

1 We thank the reviewers for their helpful feedback.

## 2 **Motivation for variational inference and concern on dataset sizes**

3 We did not clearly describe the computational challenge when using the (inverse) Wishart process. In particular, all  
4 reviewers state that our dataset sizes are manageable with exact inference; this is incorrect. Exact inference with the  
5 (inverse) Wishart process model is *intractable for even small dataset sizes*. As far as we know, an expression for the  
6 (exact) posterior distribution does not even appear anywhere in the literature. The only options for inference are to set  
7 up an MCMC procedure (as done in previous work) or to take the variational approach we now propose. We made this  
8 lack of clarity worse in the introduction (ll. 30–31) and in Sec. 3 (ll. 81–83) where we discuss “exact posterior inference  
9 on the underlying Gaussian processes.” This was unclear wording: the aforementioned MCMC routines require an  
10 exact GP posterior inference computation (costing  $O(N^3)$  on each of the  $O(D^2)$  GPs) to be performed *during each*  
11 *step of their iterative algorithms!* So we clarify that our variational approach dramatically reduces this *per iteration*  
12 computational cost (while also having an easier implementation via black-box, gradient-based inference and not losing  
13 competitive performance). We will certainly clarify this in the paper.

14 All reviewers considered the datasets in Sec. 5 to be small, however, note that there are no previous examples in the  
15 literature of inference with the Wishart process scaling to this size (of both  $N$  and  $D$ ). We reiterate: our experiments  
16 are the largest scale we’re aware of in any work using the (inverse) Wishart process. We would like to point out the  
17 inherently overparameterized nature of this problem: We are estimating a sequence of  $N$  covariance matrices, each of  
18  $O(D^2)$  size, from a *single example*. So note that even  $D = 20$  is NOT “small”. This will always be a difficult task even  
19 for small  $N$  and  $D$ , and the problem will always be overparameterized no matter how big your dataset gets. With that  
20 said, there is no reason to believe that our variational approach would not scale to larger dataset sizes.

21 Finally, we also note that one additional benefit of the variational approach is a natural way to implement *online*  
22 *inference*. We can also mention this in the paper, though this is certainly not a point we’re trying to emphasize.

## 23 **Motivation for predicting the covariance matrices**

24 R1 & R2 requested motivation for why one would want to predict the covariances  $\Sigma_n$ . The prominent application  
25 requiring such predictions is in the construction of optimal financial trading portfolios that minimize risk ( $\rho$ ), where  
26 returns are maximized based on a separate model for  $Y_n$  (not addressed in this paper) and a model for the *covariance of*  
27 *the residuals* is used to penalize risky (i.e., volatile) assets, which is the part we are tackling in our paper. It is for this  
28 reason that we present in the context of financial applications. Moreover, these predictions of “volatility” or “risk” are  
29 desired throughout finance and beyond. R4 specifically requested motivation beyond finance: such models have been  
30 used to analyze the spread of disease incidence ( $\rho$ ) and XXX ( $\rho$ ). We will make these clarifications and additions.

31 R2 similarly questions why (large) covariance matrices would be useful without sparsity. Note that the *full* covariance  
32 matrix is required to construct an optimal financial portfolio. Sparse approximations are often imposed for computational  
33 reasons, but you will *always* construct a suboptimal portfolio with a sparse approximation. This would certainly never  
34 be desired if you can handle the computational burden.

## 35 **Originality/novelty**

36 All reviewers expressed concerns that our work appears too close to previous techniques. We re-emphasize that our work  
37 demonstrates the first alternative to MCMC for inference in the Wishart process model (improving on its computational  
38 efficiency and ease of implementation), and the first experiments scaling inference to the size of the datasets in Sec. 5.

## 39 **Other points**

40 R1 points out the apparent contradiction of “low-variance” MC gradient estimates (through reparameterization) vs. the  
41 study in Sec. 4. We should have instead said: “while such gradient estimates typically have low variance, the particular  
42 form of the Wishart process likelihood introduces computational instability that renders these estimates useless.”

43 R1 & R4 requested some discussion of wall clock runtime of the methods. We will add these to the paper.

44 We are pleased that R4 recognized the significance of our identification and resolution of the computational issues  
45 involved with applying black-box variational inference to the Wishart process case. These techniques are often applied  
46 indiscriminately in practice, and it is our hope this study (the likes of which are rarely explored in the literature) will  
47 provide a useful warning to practitioners.

48 We will follow R1’s suggestion to move the details of Sec. 3 to the supplement (R2 & R4 also point out the density/clutter  
49 and notation, which we will address); this will also free up enough space to add all the proposed clarifications.