

# Trust Your Robots! Predictive Uncertainty Estimation of Neural Networks with Sparse Gaussian Processes

Jongseok Lee<sup>1,2</sup> Jianxiang Feng<sup>1,3</sup> Matthias Humt<sup>1,3</sup> Marcus G. Müller<sup>1,4</sup> Rudolph Triebel<sup>1,3</sup>

<sup>1</sup>Institute of Robotics and Mechatronics, German Aerospace Center (DLR)

<sup>2</sup>High Performance Humanoid Technologies, Karlsruhe Institute of Technology (KIT)

<sup>3</sup>Chair of Computer Vision and Artificial Intelligence, Technical University of Munich (TUM)

<sup>4</sup>Autonomous Systems Laboratory, ETH Zürich (ETHZ)

**Abstract:** This paper presents a probabilistic framework to obtain both reliable and fast uncertainty estimates for predictions with Deep Neural Networks (DNNs). Our main contribution is a practical and principled combination of DNNs with sparse Gaussian Processes (GPs). We prove theoretically that DNNs can be seen as a special case of sparse GPs, namely mixtures of GP experts (MoE-GP), and we devise a learning algorithm that brings the derived theory into practice. In experiments from two different robotic tasks – inverse dynamics of a manipulator and object detection on a micro-aerial vehicle (MAV) – we show the effectiveness of our approach in terms of predictive uncertainty, improved scalability, and run-time efficiency on a Jetson TX2. We thus argue that our approach can pave the way towards reliable and fast robot learning systems with uncertainty awareness.

**Keywords:** Robotic Introspection, Bayesian Deep Learning, Gaussian Processes

## 1 Introduction

This work aims to provide an algorithm that can estimate uncertainty of DNN predictions reliably and fast, and at the same time, is suited for integration into a large range of robotic systems. Generic solutions to this problem are crucial for safe robot learning and introspection [1, 2], giving the robots with an ability to assess their own failure probabilities and to alter their behaviors towards safety. While the state of the art is advancing [3, 4, 5], we still face practical difficulties in two main domains. First, the current methods are often not efficient at test-time, e.g. for a single input, the methods require multiple predictions from several copies of a model [4] or samples from the model’s distribution [3]. And second, in general these methods do not provide reliable uncertainty estimates when compared to GPs - often known as the golden standard of probabilistic machine learning [6, 7].

Therefore, we propose to estimate the *predictive uncertainty* of DNNs using MoE-GPs [8, 9] - a sparse variant of GPs that divides the input space into smaller local regions using a *gating function*, where individual GPs called *experts* learn and make predictions (see figure 1). To do so, we provide both theoretic foundations and a practical learning algorithm. First, we formally derive a connection between DNNs and MoE-GPs. As a result, we reveal how MoE-GPs with a DNN-based kernel [10] can provably approximate uncertainty in DNNs. Moreover, we devise a learning algorithm that brings the derived theory into practice. Our solution involves a gating function that strictly divides the input data into smaller subsets by performing clustering in kernel space, and we propose the concept of a patchwork prior, mitigating the problem of discontinuity between local GP experts at their boundaries. For efficiency, we further propose to exploit active learning and model compression techniques.

Our approach has several favorable features for many classification and regressions tasks in robotics. At run-time, it maintains the predictive power of a DNN, and its uncertainty estimates do not require combining multiple predictions of DNNs. Moreover, we inherit the benefits of MoE-GPs for uncertainty estimates. These include an improved scalability when compared to a GP with the neural tangent kernel [10] (NTK), as well as natural support for distributed training. We contend that such probabilistic methods must run efficiently on a robot and as applicable as the sparse GPs. Therefore, we not only provide ablation studies and evaluations against the state-of-the-art (SOTA), but we also show that our method (i) scales to approximately 2 million data points for the task of learning inverse

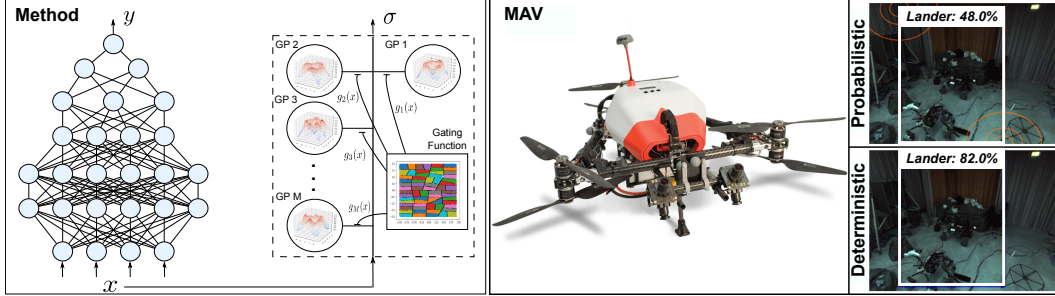


Figure 1: Our method (left) computes uncertainty  $\sigma$  of Neural Network predictions  $y$  using MoE-GPs with the Neural Tangent Kernels (NTK). With this, a MAV [11] (right) performs probabilistic object detection at an interactive frame rate, reducing overconfidence of an object detector. Only a feedforward network is shown but our method also applies to convolution and recurrent DNNs.

dynamics, and (ii) on a Jetson TX2 of a MAV, runs more than 12 times faster than the widely used MC-dropout [3], while performing object detection within a scenario for planetary exploration.

**Contributions** In summary, our main contribution is a novel method for estimating predictive uncertainty of DNNs with sparse GPs (section 3.1), backed up by (i) a theoretical connection between DNNs and MoE-GPs (section 3.2), (ii) a learning algorithm for its applicability in practice (section 3.3), and (iii) an exhaustive empirical evaluation that shows the benefits of our approach (section 4).

## 2 Related Work

Many researchers explored the idea of robotic introspection [2], i.e. robots that reason about “*when they don’t know*” in order to avoid catastrophic failures. So far, many proposed methods used the robots’ past experiences of failures [12, 13, 14], or detect the inconsistencies within the predictors [15, 16]. These works provide strong evidence that introspection improves the reliability of the robots. Other methods are based on uncertainty estimation [1], where learning algorithms such as DNNs use probability theory to express their own reliability [17, 18, 19]. Commonly, these works use a technique called MC-dropout [3]. However, the intense computational burden in estimating uncertainty via sampling or ensembles [20] limits their applicability in real-time systems [11, 21, 22].

Several techniques without sampling or ensembles have been presented for run-time efficiency. These methods often use either Dirichlet distributions or functional Bayesian Neural Networks [23, 24]. Unfortunately, the applicability of Dirichlet distributions is limited to classification, while functional Bayesian Neural Networks face hurdles in scaling to large datasets without approximations [24]. We recommend the recent works of He et al. [25] and Adlam et al. [26], who use NTK-based GPs for uncertainty estimation of DNNs and provide useful insights. However, we focus on a technique that can be readily deployed on a real-time system such as a MAV. Hence, we take a different path to ensembles of infinite width DNNs [25, 26]. We find uncertainty propagation techniques [27, 28, 29] to be practical as these methods are training-free, i.e. assume existing and already well-trained DNNs.

In addition, our theoretic contribution (section 3.2) extends on how DNNs are related to GPs (pioneered by Neal [30] and others [31, 10, 32]), where we derive a relationship between DNNs and MoE-GPs as opposed to a single GP. Theoretic insights that connect DNNs and sparse GPs may pave the way towards their application, as GPs alone do not scale to big data required by DNNs. Lastly, we extend previous works on MoE-GPs (section 3.3). MoE-GPs are efficient but suffer from non-smooth uncertainty estimates [33]. Park and Apley [34] show an elegant solution called patchwork kriging. By augmenting data at the boundaries between GP experts, the experts are forced to produce similar predictions at their boundaries. Yet, as patchwork kriging is limited to low-dimensional inputs [34], we propose several techniques that extend it to higher dimensions, like images (assuming the NTK).

## 3 Materials and Methods

We describe a method that addresses the problem of quantifying uncertainty estimates in DNN predictions, without sampling or ensembles. First, we introduce the main concept of estimating predictive uncertainty with sparse GPs, followed by theoretic results and our learning algorithm.

### 3.1 Main Idea: Fast Uncertainty Estimation with Sparse Gaussian Processes

Figure 1 visualizes our approach. The proposed predictive model uses an existing, and already well-trained DNN for accurate predictions. At the same time, the method offers the possibility to obtain the predictors’ reliable uncertainty in a closed form solution, using a DNN-based MoE-GPs. Intuitively, we can linearize and transform DNNs around a mode to obtain a linear model. As any linear models are GP regressors from a function space view, we can get a GP-based representation of DNNs with a DNN-based kernel, the NTK. We describe below how the obtained GPs can estimate the uncertainty of DNNs, but cannot replace the DNN predictions. Moreover, as full GPs do not scale to big data, we choose to use MoE-GPs, which divides the data into smaller subsets and form many smaller GPs per these subsets. This results in the proposed combination of DNNs and sparse GPs.

**Fundamentals** First, we introduce our notation. Considering a supervised learning task on input-output pairs  $\mathcal{D} = \{\mathcal{X}, \mathcal{Y}\} = \{(x_i, y_i)\}_{i=1}^N$ , where  $x_i \in \mathbb{R}^D$ ,  $y_i \in \mathbb{R}^K$ , we describe a DNN as a parametrized function  $f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^K$ , where  $\theta \in \mathbb{R}^P$ . Here, learning typically seeks to obtain an empirical risk minimizer of the loss function, i.e.  $\min_\theta \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \mathcal{L}(f_\theta(x), y) + \frac{\delta}{2} \theta^T \theta$  where  $\delta$  is an  $L_2$ -regularizer, and mini-batches  $\mathcal{B} \subset \mathcal{D}$  are used to find a local maximum-a-posteriori (MAP) solution  $\hat{\theta}$ . We assume a twice differentiable and strictly convex loss function  $\mathcal{L}$ , e.g. mean squared error (MSE), and piece-wise linear activations in  $f_\theta$  (e.g. RELU). For a clear exposition, we drop the indices  $i$ .

To set the scene for the paper, the Neural Linear Models (NLMs) [27, 35, 36, 32] are defined, which can estimate a DNNs’ predictive uncertainty. To do so, consider a transformed dataset  $\tilde{\mathcal{D}} = \{\mathcal{X}, \tilde{\mathcal{Y}}\}$  with the pseudo-output  $\tilde{y} := J_f(x)\hat{\theta} - H_{\mathcal{L}}(x, y)^{-1} R_{\mathcal{L}}(x, y)$ . This transformation uses the model derivatives of the underlying DNN around  $\hat{\theta}$ , namely the Jacobian  $J_f(x) := \partial f_\theta(x) / \partial \theta^T \in \mathbb{R}^{K \times P}$ , the Hessian  $H_{\mathcal{L}}(x, y) := \partial^2 \mathcal{L}(f_\theta(x), y) / \partial f_\theta(x)^T \partial f_\theta(x) \in \mathbb{R}^{K \times K}$  and the residuals  $R_{\mathcal{L}}(x, y) := \partial \mathcal{L}(f_\theta(x), y) / \partial f_\theta(x) \in \mathbb{R}^K$ . Assuming a white noise and an uninformative prior, the NLMs are:

$$\tilde{y} = J_f(x)\theta + \epsilon \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, H_{\mathcal{L}}(x, y)^{-1}) \quad \text{and} \quad p(\theta) \sim \mathcal{N}(0, \delta^{-1}I). \quad (1)$$

The NLMs can be thought as a Bayesian linear model with learned features from a DNN, which is obtained via a linearization around the DNNs’ last layer [27]. Thus, the isotropic prior is governed by the parameter  $\delta$  that corresponds to the  $L_2$  regularizer [37] while the white noise is governed by the term  $H_{\mathcal{L}}^{-1} R_{\mathcal{L}}$  in the described data transformation, given that the DNN fit the train data well, e.g.  $R_{\mathcal{L}} \approx 0$ . The term  $H_{\mathcal{L}}$  is independent of the input if the 2nd derivative of a loss is a constant, e.g. MSE. Moreover, for commonly used loss functions such as cross entropy and MSE, the covariance  $\Sigma$  of DNN predictions  $f_\theta$  is equal to that of the NLMs’  $\tilde{y}$  [27, 32]. This means that for the same input  $x$ , the NLMs can estimate the predictive uncertainty of DNNs. Yet,  $\tilde{y}$  cannot be used to replace  $f_\theta$ , due to their differences in  $f_\theta$  and  $\tilde{y}$ . In the Appendix, we provide concrete examples and their derivations.

**Main idea** Let’s now apply a kernel trick [38] to the NLMs, so that we can obtain an effective combination of DNNs with GPs. To explain, we can further map the NLM to a function space, as opposed to working in the weight space. Doing so, a well known function space formulation of linear models [38] turns the NLM into a GP with the kernel  $K = \frac{1}{\delta} J_f(x)^T J_f(x)$ , which is also known as the *Network Tangent Kernel* (NTK). Thus, GPs with NTK can be used to estimate the predictive uncertainty of DNNs, as their equivalent NLMs can. Here, what motivates the formulation is the idea of deep kernel machines [39], i.e. we learn the kernel representations as oppose to hand designing the features and the kernels, and in our case, the kernel representations are the tangents of DNNs’ Jacobians. Importantly, doing so, we get a training-free combination of DNNs with GPs, as we keep DNN predictions the same as the MAP, while estimating their uncertainty using GPs with NTK.

Unfortunately, such combinations of DNNs with GPs are restricted by a prominent weakness of standard GPs [40]: the cubic time complexity  $O(N^3)$  that grows with the dataset size  $N$ . So, the computational costs are prohibitive for large scale problems that DNNs assume (typically referred to  $N > 10000$  in [40, 38]). While we leave the theoretic motivation to section 3.2, we thus propose MoE-GPs with NTK, in order to advance the scalability of the proposed combination. A MoE-GP consists of  $M$  experts of GPs or learners  $\mathcal{F} = \{\tilde{f}_{\text{GP}_1}, \dots, \tilde{f}_{\text{GP}_M}\}$ , and a gating function  $g : \mathbb{R}^D \rightarrow \Delta^{M-1}$  that maps any input  $x$  to  $g(x) = [g_1(x), \dots, g_M(x)]$  [8, 9]. Each expert of the MoE learns and predicts within a subset of the input domain, and a gating function generates these subsets.

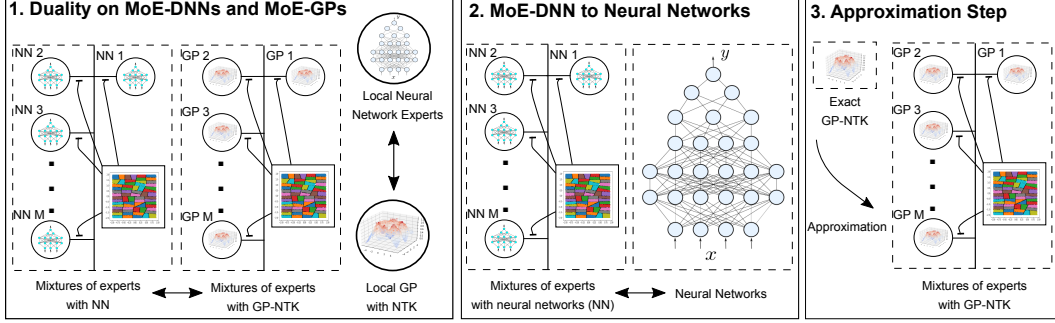


Figure 2: We visualize theoretic results as a roadmap of the proof path. Detailed theorems and their proofs are in the Appendix. (1) Assuming a strict division of data, a duality between mixtures of experts (MoE) with DNNs and GPs are established first. (2) To expand the applicability of the results beyond MoE-DNN, we show that the input-prediction relationships of a single DNN and a MoE-DNN are equivalent with an assumption on the identical local DNN experts. (3) This results in a provable approximation error between the proposed MoE-GPs and the true GPs with the NTK.

Specifically, we choose a gating function  $g_m(\mathbf{x}) = 1$  in just one coordinate for each input [41]. Such strict partition enables the use of a local model for each input [34], avoiding the combinations of multiple predictions of a model for each input, e.g. model averaging as in Meier et al. [42]. The subscripts  $m = 1, 2, \dots, M$  denote the  $m^{\text{th}}$  expert throughout the paper. Then, we write a MoE-GP as,

$$\tilde{\mathbf{y}} = \sum_{m=1}^M g_m(\mathbf{x}) \tilde{f}_{\text{GP}_m}(\mathbf{x}) + \epsilon_m \quad \text{with} \quad \tilde{f}_{\text{GP}_m}(\mathbf{x}) \sim \text{GP}(\mathbf{0}, \frac{1}{\delta_m} \mathbf{J}_{f_m}(\mathbf{x})^T \mathbf{J}_{f_m}(\mathbf{x})). \quad (2)$$

As depicted in figure 1, the gating function  $g_m(\mathbf{x})$  assigns the input data  $\mathbf{x}$  to the  $m^{\text{th}}$  subset, and individual GPs  $f_{\text{GP}_m}(\mathbf{x})$  learn and make predictions for the assigned data within the  $m^{\text{th}}$  subset. Consequently, the generative model and the predictions  $\mathcal{N}(\tilde{\mathbf{y}}_m^*, \Sigma_m(\mathbf{x}^*))$  on the new test datum  $\mathbf{x}^*$  are:

$$\begin{bmatrix} \tilde{\mathbf{y}}_m \\ \tilde{f}_{\text{GP}_m} \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K}_m + \sigma_{0,m} \mathbf{I} & \mathbf{k}_{m,*} \\ \mathbf{k}_{m,*}^T & \mathbf{k}_{m,**} \end{bmatrix}\right) \quad \text{and} \quad \begin{aligned} \tilde{\mathbf{y}}_m^* &= \mathbf{k}_{m,*}^T (\mathbf{K}_m + \sigma_{0,m} \mathbf{I})^{-1} \tilde{\mathbf{y}}_m, \\ \Sigma_m &= \mathbf{k}_{m,**} - \mathbf{k}_{m,*}^T (\mathbf{K}_m + \sigma_{0,m} \mathbf{I})^{-1} \mathbf{k}_{m,*} + \sigma_{0,m}, \end{aligned}$$

where,  $\mathbf{K}_m = \mathbf{K}_m(\mathcal{X}, \mathcal{X})$ ,  $\mathbf{k}_{m,*} = \mathbf{k}_m(\mathcal{X}, \mathbf{x}^*)$  and  $\mathbf{k}_{m,**} = \mathbf{k}_m(\mathbf{x}^*, \mathbf{x}^*)$ . This posterior predictive distribution  $\Sigma_m$  is an indicator of total uncertainty, i.e. the kernel function captures the model uncertainty while the constant term  $\sigma_{0,m} \mathbf{I}$  corresponds to the aleatoric uncertainty [43]. In GPs, as we do not have access to exact function values, the aleatoric uncertainty often relies on the assumption of additive i.i.d white noise, and in MoE-GPs, the terms  $\sigma_{0,m} \mathbf{I}$  are inferred from the data per subset, leading to a non-stationary kernel. This is achieved by optimizing the marginal likelihood [38].

The benefits of the MoE-GP formulation are two-folds. First, as shown above, the generative model of a GP is defined per subset, and therefore, the GP experts are smaller than the one on the entire dataset. This improves the computational complexity of GPs [34]. Second, the covariance computations are in a closed form. Thus, we neither need sampling nor ensemble [3, 4]. Using the Lanczos approximation [44], which approximates the matrix inversion with the multiplications (suited for GPUs), we can compute the uncertainty estimates from GPs, with a constant time complexity [44].

For classification, we leverage the framework of Lu et al. [45], which can perform classification with the regression models via confidence calibration [5]. While we refer to Lu et al. [45] for the details, this steps enables us to obtain the classification uncertainty in closed form, i.e. for a class  $c$ :

$$p(c|z_m) = \text{softmax}\left(\frac{z_m}{T_m}\right) \quad \text{with} \quad T_m = \sqrt{1 + \lambda_{m,0} \Sigma_m(\mathbf{x}^*)}, \quad (3)$$

where  $z_m$  is the activation,  $T_m$  is the temperature scaling factor, which is a function of a scaling constant  $\lambda_{m,0}$ , and GP regression uncertainty  $\Sigma_m(\mathbf{x}^*)$ . Intuitively, high prediction uncertainty will reduce the corresponding  $p(c|z_m)$  by increasing the  $T_m$ , thereby calibrating the confidences [5, 46, 45].

### 3.2 On Theory: Neural Networks as Sparse Gaussian Processes

So far, we have outlined our main idea - using the function space view of NLMs and dividing the input space into smaller subsets, we can form smaller GPs per division. These GPs then estimate the

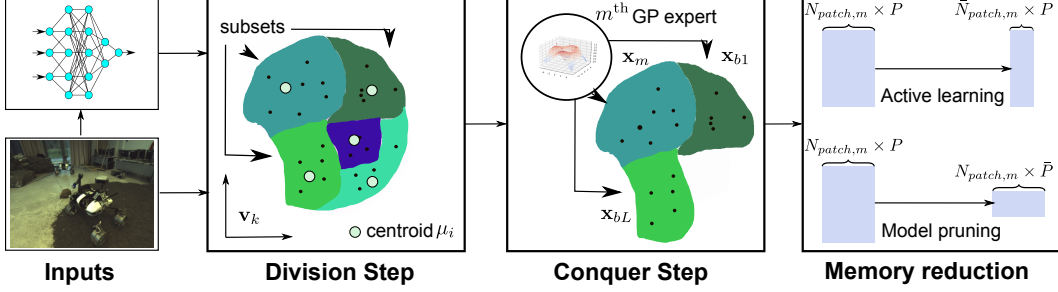


Figure 3: The proposed algorithm involves (i) the division of input data using Jacobians of already trained DNNs, (ii) training of GP experts within their partitions and neighboring regimes, and (iii) combining active learning and model pruning techniques to reduce space complexity.

predictive uncertainty of DNNs. Now, we discuss the theoretical foundations on how DNNs can be cast as MoE-GPs with NTK. Due to space constraints, we refer the reader to the Appendix for in depth treatment, where we provide various background materials, our theorems and their proofs. In section 3.3 we explain our learning algorithm while here, we briefly summarize the main results.

**Main result** Figure 2 summarizes our findings. Consider a mixtures of experts model, where the experts are DNNs and a gating function strictly divides the input space, i.e only one local DNN expert per division (MoE-DNN). Then, instead of the MAP parameter estimates  $\hat{\theta}$ , we consider Bayesian DNNs with Gaussian distributions  $p(\theta|\mathcal{D})$ , i.e representing DNNs as probability distributions over their parameters (obtained using [47, 48, 49, 50]). As a first step, as a specific instance of Khan et al. [32], we show that the DNN experts have mathematically equivalent Gaussian distributions with GPs using NTK. We further show that MoE-DNNs and MoE-GPs have equivalent Gaussian distributions, establishing the duality of the two models in a Bayesian sense. Main insight to the former is the probabilistic independence between the local experts due to a strict division of the input space.

Next, we establish a connection between a DNN, and a MoE-GP as shown in figure 2, in order to increase the applicability of our theoretic results. To do so, we point out that the input-prediction relationships of a single DNN and a MoE-DNN are equivalent, if all the local DNN experts are the same as the single DNN. Under the given conditions, we derive that a single DNN can be cast as a MoE-GP, with an assumption that the data is stationary. The resulting step yields the approximation error of  $\|\mathbf{K}(\mathcal{X}, \mathcal{X}) - \mathbf{K}_{true}(\mathcal{X}, \mathcal{X})\|_F^2 = \sum_{ij} \mathbf{K}(\mathbf{x}_i, \mathbf{x}_j)^2 - \sum_{m=1}^M \sum_{ij \in \mathcal{B}_m} \mathbf{K}(\mathbf{x}_i, \mathbf{x}_j)^2$  where  $\mathbf{K}_1 = \mathbf{K}_2 = \dots = \mathbf{K}_M$  for all  $M$  GP experts (while it is still possible to keep different hyper-parameters  $\delta_m$  and  $\sigma_m$ ), and  $\mathbf{K}_{true}$  is the kernel matrix of the true DNN-equivalent GP with NTK. This means that less approximation error occurs when less correlated data points are assumed to be independent by MoE-GP, while strongly correlated points by NTK are within the same GP experts.

### 3.3 From Theory to Practice: A Practical Learning Algorithm

Now, we attempt to bring our theoretical results into a practical tool for robotics. As shown in figure 3 our solution involves the division of data into several partitions. Then, in a conquer step, we train GP experts within the partitions. Moreover, we propose concepts to reduce the memory complexity.

**Division step** We design a gating function that divides the input data into smaller regimes s.t. highly correlated data points are grouped together whilst the points with weak correlations are separated apart. To do so, we perform NTK PCA, i.e. kernel PCA that uses NTK to reduce the dimensions of the data. Then, K-means is applied to form the clusters. Concretely, NTK PCA reduces the input data  $\mathbf{x}_i$  for  $i = 1, \dots, N$  into the low dimensional principal components  $\mathbf{v}_k = \sum_{i=1}^N \alpha_{ik} \mathbf{J}_f(\mathbf{x})^T \mathbf{J}_f(\mathbf{x}_i)$ , where  $\alpha$  is the eigenvector of  $\widetilde{\mathbf{K}} \alpha_k = \lambda_k \alpha_k$  with a centered kernel matrix  $\widetilde{\mathbf{K}} = \mathbf{K} - 2\mathbf{1}_N \mathbf{K} + \mathbf{1}_N \mathbf{K} \mathbf{1}_N$  for the matrix of ones  $\mathbf{1}_N \in \mathbb{R}^{N \times N}$ . Then, the centroid  $\mu_i$  and division labels  $c$  are computed, iterating:

$$c_i = \arg \min \sum_i^M \sum_v \|v_i - \mu_i\|^2 \quad \text{and} \quad \mu_i = \frac{\sum_i^M \mathbf{1}\{c_i = i\} v_i}{\sum_i^M \mathbf{1}\{c_i = i\}}. \quad (4)$$

The proposed technique gives a solution for Kernel K-means [51], which minimizes the derived error bound  $\|\mathbf{K}(\mathcal{X}, \mathcal{X}) - \mathbf{K}_{true}(\mathcal{X}, \mathcal{X})\|_F^2$  (section 3.2) with a balancing normalization [52]. As kernel PCA

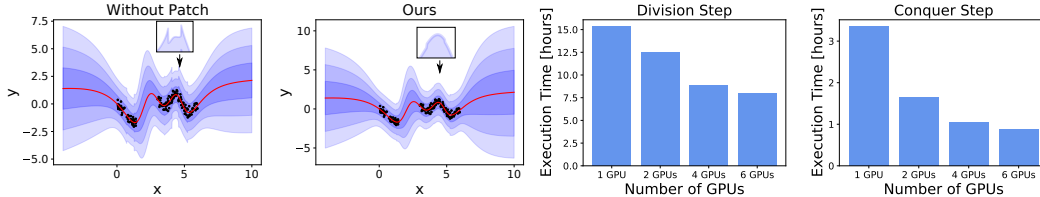


Figure 4: (Left) We visualize predictive uncertainty of DNNs with and without the proposed patchwork prior. The black dots are train data points, and the red line shows the predictions of a deterministic network. Blue shades show up to three standard deviations. Without the patch prior, sharp peaks are observed in uncertainty estimates. (Right) Training time for approximately 2 million data points is shown for a varying the number of GPUs. We learned 512 GP experts per joint torques.

does not scale to large datasets, we use 2-step kernel K-means [53], which uses randomly sampled and less data to compute the cluster centroids. We refer to Chitta et al. [53] for more details.

**Conquer step** Next, we form the individual GP experts per divided local regime. A naive strategy is to form a single GP per regime. Unfortunately, such a model suffers from discontinuity in predictions at the boundaries [33]. For example, input points at the cluster boundaries may yield different predictions from the surrounding GPs, causing sharp peaks in predictions. Therefore, we propose:

$$\mathbf{K}_{m,\text{patch}} = \left[ \mathbf{J}_{f_m}(\mathbf{x}_m), \mathbf{J}_{f_{m b_1}}(\mathbf{x}_{b_1}), \dots, \mathbf{J}_{f_{m b_L}}(\mathbf{x}_{b_L}) \right]^T \left[ \mathbf{J}_{f_m}(\mathbf{x}_m), \mathbf{J}_{f_{m b_1}}(\mathbf{x}_{b_1}), \dots, \mathbf{J}_{f_{m b_L}}(\mathbf{x}_{b_L}) \right] \quad , \quad (5)$$

for the  $m^{\text{th}}$  GP expert. We include the neighboring GP experts  $\mathbf{J}_{f_{m b_1}}, \dots, \mathbf{J}_{f_{m b_L}}$  into the prior of the  $m^{\text{th}}$  expert. This includes the information of neighboring GPs into the prior of the  $m^{\text{th}}$  GP expert. Intuitively, the discontinuities occur at the shared boundary regimes of two GP experts. Our GP prior removes such boundary regimes locally by taking into account the neighboring regimes.

**Active uncertainty learning** The complexity of each GP expert is bounded to  $O(N_{\text{patch},m}^3)$ , where the number of data points  $N_{\text{patch},m}$  includes the data points of neighboring GPs:  $N_{\text{patch},m} = N_m + N_{b_1} + \dots + N_{b_k}$ , resulting in added costs  $N_m < N_{\text{patch},m}$ . To mitigate, we perform active learning on the neighboring GP experts, and thus choose fewer, but the most informative points, e.g. IVMs [54, 55]. Intuitively, as the neighboring data are only included for reducing the discontinuity problem, we may select fewer data points. Concretely, we use the following steps: (i) for the neighboring GPs of the  $m^{\text{th}}$  expert, we draw an initial subset and train the GPs. (ii) Using the trained GP models, we rank the remaining data (stored as a pool) by uncertainty (queries generated by uncertainty sampling). (iii) the GP expert is then updated with the most uncertain point. These steps repeat until  $N_{b_k}$  is reduced to a desired, smaller value. In this way, while forming a patchwork prior, the neighboring GPs actively choose the neighboring data they want to learn from. As a result, we obtain:  $\bar{N}_{\text{patch},m} < N_{\text{patch},m}$ .

**Gaussian process compression** We further reduce the complexity of the algorithm by exploiting model compression techniques. Note that the Jacobians  $\mathbf{J}_f(\mathbf{x})$  are  $N_{\text{patch},m}$  by  $P$  matrices where  $P$  is the total number of parameters in DNNs. As  $\mathbf{J}_f(\mathbf{x})$  represents the sensitivity of each parameter to the output, the elements of  $\mathbf{J}_f(\mathbf{x})$  that belong to an unimportant DNN parameter can be removed. To do so, we propose a two-staged compression technique: (i) rank the DNN parameters by their importance using existing pruning methods and remove the corresponding elements of  $\mathbf{J}_f(\mathbf{x})$  for all  $m$ . (ii) Per GP expert, rank the elements of  $\mathbf{J}_{f_m}(\mathbf{x})$  by a metric  $|\sum \mathbf{J}_{f_m}(\mathbf{x})|$ , as smaller absolute values contribute less to the kernel. The first step is targeted at removing redundant parameters in DNNs (for the division step) while the second step is targeted at each individual expert, resulting in  $\bar{P} < P$ .

## 4 Experiments and Evaluations

We provide 5 sets of experiments to not only validate the method, but also to show the benefits of the proposed formulation namely performance, scalability and run-time. Implementation details can be found in the Appendix, and importantly, the accompanying video shows the experiments with a MAV.

**Toy illustration** With a toy regression on the Snelson dataset, we show (i) on how MoE-GPs can estimate the predictive uncertainty of DNNs, and (ii) the ability of the patchwork prior to produce smoother uncertainty estimates. For this, we train a single hidden layer Multi-layer perceptron

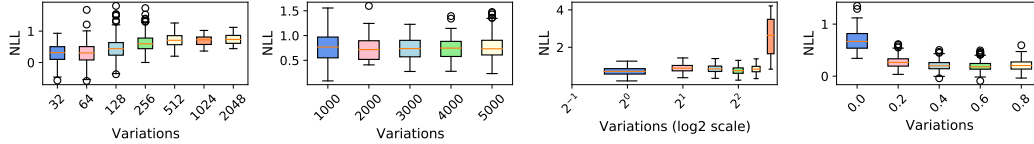


Figure 5: The effects of each hyperparameter is shown with normalized test NLL, by varying them in different steps (Variations), and fixing the others to default settings. Lower the better.

Table 1: The results of inverse dynamic experiments are reported using NLL. Lower the better.

Train	Test	Ensemble	MC-dropout	rMC-dropout	Laplace	rLaplace	Ours
sarcos	sarcos	$1.261 \pm 0.105$	$1.398 \pm 0.012$	$1.398 \pm 0.012$	$1.392 \pm 0.014$	$1.392 \pm 0.014$	<b><math>1.035 \pm 0.039</math></b>
sarcos	kuka1	$17.86 \pm 2.995$	<b><math>16.21 \pm 0.866</math></b>	$20.58 \pm 1.144$	$25.88 \pm 1.354$	$24.55 \pm 2.746$	$21.53 \pm 2.268$
sarcos	kuka2	$17.50 \pm 2.895$	<b><math>15.63 \pm 0.858</math></b>	$20.17 \pm 0.988$	$25.59 \pm 1.079$	$24.42 \pm 2.257$	$20.92 \pm 2.236$
sarcos	kukasim	<b><math>23.06 \pm 2.649</math></b>	$51.09 \pm 11.14$	$61.40 \pm 12.921$	$77.13 \pm 15.40$	$74.03 \pm 13.23$	$68.95 \pm 4.003$
kuka1	kuka1	$2.013 \pm 0.020$	$1.611 \pm 0.007$	$1.620 \pm 0.006$	$1.676 \pm 0.013$	$1.719 \pm 0.071$	<b><math>1.347 \pm 0.013</math></b>
kuka1	kuka2	$2.085 \pm 0.005$	$1.349 \pm 0.019$	$1.330 \pm 0.006$	<b><math>1.310 \pm 0.005</math></b>	<b><math>1.310 \pm 0.008</math></b>	$1.315 \pm 0.003$
kuka1	kukasim	$60.44 \pm 1.108$	$42.14 \pm 2.869$	$48.76 \pm 0.566$	$60.37 \pm 0.799$	$59.74 \pm 1.571$	<b><math>1.348 \pm 0.012</math></b>
kuka1	sarcos	$8.128 \pm 0.169$	$22.24 \pm 4.288$	$42.36 \pm 2.398$	$122.68 \pm 13.32$	$1837 \pm 2431.3$	<b><math>1.356 \pm 0.014</math></b>
kuka2	kuka2	$2.042 \pm 0.009$	$1.443 \pm 0.005$	$1.423 \pm 0.006$	$1.429 \pm 0.006$	$1.423 \pm 0.006$	<b><math>1.354 \pm 0.013</math></b>
kuka2	kukasim	$63.07 \pm 0.741$	$38.94 \pm 0.561$	$48.93 \pm 0.779$	$60.36 \pm 1.005$	$59.05 \pm 1.158$	<b><math>1.333 \pm 0.012</math></b>
kuka2	sarcos	$8.509 \pm 0.461$	$18.24 \pm 1.599$	$53.24 \pm 6.287$	$141.7 \pm 17.14$	$211.3 \pm 180.2$	<b><math>1.355 \pm 0.014</math></b>
kuka2	kuka1	$2.106 \pm 0.010$	$1.461 \pm 0.004$	$1.395 \pm 0.004$	$1.373 \pm 0.005$	$1.374 \pm 0.004$	<b><math>1.331 \pm 0.013</math></b>

(MLP) with 200 units and a tanh activation. Evaluating for out-of-distribution (OOD) data (points far from train data), and domain-shift (DS) data (removed in-between points), as shown in figure 4, the produced uncertainty estimates are high as test points move away from the train data. Moreover, when DNN predictions match the train data, the produced uncertainty estimates are close to the data noise. We also observe that the patchwork prior mitigates the discontinuity problem of local GPs.

**On hyperparameters** Next, we examine (i) the influences of hyperparameter choices, and (ii) provide a simple tuning recipe. For this, we examine a regression task using a 5-layered, 200 units MLP on SARCOS robot arm data-set [38]. Here, we vary each hyperparameter while fixing the remaining ones, and we compute the Negative Log Likelihood (NLL). The results are shown in figure 5, with 4 hyperparameters namely, (i) the number of GP experts, (ii) the subset sizes for clustering, (iii) the pruning level, and (v) the size of the informative data points. We observe the following: (i) the NLL is proportional to the number of GP experts, (ii) the subset size do not significantly influence the NLL, (iii) pruning steps have a tipping point where the NLL increases, and (iv) the active learning parameter has also a tipping point, where the NLL decrease is marginal. In contrary, by the design, the computational complexity increases with (i) lower number of experts and (ii) the size of actively selected points, and decreases with the pruning steps. We find these results confirm our intuition. For example, keeping 2048 GP experts for SARCOS results in 21 data-points per experts on average. This setting is inferior to keeping 32 GP experts, which has 1390 data points on average. Reflecting this, our strategy is to assign more data points for each GP expert, while selecting only the needed active learning points. The pruning levels can then be varied within a range that does not deteriorate the performance in order to reduce the computational complexity of the method.

**Comparison study** Now, we evaluate the performance of our method using the datasets namely SARCOS, KUKA1 and KUKA2 [42]. The baselines are MC-dropout [3], Laplace Approximation [47, 49], and their fast variants: rMC-dropout [28] and rLaplace [27, 29, 49]. The fast variants are our main competitors: principled Bayesian approaches that are training-free, i.e. works with pre-trained DNNs without modifications to the training procedures. We find decoupling DNN training to uncertainty estimates crucial, as it ensures comparisons of uncertainty estimates only by keeping the accuracy of the predictors constant amongst the baselines (in this case, the same pretrained 5-layered MLP across the baselines). We also include the ensembles [4]. Lastly, we adopt zero-shot cross-dataset transfer to systematically evaluate in-domain, OOD and DS scenarios, following the recent insights, that is, we train on a dataset, and evaluate uncertainty estimates for all other datasets.

The results are in table 1, where we averaged over 3 random initializations and report NLL as a metric (a standard for regression tasks). We find that our method often outperforms other baselines

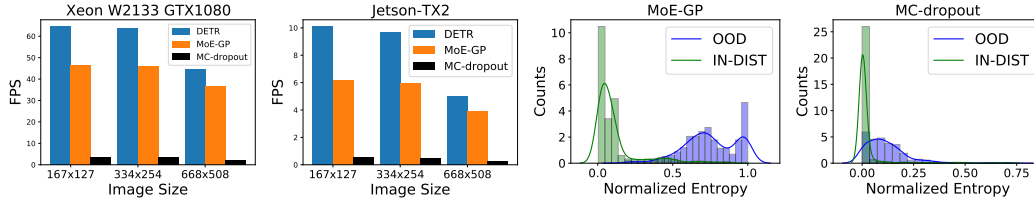


Figure 6: The run-time analysis (left) shows efficiency of our approach when compared to MC-dropout [3] with 20 samples. Any difference to a deterministic DNN can also be seen. The entropy histograms (right) show a novelty detection performance. The memory requirement is 8GB.

significantly. MC-dropout dominates in 3 experiments. Yet, MC-dropout requires sampling for uncertainty estimates - one of the core problems addressed in our work. We note again, that all the compared methods are training-free, making the comparisons meaningful. Importantly, rLaplace uses the NLM from section 3 to infer Gaussian uncertainty in weight space, in contrast to function space inference with GPs. Thus, the comparison of rLaplace and ours instantiates the comparison of inference between weight space and function space views. We believe that our results show the relevance of GP based uncertainty estimates for DNNs over uncertainty propagation methods.

**On scalability** Here, we show that (i) our approach scales to KUKA SIM dataset (contains 1984950 data points), and (ii) supports distributed training. Scalability is relevant, as DNNs operate in the regime of big data, and an exact GP does not scale to such settings. Moreover, a key benefit of MoE-GP over other sparse GPs lies in its distributed, local nature, and we can exploit it in practice. The results are depicted in figure 4, where we plot the training time in hours against the number of used GPUs. Using 512 GP experts and 100 iterations for MLL optimization, we find that our method scales to approximately 2 million data points within a day, and less than 8 hours when using 6 GPUs. This experiment shows an empirical evidence that our method scales to 2 million data points for the task of inverse dynamic learning, and validates the claim that MoE-GPs support distributed training.

**On a real robot** To validate the applicability of our method to a robot and evaluate the run-time on the hardware, we perform experiments on a MAV within the context of an on-going space demo mission for future planetary exploration [56]. For this, we train an object detector, DETR [57] with EfficientNet backbone. For testing, we create two scenarios: (i) a carry test, used for training and testing with 1008 manual labels, and (ii) flight tests with different configurations of the objects (e.g. space rover and lander). The later creates OOD samples, with a slight domain shift in data distribution. In field robotics, the assumption that the test set comes from the same training data distribution is routinely violated due to the changes in the environments, and we attempt to create such effects. Moreover, we also evaluate the complexity of our method on a Jetson-TX2, which is used on our MAV. The quantitative results are shown in figure 6, which shows that (i) our method can be fast on a Jetson-TX2, and (ii) can clearly separate OOD samples within this scenario. The accompanying video explains this setup and the qualitative results. More details can be found in the Appendix.

## 5 Discussion

If the computational complexity of GPs can be tamed in practice, our work shows that the predictive uncertainty of DNNs can be obtained from the GP formulation. By advancing the scalability of a full NTK, the advantages of our approach are demonstrated. Yet, the applicability to larger data regimes, e.g. imagenet, is confined to that of the sparse GPs, e.g. the non-parametric models require an access to the training data at run-time, which can require more memory than the parametric models such as DNNs. Sparse GPs may also face struggles when the output dimension is larger. The alternatives can be the weight space-based methods [58] such as Sharma et al. [59]. Here, various approximations have been devised so far, in order to cope with the high dimensional weight space. For robotics though, when the availability of data is limited, the dimensions of data space can be smaller than the weight space, and thus, our work can provide a reference that the Bayesian non-parametric can be a relevant tool in addressing the current challenges of the uncertainty estimation for neural networks.

In future works, we plan to apply the proposed methodology in a field robotics setting, and study how the uncertainty estimates can be used to improve the robustness of the robotic systems.



## Acknowledgments

We thank the anonymous reviewers and area chairs for their thoughtful feedback. Special thanks to Klaus Strobl, Maximilian Denninger and Antonin Raffin for proof reading the paper, and also Seok Joon Kim for the support during his internship at DLR. We also would like to acknowledge the support of Helmholtz Association, the project ARCHES (contract number ZT-0033), the Initiative and Networking Fund (INF) under the Helmholtz AI platform grant agreement (ID ZT-I-PF-5-1) and lastly, the EU Horizon 2020 project RIMA under the grant agreement number 824990. Jianxiang Feng is supported by the Munich School for Data Science (MUDS) and Rudolph Triebel is a member of MUDS.

## References

- [1] S. Thrun. Probabilistic robotics. *Communications of the ACM*, 45(3):52–57, 2002.
- [2] H. Grimmitt, R. Triebel, R. Paul, and I. Posner. Introspective Classification for Robot Perception. *International Journal of Robotics Research (IJRR)*, 2015.
- [3] Y. Gal and Z. Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059, 2016.
- [4] B. Lakshminarayanan, A. Pritzel, and C. Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in neural information processing systems*, pages 6402–6413, 2017.
- [5] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On calibration of modern neural networks. In *International Conference on Machine Learning*, pages 1321–1330, 2017.
- [6] S. Sun, G. Zhang, J. Shi, and R. Grosse. Functional variational bayesian neural networks. In *International Conference on Learning Representations*, 2019.
- [7] R. Paul, R. Triebel, D. Rus, and P. Newman. Semantic categorization of outdoor scenes with uncertainty estimates using multi-class gaussian process classification. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 2404–2410. IEEE, 2012.
- [8] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, and G. E. Hinton. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991.
- [9] V. Tresp. Mixtures of gaussian processes. In *Advances in neural information processing systems*, pages 654–660, 2001.
- [10] A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018.
- [11] P. Lutz, M. G. Müller, M. Maier, S. Stoneman, T. Tomić, I. von Bargaen, M. J. Schuster, F. Steidle, A. Wedler, W. Stürzl, and R. Triebel. ARDEA - An MAV with skills for future planetary missions. *Journal of Field Robotics (JFR)*, 2019.
- [12] S. Daftry, S. Zeng, J. A. Bagnell, and M. Hebert. Introspective perception: Learning to predict failures in vision systems. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1743–1750. IEEE, 2016.
- [13] M. S. Ramanagopal, C. Anderson, R. Vasudevan, and M. Johnson-Roberson. Failing to learn: autonomously identifying perception failures for self-driving cars. *IEEE Robotics and Automation Letters*, 3(4):3860–3867, 2018.
- [14] C. Gurău, D. Rao, C. H. Tong, and I. Posner. Learn from experience: probabilistic prediction of perception performance to avoid failure. *The International Journal of Robotics Research*, 37(9): 981–995, 2018.
- [15] C. Richter and N. Roy. Safe visual navigation via deep learning and novelty detection. In *Robotics, Science and Systems (RSS)*, 2017.

- [16] K. Wong, S. Wang, M. Ren, M. Liang, and R. Urtasun. Identifying unknown instances for autonomous driving. In L. P. Kaelbling, D. Kragic, and K. Sugiura, editors, *Proceedings of the Conference on Robot Learning*, volume 100 of *Proceedings of Machine Learning Research*, pages 384–393. PMLR, 30 Oct–01 Nov 2020.
- [17] J. Feng, M. Durner, Z.-C. Marton, F. Balint-Benczedi, and R. Triebel. Introspective robot perception using smoothed predictions from bayesian neural networks. In *2019 International Symposium on Robotics Research (ISRR)*, 2019.
- [18] A. Harakeh, M. Smart, and S. L. Waslander. Bayesod: A bayesian approach for uncertainty estimation in deep object detectors. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 87–93. IEEE, 2020.
- [19] N. Sünderhauf, F. Dayoub, D. Hall, J. Skinner, H. Zhang, G. Carneiro, and P. Corke. A probabilistic challenge for object detection. *Nature Machine Intelligence*, 1(9):443–443, 2019.
- [20] J. Gawlikowski, C. R. N. Tassi, M. Ali, J. Lee, M. Humt, J. Feng, A. Kruspe, R. Triebel, P. Jung, R. Roscher, et al. A survey of uncertainty in deep neural networks. *arXiv preprint arXiv:2107.03342*, 2021.
- [21] J. Lee, R. Balachandran, Y. S. Sarkisov, M. De Stefano, A. Coelho, K. Shinde, M. J. Kim, R. Triebel, and K. Kondak. Visual-inertial telepresence for aerial manipulation. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1222–1229. IEEE, 2020.
- [22] J. Lee, T. Muskardin, C. R. Pacz, P. Oettershagen, T. Stastny, I. Sa, R. Siegwart, and K. Kondak. Towards autonomous stratospheric flight: A generic global system identification framework for fixed-wing platforms. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 6233–6240. IEEE, 2018.
- [23] M. Sensoy, L. Kaplan, and M. Kandemir. Evidential deep learning to quantify classification uncertainty. In *Advances in Neural Information Processing Systems*, pages 3179–3189, 2018.
- [24] E. D. Carvalho, R. Clark, A. Nicastro, and P. H. Kelly. Scalable uncertainty for computer vision with functional variational inference. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12003–12013, 2020.
- [25] B. He, B. Lakshminarayanan, and Y. W. Teh. Bayesian deep ensembles via the neural tangent kernel. *arXiv preprint arXiv:2007.05864*, 2020.
- [26] B. Adlam, J. Lee, L. Xiao, J. Pennington, and J. Snoek. Exploring the uncertainty properties of neural networks’ implicit priors in the infinite-width limit. In *International Conference on Learning Representations*, 2021.
- [27] D. J. MacKay. Information-based objective functions for active data selection. *Neural computation*, 4(4):590–604, 1992.
- [28] J. Postels, F. Ferroni, H. Coskun, N. Navab, and F. Tombari. Sampling-free epistemic uncertainty estimation using approximated variance propagation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2931–2940, 2019.
- [29] A. Y. Foong, Y. Li, J. M. Hernández-Lobato, and R. E. Turner. ‘in-between’ uncertainty in bayesian neural networks. *arXiv preprint arXiv:1906.11537*, 2019.
- [30] R. M. Neal. Priors for infinite networks. In *Bayesian Learning for Neural Networks*, pages 29–53. Springer, 1996.
- [31] J. Lee, Y. Bahri, R. Novak, S. Schoenholz, J. Pennington, and J. Sohl-dickstein. Deep neural networks as gaussian processes. In *International Conference on Learning Representations*, 2018.
- [32] M. E. E. Khan, A. Immer, E. Abedi, and M. Korzepa. Approximate inference turns deep networks into gaussian processes. In *Advances in Neural Information Processing Systems*, pages 3088–3098, 2019.

- [33] C. Park and J. Z. Huang. Efficient Computation of Gaussian Process Regression for Large Spatial Data Sets by Patching Local Gaussian Processes. *Journal of Machine Learning Research*, 17(174):1–29, 2016.
- [34] C. Park and D. Apley. Patchwork kriging for large-scale gaussian process regression. *The Journal of Machine Learning Research*, 19(1):269–311, 2018.
- [35] S. W. Ober and C. E. Rasmussen. Benchmarking the neural linear model for regression. *Third Symposium on Advances in Approximate Bayesian Inference*, 2020.
- [36] J. Watson, J. A. Lin, P. Klink, and J. Peters. Neural linear models with functional gaussian process priors. In *Third Symposium on Advances in Approximate Bayesian Inference*, 2020.
- [37] C. M. Bishop. *Pattern recognition and machine learning*. springer, 2006.
- [38] C. Rasmussen and C. Williams. *Gaussian Processes for Machine Learning*. Adaptive Computation and Machine Learning. MIT Press, Cambridge, MA, USA, Jan. 2006.
- [39] A. G. Wilson, Z. Hu, R. Salakhutdinov, and E. P. Xing. Deep kernel learning. In *Artificial intelligence and statistics*, pages 370–378. PMLR, 2016.
- [40] H. Liu, Y.-S. Ong, X. Shen, and J. Cai. When gaussian process meets big data: A review of scalable gps. *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [41] S. Gross, M. Ranzato, and A. Szlam. Hard mixtures of experts for large scale weakly supervised vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6865–6873, 2017.
- [42] F. Meier, P. Hennig, and S. Schaal. Incremental local gaussian regression. In *Advances in Neural Information Processing Systems*, pages 972–980, 2014.
- [43] E. Hüllermeier and W. Waegeman. Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods. *Machine Learning*, 110(3):457–506, 2021.
- [44] G. Pleiss, J. Gardner, K. Weinberger, and A. G. Wilson. Constant-time predictive distributions for gaussian processes. In *International Conference on Machine Learning*, pages 4114–4123. PMLR, 2018.
- [45] Z. Lu, E. Ie, and F. Sha. Uncertainty estimation with infinitesimal jackknife, its distribution and mean-field approximation. *CoRR*, abs/2006.07584, 2020.
- [46] J. Wenger, H. Kjellström, and R. Triebel. Non-parametric calibration for classification. In *International Conference on Artificial Intelligence and Statistics*, pages 178–190. PMLR, 2020.
- [47] H. Ritter, A. Botev, and D. Barber. A scalable laplace approximation for neural networks. In *6th International Conference on Learning Representations, ICLR 2018-Conference Track Proceedings*, volume 6. International Conference on Representation Learning, 2018.
- [48] M. Humt, J. Lee, and R. Triebel. Bayesian optimization meets laplace approximation for robotic introspection. *arXiv preprint arXiv:2010.16141*, 2020.
- [49] J. Lee, M. Humt, J. Feng, and R. Triebel. Estimating model uncertainty of neural networks in sparse information form. In H. D. III and A. Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 5702–5713. PMLR, 13–18 Jul 2020.
- [50] K. Shinde, J. Lee, M. Humt, A. Sezgin, and R. Triebel. Learning multiplicative interactions with bayesian neural networks for visual-inertial odometry. *arXiv preprint arXiv:2007.07630*, 2020.
- [51] C. Ding and X. He. K-means clustering via principal component analysis. In *Proceedings of the twenty-first international conference on Machine learning*, page 29, 2004.

- [52] I. S. Dhillon, Y. Guan, and B. Kulis. Kernel k-means: spectral clustering and normalized cuts. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 551–556, 2004.
- [53] R. Chitta, R. Jin, T. C. Havens, and A. K. Jain. Approximate kernel k-means: Solution to large scale kernel clustering. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 895–903, 2011.
- [54] N. Lawrence, M. Seeger, and R. Herbrich. Fast sparse gaussian process methods: The informative vector machine. In *Proceedings of the 16th annual conference on neural information processing systems*, number CONF, pages 609–616, 2003.
- [55] R. Triebel, H. Grimmert, R. Paul, and I. Posner. Driven learning for driving: How introspection improves semantic mapping. In *Robotics Research*, pages 449–465. Springer, 2016.
- [56] M. J. Schuster, M. G. Müller, S. G. Brunner, H. Lehner, P. Lehner, R. Sakagami, A. Dömel, L. Meyer, B. Vodermayr, R. Giubilato, et al. The arches space-analogue demonstration mission: towards heterogeneous teams of autonomous robots for collaborative scientific sampling in planetary exploration. *IEEE Robotics and Automation Letters*, 5(4):5315–5322, 2020.
- [57] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, pages 213–229. Springer, 2020.
- [58] D. Madras, J. Atwood, and A. D’Amour. Detecting extrapolation with local ensembles. *arXiv preprint arXiv:1910.09573*, 2019.
- [59] A. Sharma, N. Azizan, and M. Pavone. Sketching curvature for efficient out-of-distribution detection for deep neural networks. *The Conference on Uncertainty in Artificial Intelligence (UAI)*, 2021.