

1 We thank the reviewers for their helpful feedback and suggestions that will significantly improve the final version of the
2 paper. The reviewers agree that the paper is well-written (R1, R3) and offers an original and significant contribution
3 to the machine learning community (R1, R2, R3). In particular, the paper proposes an efficient framework for Cox
4 processes that overcomes the limitations of the current state-of-the-art schemes (R1). The proposed technique is free
5 from the curse of dimensionality, preserves most of the dependencies in the model and does not rely on assumptions on
6 the covariance function (R3). Furthermore, the algorithm is scalable and offers competitive experimental results (R2).

7 We hope our detailed response below will further highlight the paper’s quality and originality and persuade them to
8 increase their overall scores. For this, we have attached the tag *Q5* to some of our responses, indicating that we address
9 Question 5: Improvements suggested by the reviewers that may yield to a score increase.

10 **R#1:** (1) *Factorization between the latent process and the latent locations.* We can relax this assumption by considering
11 a PPP with intensity $\int_{\mathcal{X}} \lambda^* \sigma(-f(\mathbf{x})) d\mathbf{x}$ as the joint variational distribution $q^*(M, \{\mathbf{y}_m\}_{m=1}^M)$, which is indeed the true
12 posterior [8]. In fact, we have implemented this approach (post-submission) and found out that while fully capturing
13 model dependencies, it introduces a significant computational burden due to sampling from the full approximate posterior
14 in the computation of T_3 . Empirically, this full posterior yields better uncertainty quantification but comparable point-
15 performance metrics to those reported in the paper. Increasing the computational efficiency of this new approach
16 remains an interesting research direction.

17 (2) *Q5: Higher dimensions.* We thank reviewer for suggesting testing our algorithm on higher-dimensional data. While
18 VBPP [19] does not currently support $D > 2$, we run our algorithm on the spatio-temporal Taxi dataset and found it to
19 outperform MFVB [10] both in terms of performance metrics and uncertainty quantification, e.g. $\ell_{test}[\times 10^7] : -31.26$
20 vs -42.97 . We will add this comparison in the additional page of the final version.

21 **R#2:** (1) *Clarity on marginal likelihood being optimized.* This corresponds to integrating out all latent variables in
22 Eq. 5 (after including the augmented GP prior), which is analytically intractable. However, we will show its relationship
23 to the ELBO explicitly in the final version. Many thanks for the suggestion.

24 (2) *Q5: Optimal structure of $q(\mathbf{y})$.* As mentioned above, the optimal joint distribution $q^*(M, \{\mathbf{y}_m\}_{m=1}^M)$ is a PPP with
25 intensity $\int_{\mathcal{X}} \lambda^* \sigma(-f(\mathbf{x})) d\mathbf{x}$, which we found to have comparable point-performance metrics. Critically, this fully
26 structured posterior significantly increases the computational cost. The mixture of truncated Gaussians provides a
27 flexible and computationally advantageous alternative, while satisfying the constraint of being within the domain of
28 interest. See R#1 (1) above for more details.

29 (3) *Integral in line 157:* We estimate this using Monte Carlo. As described in lines 159–163, the key to our approach,
30 which distinguishes it from previous work, is that this integral does not need to be estimated accurately, as we only
31 require it during optimization and, therefore, the quality of the posterior intensity does not depend directly on how
32 accurate this estimation is.

33 (4) *Q5: Approximate ELBO due to Stirling’s approximation.* The reviewer is correct in pointing out that the ELBO
34 claim would need to be relaxed due to the use of this approximation. However, we have found out that this term appears
35 with opposite signs in T_2 and T_4 and thus cancels out. We will clarify this in the final version but thank the reviewer for
36 the insightful comment.

37 (5) *Stochastic optimization may obfuscate results.* We first clarify that the CPU times and performances are directly
38 comparable across all methods. Our results only include one source of stochasticity due to noisy gradient estimates
39 arising from MC sampling. However, while we mention the possibility to use stochastic optimization techniques in
40 lines 179–183, we refer to the use of a second source of stochasticity due to mini-batch optimization. None of our
41 experiments actually exploit this. We will clarify this in the final version.

42 (6 & 7) *Minor edits:* Many thanks for your suggestions, we will include them in the final version and the supplement.

43 **R#3:** (1) *Standard VI.* We would like to highlight how, even though the variational inference scheme follows from
44 standard arguments, by exploiting the structure of the model and the approximate posterior we increase the algorithm
45 efficiency and avoid high-variance gradient estimates. Applying black-box variational inference naïvely to a structured
46 posterior would require sampling \mathbf{f} , λ^* , M and $\{\mathbf{y}_m\}_{m=1}^M$ thus slowing down the algorithm while leading to poor
47 convergence.

48 (2) *Q5: Clarify how superposition view helps.* The model double intractability arises first from the estimation of the
49 integral of $\lambda(\mathbf{x})$ in Eq. (1) and second from the standard posterior estimation which requires the computation of the
50 marginal likelihood. The augmentation scheme helps us with the first intractability. By superimposing two PPP with
51 opposite intensities we obtain an homogeneous PPP and thus avoid the integration of the GP over \mathcal{X} . Instead, we only
52 need to compute the measure of the input space $\int_{\mathcal{X}} d\mathbf{x}$, see Eq. (4). We will expand on this in the final version.