

---

# Incorporating Prior Knowledge into Neural Networks through an Implicit Composite Kernel

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 It is challenging to guide neural network (NN) learning with prior knowledge.  
2 In contrast, many known properties, such as spatial smoothness or seasonality,  
3 are straightforward to model by choosing an appropriate kernel in a Gaussian  
4 process (GP). Many deep learning applications could be enhanced by modeling  
5 such known properties. For example, convolutional neural networks (CNNs) are  
6 frequently used in remote sensing, which is subject to strong seasonal effects. We  
7 propose to blend the strengths of deep learning and the clear modeling capabilities  
8 of GPs by using a composite kernel that combines a kernel implicitly defined by a  
9 neural network with a second kernel function chosen to model known properties  
10 (e.g., seasonality). Then, we approximate the resultant GP by combining a deep  
11 network and an efficient mapping based on the Nyström approximation, which we  
12 call Implicit Composite Kernel (ICK). ICK is flexible and can be used to include  
13 prior information in neural networks in many applications. We demonstrate the  
14 strength of our framework by showing its superior performance and flexibility  
15 on both synthetic and real-world data sets. The code is available at: [https://anonymous.4open.science/r/ICK\\_NNGP-17C5/](https://anonymous.4open.science/r/ICK_NNGP-17C5/).  
16

## 17 1 Introduction

18 In complex regression tasks, input data often contains *multiple sources of information*. These sources  
19 can be presented in both high-dimensional (e.g. images, audios, texts, etc.) and low-dimensional  
20 (e.g. timestamps, spatial locations, etc.) forms. A common approach to learn from high-dimensional  
21 information is to use neural networks (NNs) [21, 33], as NNs are powerful enough to capture the  
22 relationship between complex high-dimensional data and target variables of interest. In many areas,  
23 NNs are standard practice, such as the dominance of Convolutional Neural Networks (CNNs) for  
24 image analysis [26, 61, 62]. In contrast, for low-dimensional information, we usually have some  
25 prior knowledge on how the information relates to the predictions. As a concrete example, consider  
26 a remote sensing problem where we predict ground measurements from satellite imagery with  
27 associated timestamps. *A priori*, we expect the ground measurements to vary periodically with  
28 respect to time between summer and winter due to seasonal effects. We would typically use a CNN  
29 to capture the complex relationship between the imagery and the ground measurements. In this case,  
30 we want to guide the learning of the CNN with our prior knowledge about the seasonality. This is  
31 challenging because knowledge represented in NNs pertains mainly to correlation between network  
32 units instead of quantifiable statements [36].

33 Conversely, Gaussian processes (GPs) have been used historically to incorporate relevant prior beliefs  
34 by specifying the appropriate form of its kernel function (or covariance function) [2, 54]. One  
35 approach to modeling multiple sources of information is to assign a relevant kernel function to each  
36 source of information respectively and combine them through addition or multiplication, resulting in a

37 *composite kernel function* [14]. This formulation means that specifying a kernel to match prior beliefs  
38 on one source of information is straightforward. Such composite kernel learning techniques are  
39 extensively used in many application areas such as multi-media data [40], neuroimaging [60], spatial  
40 data analysis, and environmental data analysis [28, 44]. In view of the clear modeling capabilities of  
41 GP, it is desirable to examine how a NN could be imbued with the same modeling ease.

42 In recent years, researchers have come up with a variety of methods to incorporate prior knowledge  
43 into NNs. These efforts can be broken into many categories, such as those that add prior information  
44 through loss terms like physics-informed NNs [32, 41]. Here, we focus on the major category of those  
45 methods that build integrated models of NNs and GPs with various structures [50, 57, 58]. Related  
46 to our proposed methodology, Pearce et al. [43] exploited the fact that a Bayesian neural network  
47 (BNN) approximates a GP to construct additive and multiplicative kernels, but they were limited  
48 to specific predefined kernels. Matsubara et al. [38] then resolved this limitation by constructing  
49 priors of BNN parameters based on the ridgelet transform and its dual, but they did not explicitly  
50 show how their approach works for data with multiple sources of information. To our knowledge,  
51 none of these existing approaches allows a modeler to choose any appropriate kernel of known prior  
52 information from multiple sources. We address this limitation by presenting a simple yet novel  
53 Implicit Composite Kernel (ICK) framework, which processes high-dimensional information using a  
54 kernel implicitly defined by a neural network and low-dimensional information using a chosen kernel  
55 function. The low-dimensional kernels are mapped into the neural network framework to create a  
56 straightforward and simple-to-learn implementation. Our key results and contributions are:

- 57 • We analytically show our ICK framework, under reasonable assumptions, is approximately  
58 equivalent to a Gaussian process regression (GPR) model with a composite kernel *a priori*.
- 59 • We demonstrate that our ICK framework yields better performance on both prediction and  
60 forecasting tasks even with very limited data.
- 61 • We compare to joint deep learning models, such as a neural network-random forest joint  
62 model, to show that ICK can flexibly capture the patterns of the low-dimensional information  
63 without deliberately designing a pre-processing procedure or complex NN structures.

64 Based on these contributions, we believe ICK will be useful in the context of learning from complex  
65 *hybrid* data with prior knowledge, especially in the field of remote sensing and spatial statistics.

## 66 2 Related Work

67 **Equivalence between NNs and GPs** The equivalence between GPs and randomly initialized single-  
68 layer NNs with infinite width was first shown by Neal [42]. With the development of modern deep  
69 learning, researchers further extended this relationship to deep networks [34, 39] and convolutional  
70 neural networks (CNNs) [17]. This relationship is crucial for proving the resemblance between GPR  
71 and our ICK framework, which will be discussed in Section 4.1.

72 **NNs with prior knowledge** As mentioned before, one approach to equip NNs with prior knowledge  
73 is to modify the loss function. For example, Lagaris et al. [32] solved differential equations (DEs)  
74 using NNs by setting the loss to be a function whose derivative satisfies the DE conditions. Another  
75 approach is to build integrated models of NNs and GPs. For example, Wilson et al. [58] implemented  
76 a regression network with GP priors over latent variables and made inference by approximating  
77 the posterior using Variational Bayes or sampling from the posterior using Gibbs sampling scheme.  
78 Garnelo et al. [16] introduced a class of neural latent variable models called Neural Processes (NPs)  
79 which are capable of learning efficiently from the data and adapting rapidly to new observations. Zhu  
80 et al. [10] proposed NeuralEF which can accurately approximate kernel functions by using a series  
81 of objective functions parameterized by NNs under the principle of eigen-decomposition.

82 **GP with composite kernels** Composite kernel GPs are widely used in both machine learning  
83 [14, 54] and geostatistical modeling [9, 18]. GPR in geostatistical modeling is also known as *kriging*  
84 [27, 31], which serves as a surrogate model to replace expensive function evaluations. The inputs for  
85 a composite GP are usually low-dimensional (e.g. spatial distance) as GPs do not scale well with the  
86 number of samples for high-dimensional inputs [4, 5]. To overcome this issue, Pearce et al. [43] and  
87 Matsubara et al. [38] developed BNN analogue for composite GPs. Similar to these studies, our ICK  
88 framework can also be viewed as a simulation for composite GPs.

89 **Approximation methods for GP** For large data sets, approximation methods are needed as exact  
90 kernel learning and inference scales  $\mathcal{O}(N^3)$ . Nyström low-rank matrix approximation [12, 53] and  
91 Random Fourier Features [45, 46] are two of the most commonly used methods. A common technique  
92 is to choose inducing points as pseudo-inputs to efficiently approximate the full kernel matrix [49, 23].  
93 Our work is inspired by these approximation techniques and we use them as *transformation functions*  
94 to map the kernel matrix into latent space representations in Section 4.2.

## 95 3 Background

96 Before elaborating on the details of our ICK framework, we introduce our notation, briefly go over  
97 the concepts of composite GPs, and describe the relationship between GPs and NNs.

### 98 3.1 Problem Setup

99 To formalize the problem, we have a training data set which contains  $N$  data points  $\mathbf{X} =$   
100  $[\mathbf{x}_i]_{i=1}^N = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$  and the corresponding labels of these data points are  $\mathbf{y} = [y_i]_{i=1}^N =$   
101  $[y_1, y_2, \dots, y_N]^T$  where  $y_i \in \mathbb{R}$ . Each data point  $\mathbf{x}_i = \{x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(M)}\}$  is composed of informa-  
102 tion from  $M$  different sources where the  $m^{\text{th}}$  source of information of the  $i^{\text{th}}$  data point is denoted as  
103  $x_i^{(m)} \in \mathbb{R}^{D_m}$ . Our goal is to learn a function  $\hat{y}_i = f(\mathbf{x}_i) : \{\mathbb{R}^{D_1}, \mathbb{R}^{D_2}, \dots, \mathbb{R}^{D_M}\} \rightarrow \mathbb{R}$  which takes  
104 in a data point  $\mathbf{x}_i$  and outputs a predicted value  $\hat{y}_i$ .

### 105 3.2 Composite GPs

106 A Gaussian process (GP) describes a distribution over functions [54]. A key property of GP is that  
107 it can be completely defined by a mean function  $\mu(\mathbf{x})$  and a kernel function  $K(\mathbf{x}, \mathbf{x}')$ . The mean  
108 function  $\mu(\mathbf{x})$  is often assumed to be zero for simplicity. In that case, the outcome function is

$$f(\mathbf{x}) \sim \mathcal{GP}(0, K(\mathbf{x}, \mathbf{x}')). \quad (1)$$

109 Any finite subset of these random variables has a multivariate Gaussian distribution with mean  $\mathbf{0}$   
110 and kernel matrix  $\mathbf{K}$  whose entries can be calculated as  $\mathbf{K}_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$ . In many situations, the  
111 full kernel function is built by a composite kernel by combining simple kernels through addition  
112  $K_{\text{comp}}(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}') + K_2(\mathbf{x}, \mathbf{x}')$  or multiplication  $K_{\text{comp}}(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}')K_2(\mathbf{x}, \mathbf{x}')$   
113 [14]. A useful property that ICK exploits is that  $K_1$  and  $K_2$  can take different subparts of  $\mathbf{x}$  as  
114 their inputs. For example,  $K_{\text{comp}}(\mathbf{x}, \mathbf{x}') = K_1(x^{(1)}, x^{(1)'}) + K_2(x^{(2)}, x^{(2)'})$  or  $K_{\text{comp}}(\mathbf{x}, \mathbf{x}') =$   
115  $K_1(x^{(1)}, x^{(1)'}) K_2(x^{(2)}, x^{(2)'})$ . Other methods such as functional mapping are also valid if the  
116 resulting kernel matrix  $\mathbf{K}$  is positive semidefinite (PSD) for all possible choices of data set  $\mathbf{X}$  [47].

### 117 3.3 Correspondence between GPs and NNs

118 Neal [42] proved that a single-hidden layer network with infinite width is *exactly equivalent* to a GP  
119 over data indices  $i = 1, 2, \dots, N$  under the assumption that the weight and bias parameters of the  
120 hidden layer are i.i.d. Gaussian with zero mean. Lee et al. [34] and Garriga et al. [17] then extended  
121 this statement to deep neural networks and convolutional neural networks (CNNs), respectively.  
122 However, if the width (or the number of channels) of a network is finite, then these results state that  
123 the network *approximately* converges to a GP with zero mean as in the following lemma.

124 **Lemma 1** *Let  $\mathbf{z} = f_{\text{NN}}(x^{(1)}) : \mathbb{R}^{D_1} \rightarrow \mathbb{R}^p$  be the latent representation extracted from  $x^{(1)}$  where*  
125  *$p$  is the dimension of the extracted representation and  $f_{\text{NN}}$  is a neural network with finite width*  
126 *and zero-mean i.i.d. parameters. The  $k^{\text{th}}$  entry of this representation will approximately follow a*  
127 *NN-implied GP*

$$z_k = f_{\text{NN}}(x^{(1)})_k \sim \mathcal{GP}_{\text{approx}}(0, K_{\text{NN}}(x^{(1)}, x^{(1)'})) \quad (2)$$

128 That is to say, the  $k^{\text{th}}$  component  $z_k$  of the representation extracted by the network has zero mean  
129  $\mathbb{E}_{p(\theta^{(1)})} [z_{ik}^{(1)}] = 0$  for all  $i = 1, 2, \dots, N$  where  $\theta$  represents the network parameters. The co-  
130 variance between  $z_{ik}^{(1)}$  and  $z_{jk}^{(1)}$  for *different* data indices  $i, j = 1, 2, \dots, N$  can be approximated as

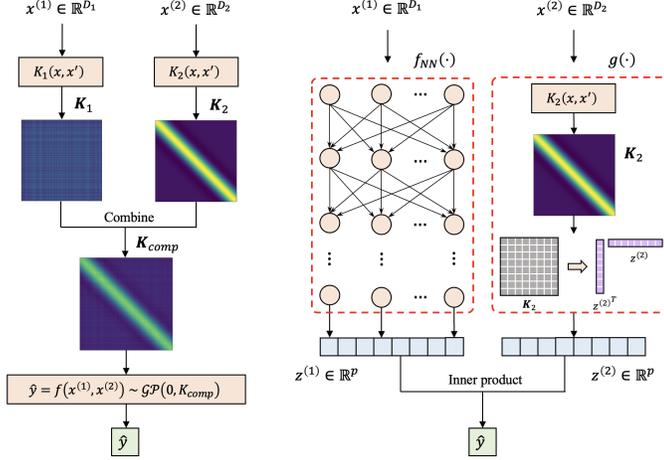


Figure 1: Given data containing 2 sources of information  $x^{(1)}$  and  $x^{(2)}$ , we can process the data using either **(Left)** a composite Gaussian process regression (GPR) model or **(Right)** our ICK framework where  $x^{(1)}$  is processed with a neural network  $f_{\text{NN}}(\cdot)$  and  $x^{(2)}$  is processed with  $g(\cdot)$  where  $g(\cdot)$  consists of a kernel function  $K_2$  and some transformation which maps the kernel matrix  $K_2$  into the latent space.

131  $\text{cov}(z_{ik}^{(1)}, z_{jk}^{(1)}) = \mathbb{E}_{p(\theta^{(1)})} [z_{ik}^{(1)} z_{jk}^{(1)}] \approx K_{\text{NN}}(x_i^{(1)}, x_j^{(1)})$  where  $x_i^{(1)}$  and  $x_j^{(1)}$  are the correspond-  
 132 ing inputs for the network and  $K_{\text{NN}}$  is the kernel function implied by the network.

#### 133 4 Implicit Composite Kernel (ICK) Framework

134 We show the structure of a composite GPR model and our ICK framework in Figure 1. To make  
 135 the illustration clear, we limit ourselves to data with information from 2 different sources  $x =$   
 136  $\{x^{(1)}, x^{(2)}\}$  where  $x^{(1)}$  is high-dimensional and  $x^{(2)}$  is low-dimensional (i.e.  $D_1 \gg D_2$ ) with some  
 137 known relationship with the target  $y$ . We are inspired by composite GPR, which computes 2 different  
 138 kernel matrices  $K_1$  and  $K_2$  and then combines them into a single composite kernel matrix  $K_{\text{comp}}$ .  
 139 However, as discussed before, it is more suitable to use a NN to learn from the high dimensional  
 140 information  $x^{(1)}$ . In our ICK framework, we process  $x^{(1)}$  with a NN  $f_{\text{NN}}(\cdot) : \mathbb{R}^{D_1} \rightarrow \mathbb{R}^p$  and  $x^{(2)}$   
 141 with a mapping  $g(\cdot) : \mathbb{R}^{D_2} \rightarrow \mathbb{R}^p$  which consists of a kernel function  $K_2$  followed by a kernel-  
 142 to-latent-space transformation (described in Section 4.2), resulting in two latent representations  
 143  $z^{(1)}, z^{(2)} \in \mathbb{R}^p$ . Then, we make a prediction  $\hat{y}$  by doing an *inner product* between these two  
 144 representations  $\hat{y} = f_{\text{NN}}(x^{(1)}) \cdot g(x^{(2)})$ . Finally, the parameters of both the NN and the kernel  
 145 function are learned via gradient-based optimization methods [3, 30].

146 In the sections below, we first analytically show that our ICK framework is *approximately* equivalent  
 147 to a composite GPR model *a priori* using a multiplicative kernel between the kernel implicitly defined  
 148 by the NN on  $x^{(1)}$  and the chosen kernel on  $x^{(2)}$ . **This theory is used to motivate the model form.**  
 149 **The model will deviate from the GP solution after learning, but we note that recent work suggests**  
 150 **that the predictions from a trained NN may not vary too much from its GP equivalent [34].** We then  
 151 show how we implement the kernel-to-latent-space transformation in detail. Here, we note that we  
 152 apply ICK for multiplicative kernels but note that an additive kernel may be constructed using the  
 153 methods of Pearce et al. [43].

#### 154 4.1 Resemblance between Composite GPR and ICK

155 We will analytically prove the following theorem for data with information from 2 different sources  
 156  $x = \{x^{(1)}, x^{(2)}\}$  for clarity, and we note this theorem can be straightforwardly extended to  $M > 2$ .

157 **Theorem 1** Let  $f_{\text{NN}} : \mathbb{R}^{D_1} \rightarrow \mathbb{R}^p$  be a NN function with random weights and  $g : \mathbb{R}^{D_2} \rightarrow \mathbb{R}^p$  be a  
 158 mapping function, and define an inner product  $\hat{y}$  between the representations  $z^{(1)} = f_{\text{NN}}(x^{(1)})$  and  
 159  $z^{(2)} = g(x^{(2)})$ . Then this inner product approximately follows a composite GPR model

$$\hat{y} = f_{\text{ICK}}(x^{(1)}, x^{(2)}) = f_{\text{NN}}(x^{(1)}) \cdot g(x^{(2)}) \sim \mathcal{GP}_{\text{approx}}(0, K_{\text{NN}}^{(1)} K^{(2)}), \quad (3)$$

160 if  $g$  includes the following deterministic kernel-to-latent-space transformation

$$K^{(2)}(x_i^{(2)}, x_j^{(2)}) \approx z_i^{(2)T} z_j^{(2)} = g(x_i^{(2)})^T g(x_j^{(2)}), \quad (4)$$

161 where  $K_{NN}^{(1)}$  is a NN-implied kernel and  $K^{(2)}$  is any valid kernel of our choice.

162 To prove Theorem 1, we first make the following assumption.

163 **Assumption 1** For latent representations  $\mathbf{z}_i^{(m)}$  and  $\mathbf{z}_j^{(m)}$  extracted from different data points  $\mathbf{x}_i$  and  
 164  $\mathbf{x}_j$  where  $i \neq j$  and  $m \in \{1, 2\}$ , the interactions between different entries of  $\mathbf{z}_i^{(m)}$  and  $\mathbf{z}_j^{(m)}$  can be  
 165 reasonably ignored. In other words, let  $\theta^{(m)}$  be the parameters of the network or the kernel function  
 166 which takes in  $x^{(m)}$  and outputs  $\mathbf{z}^{(m)}$ , we have  $\mathbb{E}_{p(\theta^{(m)})} [z_{ik}^{(m)} z_{jl}^{(m)}] = 0$  for all  $k \neq l$ .

167 With Assumption 1 and Lemma 1, let  $\Theta = \{\theta^{(1)}, \theta^{(2)}\}$  represent the parameters of the ICK frame-  
 168 work, we can calculate the covariance between  $\hat{y}_i$  and  $\hat{y}_j$  for different data indices  $i \neq j$  as follows

$$\text{cov}(\hat{y}_i, \hat{y}_j) = \mathbb{E}_{p(\Theta)}[\hat{y}_i \hat{y}_j] - \mathbb{E}_{p(\Theta)}[\hat{y}_i] \mathbb{E}_{p(\Theta)}[\hat{y}_j] \quad (5)$$

$$= \mathbb{E}_{p(\Theta)} \left[ \left( \sum_{k=1}^p z_{ik}^{(1)} z_{ik}^{(2)} \right) \left( \sum_{k=1}^p z_{jk}^{(1)} z_{jk}^{(2)} \right) \right] \quad (6)$$

$$= \mathbb{E}_{p(\Theta)} \left[ \sum_{k=1}^p \sum_{l=1}^p z_{ik}^{(1)} z_{jl}^{(1)} z_{ik}^{(2)} z_{jl}^{(2)} \right] \quad (7)$$

$$= \mathbb{E}_{p(\Theta)} \left[ \sum_{k=1}^p z_{ik}^{(1)} z_{jk}^{(1)} z_{ik}^{(2)} z_{jk}^{(2)} \right] \quad (8)$$

$$= \sum_{k=1}^p \mathbb{E}_{p(\theta^{(1)})} [z_{ik}^{(1)} z_{jk}^{(1)}] \mathbb{E}_{p(\theta^{(2)})} [z_{ik}^{(2)} z_{jk}^{(2)}] \quad (9)$$

$$\approx K_{NN}^{(1)}(x_i^{(1)}, x_j^{(1)}) \sum_{k=1}^p \mathbb{E}_{p(\theta^{(2)})} [z_{ik}^{(2)} z_{jk}^{(2)}]. \quad (10)$$

169 Here, from Equation 5 to Equation 6, we use the statement  $\mathbb{E}_{p(\theta^{(1)})} [z_{ik}^{(1)}] = 0$  from Lemma 1 and the  
 170 independence between  $\theta^{(1)}$  and  $\theta^{(2)}$ , which leads to  $\mathbb{E}_{p(\Theta)}[\hat{y}_i] = \mathbb{E}_{p(\Theta)}[\hat{y}_j] = 0$ . From Equation 7 to  
 171 Equation 8, we get rid of all the cross terms under Assumption 1. From Equation 8 to Equation 9, we  
 172 again make use of the independence between  $\theta^{(1)}$  and  $\theta^{(2)}$ . From Equation 9 to Equation 10, we use  
 173 the statement  $\mathbb{E}_{p(\theta^{(1)})} [z_{ik}^{(1)} z_{jk}^{(1)}] \approx K_{NN}^{(1)}(x_i^{(1)}, x_j^{(1)})$  from Lemma 1. If the kernel-to-latent-space  
 174 transformation in  $g(\cdot)$  is *deterministic*, we can remove the expectation sign from the summation term  
 175 in Equation 10 and the covariance can be further expressed as

$$\text{cov}(\hat{y}_i, \hat{y}_j) \approx K_{NN}^{(1)}(x_i^{(1)}, x_j^{(1)}) \left( \mathbf{z}_i^{(2)T} \mathbf{z}_j^{(2)} \right) = K_{NN}^{(1)}(x_i^{(1)}, x_j^{(1)}) K^{(2)}(x_i^{(2)}, x_j^{(2)}), \quad (11)$$

176 which means that  $\hat{y}$  approximately follows a GP with composite kernel  $K_{\text{comp}}(\mathbf{x}_i, \mathbf{x}_j) =$   
 177  $K_{NN}^{(1)}(x_i^{(1)}, x_j^{(1)}) K^{(2)}(x_i^{(2)}, x_j^{(2)})$  *a priori*. This completes our proof of Theorem 1.

## 178 4.2 Kernel-to-latent-space Transformation

179 We now show how we can construct an appropriate mapping  $g(\cdot)$  that approximately satisfies the  
 180 assumed form of (4) and is used in the derivation of ICK from (10) to (11). Here we adopt two  
 181 methods, Nyström approximation and Random Fourier Features (RFF), to map the kernel matrix into  
 182 the latent space. Below, we give the formulations and results for the Nyström method, and give the  
 183 methods and results for RFF in Appendix B. According to Yang et al. [59], the Nyström method  
 184 will yield much better performance than RFF if there exists a large gap in the eigen-spectrum of the  
 185 kernel matrix. In our applications, we also observe a large eigen-gap (see details in Appendix C)  
 186 and Nyström method does generalize much better than RFF. We name our framework with Nyström  
 187 method and random Fourier Features ICKy and ICKr, respectively.

### 188 4.2.1 Nyström Method

189 The main idea of Nyström method [53] is to approximate the kernel matrix  $\mathbf{K} \in \mathbb{R}^{N \times N}$  with a much  
 190 smaller low-rank matrix  $\mathbf{K}_q \in \mathbb{R}^{q \times q}$  where  $q \ll N$  so both the computational and space complexity  
 191 of kernel learning can be significantly reduced

$$\mathbf{K} \approx \hat{\mathbf{K}} = \mathbf{K}_{nq} \mathbf{K}_q^{-1} \mathbf{K}_{nq}^T. \quad (12)$$

192 The entries of  $\mathbf{K}_q$  and  $\mathbf{K}_{nq}$  can be calculated as  $(\mathbf{K}_q)_{ij} = K(\hat{x}_i, \hat{x}_j), i, j \in \{1, 2, \dots, q\}$  and  
 193  $(\mathbf{K}_{nq})_{ij} = K(x_i, \hat{x}_j), i \in \{1, 2, \dots, N\}, j \in \{1, 2, \dots, q\}$ , respectively.  $x$  represents the original  
 194 data points and  $\hat{x}$  represents pre-defined inducing points (or pseudo-inputs [49]). In our study, these  
 195 inducing points are chosen by defining an evenly spaced vector over the range of original data points.  
 196 By performing Cholesky decomposition  $\mathbf{K}_q^{-1} = \mathbf{U}^T \mathbf{U}$ , where  $\mathbf{U} \in \mathbb{R}^{q \times q}$ ,  $\hat{\mathbf{K}}$  is

$$\hat{\mathbf{K}} = \mathbf{K}_{nq} \mathbf{K}_q^{-1} \mathbf{K}_{nq}^T = \mathbf{K}_{nq} \mathbf{U}^T \mathbf{U} \mathbf{K}_{nq}^T = \left( \mathbf{U} \mathbf{K}_{nq}^T \right)^T \left( \mathbf{U} \mathbf{K}_{nq}^T \right). \quad (13)$$

197 Therefore, if we set the number of inducing points to be  $q = p$ , then we can use  $\mathbf{z}_i \triangleq \mathbf{U} \left( \mathbf{K}_{np}^T \right)_{:,i}$   
 198 as a kernel-to-latent-space transformation because each element in  $\mathbf{K}$  **approximately satisfies** (4)  
 199 as stated in Theorem 1:  $K(x_i, x_j) = K_{ij} \approx \hat{K}_{ij} = \mathbf{z}_i^T \mathbf{z}_j$ . Conveniently, modern deep learning  
 200 frameworks can propagate gradients through the Cholesky operation, making it straightforward to  
 201 update the kernel parameters with gradient methods. Note that as we increase the number of inducing  
 202 points  $p$ , the approximation error between  $\mathbf{K}$  and  $\hat{\mathbf{K}}$  decreases. However, it is not recommended to  
 203 set  $p$  very large as updating the Cholesky decomposition requires  $\mathcal{O}(p^3)$ . **The empirical impact of**  
 204  **$p$  on computational time and performance is shown in Appendix E. In our experiments, only mild**  
 205 **values of  $p$  are necessary and the impact on computational is relatively small.**

## 206 5 Experimental Results

207 We evaluate ICKy on 3 different data sets: a synthetic data set, a remote sensing data set, and a data  
 208 set obtained from UCI Machine Learning Repository [13]. Note that in all the 3 experiments, our  
 209 ICKy framework only consists of 2 kernels (i.e.  $M = 2$ ), one NN-implied kernel and one chosen  
 210 kernel function with trainable parameters. To verify that ICKy can work with more than 2 kernels,  
 211 we create another synthetic data set with 3 kernels and show the corresponding results in Appendix  
 212 A. In addition, the experimental results for ICKr is provided in Appendix B. All experiments are  
 213 conducted on a computer cluster equipped with a GeForce RTX 2080 Ti GPU. **The implementation**  
 214 **details of all the experiments in this section are provided in Appendix G.**

### 215 5.1 Synthetic Data

216 To verify that ICKy can simulate a *multiplicative kernel*, we create a synthetic data set  $y \sim$   
 217  $\mathcal{GP}(0, K_1 K_2)$  containing 3000 data points where  $x^{(1)} \in [0, 1]$  is the input for the linear kernel  
 218  $K_1$  and  $x^{(2)} \in [0, 2]$  is the input for the *spectral mixture* kernel [56]  $K_2$  with 2 components. With  
 219 ICKy, we process  $x^{(1)}$  with a single-hidden-layer NN and  $x^{(2)}$  with a spectral mixture kernel function.  
 220 We evaluate ICKy on both a *prediction* task (where we first randomly shuffle the data points and do a  
 221 50:50 train-test split) and a *forecasting* task (where we use only the data points with  $x^{(2)} < 0.6$  for  
 222 training and the rest for testing).

223 We then compare ICKy with two models: a plain multi-layer perceptron (MLP) applied to the  
 224 concatenated features and a novel multi-layer perceptron-random forest (MLP-RF) joint model  
 225 employed by Zheng et al. [61] where MLP learns from  $x^{(1)}$  and RF learns from  $x^{(2)}$ . We believe  
 226 MLP-RF serves as a good benchmark model as it is a joint model with similar architecture to our  
 227 ICKy framework. To see how ICKy simulates the spectral mixture kernel, we plot only  $x^{(2)}$  against  
 228 the predicted value of  $y$  as shown in Figure 2. As can be seen from the figure, in the prediction task  
 229 (top row), plain MLP only captures the linear trend. MLP-RF only captures the mean of the spectral  
 230 mixture components. In contrast, our ICKy framework captures both the mean and the variance of the  
 231 spectral mixture kernel. In the forecasting task (bottom row), ICKy also outperforms plain MLP and  
 232 MLP-RF as it approximately captures the rising trend in the range of  $x^{(2)} \in [0.6, 1]$ . When  $x^{(2)} > 1$ ,  
 233 ICKy is unable to confidently extrapolate, so it starts to "*fail gracefully*," **by which we mean that the**  
 234 **prediction reverts to the mean of the prior (e.g., no observed information case.), just as would be**  
 235 **expected in a GP. However, we do not evaluate the posterior distribution to get a full sense of the**  
 236 **posterior uncertainty.**

237 We also test plain MLP, MLP-RF, and ICKy on the prediction task using different number of training  
 238 samples. As displayed in Figure 3, ICKy yields the smallest error among all the 3 frameworks even  
 239 with very limited data. In addition, to test the robustness of ICKy, we conduct the same experiments  
 240 on another synthetic data set in Appendix D to confirm that ICKy can simulate an *additive kernel*.

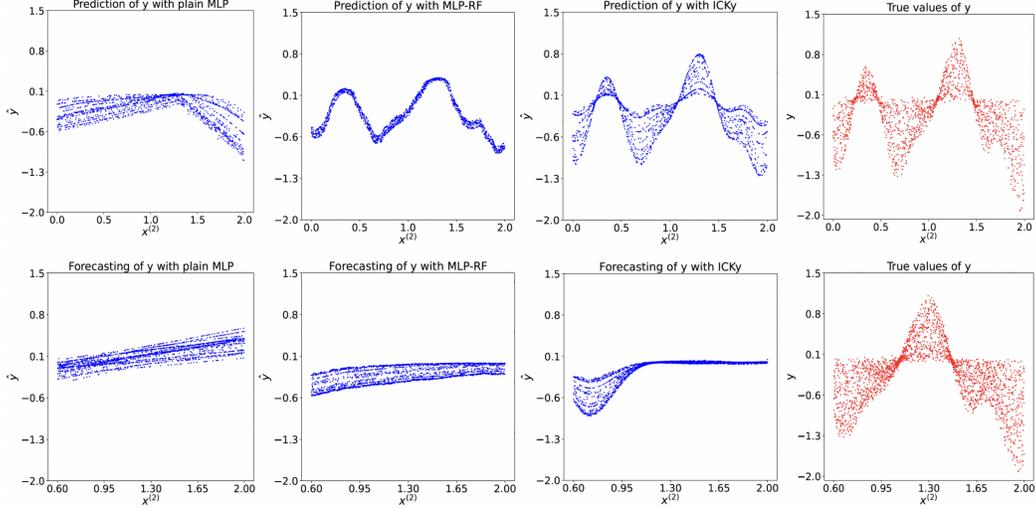


Figure 2: Prediction (top row) and forecasting (bottom row) of  $y \sim \mathcal{GP}(0, K_1 K_2)$ , where  $x^{(1)}$  is input to a linear kernel  $K_1$  and  $x^{(2)}$  is input to a spectral mixture kernel  $K_2$ . We plot  $x^{(2)}$  against the predicted  $y$ . We show results from a plain MLP (left column), MLP-RF (middle left column), and ICKy framework (middle right column), and we compare to the true values of  $y$  (right column).

## 241 5.2 Remote Sensing Data

242 We believe ICKy will be particularly useful for remote sensing applications. In this experiment, we  
 243 collect remote sensing data from 51 air quality monitoring (AQM) stations located in the National Capital  
 244 Territory (NCT) of Delhi and its satellite cities over the period from January 1, 2018 to June 30, 2020 (see  
 245 Appendix F for notes on data availability). Each data point  $\mathbf{x} = \{x, t\}$  contains 2 sources of information:  
 246 a three-band natural color (red-blue-green) satellite image  $x$  as the high-dimensional information and  
 247 the corresponding timestamp as the low-dimensional information. Note that we convert the timestamps  
 248 into numerical values  $t$  (where the day 2018-01-01 corresponds to  $t = 0$ ) before feeding them into the  
 249 models. Our goal is to predict the ground-level  $\text{PM}_{2.5}$  concentration  $\hat{y} = f(x, t)$  using both sources of in-  
 250 formation.  
 251  
 252  
 253  
 254  
 255  
 256  
 257  
 258

259 We split the train and test data set based on  $t$ . Specifically, we use all the data points with  $t \geq 500$  for  
 260 testing and the rest for training. As  $\text{PM}_{2.5}$  varies with time on a yearly basis, we use an *exponential-*  
 261 *sine-squared* kernel with a period of  $T = 365$  (days) to process the low-dimensional information  
 262  $t$ . The satellite images are processed with a CNN. Figure 4 shows the true versus the forecasted  
 263  $\text{PM}_{2.5}$  values by both ICKy and 2 benchmarks: a Convolutional Neural Network-Random Forest  
 264 (CNN-RF) joint model [61, 62] (similar to the MLP-RF model in Section 5.1, where RF learns the  
 265 temporal variation of  $\text{PM}_{2.5}$  and CNN captures the spatial variation of  $\text{PM}_{2.5}$  from satellite images)  
 266 and a *carefully designed* CNN-RF model that maps  $t$  into two new features,  $\sin(2\pi t/365)$  and  
 267  $\cos(2\pi t/365)$ , to explicitly model seasonality. As can be seen, ICKy outperforms both benchmarks  
 268 with the highest correlation coefficients and the lowest errors on the forecasting task. Specifically,  
 269 regular CNN-RF joint model fails to forecast  $\text{PM}_{2.5}$  as shown in Figure 4a. After including seasonality,  
 270 CNN-RF performs significantly better as shown in Figure 4b, but the forecasted  $\text{PM}_{2.5}$  values are still  
 271 less smooth than those from ICKy (Figure 4c) due to the discontinuous nature of the RF regressor  
 272 [6, 19]. We also visualize these results in the form of time series in Appendix F

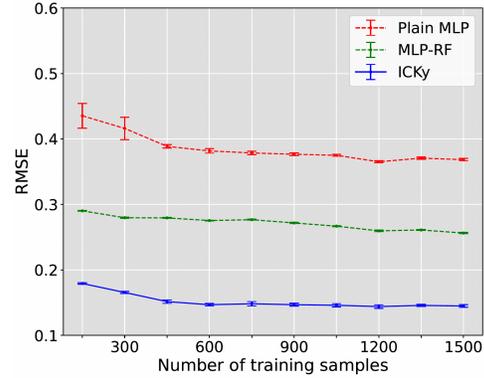


Figure 3: Prediction error of plain MLP, MLP-RF, and ICKy with different amount of training data generated by  $y = \mathcal{GP}(0, K_1 K_2)$ .

Table 1: Correlation and error statistics of ICKy and other joint deep models with both convolutional and attention-based architectures on the  $PM_{2.5}$  forecasting task

	Spearman R $\uparrow$	Pearson R $\uparrow$	RMSE $\downarrow$	MAE $\downarrow$
CNN-RF [61]	0.48	0.32	70.09	54.28
ViT-RF [11]	0.42	0.32	70.58	55.15
Seasonal CNN-RF	0.65	0.73	52.60	39.25
Seasonal ViT-RF	0.66	0.74	49.92	36.22
Seasonal DeepViT-RF [63]	0.68	0.76	48.87	35.32
Seasonal MAE-ViT-RF [22]	0.68	0.76	48.43	34.92
CNN-ICKy	<b>0.70</b>	<b>0.77</b>	<b>47.15</b>	<b>32.84</b>
ViT-ICKy	0.66	<b>0.77</b>	47.17	34.00
DeepViT-ICKy	0.62	0.73	48.68	34.10

273 We note that the inner product operation in ICK is similar in mathematical structure to attention-based  
 274 mechanisms [51] popular in many deep learning frameworks. Therefore, we compare ICKy with  
 275 4 attention-based benchmarks that we constructed based off of a Vision Transformer (ViT) [11]  
 276 architecture, including ViT-RF, Seasonal ViT-RF, Seasonal DeepViT-RF [63], and Seasonal MAE-  
 277 ViT-RF where ViT is pre-trained by a Masked Autoencoder [22], as displayed in Table 1. These use  
 278 the same RF and sinusoidal mappings as described previously to input the temporal information into  
 279 the model. We note that we are unaware of Vision Transformers being used in this manner, and that  
 280 all these models represent novel formulations. It can be observed that standard ViT-RF model fails to  
 281 forecast  $PM_{2.5}$  without seasonality incorporated, just as in CNN-RF. After introducing seasonality  
 282 by mapping  $t$  into sinusoidal features, ViT-based joint models yield higher correlation and smaller  
 283 error than CNN-RF but still underperform CNN-based ICKy. We also considered using a ViT-based  
 284 architecture for the CNN ICKy, and observed similar performance in these ICKy variants.

### 285 5.3 UCI Machine Learning Repository Data

286 To see if our ICKy framework generalizes to other domains, we acquire another data set containing  
 287 the normalized productivity and corresponding features of garment workers from the UCI Machine  
 288 Learning Repository. Imran et al. [1] employed a dense MLP with 2 hidden layers to predict the  
 289 worker productivity with collected features such as date, team number, targeted productivity, etc. To  
 290 test our ICKy framework, we separate out the *date* and use it as the low-dimensional information.  
 291 The rest of the features (excluding the temporal information) are then concatenated together to  
 292 serve as the high-dimensional information. Observing that the *daily averaged* worker productivity  
 293 has an approximate *monthly trend*, we again use an *exponential-sine-squared* kernel. The network  
 294 architecture of ICKy is the same as that of the two-hidden-layer MLP benchmark. To demonstrate  
 295 the strength of ICKy compared to other methods that equip NNs with GPs, we also add 2 additional  
 296 benchmarks here: a Gaussian Neural Process (GNP) [7] and an Attentive Gaussian Neural Process  
 297 (AGNP) [29]. Based on the results shown in Table 2, ICKy has the best performance when the  
 298 period parameter of the kernel is set to be  $T = 30$  (days) and it outperforms both MLP and NP  
 299 benchmarks by almost **one order of magnitude**. When we set  $T = 2$  or  $T = 7$ , this improvement is  
 300 less significant, which aligns with our initial observation that the daily averaged productivity has a  
 301 monthly seasonal trend. It is also worth noting that the GNP benchmarks here yield larger errors than  
 302 MLPs. A possible explanation is that GNP does not allow explicit assignment of a stationary kernel  
 303 (as the kernel models a posterior covariance) so it is hard for GNP to identify specific patterns in the  
 304 data such as seasonality without being given the pattern *a priori*.

Table 2: Prediction error of actual worker productivity on the test data set with ICKy and other benchmark models (MLPs and NPs)

	MSE $\downarrow$ ( $\ast 10^{-3}$ )	MAE $\downarrow$ ( $\ast 10^{-2}$ )	MAPE $\downarrow$
MLP [1]	20.16 $\pm$ 1.26	9.93 $\pm$ 0.36	17.30 $\pm$ 0.82
Cyclic MLP	20.97 $\pm$ 1.98	10.16 $\pm$ 0.77	17.48 $\pm$ 1.37
GNP [7, 37]	57.25 $\pm$ 4.31	19.39 $\pm$ 0.94	29.58 $\pm$ 1.63
AGNP [29]	43.11 $\pm$ 5.95	14.38 $\pm$ 0.88	22.59 $\pm$ 1.42
ICKy, $T = 2$	3.43 $\pm$ 1.42	4.85 $\pm$ 1.00	6.74 $\pm$ 1.37
ICKy, $T = 7$	0.44 $\pm$ 0.13	1.43 $\pm$ 0.15	2.22 $\pm$ 0.21
ICKy, $T = 30$	<b>0.31 <math>\pm</math> 0.09</b>	<b>1.17 <math>\pm</math> 0.14</b>	<b>1.79 <math>\pm</math> 0.22</b>

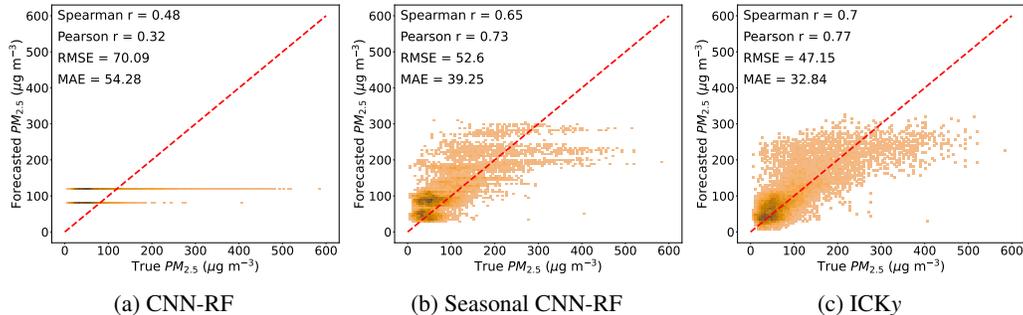


Figure 4: Density plots of the true  $PM_{2.5}$  concentrations against the **forecasted**  $PM_{2.5}$  concentrations for  $t \geq 500$  using (a) a CNN-RF joint model [61, 62], (b) a CNN-RF joint model with seasonality incorporated, and (c) our ICKy framework

## 6 Discussion

**Efficiency, Flexibility, and Generalization** Compared to exact composite GP models which scale  $\mathcal{O}(N^3)$ , the training process of our ICK framework is more efficient as it leverages standard back-propagation to learn both the parameters of NN and the kernel function. In addition, the network architecture of ICK can be very simple, as can be seen in all 3 experiments of ours, which further reduces its time and space complexity. Besides efficiency, our ICK framework is more flexible compared to other joint models (i.e. BNNs and CNN-RF). To be specific, the BNNs implemented by Pearce et al. [43] cannot simulate complicated kernels such as the spectral mixture kernel we use in Section 5.1. The CNN-RF joint model implemented by Zheng et al. [61] requires us to carefully design the input pre-processing procedure. Also, ICK generalizes well to unseen data even with very limited training samples. **There is a potential concern that ICKy may run into computational challenges when a large number of inducing points are required. This was not a problem in our experiments, but in large scale models this could be tackled by considering conjugate gradient methods, which have been recently popular in GP inference [15].**

**Limitations** A major limitation of ICK lies in our method of combining latent representations as the nature of inner product (i.e. the effect of multiplying small numbers) may cause *vanishing gradient* problems when we have a large number of sources of information (i.e.  $M$  is large). Furthermore, this paper only discusses the theoretical relationship between ICK and composite GPR *a priori*. This relationship will not exactly hold true *a posteriori*, although empirical results [34] and theoretical results [24] in slightly different contexts suggest that they may be close. Future work will evaluate this gap by exploring Bayesian neural networks and *a posteriori* properties.

**Broader Impacts** We believe our framework is extensively applicable to regression problems in many fields of study involving high-dimensional data and multiple sources of information with perceptible trends, such as remote sensing, spatial statistics, or clinical diagnosis.

## 7 Conclusion

This paper presents a novel yet surprisingly simple Implicit Composite Kernel (ICK) framework to learn from *hybrid* data containing both high-dimensional information and low-dimensional information with prior knowledge. We first analytically show the resemblance between ICK and composite GPR models and then conduct experiments using both synthetic and real-world data. It appears that ICK outperforms various benchmark models in our experiments with lowest prediction errors and highest correlations even with very limited data. Overall, we show that our ICK framework is exceptionally powerful when learning from *hybrid* data with our prior knowledge incorporated, and we hope our work can inspire more future research on joint machine learning models, enhancing their performance, efficiency, flexibility, and generalization capability.

339 **References**

- 340 [1] Abdullah Al Imran, Md Nur Amin, Md Rifatul Islam Rifat, and Shamprikta Mehreen. Deep  
341 neural network approach for predicting the productivity of garment employees. In *2019 6th*  
342 *International Conference on Control, Decision and Information Technologies (CoDIT)*, pages  
343 1402–1407. IEEE, 2019.
- 344 [2] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*,  
345 volume 4. Springer, 2006.
- 346 [3] Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine  
347 learning. *Siam Review*, 60(2):223–311, 2018.
- 348 [4] Mohamed A Bouhlel and Joaquim RRA Martins. Gradient-enhanced kriging for high-  
349 dimensional problems. *Engineering with Computers*, 35(1):157–173, 2019.
- 350 [5] Mohamed Amine Bouhlel, Nathalie Bartoli, Abdelkader Otsmane, and Joseph Morlier. Improv-  
351 ing kriging surrogates of high-dimensional design models by partial least squares dimension  
352 reduction. *Structural and Multidisciplinary Optimization*, 53(5):935–952, 2016.
- 353 [6] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- 354 [7] Wessel P Bruinsma, James Requeima, Andrew YK Foong, Jonathan Gordon, and Richard E  
355 Turner. The gaussian neural process. *arXiv preprint arXiv:2101.03606*, 2021.
- 356 [8] Nidhan Choudhuri, Subhashis Ghosal, and Anindya Roy. Nonparametric binary regression  
357 using a gaussian process prior. *Statistical Methodology*, 4(2):227–243, 2007.
- 358 [9] Abhirup Datta, Sudipto Banerjee, Andrew O Finley, and Alan E Gelfand. Hierarchical nearest-  
359 neighbor gaussian process models for large geostatistical datasets. *Journal of the American*  
360 *Statistical Association*, 111(514):800–812, 2016.
- 361 [10] Zhijie Deng, Jiaxin Shi, and Jun Zhu. Neuref: Deconstructing kernels by deep neural networks.  
362 *arXiv preprint arXiv:2205.00165*, 2022.
- 363 [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
364 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.  
365 An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*  
366 *arXiv:2010.11929*, 2020.
- 367 [12] Petros Drineas, Michael W Mahoney, and Nello Cristianini. On the nyström method for  
368 approximating a gram matrix for improved kernel-based learning. *journal of machine learning*  
369 *research*, 6(12), 2005.
- 370 [13] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
- 371 [14] David Duvenaud. *Automatic model construction with Gaussian processes*. PhD thesis, Univer-  
372 sity of Cambridge, 2014.
- 373 [15] Jacob Gardner, Geoff Pleiss, Kilian Q Weinberger, David Bindel, and Andrew G Wilson.  
374 Gpytorch: Blackbox matrix-matrix gaussian process inference with gpu acceleration. *Advances*  
375 *in neural information processing systems*, 31, 2018.
- 376 [16] Marta Garnelo, Jonathan Schwarz, Dan Rosenbaum, Fabio Viola, Danilo J Rezende, SM Eslami,  
377 and Yee Whye Teh. Neural processes. *arXiv preprint arXiv:1807.01622*, 2018.
- 378 [17] Adrià Garriga-Alonso, Carl Edward Rasmussen, and Laurence Aitchison. Deep convolutional  
379 networks as shallow gaussian processes. *arXiv preprint arXiv:1808.05587*, 2018.
- 380 [18] Alan E Gelfand and Erin M Schliep. Spatial statistics and gaussian processes: A beautiful  
381 marriage. *Spatial Statistics*, 18:86–104, 2016.
- 382 [19] Pierre Geurts, Damien Ernst, and Louis Wehenkel. Extremely randomized trees. *Machine*  
383 *learning*, 63(1):3–42, 2006.
- 384 [20] Mark Girolami and Simon Rogers. Variational bayesian multinomial probit regression with  
385 gaussian process priors. *Neural Computation*, 18(8):1790–1817, 2006.
- 386 [21] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- 387 [22] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked  
388 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on*  
389 *Computer Vision and Pattern Recognition*, pages 16000–16009, 2022.

- 390 [23] James Hensman, Nicolo Fusi, and Neil D Lawrence. Gaussian processes for big data. *arXiv*  
391 *preprint arXiv:1309.6835*, 2013.
- 392 [24] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and  
393 generalization in neural networks. *Advances in neural information processing systems*, 31,  
394 2018.
- 395 [25] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks.  
396 *Advances in neural information processing systems*, 28, 2015.
- 397 [26] Ziyang Jiang, Tongshu Zheng, Mike Bergin, and David Carlson. Improving spatial variation  
398 of ground-level pm<sub>2.5</sub> prediction with contrastive learning from satellite imagery. *Science of*  
399 *Remote Sensing*, page 100052, 2022.
- 400 [27] Andre G Journel and Charles J Huijbregts. *Mining geostatistics*. The Blackburn Press, 1976.
- 401 [28] Hyoung-Moon Kim, Bani K Mallick, and Chris C Holmes. Analyzing nonstationary spatial  
402 data using piecewise gaussian processes. *Journal of the American Statistical Association*,  
403 100(470):653–668, 2005.
- 404 [29] Hyunjik Kim, Andriy Mnih, Jonathan Schwarz, Marta Garnelo, Ali Eslami, Dan Rosenbaum,  
405 Oriol Vinyals, and Yee Whye Teh. Attentive neural processes. *arXiv preprint arXiv:1901.05761*,  
406 2019.
- 407 [30] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
408 *arXiv:1412.6980*, 2014.
- 409 [31] Daniel G Krige. A statistical approach to some basic mine valuation problems on the witwa-  
410 tersrand. *Journal of the Southern African Institute of Mining and Metallurgy*, 52(6):119–139,  
411 1951.
- 412 [32] Isaac E Lagaris, Aristidis Likas, and Dimitrios I Fotiadis. Artificial neural networks for solving  
413 ordinary and partial differential equations. *IEEE transactions on neural networks*, 9(5):987–  
414 1000, 1998.
- 415 [33] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444,  
416 2015.
- 417 [34] Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S Schoenholz, Jeffrey Pennington,  
418 and Jascha Sohl-Dickstein. Deep neural networks as gaussian processes. *arXiv preprint*  
419 *arXiv:1711.00165*, 2017.
- 420 [35] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous  
421 relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- 422 [36] Gary Marcus. Deep learning: A critical appraisal. *arXiv preprint arXiv:1801.00631*, 2018.
- 423 [37] Stratis Markou, James Requeima, Wessel Bruinsma, and Richard Turner. Efficient gaussian  
424 neural processes for regression. *arXiv preprint arXiv:2108.09676*, 2021.
- 425 [38] Takuo Matsubara, Chris J Oates, and François-Xavier Briol. The ridgelet prior: A covari-  
426 ance function approach to prior specification for bayesian neural networks. *arXiv preprint*  
427 *arXiv:2010.08488*, 2020.
- 428 [39] Alexander G de G Matthews, Mark Rowland, Jiri Hron, Richard E Turner, and Zoubin  
429 Ghahramani. Gaussian process behaviour in wide deep neural networks. *arXiv preprint*  
430 *arXiv:1804.11271*, 2018.
- 431 [40] Brian McFee, Gert Lanckriet, and Tony Jebara. Learning multi-modal similarity. *Journal of*  
432 *machine learning research*, 12(2), 2011.
- 433 [41] Ben Moseley, Andrew Markham, and Tarje Nissen-Meyer. Solving the wave equation with  
434 physics-informed deep learning. *arXiv preprint arXiv:2006.11894*, 2020.
- 435 [42] Radford M Neal. Priors for infinite networks. In *Bayesian Learning for Neural Networks*, pages  
436 29–53. Springer, 1996.
- 437 [43] Tim Pearce, Russell Tsuchida, Mohamed Zaki, Alexandra Brintrup, and Andy Neely. Expressive  
438 priors in bayesian neural networks: Kernel combinations and periodic functions. In *Uncertainty*  
439 *in artificial intelligence*, pages 134–144. PMLR, 2020.

- 440 [44] Dejan Petelin, Alexandra Grancharova, and Juš Kocijan. Evolving gaussian process models  
441 for prediction of ozone concentration in the air. *Simulation modelling practice and theory*,  
442 33:68–80, 2013.
- 443 [45] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. *Advances*  
444 *in neural information processing systems*, 20, 2007.
- 445 [46] Ali Rahimi and Benjamin Recht. Weighted sums of random kitchen sinks: Replacing minimiza-  
446 tion with randomization in learning. *Advances in neural information processing systems*, 21,  
447 2008.
- 448 [47] John Shawe-Taylor, Nello Cristianini, et al. *Kernel methods for pattern analysis*. Cambridge  
449 university press, 2004.
- 450 [48] Tao Shi and Steve Horvath. Unsupervised learning with random forest predictors. *Journal of*  
451 *Computational and Graphical Statistics*, 15(1):118–138, 2006.
- 452 [49] Edward Snelson and Zoubin Ghahramani. Sparse gaussian processes using pseudo-inputs.  
453 *Advances in neural information processing systems*, 18, 2005.
- 454 [50] Mark Van der Wilk, Carl Edward Rasmussen, and James Hensman. Convolutional gaussian  
455 processes. *Advances in Neural Information Processing Systems*, 30, 2017.
- 456 [51] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
457 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*  
458 *processing systems*, 30, 2017.
- 459 [52] Hao Wang, Hao He, and Dina Katabi. Continuously indexed domain adaptation. *International*  
460 *Conference on Machine Learning*, 2020.
- 461 [53] Christopher Williams and Matthias Seeger. Using the nyström method to speed up kernel  
462 machines. *Advances in neural information processing systems*, 13, 2000.
- 463 [54] Christopher K Williams and Carl Edward Rasmussen. *Gaussian processes for machine learning*,  
464 volume 2. MIT press Cambridge, MA, 2006.
- 465 [55] Christopher KI Williams and David Barber. Bayesian classification with gaussian processes.  
466 *IEEE Transactions on pattern analysis and machine intelligence*, 20(12):1342–1351, 1998.
- 467 [56] Andrew Wilson and Ryan Adams. Gaussian process kernels for pattern discovery and extrapo-  
468 lation. In *International conference on machine learning*, pages 1067–1075. PMLR, 2013.
- 469 [57] Andrew Gordon Wilson, Zhiting Hu, Ruslan Salakhutdinov, and Eric P Xing. Deep kernel  
470 learning. In *Artificial intelligence and statistics*, pages 370–378. PMLR, 2016.
- 471 [58] Andrew Gordon Wilson, David A Knowles, and Zoubin Ghahramani. Gaussian process  
472 regression networks. *arXiv preprint arXiv:1110.4411*, 2011.
- 473 [59] Tianbao Yang, Yu-Feng Li, Mehrdad Mahdavi, Rong Jin, and Zhi-Hua Zhou. Nyström method  
474 vs random fourier features: A theoretical and empirical comparison. *Advances in neural*  
475 *information processing systems*, 25, 2012.
- 476 [60] Daoqiang Zhang, Yaping Wang, Luping Zhou, Hong Yuan, Dinggang Shen, Alzheimer’s  
477 Disease Neuroimaging Initiative, et al. Multimodal classification of alzheimer’s disease and  
478 mild cognitive impairment. *Neuroimage*, 55(3):856–867, 2011.
- 479 [61] Tongshu Zheng, Michael Bergin, Guoyin Wang, and David Carlson. Local pm2. 5 hotspot  
480 detector at 300 m resolution: A random forest–convolutional neural network joint model jointly  
481 trained on satellite images and meteorology. *Remote Sensing*, 13(7):1356, 2021.
- 482 [62] Tongshu Zheng, Michael H Bergin, Shijia Hu, Joshua Miller, and David E Carlson. Estimating  
483 ground-level pm2. 5 using micro-satellite images by a convolutional neural network and random  
484 forest approach. *Atmospheric Environment*, 230:117451, 2020.
- 485 [63] Daquan Zhou, Bingyi Kang, Xiaojie Jin, Linjie Yang, Xiaochen Lian, Zihang Jiang, Qibin Hou,  
486 and Jiashi Feng. Deepvit: Towards deeper vision transformer. *arXiv preprint arXiv:2103.11886*,  
487 2021.

488 **Checklist**

- 489 1. For all authors...
- 490 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
491 contributions and scope? [Yes] See Section 1.
- 492 (b) Did you describe the limitations of your work? [Yes] See Section 6
- 493 (c) Did you discuss any potential negative societal impacts of your work? [No] We are  
494 unaware of direct negative societal impacts.
- 495 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
496 them? [Yes]
- 497 2. If you are including theoretical results...
- 498 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Section 4.1
- 499 (b) Did you include complete proofs of all theoretical results? [Yes] See Section 4.1
- 500 3. If you ran experiments...
- 501 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
502 mental results (either in the supplemental material or as a URL)? [Yes] The code, data,  
503 and instructions to reproduce the experimental results are provided via the anonymous  
504 repository URL in the abstract.
- 505 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
506 were chosen)? [Yes] See Section 5.
- 507 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
508 ments multiple times)? [Yes] See Section 5.1 and 5.3.
- 509 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
510 of GPUs, internal cluster, or cloud provider)? [Yes] See section 5
- 511 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 512 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 513 (b) Did you mention the license of the assets? [Yes] See Appendix I
- 514 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
515 We have included code, with the link in the abstract to an anonymous repository. This  
516 will be switched to a user-maintained repository if accepted.
- 517 (d) Did you discuss whether and how consent was obtained from people whose data you're  
518 using/curating? [N/A] We do not release data.
- 519 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
520 information or offensive content? [Yes] None of our datasets contain personally  
521 identifiable information or offensive content.
- 522 5. If you used crowdsourcing or conducted research with human subjects...
- 523 (a) Did you include the full text of instructions given to participants and screenshots, if  
524 applicable? [N/A]
- 525 (b) Did you describe any potential participant risks, with links to Institutional Review  
526 Board (IRB) approvals, if applicable? [N/A]
- 527 (c) Did you include the estimated hourly wage paid to participants and the total amount  
528 spent on participant compensation? [N/A]

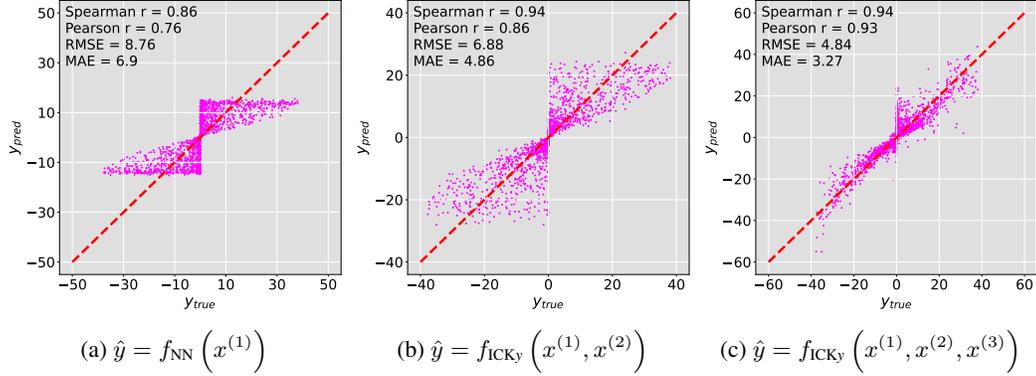


Figure A1: Scatter plots of the true values of  $y$  against the predicted values of  $y$  using our ICKy framework with (a) one source, (b) 2 sources, and (c) 3 sources of information

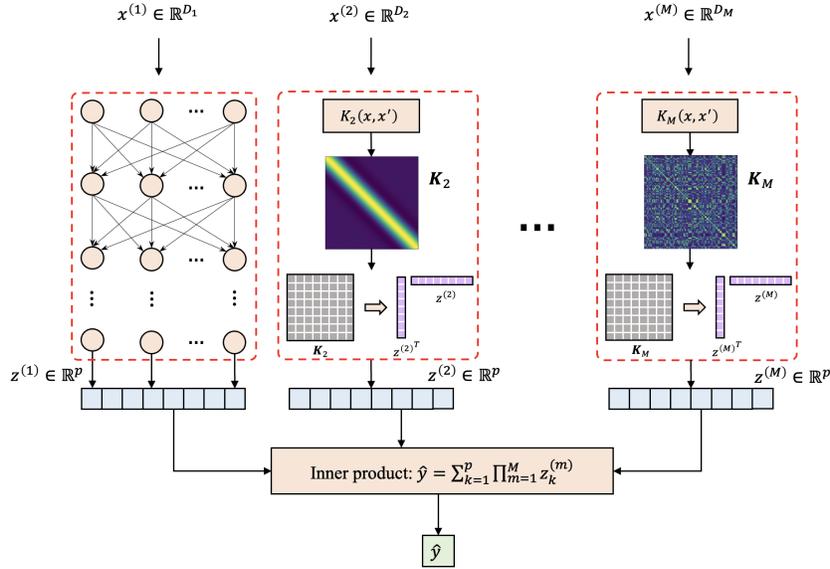


Figure A2: Given data containing  $M$  sources of information  $\mathbf{x} = \{x^{(1)}, x^{(2)}, \dots, x^{(M)}\}$ , we can process the data using our ICK framework where high-dimensional information (e.g.  $x^{(1)}$  in the figure) is processed using a neural network and low-dimensional information (e.g.  $x^{(2)}$  in the figure) is processed using a kernel function followed by Nyström or RFF transformation.

## 529 A ICK with More Than Two Kernels

530 Besides the visualization presented in Figure 1, we also show our ICK framework for processing data  
 531  $\mathbf{x} = \{x^{(1)}, x^{(2)}, \dots, x^{(M)}\}$  with  $M > 2$  sources of information in Figure A2. Here  $K^{(2)}, \dots, K^{(M)}$   
 532 represent different types of kernels with trainable parameters. The final prediction is calculated by a  
 533 chained inner product of all extracted representations  $\hat{y} = \sum_{k=1}^p \prod_{m=1}^M z_k^{(m)}$ .

534 To confirm that ICKy can work with more than 2 kernels, we construct another synthetic data set  
 535 containing 3000 data points in total. Each input  $\mathbf{x} = \{x^{(1)}, x^{(2)}, x^{(3)}\}$  has 3 sources of information.  
 536 The output  $y$  is generated by  $y = x^{(3)} \tanh(2x^{(1)} \cos^2(\pi x^{(2)}/50)) + \epsilon$  where  $\epsilon$  is a Gaussian noise  
 537 term. We process  $x^{(1)}$  with a small single-hidden-layer NN,  $x^{(2)}$  with an *exponential sine squared*  
 538 kernel, and  $x^{(3)}$  with a *radial-basis function* (RBF) kernel. Figure A1 shows the prediction results  
 539 as we progressively add more sources of information into our ICKy framework with corresponding  
 540 kernel functions. It can be observed that ICKy yields both smallest error and highest correlation

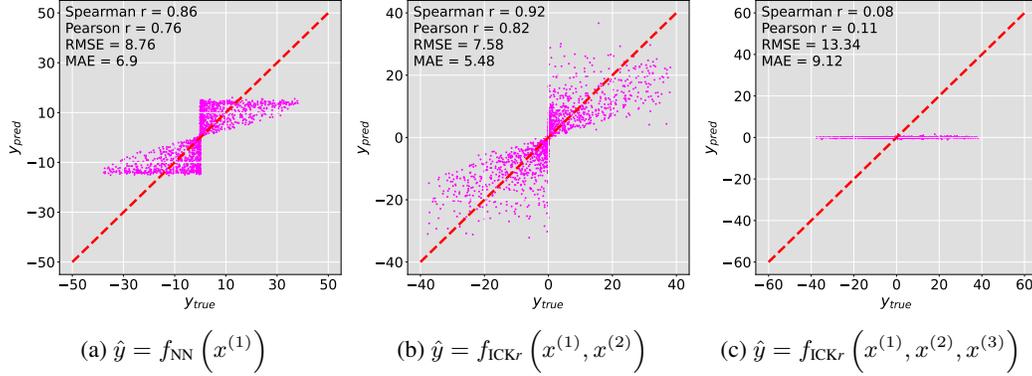


Figure B1: Scatter plots of the true values of  $y$  against the predicted values of  $y$  using our ICKr framework with (a) one source of information, (b) 2 sources of information, and (c) 3 sources of information. Note that here we use RFF for kernel-to-latent-space transformation.

541 with information from all 3 different sources. Hence, ICKy works well with the  $M = 3$  case and the  
 542 regression performance is improved as we add in more information related to the target.

## 543 B Random Fourier Features

### 544 B.1 Methodology

545 Random Fourier Features (RFF) is another popular approximation method used for kernel learning  
 546 [45]. Unlike the Nyström method which approximates the entire kernel matrix, RFF directly ap-  
 547 proximates the kernel function  $K$  using some randomized feature mapping  $\phi: \mathbb{R}^{D_m} \rightarrow \mathbb{R}^{2d_m}$  such  
 548 that  $K(x_i^{(m)}, x_j^{(m)}) \approx \phi(x_i^{(m)})^T \phi(x_j^{(m)})$ . To obtain the feature mapping  $\phi$ , based on Bochner's  
 549 theorem, we first compute the Fourier transform  $p$  of kernel  $K$

$$p(\omega) = \frac{1}{(2\pi)^{D_m}} \int_{-\infty}^{+\infty} e^{-j\omega^T \delta} K(\delta) d\delta, \quad (14)$$

550 where  $\delta = x_i^{(m)} - x_j^{(m)}$ . Then we draw  $d_m$  i.i.d. samples  $\omega_1, \omega_2, \dots, \omega_{d_m}$  from  $p(\omega)$  and construct  
 551 the feature mapping  $\phi$  as follows

$$\phi(x^{(m)}) \equiv d_m^{-1/2} \left[ \cos(\omega_1^T x^{(m)}), \dots, \cos(\omega_{d_m}^T x^{(m)}), \sin(\omega_1^T x^{(m)}), \dots, \sin(\omega_{d_m}^T x^{(m)}) \right]. \quad (15)$$

552 Since  $\phi(x^{(m)}) \in \mathbb{R}^{2d_m}$ , we need to set  $d_m = p/2$  when using RFF as a kernel-to-latent-space  
 553 transformation. In addition, since RFF involves sampling from a distribution, the kernel parameters  
 554 are thus not directly differentiable and we need to apply a reparameterization trick [35] to learn those  
 555 parameters.

## 556 B.2 Experimental Results

### 557 B.2.1 Synthetic Data

558 We use the same toy data set where each data point  $x = \{x^{(1)}, x^{(2)}, x^{(3)}\}$  contains 3 sources of  
 559 information as described in Appendix A. Also, we use the same types of kernels as those in ICKy as  
 560 discussed in Appendix A. The only difference here is that we use RFF instead of Nyström method to  
 561 transform the kernel matrix into the latent space in ICKr framework.

562 The results are displayed in Figure B1. It can be observed that when we add in only the side  
 563 information  $x^{(2)}$  along with the *exponential sine squared* kernel, both the correlation and the predictive  
 564 performance are improved (though not as good as the results from ICKy as shown in Figure A1).

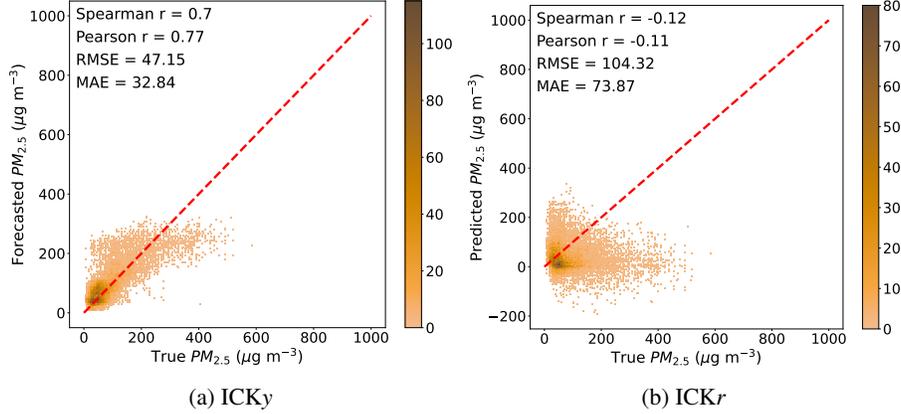


Figure B2: Density plots of the true  $PM_{2.5}$  concentrations against the forecasted  $PM_{2.5}$  concentrations for  $t \geq 500$  using our ICK framework with (a) ICKy and (b) ICKr

565 However, after we further include  $x^{(3)}$  with the *RBF* kernel, we realize that the parameters of ICKr  
 566 become very hard to optimize and it fails to make valid predictions and starts to guess randomly  
 567 around zero.

### 568 B.2.2 Remote Sensing Data

569 We also try ICKr on the *forecasting* task using the remote sensing data (see Section 5.2) and compare  
 570 the results with those from ICKy. Each data point  $x = \{x, t\}$  contains a satellite image  $x$  as the  
 571 high-dimensional information and its corresponding timestamp  $t$  as the low-dimensional information.  
 572 The satellite images are processed with a two-layer CNN and the timestamps are processed with an  
 573 *exponential-sine-squared* kernel with a period of  $T = 365$  (days). As can be observed from Figure  
 574 B2, ICKr yields much higher error compared to ICKy.

## 575 C Estimated Kernel Matrix and its Eigen-spectrum

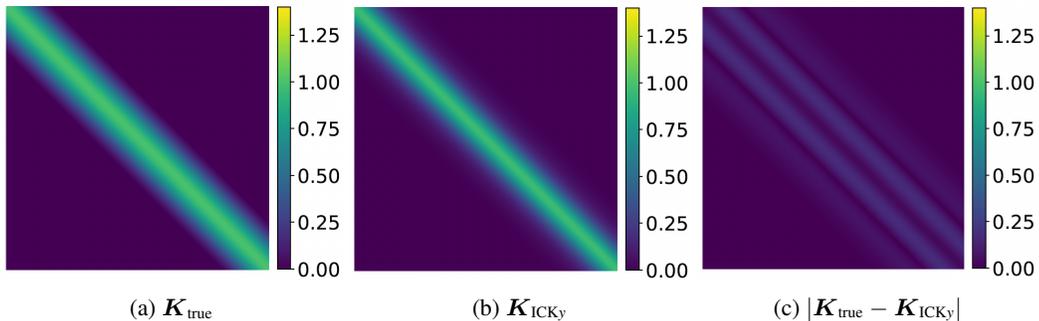


Figure C1: Visualization of (a) True matrix (b) estimated matrix by our ICKy framework, and (c) absolute difference between the true and estimated matrix for the spectral mixture kernel

576 We first examine whether ICKy can retrieve the spectral mixture kernel in the prediction task. After  
 577 fitting the parameters of the spectral mixture kernel in ICKy, we compute the kernel matrix  $K_{ICKy}$   
 578 using these learned parameters and compare it with the true kernel matrix  $K_{true}$  by calculating the  
 579 absolute difference between them as displayed in Figure C1. As can be observed,  $K_{ICKy}$  and  $K_{true}$   
 580 are similar and their absolute difference is relatively small, indicating that ICKy can approximately  
 581 retrieve the spectral mixture kernel.

582 Yang et al. [59] studied the fundamental difference  
 583 between Nyström method and Random Fourier Fea-  
 584 tures (RFF). They conclude that Nyström-method-  
 585 based approaches can yield much better generaliza-  
 586 tion error bound than RFF-based approaches if there  
 587 exists a large gap in the eigen-spectrum of the kernel  
 588 matrix. This phenomenon is mainly caused by how  
 589 these two methods construct their basis functions. In  
 590 particular, the basis functions used by RFF are sam-  
 591 pled from a Gaussian distribution that is independent  
 592 from the training examples, while the basis functions  
 593 used by the Nyström method are sampled from the  
 594 training samples so they are data-dependent. In our  
 595 synthetic data experiments, we train our ICK frame-  
 596 work using a batch size of 50. The eigenvalues of the  
 597 kernel matrices computed from the first 4 batches of  
 598 the synthetic data set are displayed in Figure C2. It  
 599 can be observed that the first few eigenvalues of the  
 600 kernel matrix are much larger than the remaining eigenvalues. Namely, there exists a large gap in the  
 601 eigen-spectrum of the kernel matrix, which helps explain why ICK<sub>y</sub> has a much better performance  
 602 than ICK<sub>r</sub>.

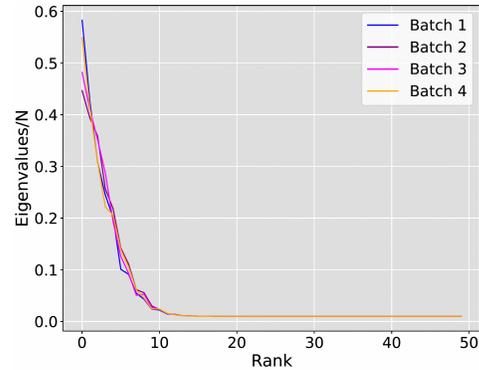


Figure C2: Eigenvalues of the kernel matrix computed from the first 4 batches of training data

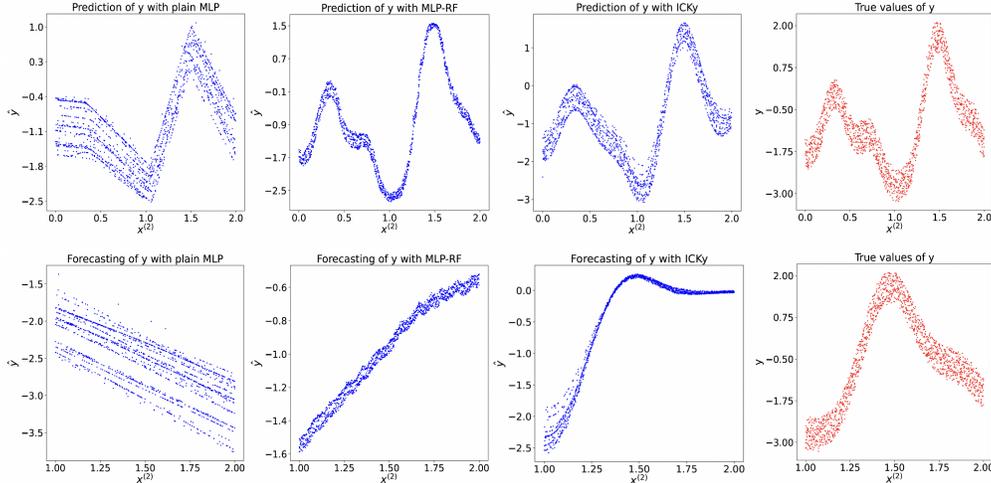


Figure D1: Prediction (top row) and forecasting (bottom row) of  $y \sim \mathcal{GP}(0, K_1 + K_2)$ , where  $x^{(1)}$  is input to a linear kernel  $K_1$  and  $x^{(2)}$  is input to a spectral mixture kernel  $K_2$ . Here we only plot  $x^{(2)}$  against the predicted  $y$ . We implement 3 types of models: plain MLP (left column), MLP-RF (middle left column), and our ICKy framework (middle right column), and we compare the results with the true values of  $y$  (right column).

## 603 D Simulation of an Additive Kernel

604 While ICK is designed to capture multiplicative kernels,  
 605 we evaluated how well it could capture additive  
 606 kernels. We conduct experiments using another  
 607 synthetic data set generated by an additive kernel  
 608  $y \sim \mathcal{GP}(0, K_1 + K_2)$  with the same training settings.  
 609 As shown in Figure D1, ICKy again outperforms  
 610 plain MLP and MLP-RF in both the prediction and  
 611 the forecasting tasks. Moreover, we again test plain  
 612 MLP, MLP-RF, and ICKy on the prediction task using  
 613 different number of training samples. As displayed  
 614 in Figure D2, ICKy yields the smallest error among  
 615 all the 3 frameworks. Also, the performance gap  
 616 between ICKy and the other 2 benchmark models  
 617 shrink as we feed in more training data. Therefore,  
 618 we conclude ICKy is robust enough to simulate both  
 619 additive and multiplicative kernels.

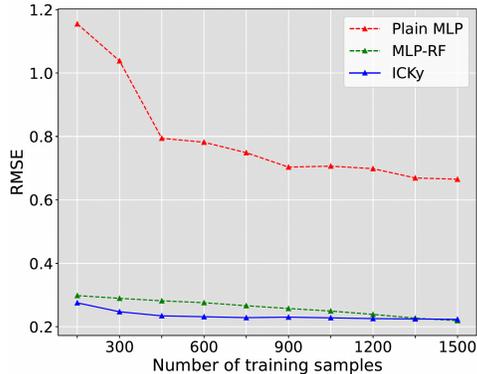


Figure D2: Prediction error of plain MLP, MLP-RF, and ICKy with different amount of training data generated by  $y \sim \mathcal{GP}(0, K_1 + K_2)$

## 620 E Number of Inducing Points

621 As discussed in Section 4.2.1, as we increase the number of inducing points  $p$ , we expect the  
 622 approximation error between the true kernel matrix  $K$  and the approximated kernel matrix  $\hat{K}$  to  
 623 decrease. Here, we empirically show how the value of  $p$  impacts our predictions. In Figure E1a, we  
 624 plot the prediction error of  $\hat{y} = f_{\text{ICKy}}(x^{(1)}, x^{(2)}, x^{(3)})$  against the number of inducing points using  
 625 the synthetic data generated in Appendix A. As can be observed, the prediction error drops sharply as  
 626 we raise  $p$  from a small value (e.g.  $p = 2$ ). When  $p$  is relatively large, increasing  $p$  yields smaller  
 627 improvement on the predictions. Additionally, in Figure E1b, we plot the total training time against  
 628  $p$ . The total training time is dependent on how long a single iteration takes and the total number of  
 629 epochs required. We note that once  $p > 80$  the training time is relatively flat, which is due to the  
 630 fact that the total computation in the Cholesky is less than the computation in the neural network.  
 631 Interestingly, it appears that when  $p$  is very small, ICKy takes longer to converge due to the need for  
 632 many more epochs. As we increase  $p$ , the training time goes down and then goes up again due to the

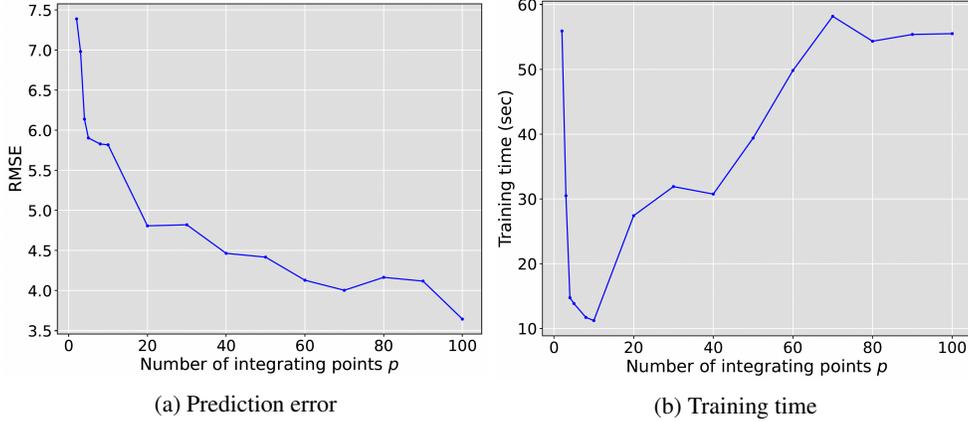


Figure E1: Plots of (a) prediction error and (b) training time of  $\hat{y} = f_{\text{ICKy}}(x^{(1)}, x^{(2)}, x^{(3)})$  against the number of inducing points  $p$

633 computational complexity, i.e.  $\mathcal{O}(p^3)$ , of the Cholesky decomposition. Based on these observations,  
 634 we are overly concerned about the computational complexity for reasonable values of  $p$ .

## 635 F Visualization of Remote Sensing Data as Time Series

636 To better illustrate the results in Section 5.2, we visualize those results in the form of time series. As  
 637 shown in Figure F1, the plain CNN-RF model does not work as it tends to forecast constant  $\text{PM}_{2.5}$   
 638 values. In contrast, both the seasonal CNN-RF model and our ICKy framework captures the overall  
 639 trend of the true daily averaged  $\text{PM}_{2.5}$  values, but the forecasted values by ICKy are smoother and  
 640 yield smaller error.

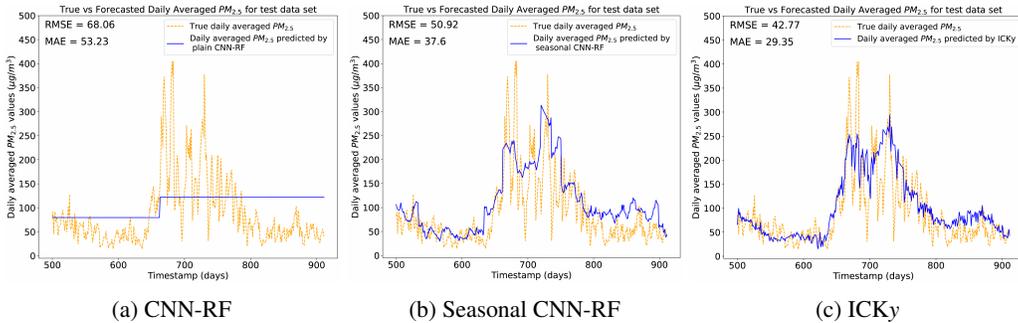


Figure F1: Time series visualization of the true against the forecasted *daily averaged*  $\text{PM}_{2.5}$  concentrations for  $t \geq 500$  using (a) a CNN-RF joint model [61, 62], (b) a CNN-RF joint model with seasonality incorporated, and (c) our ICKy framework

## 641 G Experimental Details

### 642 G.1 Synthetic Data

643 We use the GPytorch package [15] to generate the synthetic data. Before feeding  $x^{(1)}$  into MLP, we  
 644 first map  $x^{(1)}$  into higher dimension using an unsupervised algorithm called Totally Random Trees  
 645 Embedding [48]. All the MLP structures in this experiment (including those in MLP-RF and ICKy)  
 646 contain one single fully connected (FC) layer of width 1000, which serves as a simple benchmark  
 647 since a one-hidden-layer MLP can only capture linear relationship between the input and output. For  
 648 model training, we optimize a Mean Squared Error (MSE) objective using Adam optimizer [30] with  
 649 a weight decay of 0.1.

Table 3: Model architecture and training details for remote sensing data experiment in Section 5.2

	Backbone architecture details	Output FC layers dimension	Optimizer
CNN-RF	# Conv blocks = 1, # Channels = 16, Kernel size = 3, Stride = 1	1000, 512, 512, 1	Adam $\beta_1 = 0.9$ $\beta_2 = 0.999$
ViT-RF	# Transformer blocks = 6, # Attention heads = 8, Dropout ratio = 0.1	1000, 512, 512, 1	Adam $\beta_1 = 0.9$ $\beta_2 = 0.999$
DeepViT-RF	# Transformer blocks = 6, # Attention heads = 8, Dropout ratio = 0.1	1000, 512, 512, 1	Adam $\beta_1 = 0.9$ $\beta_2 = 0.999$
MAE-ViT-RF	# Transformer blocks = 6, # Attention heads = 8, Dropout ratio = 0.1, Masking ratio = 0.75	1000, 512, 512, 1	Adam $\beta_1 = 0.9$ $\beta_2 = 0.999$
CNN-ICKy	# Conv blocks = 1, # Channels = 16, Kernel size = 3, Stride = 1	1000, 512, $p$	SGD momentum = 0.9
ViT-ICKy	# Transformer blocks = 6, # Attention heads = 8, Dropout ratio = 0.1	1000, 512, $p$	SGD momentum = 0.9
DeepViT-ICKy	# Transformer blocks = 6, # Attention heads = 8, Dropout ratio = 0.1	1000, 512, $p$	SGD momentum = 0.9

650 **G.2 Remote Sensing Data**

651 The model architecture and training details are listed in Table 3. Here  $p$  denotes the length of latent  
652 representations  $z$  as discussed in Section 4. Note that we use stochastic gradient descent (SGD)  
653 optimizer with a momentum of 0.9 for ICKy as we realize that SGD helps ICKy find a local minimum  
654 on the objective more efficiently. We use MSE objective for ICKy and all benchmark models in this  
655 experiment.

656 **G.3 UCI Machine Learning Repository Data**

657 The MLPs (including the MLP part in ICKy) in this experiment share the same structure as the one  
658 used in [1], which consist of 3 hidden layers of width 128, 32, and 32, respectively. For plain MLP,  
659 cyclic MLP, and ICKy, we use the mean absolute error (MAE) objective to put less weight on the  
660 outliers and thus enhance the model performance. For GNP and AGNP, we maximize a biased Monte  
661 Carlo estimate of the log-likelihood objective as discussed in [37]. All these objectives are optimized  
662 by an Adam optimizer with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ .

663 **H Adapting ICK for Classification**

664 While regression tasks are the primary motivation for this paper, there are many ways to adapt GPR  
665 for classification tasks. For example, a binary classification model can be created by using a sigmoid  
666 [55] or probit link [8] on the output of the GP. Succinctly, given a function  $f(\mathbf{x}) \sim \mathcal{GP}(0, K(\mathbf{x}, \mathbf{x}'))$ ,  
667 the binary outcome probability is given as  $p(y = 1|f(\mathbf{x})) = \sigma(f(\mathbf{x}))$ . Likewise, a multiple  
668 classification model can be constructed by using a multi-output GP (or multiple GPs) and putting  
669 the outputs through a softmax function [55] or multinomial probit link [20]. This strategy can be  
670 summarized by calculating  $C$  different functions  $f_c(\mathbf{x}) \sim \mathcal{GP}(0, K(\mathbf{x}, \mathbf{x}'))$  for  $c = 1, \dots, C$ , where  
671  $C$  is the number of classes, and then calculating the class probabilities through a link function,  
672  $p(y|x) = \text{softmax}([f_1(x), f_2(x), \dots, f_C(x)])$ .

673 This same logic can be used to construct a multiple classification model from ICKy. Succinctly, let  
674  $r_c = f_{NN,c}(x^{(1)}) \odot z_c^{(2)}$ , where  $f_{NN,c}$  denotes a neural network specific to the  $c^{th}$  class and  $z_c^{(2)}$   
675 represents the Nyström approximation specific to the kernel for the  $c^{th}$  class. We note that often in

676 a multi-output case the kernel parameters are shared, and so  $z_c^{(2)}$  would be an identical vector for  
677 each class. Then, the output probabilities for a data sample as  $p(y|x) = \text{softmax}([r_1, \dots, r_C])$ . This  
678 framework is learned with a cross-entropy loss.

679 To provide *proof-of-concept* of this multiple classification strategy, we implemented this model on a  
680 version of Rotating MNIST. In this task, a dataset was created by rotating each image in the dataset  
681 by a uniform random value  $\phi \in [0, 2\pi)$ , thus creating a dataset with 60,000 images each with an  
682 associated rotation covariate  $\phi$ . We implemented the above multiple classification model with a  
683 periodic kernel over the rotation angle. This strategy yielded an accuracy of 92.3% on the validation  
684 data. This is lower than methods such as spatial transformers [25] that report accuracy greater than  
685 99%. However, those models explicitly use the fact that the information is simply rotated, whereas  
686 ICK is modeling a smooth transformation in the prediction function as a function of angle. This  
687 ICK classification model is much closer in concept to the way Rotating MNIST is used to evaluate  
688 unsupervised domain adaptation. While the evaluation strategy is different than our random validation  
689 set, the state-of-the-art accuracy on unsupervised domain adaption is 87.1% [52]. Due to the lack of  
690 complete and fair comparisons, we are not claiming that ICKy is state-of-the-art for classification, but  
691 ICKy’s classification model does seem reasonable and viable based upon this result.

## 692 I Generation, Accessibility, and Restrictions of the Data

693 The synthetic data  $y \sim \mathcal{GP}(0, K_1 K_2)$  in Section 5.1 and  $y \sim \mathcal{GP}(0, K_1 + K_2)$  in Appendix D  
694 are generated using the GPyTorch package. The remote sensing data in Section 5.2 is downloaded  
695 using PlanetScope API whose content is protected by copyright and/or other intellectual property  
696 laws. To access the data on PlanetScope, the purchase of an end-user license is required. When this  
697 manuscript is accepted, we will provide the codes we used to acquire the data. The UCI machine  
698 learning repository data we use in Section 5.3 has an open access license, meaning that the data is  
699 freely available online.