

---

# Meta-Learning to Perform Bayesian Inference in a single Forward Propagation

---

**Samuel G. Müller**  
University of Freiburg  
muellesa@cs.uni-freiburg.de

**Noah Hollmann**  
Charité Berlin  
noah.hollmann@charite.de

**Sebastian Pineda Arango**  
University of Freiburg  
pineda@cs.uni-freiburg.de

**Josif Grabocka**  
University of Freiburg  
grabocka@informatik.uni-freiburg.de

**Frank Hutter**  
University of Freiburg &  
Bosch Center for Artificial Intelligence, Germany  
fh@cs.uni-freiburg.de

## Abstract

Currently, it is hard to reap the benefits of deep learning for Bayesian methods. We present *Prior-Data Fitted Networks (PFNs)*, a method that allows to employ large-scale machine learning techniques to approximate a large set of posteriors. The only requirement for PFNs is the ability to sample from a prior distribution over supervised learning tasks (or functions). The method repeatedly draws a task (or function) from this prior, draws a set of data points and their labels from it, masks one of the labels and learns to make probabilistic predictions for it based on the set-valued input of the rest of the data points. Presented with samples from a new supervised learning task as input, it can then make probabilistic predictions for arbitrary other data points in a single forward propagation, effectively having learned to perform Bayesian inference. We demonstrate that PFNs can near-perfectly mimic Gaussian processes and also enable efficient Bayesian inference for intractable problems, with over 200-fold speedups in multiple setups compared to current methods. We obtain strong results in such diverse areas as Gaussian process regression and Bayesian neural networks, demonstrating the generality of PFNs.

## 1 Introduction

In the last decade, supervised machine learning (ML) methods using deep learning architectures have made substantial progress on machine learning tasks where a large amount of training data is available (Vaswani et al., 2017; He et al., 2016; Krizhevsky et al., 2012). A very important problem in ML is thus to transfer these successes to smaller-scale setups with less data available. In this paper, we propose a way to build models that approximate posteriors with flexible and replaceable priors using deep learning models.

While the success of deep learning on large datasets can be attributed to the capacity of neural networks to approximate any function, there is still a need for encoding prior knowledge through model architecture (e.g. Convolutional Neural Networks (LeCun et al., 1989)) or regularizers (e.g. data augmentations (Hendrycks et al., 2019; Cubuk et al., 2020)). Otherwise, the no free lunch theorems show that there are no good methods to solve the class of prediction problems (Wolpert & Macready, 1997). In particular, this is true for ML applications to small datasets; thus, a large number of specialized algorithms have been developed for different tasks (LeCun et al., 1989; Kadra et al., 2021). Encoding prior information into such a model can, however, be challenging.

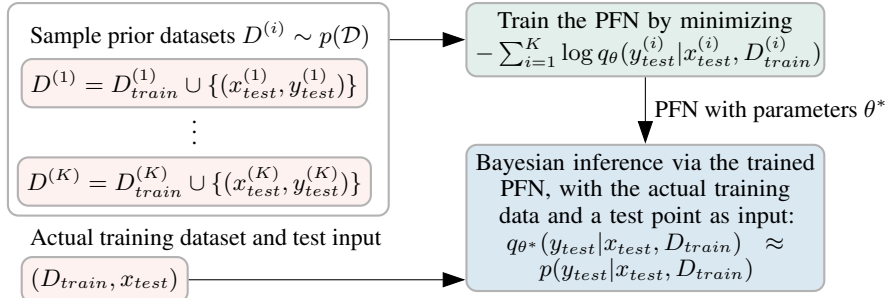


Figure 1: A visualization of Prior-Data Fitted Networks (PFNs). We sample (meta-train) datasets from a prior and fit a PFN on hold-out examples of these datasets. Given an actual (meta-test) dataset, we feed it and a test point to the PFN and obtain an approximation to Bayesian inference in a single forward propagation.

A well-defined way to bias a model is to use Bayesian inference. The foundation of Bayesian inference is an assumption on the distribution of the data to appear in a real world application. This assumption gives rise to a prior belief over the probability for the data to follow a particular model. One might, e.g., implement a prior, encoding the data as created by a neural network (Bayesian Neural Networks, (MacKay, 1992)), by a polynomial or the likelihood of a Gaussian mixture. In practice, however, it is often very difficult or intractable to retrieve the Posterior Predictive Distribution (PPD) for a given prior exactly.

In this paper, we use large-scale neural networks to approximate Bayesian models. We outline our method of Prior-data Fitted Networks (PFNs) in Figure 1. We assume a given representative prior distribution over supervised learning tasks (or functions), which provides our inductive bias. To train a PFN, we repeatedly sample a meta-train task (or function) from this given prior, draw a set of data points and their labels from it, mask one of the labels and learn to make probabilistic predictions for it based on the set-valued input of the rest of the data points. Given an actual (meta-test) dataset, we feed it and a test point as inputs to the PFN and the network outputs its prediction for the test point, conditioned on the dataset. As we will demonstrate, this probabilistic prediction approximates exact Bayesian posterior prediction.

Our PFNs thus allow us to approximate the posterior predictive distribution for *any* prior from which we are able to sample data. This is a very weak requirement compared to the standard assumption that unnormalized posterior probabilities can be computed, which is made by currently common methods, such as Markov Chain Monte Carlo (MCMC) Neal (1996); Andrieu et al. (2003); Welling & Teh (2011); Hoffman et al. (2014) and Variational Inference (VI) (Jordan et al., 1999; Wainwright & Jordan, 2008; Hoffman et al., 2013). This weak requirement of being able to sample allows a simple approximation of a large set of priors, including priors that are very hard to approximate with currently available tools. We make the following contributions:

- We introduce the technique of Prior-data Fitted Networks (PFNs);
- We show that via meta learning on data sets sampled from a dataset prior, PFNs yield an approximation technique for the posterior predictive distribution (PPD);
- We present architectural changes to successfully use Transformers for PPD approximation, including a novel predictive distribution for regression tasks;
- We demonstrate the successful application of PFNs of in approximating the PPD for Gaussian processes and Bayesian neural networks.

## 2 Background on Bayesian Inference

Bayesian inference relies on priors for real world applications. By combining the knowledge from a general prior distribution with the data from a particular task, one obtains a posterior distribution, which can be used to predict outputs for new data points from the task at hand.

In this work, we are interested in supervised learning problems, i.e., we model the relationship between inputs  $x$  and outputs  $y$ . Generally, we consider predictions based on datasets of arbitrary size  $n$ ,  $\mathcal{D} = \{(x_i, y_i) | i = 1, \dots, n\}$ , where  $y_i$  is the *output* for  $x_i$ .

There are two kinds of Bayesian models commonly used:

**i)** In *parametric models*, e.g. Bayesian Neural Networks (BNNs), we model a distribution  $p(t)$  over functions or tasks explicitly, where  $t$  is the random variable representing the task. In parametric models,  $t$  can be represented explicitly, e.g., as a vector, and either represents a conditional distribution  $p(y|x, t)$  or a joint distribution  $p(x, y|t)$  over the data. Thus, we can generate a prior over datasets  $p(\mathcal{D}) = p(t) \prod_{i=1}^n p(x_i, y_i|t)$ , where in the case of the conditional definition  $p(x_i, y_i|t)$  factors into  $p(y_i|x_i, t)$  and a simple task-independent distribution  $p(x_i)$ .

**ii)** *Non-parametric models*, e.g., Gaussian Processes (GPs), can be viewed similarly to the above, but here we cannot explicitly instantiate  $t$ . Thus, we directly model, that means without the indirection over the task  $t$ , both the prior on data sets  $p(\mathcal{D})$ , as well as the PPD  $p(y|x, \mathcal{D})$ .

For both models, the prior distribution (PD)  $p(\mathcal{D})$  is easily accessible. For all cases we consider, it is simple and cheap to sample from  $p(\mathcal{D})$ , fulfilling the only requirement our method has. In practice, using Bayesian inference has the advantages that (i) it has a theoretical foundation that makes it valid in setups which the prior  $p(t)$  fits, (ii) it can thus better account for the actual likelihood of different events; (iii) it is well calibrated and (iv) it is very interpretable as the prior describes the expectation of the model.

### 3 PPD Approximation with Prior-Data Fitting

Let us consider a parameterized model  $q_\theta$  that can accept a set  $\mathcal{D} = \{(x_i, y_i) | i = 1, \dots, n\}$ , as well as a query  $x$  as input, and which predicts a distribution over possible values of  $y$  for the query  $x$ . Many current neural network models can be used as such a model; in this paper, we use a variant of Transformers (Vaswani et al., 2017) as they are a reliable and powerful architecture.

We train this model by cross-entropy over samples drawn from the prior. Our proposed loss, the *Data Prior Negative Log-Likelihood (Prior-Data NLL)*  $\ell$  is defined as

$$\ell_\theta = -\mathbb{E}_{D \sim p(x, y) \sim p(\mathcal{D})} [\log q_\theta(y|x, \mathcal{D})], \quad (1)$$

where  $D \cup \{x, y\}$  is simply a dataset of size  $|D| + 1$  sampled from  $p(\mathcal{D})$ .

We minimize the Prior-Data NLL using a stochastic approximation of  $\ell_\theta$ , as we show in Algorithm 1. That means, we draw many samples of datasets from our prior and fit our model to predict a hold-out example correctly for these.

---

**Algorithm 1:** Training the model by Fitting Prior-Data

---

**Result:** A model  $q_\theta$  that will approximate the PPD

Define data prior  $p(\mathcal{D})$ ; initialize neural network  $q_\theta$ ;

**while** not converged **do**

Sample  $D \cup \{(x_i, y_i)\}_{i=1}^m \sim p(\mathcal{D})$ ;

Compute stochastic loss approximation  $\bar{\ell}_\theta = \sum_{i=1}^m (-\log q_\theta(y_i|x_i, D))$ ;

Update model  $\theta$  with gradient  $\nabla_\theta \bar{\ell}_\theta$ ;

**end**

---

We can show that minimizing  $\ell$  yields an approximation to the true PPD in a similar way as in Gordon et al. (2019).

**Theorem 1.** *The proposed objective  $\ell_\theta$  is equal to the expectation of the cross-entropy  $\mathbb{E}_{x, D} [H(p(\cdot|x, \mathcal{D}), q_\theta(\cdot|x, \mathcal{D}))]$  between the PPD and its approximation.*

*Proof.* The above can be shown with a simple derivation. We mark transformations for easy reading.

$$\ell_\theta = - \int_{D,x,y} p(x, y, D) \log q_\theta(y|x, D) \tag{2}$$

$$= - \int_{D,x} p(x, D) \int_y p(y|x, D) \log q_\theta(y|x, D) \tag{3}$$

$$= \int_{D,x} p(x, D) \mathbf{H}(p(\cdot|x, D), q_\theta(\cdot|x, D)) \tag{4}$$

$$= \mathbb{E}_{x,D}[\mathbf{H}(p(\cdot|x, D), q_\theta(\cdot|x, D))] \tag{5}$$

□

We see that by optimizing  $\ell_\theta$ , we optimize for similarity in terms of cross-entropy between the model and the PPD across the support of  $p(x, D)$ .

The proof above implies a minimization in terms of *KL-Divergence*, too.

**Corollary 1.1.** *The loss  $\ell_\theta$  is up to a constant, equal to the expected KL-Divergence of  $p(\cdot|x, D)$  with  $q_\theta(\cdot|x, D)$  over prior data  $x, D$ .*

For the proof, please refer to Appendix A.

In the following, we consider the optimum of the given optimization problem  $\theta = \arg \min_\theta \ell_\theta$ . We consider the optimal result of our process for the case where  $q_\theta$  describes a distribution family that can exactly approximate the posterior, that means there is a  $\theta$  such that  $q_\theta = p$ .

**Corollary 1.2.** *If  $q_\theta$  is a distribution family that can exactly approximate the posterior, we have  $q_{\theta^*}(\cdot|x, D) = p(\cdot|x, D)$  for all  $x, D$  with  $p(x, D) > 0$ .*

*Proof.* We assume a  $\theta$  with  $q_\theta = p$  exists and we know that the cross-entropy is minimized for equal distributions; thus the optimum  $\theta$  fulfills  $q_{\theta^*}(\cdot|x, D) = p(\cdot|x, D)$  for all  $x, D$  with  $p(x, D) > 0$ . □

## 4 Adapting the Transformer for Bayesian Inference

We use a standard transformer (Vaswani et al., 2017) with only a few changes to its architecture to fit our problem better. Transformers were developed with sequence modelling and prediction in mind, but in our application we want to map a set  $D = \{(x_i, y_i)\}_{i=1}^K$  of data points  $(x_i, y_i)$  and a query  $x$  to a distribution  $q_\theta(y|x, D)$  over  $y$ , which in many cases should be continuous. We propose an efficient architecture that allows batch-wise training and prediction, as well as a novel regression head for continuous  $y$ .

### 4.1 An Efficient Architecture

The first step in the architecture is to encode  $x$  and  $y$ ; in all our experiments, we use simple linear layers for this, to project from the given dimension to the dimension of the transformer, but one could also use other application-specific encoding layers. Our setup should be equivariant to the ordering of the input dataset’s examples. To achieve this, we remove positional encodings and feed  $(x, y)$  pairs together to the transformer as a sum of their encodings; see Appendix D for an ablation of the impact of the equivariance to the ordering.

In practice, we feed multiple query points to the transformer for efficiency reasons; these are the only inputs for which position matters, and we thus use the output at their positions as prediction for their corresponding PPD over  $y$ . Since different query points should not depend on each other, we use an attention mask that allows attention from every  $(x, y)$  pair to every other  $(x, y)$  pair and allows attention from all query points to the  $(x, y)$  pairs. The mask  $M$  we use can be characterized by the number of input examples  $m$  and the number of query points  $n$ , we write  $M_{m,n}$ . An example for

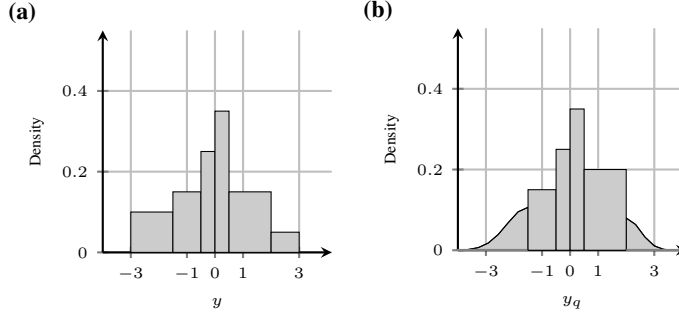


Figure 2: A visualization of the Riemann probability distribution.

such a matrix is

$$M_{3,2} = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 & 1 \end{bmatrix}. \quad (6)$$

We predict the distribution of  $y$  from the output of the last time step. During training, we fix the total size  $M = m + n$  and sample  $m$  from a distribution, in order to allow the Transformer to learn to work with different dataset sizes. Since for smaller  $m$  we have more query points, we over-sample larger  $m$ , s.t. the number of query points seen during training for each  $m$  is equal during training. That is, we sample each  $m \in \{0, \dots, M - 1\}$  with a weight  $1/n = 1/(M - m)$ .

## 4.2 Riemann Distribution

The modelling of continuous distributions is hard with neural networks. To yield strong performance in modelling PPDs, we developed a distribution that works particularly well with neural networks. Based on the knowledge that neural networks work well for classification and inspired by discretizations in distributional reinforcement learning (Bellemare et al., 2017), we created a discretized continuous distribution that we call *Riemann Distribution*, as it is based on a similar idea to approximation as Riemann integrals. PDFs of our distribution are bar plots, see Figure 2a.

The problem remaining is how to select the boundaries of the buckets. To make the problem balanced for the Transformer, we simply select buckets such that the likelihood of each bucket  $b$  has the same prior predictive probability  $p(y \in b) = 1/|\mathbf{B}|, \forall b \in \mathbf{B}$ . We approximate this by taking a large sample of  $y$  values from our prior and adjusting the buckets such that an equal number of these values fall in each bucket. The distributions described above are very powerful; e.g., in contrast to most parametric distributions they can also represent multimodal distributions. Their downside, however, is that they have finite support: the probability for  $y$  to be greater than the highest border is 0. This is a problem whenever one tries to model distributions that have unbounded support. A simple trick that worked well for us is to replace the last bar on each side with an appropriately scaled half-normal distribution. We visualize a resulting PDF in Figure 2b. For an exact mathematical definition of the distributions, we refer the reader to Appendix C. For an ablation of the Riemann distribution, see Appendix D.

## 5 Posterior Approximation Studies

In our first set of experiments, we study the capability of PFNs to perform Bayesian inference for the tractable case of Gaussian Processes (GPs) with fixed hyperparameters (where we can compare to ground truth data; Section 5.1) and the intractable cases of GPs with unknown hyperparameters (Section 5.2) and Bayesian Neural Networks (BNNs; Section 5.3).

Based on Theorem 1, we present a novel way to evaluate methods approximating PPDs and Posteriors, the Prior-Data NLL itself. We simply take the negative log-likelihood of a method on data generated from the prior over datasets. We can do this, of course, for our models, but we can do this as easily with other well-known methods. We simply treat samples from the prior as training and test sets

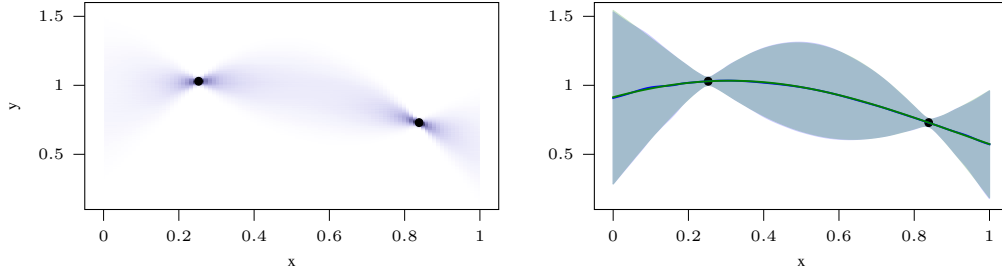


Figure 3: Left: the PFN’s probability for query points  $x$  given two evaluations of the function, highlighting the binning in its predictions. Right: The PFN’s mean (solid blue line) and 95% confidence interval (blue shaded region), alongside the exact GP mean posterior (solid green line) and 95% confidence interval (green shaded region). The two are near identical, with tiny differences. For additional comparisons, see Figure 6 in the appendix.

for these methods and evaluate the negative log-likelihood of the approximated distributions on the predicted outputs for the test sets. As we show in Corollary 1.1, the metric we evaluate is directly representative of the  $KL$  divergence between the true PPD and the approximation. There only is a constant difference. If a method only approximates the posterior and not the PPD directly, we approximate PPD with a large MC integration.

### 5.1 Gaussian Process approximation

We begin by approximating GPs as a study as to what the transformers are capable of; GPs with fixed hyperparameters are convenient for this purpose since the tractability of their PPD allows us to assess how close our approximation is. In this case, our meta-train datasets are functions  $f$  sampled from a GP with the given fixed hyperparameters. Since GPs describe conditional distributions given an input  $x$ , we sample  $n + 1$  inputs  $x_i$  uniformly at random from the unit-cube and evaluate  $y_i = f(x_i)$ . As usual, we then train our PFNs to predict  $y_{n+1}$  from the dataset  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  and the query point  $x_{n+1}$ . The trained PFN can then be used to compute the PPD for new datasets and arbitrary query points. Figure 3 visualizes the predictions of such a PFN for a small dataset with two data points and query points  $x$  ranging from 0 to 1. In Figure 3 (left), one can still see the very tiny boxes in which the Riemann distribution is discretized. In Figure 3 (right), we compare these predictions to those of a ground-truth GP with the given hyperparameters; both the mean function and the plotted 95% confidence interval are virtually indistinguishable from this ground truth. The model also learns on its own, completely automatically, to generate smooth distributions without any explicit knowledge of the positions of the bars in  $\mathbb{R}$ . In Figure 6 in the appendix we show more qualitative examples of the behavior of our approximation together with the exact posterior. Figure 4a shows that the qualitative results above generalize to much large datasets with multiple inputs. The GP Posterior achieves optimal performance by definition, and we can approximate it very closely, and better so with longer training. Additionally, there does not seem to be a trend towards worse approximation with larger datasets.

### 5.2 Approximating A GP with Hyper-Priors

A common practice in fields applying GPs in the real world is to define hyper-priors over the hyperparameters of GPs Rasmussen & Williams (2006); Snoek et al. (2012). Thus, the prior of these models considers a more diverse set of functions. The correlation between data points, the smoothness of the function and the scale of outputs can be variable. This allows the user to apply GPs to a larger set of tasks. The downside is that it is not possible to compute the PPD of the GP with hyper-priors exactly. There are two common practices to approximate it anyways and we evaluate against both. (i) Firstly, MLE-II is the most common, which is a simple special case of variational inference. Here, suitable hyperparameters are found by maximum a posteriori (MAP) estimation (e.g., Rasmussen & Williams (2006)). We use the fitting setup of BoTorch (Balandat et al., 2020) from which our Hyper-Priors for all models are inspired as well. (ii) Secondly, Markov Chain Monte Carlo (MCMC), a method to sample from any distribution for which one can calculate non-normalized probabilities, is also frequently applied to sample hyperparameters (e.g., Snoek et al. (2012)). Here

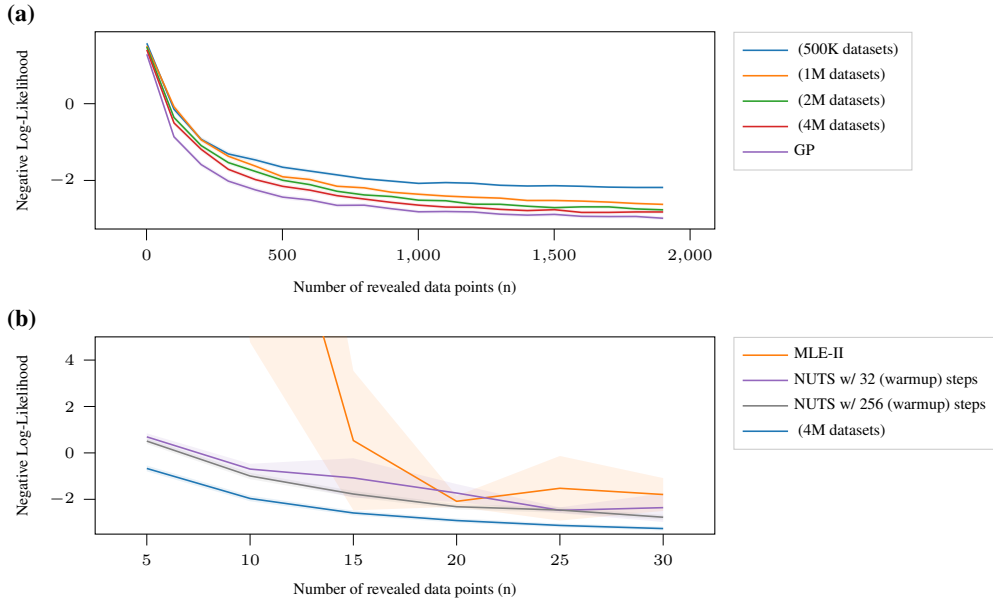


Figure 4: (a) NLL with an RBF-Kernel on different sizes of sets. Trained with Transformer. Hyperparameters: length-scale .6, noise  $1e-4$ , out scale 1., 5 inputs. (b) Plot with standard SingleTaskGP (with a different noise distribution:  $\text{Gamma}(.0001, 1.0)$ )

we compare against doing so by the NUTS method (Hoffman et al., 2014), a state-of-the-art method that uses gradients to facilitate the Hamiltonian Monte Carlo algorithm. We plot the performance of the different methods in Figure 4b. Our model can approximate the PPD on the prior much closer than both of these methods, and its inference is more than  $200\times$  faster than MLE-II and  $1000\times$  to  $8000\times$  faster than NUTS.

### 5.3 Approximating the BNN PPD

Bayesian Neural Networks (BNNs) provide an inherent ability to model uncertainty and strong regularization by inherent ensembling. Recently, approximation methods such as stochastic variational inference or stochastic gradient Markov Chain Monte Carlo methods allow for much faster convergence, which has sparked a great deal of interest in the area (Izmailov et al., 2021; Fortuin et al., 2021). However, these methods still converge too slowly for many practical applications.

In Figure 5, we show that our method can outperform today's default methods using a much smaller compute budget. The figure shows the approximation accuracy on synthetic datasets generated using the BNN prior that is used during inference by all solvers. To generate these synthetic datasets, we sampled random weights for a BNN  $BNN$  and used normally-distributed i.i.d. features  $X$  to obtain  $y := BNN(X)$  with  $\mathcal{D} := (X_{1:n}, y_{1:n})$ .

## 6 Conclusion & Future Work

We present a novel way to approximate the posterior predictive distribution very efficiently using deep neural networks. We show its capability on a set of diverse tasks and priors. This opens the possibility for a multitude of ground breaking future work. We expect especially the following to be fruitful: (i) Work on finding novel priors that are cheap to sample for particular application areas. (ii) Work on architectures that are well-fit for this task, as we simply used a slight adaption of a current model.

## References

Ron Amit and Ron Meir. Meta-learning by adjusting priors based on extended pac-bayes theory. In *International Conference on Machine Learning*, pp. 205–214. PMLR, 2018.

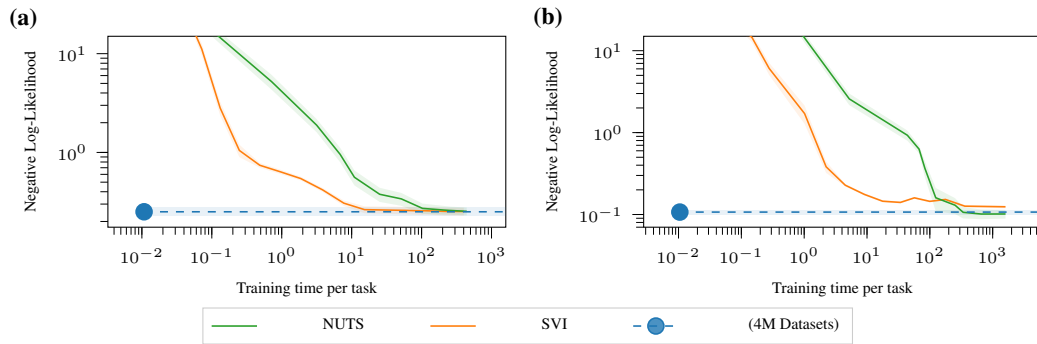


Figure 5: Training time spent per dataset at a given approximation accuracy. Negative Log Likelihood (NLL) is calculated as the mean over 100 tasks and 95% confidence intervals over tasks is indicated. Transformer and SVI are evaluated on a GPU while MCMC is evaluated on a CPU since its speed surpasses GPU speed. (a) Bayesian Neural Network prior with 3 features, 2 layers and a hidden dimensionality of 5. (b) Bayesian Neural Network prior with 8 features, 2 layers and a hidden dimensionality of 64.

Christophe Andrieu, Nando De Freitas, Arnaud Doucet, and Michael I Jordan. An introduction to mcmc for machine learning. *Machine learning*, 50(1):5–43, 2003.

Maximilian Balandat, Brian Karrer, Daniel R. Jiang, Samuel Daulton, Benjamin Letham, Andrew Gordon Wilson, and Eytan Bakshy. BoTorch: A Framework for Efficient Monte-Carlo Bayesian Optimization. In *Advances in Neural Information Processing Systems 33*, 2020.

Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement learning. In *International Conference on Machine Learning*, pp. 449–458. PMLR, 2017.

Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall, and Noah D. Goodman. Pyro: Deep Universal Probabilistic Programming. *Journal of Machine Learning Research*, 2018.

David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, Apr 2017. ISSN 1537-274X. doi: 10.1080/01621459.2017.1285773.

C Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra. Weight uncertainty in neural network. In F. Bach and D. Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning (ICML’15)*, volume 37, pp. 1613–1622. Omnipress, 2015.

Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.

Jeffrey Chan, Valerio Perrone, Jeffrey P Spence, Paul A Jenkins, Sara Mathieson, and Yun S Song. A likelihood-free inference framework for population genetic data using exchangeable neural networks. In *NeurIPS*, 2018.

Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’16*, pp. 785–794, New York, NY, USA, 2016. Association for Computing Machinery. ISBN 9781450342322. doi: 10.1145/2939672.2939785.

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.

Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15750–15758, 2021.



- Kyle Cranmer, Johann Brehmer, and Gilles Louppe. The frontier of simulation-based inference. *Proceedings of the National Academy of Sciences*, 117(48):30055–30062, 2020.
- Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pp. 702–703, 2020.
- Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *J. Mach. Learn. Res.*, 7: 1–30, December 2006. ISSN 1532-4435.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In D. Precup and Y. Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning (ICML’17)*, volume 70, pp. 1126–1135. Proceedings of Machine Learning Research, 2017a.
- C. Finn, T. Yu, T. Zhang, P. Abbeel, and S. Levine. One-shot visual imitation learning via meta-learning. In *Proceedings of the 1st Annual Conference on Robot Learning (CoRL)*, volume 78, pp. 357–368. PMLR, 2017b.
- Vincent Fortuin, Adrià Garriga-Alonso, Florian Wenzel, Gunnar Rätsch, Richard Turner, Mark van der Wilk, and Laurence Aitchison. Bayesian neural network priors revisited, 2021.
- Marta Garnelo, Dan Rosenbaum, Christopher Maddison, Tiago Ramalho, David Saxton, Murray Shanahan, Yee Whye Teh, Danilo Rezende, and SM Ali Eslami. Conditional neural processes. In *International Conference on Machine Learning*, pp. 1704–1713. PMLR, 2018a.
- Marta Garnelo, Jonathan Schwarz, Dan Rosenbaum, Fabio Viola, Danilo J Rezende, SM Eslami, and Yee Whye Teh. Neural processes. *arXiv preprint arXiv:1807.01622*, 2018b.
- P. Gijsbers, E. LeDell, S. Poirier, J. Thomas, B. Bischl, and J. Vanschoren. An open source automl benchmark. *arXiv preprint arXiv:1907.00909 [cs.LG]*, 2019. Accepted at AutoML Workshop at ICML 2019.
- Jonathan Gordon, John Bronskill, Matthias Bauer, Sebastian Nowozin, and Richard Turner. Meta-learning probabilistic inference for prediction. In *International Conference on Learning Representations*, 2019.
- Kazuyuki Hara, Daisuke Saitoh, and Hayaru Shouno. Analysis of dropout learning regarded as ensemble learning. *CoRR*, abs/1706.06859, 2017.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Dan Hendrycks, Norman Mu, Ekin Dogus Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. Augmix: A simple data processing method to improve robustness and uncertainty. In *International Conference on Learning Representations*, 2019.
- Sepp Hochreiter, A Steven Younger, and Peter R Conwell. Learning to learn using gradient descent. In *International Conference on Artificial Neural Networks*, pp. 87–94. Springer, 2001.
- Matthew D. Hoffman, David M. Blei, Chong Wang, and John Paisley. Stochastic variational inference. *Journal of Machine Learning Research*, 14(4):1303–1347, 2013.
- Matthew D Hoffman, Andrew Gelman, et al. The no-u-turn sampler: adaptively setting path lengths in hamiltonian monte carlo. *J. Mach. Learn. Res.*, 15(1):1593–1623, 2014.
- Pavel Izmailov, Sharad Vikram, Matthew D. Hoffman, and Andrew Gordon Wilson. What are bayesian neural network posteriors really like? *CoRR*, abs/2104.14421, 2021.
- Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233, 1999.

- Arlind Kadra, Marius Lindauer, Frank Hutter, and Josif Grabocka. Regularization is all you need: Simple neural nets can excel on tabular data, 2021.
- Kirthevasan Kandasamy, Willie Neiswanger, Jeff Schneider, Barnabas Poczos, and Eric Xing. Neural architecture search with bayesian optimisation and optimal transport, 2019.
- Hyunjik Kim, Andriy Mnih, Jonathan Schwarz, Marta Garnelo, Ali Eslami, Dan Rosenbaum, Oriol Vinyals, and Yee Whye Teh. Attentive neural processes. In *International Conference on Learning Representations*, 2018.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR (Poster)*, 2015.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger (eds.), *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012.
- Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation applied to handwritten zip code recognition. *Neural Computation*, 1(4):541–551, 1989. doi: 10.1162/neco.1989.1.4.541.
- Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. In *International Conference on Learning Representations*, 2018.
- Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm. *arXiv preprint arXiv:1608.04471*, 2016.
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- Jan-Matthis Lueckmann, Jan Boelts, David Greenberg, Pedro Goncalves, and Jakob Macke. Benchmarking simulation-based inference. In *International Conference on Artificial Intelligence and Statistics*, pp. 343–351. PMLR, 2021.
- David JC MacKay. A practical bayesian framework for backpropagation networks. *Neural Computation*, 4(3):448–472, 1992.
- Mahdi Pakdaman Naeni, Gregory F. Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, AAAI’15, pp. 2901–2907. AAAI Press, 2015. ISBN 0262511290.
- R. Neal. *Bayesian Learning for Neural Networks*. 1996.
- George Papamakarios and Iain Murray. Fast  $\epsilon$ -free inference of simulation models with bayesian conditional density estimation. In *Advances in neural information processing systems*, pp. 1028–1036, 2016.
- Liudmila Prokhorenkova, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush, and Andrey Gulin. Catboost: Unbiased boosting with categorical features. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, NIPS’18, pp. 6639–6649, Red Hook, NY, USA, 2018. Curran Associates Inc.
- C. Rasmussen and C. Williams. *Gaussian Processes for Machine Learning*. The MIT Press, 2006.
- Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)*. The MIT Press, 2005. ISBN 026218253X.
- Jonas Rothfuss, Vincent Fortuin, Martin Josifoski, and Andreas Krause. Pacoh: Bayes-optimal meta-learning with pac-guarantees. In *International Conference on Machine Learning*, pp. 9116–9126. PMLR, 2021.

- Binxin Ru, Xingchen Wan, Xiaowen Dong, and Michael Osborne. Interpretable neural architecture search via bayesian optimisation with weisfeiler-lehman kernels. In *International Conference on Learning Representations*, 2021.
- Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-learning with memory-augmented neural networks. In *International conference on machine learning*, pp. 1842–1850. PMLR, 2016.
- J. Snoek, H. Larochelle, and R. Adams. Practical Bayesian optimization of machine learning algorithms. In P. Bartlett, F. Pereira, C. Burges, L. Bottou, and K. Weinberger (eds.), *Proceedings of the 25th International Conference on Advances in Neural Information Processing Systems (NeurIPS’12)*, pp. 2960–2968, 2012.
- Ray J Solomonoff. The discovery of algorithmic probability. *Journal of Computer and System Sciences*, 55(1):73–88, 1997. ISSN 0022-0000. doi: <https://doi.org/10.1006/jcss.1997.1500>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- Martin J Wainwright and Michael Irwin Jordan. *Graphical models, exponential families, and variational inference*. Now Publishers Inc, 2008.
- Ben Wang and Aran Komatsuzaki. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.
- Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pp. 681–688. Citeseer, 2011.
- David H Wolpert and William G Macready. No free lunch theorems for optimization. *IEEE transactions on evolutionary computation*, 1(1):67–82, 1997.
- Jaesik Yoon, Taesup Kim, Ousmane Dia, Sungwoong Kim, Yoshua Bengio, and Sungjin Ahn. Bayesian model-agnostic meta-learning. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pp. 7343–7353, 2018.
- Sheheryar Zaidi, Arber Zela, Thomas Elsken, Chris Holmes, Frank Hutter, and Yee Whye Teh. Neural ensemble search for uncertainty estimation and dataset shift. *arXiv preprint arXiv:2006.08573*, 2020.
- Hongpeng Zhou, Minghao Yang, Jun Wang, and Wei Pan. Bayesnas: A bayesian approach for neural architecture search. *CoRR*, abs/1905.04919, 2019.
- Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8697–8710, 2018.

## A Proof of Corollary 1.1

*Proof.* This can easily be seen by considering the result of Theorem 1

$$\mathbb{E}_{x,D}[KL(p(\cdot|x, \mathcal{D}), q_\theta(\cdot|x, \mathcal{D}))] \quad (7)$$

$$= -\mathbb{E}_{x,D} \left[ \int_y p(y|x, \mathcal{D}) \log \frac{q_\theta(y|x, \mathcal{D})}{p(y|x, \mathcal{D})} \right] \quad (8)$$

$$= -\mathbb{E}_{x,D} \left[ \int_y p(y|x, \mathcal{D}) \log q_\theta(y|x, \mathcal{D}) \right] + \mathbb{E}_{x,D} \left[ \int_y p(y|x, \mathcal{D}) \log p(y|x, \mathcal{D}) \right] \quad (9)$$

$$= \mathbb{E}_{x,D}[\mathbb{H}(p(\cdot|x, \mathcal{D}), q_\theta(\cdot|x, \mathcal{D})) - \mathbb{H}(p(\cdot|x, \mathcal{D}))] \quad (10)$$

$$= \ell_\theta + C, \quad (11)$$

where  $C = -\mathbb{E}_{x,D}[\mathbb{H}(\cdot|x, \mathcal{D})]$ .  $\square$

## B Qualitative Analysis

[h]

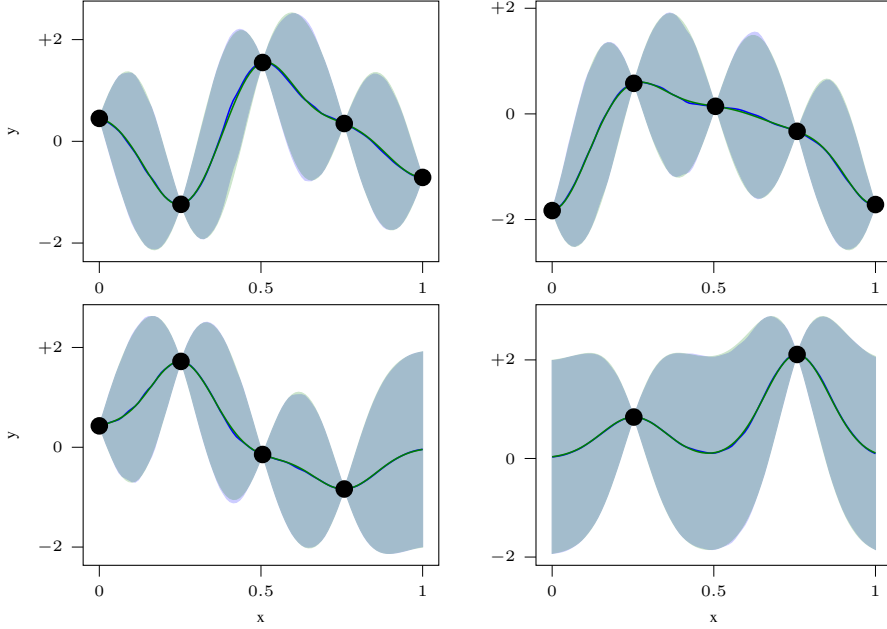


Figure 6: Green are the 95% intervals of the original GP posterior and blue are the corresponding transformer approximations. Lengthscale .1, Outscale 1., noise: 1e-4

## C Riemann Distribution Definition

Given a set of buckets  $\mathbf{B}$ , where each bucket  $b \in \mathbf{B}$  describes an interval such that  $\bigcup_{b \in \mathbf{B}} b = [a, b]$  and for any  $b, b' \in \mathbf{B}$  we have  $b \cap b' = \emptyset$ . Additionally, we define the upper-bound (lower-bound) of the lowest (highest) bucket to be  $l$  ( $f$ ). Additionally, we use a mapping  $B(\cdot)$  that maps a  $y \in \mathbb{R}$  to the unique bucket  $b$  with  $y \in b$  and the width of each bucket to be defined as  $w(b) = \max_{y \in b} y - \min_{y \in b} y$ . The model gives us a probability for each bucket  $p_b$ .

In the finite case the distribution is defined as

$$p(y) = p_{b(y)} / w(b(y)). \quad (12)$$

We normalize with the width of the bucket since  $p_b$  describes the probability  $p(y \in b)$ , but we are interested in the probability  $p(y = y)$ .

For the case of infinite support given probabilities  $p_b$  for each bucket  $b$ , we have

$$f(y_q) = \begin{cases} 2 \cdot f_N(y_q - l), & \text{if } y_q \geq l \\ 2 \cdot f_N(f - y_q), & \text{if } y_q \leq f \\ p_{b(y_q)} / w(b(y_q)), & \text{otherwise.} \end{cases} \quad (13)$$

Here, the Half-Normal distributions are scaled such that half the probability weight is in the first (last) bucket.

## D Ablation

Above, we introduce a novel way of doing regression with neural networks that departs from the tradition of using normal distributions. Additionally, we remove the default positional embeddings

from the Transformer. We show the impact of both on the Transformer in Figure 7. We can see that, while the effect of the Riemann Distribution is very pronounced, the effect of making our architecture permutation invariant only slowly increases with sequence length.

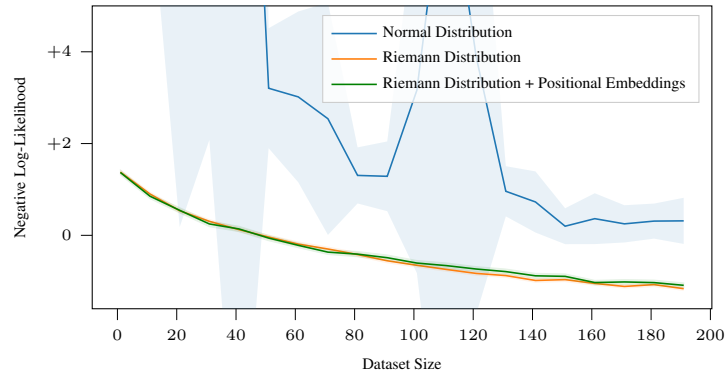


Figure 7: In this figure we show the impact of different architectural decisions on the performance of our method in fitting a Gaussian Process. We use the same Gaussian Process as in Section 5.1.