1  We thank the reviewers for the valuable feedback and address specific comments below.

2  **Motivation [R3]:** We emphasize that the major advantage of our method is fast identification of reasonable kernel
3  hyperparameters since MLL-opt is costly and has no global optima guarantees. A side benefit is that our approach
4  empirically does not result in catastrophically bad hyperparameters because collective training across different tasks
5  helps avoid bad local minima: averaging the gradient over tasks seems to help the training escape from pathological
6  solutions. Also, overparameterization in deep neural nets is known to help SGD find good optima both in terms of
7  optimization and generalization, and sometimes with theoretical support [Allen-Zhu et al., 2019, Du et al., 2019].

8  **Hierarchical attention architecture [R1 & R3]:** We would like to first clarify that our architecture consists of
9  three parts: *local transformer*, *aggregation function* and *global transformer*. The *local transformer* sees every
10 dimension independently and an aggregated representation is calculated for each dimension. The **interactions between**
11 **dimensions** are later modeled by the *global transformer*. The architecture is carefully designed for the GP problem with
12 desired properties such as versatility and permutation equivariance. We are not aware of similar existing architecture for
13 learning on datasets. Kim et al. [2019] uses attention to aggregate information of the context data and relates target
14 data to context data for prediction. The major difference is that our method uses hierarchical attention and it is able to
15 take in datasets with **different dimensionality** with **one model**. We will expand this explanation in the paper. We also
16 expect the proposed architecture to be useful for other tasks that involve learning on sets (e.g., the neural statistician).

17 **Training setup [R1, R3 & R4]:** With regard to R1's comments on generating large synthetic datasets, we would
18 like to mention that our model does not necessarily need to be trained on large synthetic datasets for it to generalize
19 to large-scale problems. In our experiments, we train a single neural model using synthetic datasets generated with
20 stationary kernels of dimensionality $2\sim15$ and 30 data points on average. (Full details are available in Appendix C.) This
21 **single trained model** is applied to **all the real-world tasks** with varying dimensionality and cardinality, yet it is able to
22 perform fairly well on those **unseen tasks**. This demonstrates the effectiveness of our proposed architecture: although
23 it is trained on small synthetic datasets, it learns to extract the underlying structure for large unseen datasets. Our model
24 is expected to work for a relatively low-dimensional GP regression task that follows the stationarity assumption.

25 **Scalability [R1 & R3]:** Full GP MLL-opt scales as $\mathcal{O}(N^3)$ using Cholesky decomposition (GPy) and $\mathcal{O}(TN^2)$ using
26 the conjugate gradient method (GPyTorch). It can be reduced to $\mathcal{O}(NM^2)$ with a sparse GP. All the methods require
27 iterative optimization, meaning the computation cost incurred in practice is the complexity multiplied by the number of
28 opt iterations. In comparison, our method scales as $\mathcal{O}(N^2 + D^2)$ and only requires one forward pass of the trained
29 neural net. In terms of actual running time, our method compares much more favorably with all the baselines above.
30 The $2,000$ datapoints used in our experiments are already considered large-scale for GP. To give a larger-scale example,
31 for *power plant* with $4,000$ data points, our method is 226 times faster than GPyTorch on GPU while 332 times faster
32 than GPy on CPU. We will include these larger-scale experiments in the revised version. It is possible to go even larger
33 for our model, but full GP MLL-opt methods either take too long on CPU or report out of memory issues on GPU.
34 Note that our model's complexity can be further reduced to $\mathcal{O}(MN + MD)$ if $M$-sparse attention is introduced. With
35 regards to R3's comments that BO is usually small-scale, it might be observed that 1) BO is small-scale in part because
36 of the cost of hyperparameter inference, and 2) AHGP is often empirically more robust in the small-data regime.

37 **Predictive variance [R1]:** It was presented in Appendix C. We will improve the presentation to make it more accessible.

38 **MLL when AHGP performs better than MLL-opt [R2]:** Upon checking, we observed that AHGP found an MLL
39 that is slightly smaller than that of MLL-opt. This illustrates that MLL-opt can potentially overfit in the low-data regime,
40 while AHGP adds extra regularization through learning to optimize over many different GP tasks.

41 **Noise variance [R1]:** In all experiments, we normalize the data input and output. The noise variance is fixed at 0.01. It
42 is possible for our method to learn the noise variance. In practice, we find the predictive performance is not sensitive to
43 the noise variance and setting it to 0.01 gives competitive performance for all baselines. Also, spectral mixture kernel is
44 flexible enough to model the noise as well (with $\mu = 0$, small $\sigma^2$). We will add the clarification to make it more clear.

45 **Neural processes (NPs) [R1]:** In terms of the problem formulation, our work targets efficient estimation of GP
46 hyperparameters as part of a full GP inference procedure, while NP serves as an end-to-end approximation to the GP
47 inference itself. We also hope to emphasize that our method tackles a richer distribution of tasks (i.e., regression tasks
48 generated from stationary GP, with **different dimensionality** and cardinality) with **a single trained model**, while NP
49 aims at solving one particular task or similar tasks (**same dimensionality**) with one neural model.

50 **BO hyperparameter transfer learning [R3]:** Thanks for the pointer. We will include them in related work. This type
51 of approach focuses on better warm-starting BO through leveraging data of previous related BO runs, such that a good
52 surrogate model can be built with very few evaluations of the target function. In comparison, our method focuses on
53 fast inference of kernel hyperparameters, and our single trained model is applicable to a wider range of GP use cases.

54 **Inducing point methods [R2]:** Thanks for the pointer. We will include the reference on efficient inducing point
55 method [Burt et al., 2019] and add experiments with an increasing number of inducing points for sparse GP.