

1 We thank the reviewers for their thoughtful comments and support. We want to emphasize that the paper provides many  
2 timely insights unified by a strong narrative around probabilistic model construction and generalization, showing the  
3 role of multimodal marginalization, neural networks priors, tempering, support, and inductive biases — with significant  
4 demonstrations, including an exhaustive empirical study of marginalization, a demonstration of the role of marginal  
5 likelihood in resolving questions around generalization, and a Bayesian resolution of double descent (tying together the  
6 opening narrative and the importance of multimodal marginalization). The material is particularly timely, given recent  
7 questions about Bayesian methods in deep learning, such as the treatment of deep ensembles as a competing approach  
8 to Bayesian neural networks, and the cold posterior experiments in Wenzel et. al (ICML 2020), which are resolved in  
9 our discussion of tempering. While papers with many types of contributions can be difficult to assess, we believe this  
10 paper makes an important and timely contribution to NeurIPS, and appreciate the strong support of reviewers.

11 Inspired by reviewer comments, we additionally evaluated the effect of deep ensembles. In the setting of Fig. 6 (c), deep  
12 ensembles achieve errors 60., 65., 68.6, 70., 70.4, 70.3, 71.2, 71.4 for widths 5, 7, 10, 15, 20, 25, 30, 50 respectively. In  
13 agreement with our Bayesian perspective of deep ensembles, they almost resolve double descent and provide similar  
14 but worse results compared to MultiSWAG. We will add the results to Figure 6.

- 15 ● **R1.** Thank you for the thoughtful and supportive review. We describe the unifying narrative and contributions above.  
16 We appreciate the feedback and we will further clarify these connections in a final version. As you say, the main  
17 contribution of the MultiSWAG part of the paper is not the algorithm but the Bayesian perspective, the exhaustive  
18 empirical study, and the resolution of double descent. We believe these are major contributions, and that the paper  
19 overall has a significant degree of novelty (which we believe should not be constrained to algorithmic innovation).
- 20 ● We will follow up with the AC on the paper you mention, but we can assure you this brief note is in no way in conflict  
21 with our submission and should not weigh in the decision.
- 22 ● L138: in some applications such as continual learning we may want to approximate the posterior over the parameters  
23 rather than the posterior predictive distribution. The Bayesian posterior predictive is not necessarily the optimal  
24 model average under model misspecification. In particular, we argue for temperature scaling in appendix E, which  
25 leads to a different predictive distribution compared to the *true* posterior.
- 26 ● Figure 3: Wasserstein distance can be easily computed from samples from the two distributions and provides a useful  
27 measure of their difference without the mode seeking or mode covering behavior of KL divergence.
- 28 ● The weights of the components in MultiSWAG should reflect the relative mass of the different modes in the posterior.  
29 We expect these masses to be similar, but estimating them empirically is an interesting direction for future work.
- 30 ● **R2.** Thank you for your thoughtful review. We think there are some misunderstandings, and hope you can consider  
31 our clarifications in your updated assessment. We would like to emphasize that the material in this paper should not  
32 be viewed as a background or tutorial on BDL, leading to the proposal of MultiSWAG. Our paper presents a novel  
33 perspective as well as truly exhaustive experiments that provide insights into BDL. The main purpose of MultiSWAG  
34 is to demonstrate properties of multi-basin marginalization, as part of the larger narrative of the paper.
- 35 ● We will add a more detailed description of MultiSWA & MultiSWAG but note that these straightforwardly ensemble  
36 independently trained SWAG and SWA models (which are described in full papers) and the point here is about  
37 multibasin marginalization. MultiSWAG has the same computational complexity as Deep Ensembles at training time,  
38 and the memory overhead comes from storing 5-20 copies of the weights for each mode (note that these weights can  
39 be stored on the disk rather than in GPU memory and take significantly less memory compared to the activations). At  
40 test time, MultiSWAG has more overhead due to samples within each basin.
- 41 ● In Fig. 3 the vague prior was used to prevent over-regularization of the network that would lead to a poor fit of the  
42 ground truth (obtained with HMC) to the data. The results do not hinge on the specific values of the hyper-parameters.
- 43 ● We want to clarify that Fig. 3 is not intended to provide a like-for-like comparison. This experiment is intended to  
44 show the importance of multibasin representations for approximate BMA integration. In particular, we do not argue  
45 against VI in general, and a multibasin ensemble of VI approximations would also support our findings. We will  
46 clarify, and include a discussion of the reference you suggested.
- 47 ● We do actually have accuracy results for Exp 2 in the Appendix, Fig. 18. We also report accuracy for double descent.
- 48 ● We include deep ensembles to our experiment on double descent as you suggested (see above).
- 49 ● **R3.** Thank you for your supportive review. We have now included a new experiment on deep ensembles performance  
50 in the context of double descent.
- 51 ● **R4.** Thanks for your supportive review. (1) We view deep ensembles as a compelling mechanism for approximate  
52 BMA integration under constraints – different from both variational methods and MCMC, which are typically  
53 combined with simple MC. In this vein, we do not view deep ensemble weights as samples from an approximate  
54 posterior, and we would not recover the exact predictive distribution by taking an infinite number of ensembles. We  
55 also provide a related discussion in Appendix C and Appendix Figure 8. We will clarify in the final version. (2)  
56 Re: Fig 3. The point here is more about multimodal vs unimodal in approximating the BMA integral than about  
57 a deficiency of variational methods. Multiple independently trained VI models (for a multi-basin posterior) could  
58 indeed perform well and also make this point. Thanks for the question. We will clarify.