}[\mathcal{F}] \right) \mathbf{t} \right) \right], \quad (131)$$

which is the characteristic function of a variable defined as:

$$[\tilde{z}[\mathcal{F}, \mathcal{H}]]_n := \left\langle \mathbf{h}_n, H^{(\ell+1)} \left[ \sigma(Z_\infty^{(\ell)}[\mathbf{f}_n]) \right] \right\rangle. \quad (132)$$

Therefore,  $z[\mathcal{F}, \mathcal{H}]$  iteratively converges in distribution to  $\tilde{z}[\mathcal{F}, \mathcal{H}]$ , as  $J_\ell \rightarrow \infty$  for every  $\ell \leq \ell$ . Since the set  $(\mathcal{F}, \mathcal{H})$  is arbitrary, we can conclude that  $(H^{(\ell+1)} \circ \sigma \circ Z_J^{(\ell)})$  also converges in distribution to  $(H^{(\ell+1)} \circ \sigma \circ Z_\infty^{(\ell)})$  as a random operator.

From the induction step, we know that the entries in  $(\sigma \circ Z_\infty^{(\ell)})$  are i.i.d. since the entries of  $\sigma \circ Z_\infty^{(\ell)}$  are also i.i.d. Therefore, we use Lemma 3.2 to show that  $\lim_{J_\ell \rightarrow \infty} (H^{(\ell+1)} \circ \sigma \circ Z_\infty^{(\ell)})$  converges in distribution to a function-valued Gaussian process with covariance function

$$\mathbf{C}^{(\ell+1)}[\mathbf{f}, \mathbf{g}] = c_{H^{(\ell+1)}}[c_{(\sigma \circ Z_\infty^{(\ell)})}] \mathbf{I}_{J_{\ell+1}} = c^{(\ell+1|\ell)}[c_\sigma[c^{(\ell)}[\mathbf{f}, \mathbf{g}]]] \mathbf{I}_{J_{\ell+1}}. \quad (133)$$

Therefore, we just proved by induction that the iterated limit  $\lim_{J_{D-1} \rightarrow \infty} \cdots \lim_{J_1 \rightarrow \infty} Z_J^{(D)}$  converges in distribution to a  $Z_\infty^{(D)} \sim \text{GP}(\mathbf{0}, c_\infty \mathbf{I}_{J_\ell})$  and this covariance function is equal to:

$$c_\infty[\mathbf{f}, \mathbf{g}] = c^{(d)}[\mathbf{f}, \mathbf{g}] = c^{(d|d-1)}[c^{(d-1|d-2)}[\cdots c^{(2|1)}[c^{(1)}[\mathbf{f}, \mathbf{g}]] \cdots]] \quad (134)$$

□

## C Experimental details

In this section, we describe the setup for our experiments. As previously mentioned, all experiments were run in a desktop machine with a 3.8GHz Intel Core i7-9800X CPU and a 24GB NVIDIA Titan RTX (TU102) GPU. More details for each experiment can be found below.

### C.1 Empirical demonstration of results

For both experiments, the input function  $f: \mathbb{T} \rightarrow \mathbb{R}$  has band-limit  $B = 3$ , with its output values  $f(x)$  sampled from a uniform distribution  $\mathcal{U}(-1, 1)$ . In other words, we can express this band-limited function as:

$$f(x) = \frac{1}{7} \sum_{s=-3}^3 f_s \sum_{s'=-3}^3 \psi_{s'} \left( x - \frac{2\pi}{7} s \right), \quad (135)$$

where each  $f_s \sim \mathcal{U}(-1, 1)$  is independent and identically distributed.

In the first experiment of Fig. 2, we construct the operator layer  $H$  under the usual formulation:

$$H[f](x): L^2(\mathbb{T}) \rightarrow L^2(\mathbb{T}) = A_k[f](x) + wf(x), \quad (136)$$

where  $w \sim \mathcal{N}(0, 1)$  and  $k$  follows the band-limited Gaussian process distribution (Section 4.1 and Appendix A.1) with band-limit  $B = 3$  and variance  $\sigma^2 = 1/7$ . Then, the operator on  $f$  is evaluated at zero  $H[f](0)$  with increasing sample sizes.

For the second experiment of Fig. 1, we construct the single-layer neural operator:

$$Z[f](x): L^2(\mathbb{T}) \rightarrow L^2(\mathbb{T}; \mathbb{R}^J) \rightarrow L^2(\mathbb{T}) = (\mathbf{w}_2^\top \circ \text{ReLU} \circ (A_k + \mathbf{w}_1))[f](x), \quad (137)$$

where  $J$  is the width of the hidden layer, and  $\mathbf{w}_2 \sim \mathcal{N}(0, 1/J)$ ,  $\mathbf{w}_1 \sim \mathcal{N}(0, 1)$ , and  $k$  follows an i.i.d. band-limited Gaussian process distribution (Section 4.1 and Appendix A.1) with band-limit  $B = 3$  and variance  $\sigma^2 = 1/7$ . For varying widths  $J \in \{1, 10, 100, 1000\}$ , we evaluate 10,000 samples of the operator on  $f$  at zero  $Z[f](0)$  and show the density of the empirical distribution using kernel density estimation (KDE) with a Gaussian kernel.

These experiments are implemented in the file `experiments/fno_limit.ipynb`.

## C.2 Regression

We consider FNOs of increasing width,  $J \in \{1, 10, 100\}$  and  $J \in \{1, 3, 10, 100, 500\}$  for the synthetic and 1D Burgers' respectively, as well as  $\infty$ -FNOs, both with increasing kernel band-limits  $B \in \{1, 5, 20\}$ . These single-layer neural operators are constructed as:

$$Z_{J,B}[f](x): L^2(\mathbb{T}) \rightarrow L^2(\mathbb{T}; \mathbb{R}^J) \rightarrow L^2(\mathbb{T}) = (\mathbf{w}_2^\top \circ \text{ReLU} \circ (A_{\mathbf{k}} + \mathbf{w}_1))[f](x), \quad (138)$$

where  $J$  is the width of the hidden layer, and  $\mathbf{w}_2 \sim \mathcal{N}(0, 1/J)$ ,  $\mathbf{w}_1 \sim \mathcal{N}(0, 1)$ , and  $\mathbf{k}$  follow an i.i.d. band-limited Gaussian process distribution (Section 4.1 and Appendix A.1) with variance  $\sigma^2 = 1/(2B + 1)$ .

The hyperparameters of the  $\infty$ -FNO are estimated using L-BFGS, while the parameters of the FNOs are optimized with Adam using a step size of 0.001. We evaluate all models using 5-fold cross-validation and report the average and standard deviation of the empirical  $L^2$  norm of the prediction error. For  $\infty$ -FNOs, we use the posterior mean as the prediction.

This experiment is implemented in the file `experiments/train.py`.

### Synthetic regression

We start by defining the ground truth Fourier neural operator (FNO) which will generate our training and test data:

$$Z_{\text{true}}[f](x): L^2(\mathbb{T}) \rightarrow L^2(\mathbb{T}) \rightarrow L^2(\mathbb{T}) = (w_2 \circ \text{ReLU} \circ (A_{\mathbf{k}} + w_1))[f](x), \quad (139)$$

where the hidden layer's width is 1 and the band-limit of  $\mathbf{k}$  is equal to 5. Next, we sample  $n = 100$  input functions  $f_i: \mathbb{T} \rightarrow \mathbb{R}$  with the same band-limit  $B = 5$  and uniformly-distributed outputs  $\mathcal{U}(-1, 1)$ , so that we have:

$$f_i(x) = \frac{1}{11} \sum_{s=-5}^5 f_{is} \sum_{s'=-5}^5 \psi_{s'} \left( x - \frac{2\pi}{11}s \right), \quad (140)$$

where each  $f_{is} \sim \mathcal{U}(-1, 1)$  is independent and identically distributed. We then compute  $Z_{\text{true}}[f_i]$  on an equally spaced grid given by  $\{-5\frac{2\pi}{11}, \dots, 5\frac{2\pi}{11}\} \subset \mathbb{R}^{11}$ .

### 1D Burgers' equation

This dataset is provided from PDEBench (Takamoto et al., 2022), which includes solutions to the 1D Burgers' equation:

$$\frac{\partial}{\partial t} u(t, x) + \frac{1}{2} \frac{\partial}{\partial x} u^2(t, x) = \frac{\nu}{\pi} \frac{\partial^2}{\partial x^2} u(t, x), \quad (141)$$

where  $x \in (0, 1)$  and  $t \in (0, 2]$  are independent variables and  $\nu$  is the diffusion coefficient.

The regression task is set up with  $\nu = 0.002$  and a collection of initial conditions  $\{u(0, \cdot) = f_i\}_{i=1}^n$  and their respective end states  $\{u(2, \cdot) = g_i\}_{i=1}^n$ . Due to memory constraints when creating the covariance matrices for  $\infty$ -FNO, we subsample the original dataset to  $n = 100$  functions and a grid size of  $m = 103$ . The original data can be downloaded at <https://darus.uni-stuttgart.de/api/access/datafile/268193>.