

2 We thank the anonymous reviewers for their enthusiasm and detailed comments on the manuscript. We summarise the  
 3 reviews as positive, and the main concerns were related to clarifying the motivation and the experiments. We agree with  
 4 the requests, and we will use the additional ninth page in the camera-ready paper for expanding the details as requested.  
 5 We start by addressing R2’s concerns as they had the lowest score.

6 **R2: (1)** The main motivation for this work is to establish understanding about the link between Matérn GP priors and  
 7 neural network activation functions. This link is *explicit* as shown in this paper. A motive for this is to allow NN models  
 8 to incorporate some of the appealing properties of GP models (*e.g.*, well-characterized uncertainties), while maintaining  
 9 the flexibility and efficiency of NNs. The choice of kernel/activation function is up to the modelling task and ‘expert  
 10 knowledge’. We merely provide a building block. The Matérn is a widely used prior, and worth adding to the NN tool set.  
 11 Stationarity encodes conservative behaviour suitable for uncertainty quantification (see R4(2)). **(2)** The references R2  
 12 provides ([1–3, 5]) are tackling a different problem, where the kernel is not used for encoding specific prior information,  
 13 but inferred from data (*cf.*, ‘automatic statistician’), and [4] is covered in this paper (limiting RBF case and NN kernel).

14 **(3)** In the experiments, we originally reported only accuracy UCI classification tasks with Matérn-3/2 activation/kernel showing the AUC  
 15 and NLPD (accounting for uncertainty). As requested, we added AUC to the results, which is in-line with the previous  
 16 results (see table), with our proposed model outperforming the baselines. Comparison to SIREN is interesting (NB: the  
 17 SIREN paper was put on arXiv *after* the submission DL), and to answer your question, we ran the experiments with it  
 18 as well (see table). On the UCI tasks, SIREN performs comparably to our method. On the OOD image classification  
 21 task, it performs clearly worse (known-class NLPD: 0.103 vs. 0.106, unknown-class NLPD: 0.896 vs. 1.62), but still  
 22 better than the baselines. Note that SIREN encodes a different type of prior (infinite smoothness, like the RBF).

UCI classification tasks with Matérn-3/2 activation/kernel showing the AUC metric. Also results for SIREN activations included.

| (10-fold cv) | n      | d  | c | SVGP      | GPDNN     | SV-DKL    | Matérn activ. | SIREN activation |           |           |
|--------------|--------|----|---|-----------|-----------|-----------|---------------|------------------|-----------|-----------|
|              |        |    |   | AUC       | AUC       | AUC       | AUC           | NLPD             | ACC       | AUC       |
| Adult        | 45222  | 14 | 2 | .893±.004 | .774±.052 | .912±.003 | .913±.004     | .314±.006        | .854±.005 | .912±.004 |
| Connect-4    | 67556  | 42 | 3 | .824±.005 | .675±.019 | .909±.013 | .913±.004     | .449±.008        | .825±.005 | .909±.003 |
| Covtype      | 581912 | 54 | 7 | .971±.001 | .943±.015 | .998±.000 | .998±.000     | .119±.002        | .957±.001 | .998±.000 |
| Diabetes     | 768    | 8  | 2 | .817±.049 | .769±.053 | .512±.095 | .838±.051     | .487±.006        | .771±.054 | .835±.054 |

23 **R1: (1)** We are glad that the question about relation to kernel feature expansions was brought up. Fourier features  
 24 (random, dense/structured, sparse) are typically leveraged for stationary kernels by projecting the GP problem on a set  
 25 of harmonic basis functions. While we share the idea of using the Fourier duality, the resulting model is spanned by  
 26 different basis functions; *e.g.*, sinusoidal FFs enforce (global) stationarity (approximation is based on Eq. (5)), while  
 27 our approach is *locally* stationary as defined by the Gaussian weights in Eq. (2), which the approximation is based on.  
 28 This discussion was left out in the interest of space, but the additional page will give us space to cover this. **(2)** We  
 29 appreciate the suggestions for improving the visualizations. We did focus on real data (with three different real-world  
 30 experiment setups) in our quantitative experiments, and the toy data examples were to give a general understanding.  
 31 Your suggestion of varying the number of hidden units and MC samples is good and easy to do. We’ll include this in the  
 32 appendix to facilitate understanding of the effect of these parameters. We’ll also run a GP on CIFAR-10 as suggested.  
 33 **(3)** We recognize your concern with the term ‘uncertainty calibration’ (used here as GPs are commonly said to have  
 34 representative uncertainties). We will replace it with the less loaded ‘uncertainty quantification’ where applicable.

35 **R3: (1)** The paper ‘Deep Neural Networks as Gaussian Processes’ (thank you for the reference) works in the scope  
 36 of our Sec. 2 (connection from activation functions to GP kernels for principled inference). We take the opposite  
 37 direction (GP prior → activation function) to encode properties of the GP prior into the NN. **(2)** Indeed, the SPDE link  
 38 for constructing Gauss–Markov random fields from stationary kernels is relevant here. The deeper-level connection is  
 39 related to the spectrum of the corresponding ‘covariance operator’ (the covariance function is formally the kernel of this  
 40 operator). We agree that there is still a lot to uncover in this space. **(3)** The computational complexity of our model is  
 41 on par with using other activation functions (say, a ReLU) in the NN models. In practice, we may converge slightly  
 42 slower (*e.g.*, 53 s vs. 42 s for the results on the ‘adult’ data set). Compared to the GP models, especially with large  
 43 data sets, we gain a considerable speed-up. **(4)** You are right that in this regard our approach can be thought of as a  
 44 simplified Bayesian neural network—but yet as one, where the Matérn prior can be directly encoded.

45 **R4: (1)** We understand your comment on not dedicating attention to the inference methods (also raised by R2). This  
 46 is true and mostly on purpose to retain focus. MC dropout is neither fast nor exciting, but does its job and unlikely  
 47 introduces complications that would raise suspicions/complications with the model. Extending beyond it is left as future  
 48 work. **(2)** In both Fig. 1–2 the number of hidden units is left small (also no ensembling or such used) to highlight the  
 49 noisiness of the corresponding NN models. In Fig. 1, the uncertainty does not always increase right outside of the data  
 50 range as it does for the GP models on the top. This property reflects the remaining suboptimality of the model. In Fig. 2,  
 51 the trained models end up in different local optima (this could probably be tuned). Mean-reversion is characteristic for  
 52 stationary models (outside data the model knows it’s uncertain and reverts to the mean) and a desired property (also  
 53 applies to the RBF fig). **(3)** The activation func. lengthscale parameter is fixed in all the NN experiments, because the  
 54 preceding layer(s) take care of scaling the inputs, which serves the same purpose. We will discuss this in the main paper.  
 55 **(4)** As discussed in Sec. 5, the only practical problems were encountered with the spiky (non-differentiable) Matérn-1/2  
 56 activation. In general, tuning the learning rate might also help prevent possible issues with convergence.