

1 Thanks to all of the reviewers for their time and effort, and both constructive and critical comments.

2 ▷ **All Reviewers.** General Comment 1 - code and reproducibility Due to file size limits we included our code as an  
3 anonymous downloadable link in Appendix J, but we need to make this more obvious, thanks. We gladly commit  
4 to making the code publicly available on paper acceptance. General Comment 2.1 - Marginal Likelihood (ML)  
5 approximation Our ML estimate follows the standard approach in EP, which is to assume that the site approximations  
6 are fixed. Proving that “at convergence, the parameters of the approximate factors can be considered fixed” may  
7 require analyzing the energy function of our model and a good idea for future work, thanks. General Comment 2.2 -  
8 estimate of ML approximation We use the ML of the approximate GP in our work and also in EP. A previous study on  
9 GP classification [Assessing Approximate Inference for Binary Gaussian Process Classification, Kuss & Rasmussen]  
10 thoroughly compared various ML approximations, and found that the ML approximation we use matches the predictive  
11 accuracy very well. General Comment 3 - novelty We thank Reviewer 1 and Reviewer 3 for acknowledging the  
12 non-triviality of the locality property and practical efficiency of our work. Locality is our central result, and guarantees  
13 that the complex nested optimization of the Wasserstein distance reduces to a relatively simple and efficient 1-d update.  
14 Reviewer 1 puts it well, that “replacing it [KL] by the L2 Wasserstein should strike the majority of researchers as an  
15 obvious desirable improvement”, and moreover that “Wasserstein distance is very hard to calculate. It is even harder to  
16 do approximate inference with it. A general procedure for approximate Bayesian inference by minimising some sort  
17 of Wasserstein distance to the posterior would be a large boon for the field.” General Comment 4 - numerical stability  
18 If the CDF is accessible, Equation (5), which forms the basis of our lookup tables, is stable and it avoids divergence.  
19 Intuitively, this is because the Gaussian function in the integration is bounded by 1 and when the integration variable is  
20 very large or small, the Gaussian function approaches 0 rapidly. If the CDF isn’t accessible, there are double integrations  
21 which can be computed in one pass. In such cases, the numerical stability may be problematic; investigating this would  
22 fairly deserve another paper.

23 ▷ **Reviewer 1.** We thank reviewer 1 for their supportive comments and helpful suggestions on *e.g.* the broader impact,  
24 which we will incorporate in the final version. “Found analytically for fewer distributions than EP” While we agree  
25 that the number of analytically tractable cases for QP is probably less than that of EP, various numerical schemes may  
26 be employed and in many cases made efficient using lookup tables; see also General Comment 4.

27 ▷ **Reviewer 2.** “No code is given” Please see General Comment 1. “not clear if the proposed method is worth it” Please  
28 see General Comment 3. “not clear how the estimate of the marginal likelihood is optimized”; “in EP one can show  
29 that at convergence, the parameters of the approximate factors can be considered fixed”; “the accuracy of the estimate  
30 of the marginal likelihood”. Please see General Comments 2.1 and 2.2.

31 ▷ **Reviewer 3.** “empirical advantages of the method are not demonstrated” We respectfully remind that in 8 out of 10  
32 tasks, our method achieves consistently better test log-likelihoods on repeated experiments, than the strong baseline that  
33 is EP. Besides, Figure 1.a and 1.b illustrate the effectiveness of our method in alleviating the over-estimation of variances  
34 of EP. “primary motivation and reason for pursuing QP over EP” Please see General Comment 3. “No analysis of fixed  
35 points is given” Theoretical analysis of fixed points is an interesting area for future work, thanks. At present we offer an  
36 empirical analysis of fixed points as given in our experiment section. “any modifications for future work for which you  
37 could prove convergence and analyze the fixed points?” Changes such as the double loop EP [Expectation Consistent  
38 Approximate Inference, Opper & Winther] and proving the property pointed out by Reviewer 2 that “at convergence,  
39 the parameters of the approximate factors can be considered fixed” are an interesting direction for future work, thanks.  
40 However, we cannot yet rule out that our method is *already* provably convergent under appropriate assumptions. “Is  
41 there any guidance for constructing the look-up tables” Please see General Comment 1.

42 ▷ **Reviewer 4.** “a lot of references to the Appendix ... the paper a little hard to digest” We will take the suggestion  
43 about bringing part of the supplement to the main paper, *e.g.* with regards to more explanation and discussion on WD.  
44 We will also use the extra page in the final version for this. “I’m not sure whether the page on the locality property is  
45 enlightening and really surprising” We respectfully disagree and are deeply disappointed that the reviewer places this  
46 comment in the list of ‘weaknesses’. Locality is central to our contribution, and much harder to show here than for  
47 EP. We ask the reviewer to kindly consider the broader relevance outlined in General Comment 3. “Note sure whether  
48 the authors intend to release code” Please see General Comment 1. “The degree of novelty is pretty small” Please see  
49 General Comment 3. “marginal likelihood and its accuracy” Please see General Comments 2.1 and 2.2. “numerical  
50 aspects of the L2 Wasserstein distance computations” Please see General Comment 4. “A discussion why values for  $p$   
51 different from 2 are not interesting to consider.” We briefly mention  $p$  different from 2 in *e.g.* line 72–73, 188–189 and  
52 Appendix B. These cases are interesting but also *even more* challenging to handle. “does not provide evidence ... better  
53 suiting in cases where EP has ‘deficiencies’” The specific shortcoming of EP is over-estimation of variances as pointed  
54 out on *e.g.* lines 13, 29-31; we will clarify this even further in the final version, thanks. We give both theoretical  
55 analysis of local updates mitigating over-estimation of variances (sec. 4.3) and empirical evidence that at convergence  
56 our method predicts better and with lower predictive variances.