

1 We thank the reviewers for their thoughtful reviews and feedback which will help us improve the paper. Overall, we  
2 found the reviews to be very positively worded, noting e.g. that the work is "comprehensive", "thorough", "impactful",  
3 "well-motivated", "clearly written", "addresses an important problem", etc. We thank the reviewers for their kind words  
4 and will address their comments below. We address common questions first and then address reviewers individually.

5 **Code release:** We would first like to share that our code and all model predictions from our experiments have been  
6 shared publicly online and are already being used by other research groups (for anonymity we will not share links here).  
7 We are also continuously refining this benchmark and have added a variational GP classifier as a last layer method.

8 *"Why does SVI perform well on MNIST but poorly on every other dataset considered?" (R2,R3):* We agree that this  
9 should be discussed in the paper. Essentially, due to a variety of reasons (initialization, lack of priors for Bayesian  
10 neural networks, posterior collapse, optimization issues) we are not aware of a SVI variant that consistently achieves  
11 competitive accuracy as well as reliable uncertainty on large-scale problems (i.e. bigger than MNIST). Indeed, even  
12 after tremendous effort and incorporating a bag of tricks from the literature (careful initialization, empirical Bayes for  
13 the prior standard deviation), we were unable to get competitive accuracy on e.g. CIFAR. We will provide more details  
14 and analysis in the paper. Apart from optimization issues, another possibility is that SVI works well for datasets with  
15 well-specified models (e.g. models on MNIST achieve more than 99% accuracy), and underperforms on datasets where  
16 model is mis-specified. See also "Bootstrap prediction and Bayesian prediction under misspecified models" (Fushiki  
17 2005) for theoretical arguments. To test this hypothesis, we will add an experiment on MNIST where we increase  
18 model mis-specification by reducing capacity (and accuracy) and compare the performance of ensemble and SVI.

19 **R1:** *"some of the baselines, seem to not be properly tuned... For example, for ImageNet and Skew Intensity 1, LL-SVI*  
20 *and Dropout underperform compared to Vanilla, ..."* In Sec A.7 we detail the tuning methodology that was used across  
21 models. However, R1 brings up a good point w.r.t. model capacity. We fixed the size of the models to have an "apples  
22 to apples" comparison, but dropout may require larger capacity models. Note, LL-SVI was trained with the rest of the  
23 model (under the ELBO) and not post-hoc which may work better. We will include these experiments in the final paper.

24 **R2:** *"In my opinion, the significance is somewhat reduced because only neural network models are used"* We empathize  
25 and agree that other models may allow for more principled uncertainty. However, given the widespread use of deep  
26 networks, we believe that understanding and benchmarking the uncertainty of this class of models is important. We  
27 struggled to find alternatives that were competitive across these tasks. Nevertheless, scaling up other models is a priority  
28 for us and we trained a variational GP on raw MNIST (incl a variety of kernels) to address this. However, we found that  
29 they were not competitive in accuracy. We will add this (and on Criteo) to the camera ready.

30 *"Should one create one's own skew and OOD datasets? Are there any principles that are important to keep in mind*  
31 *while doing that? How hard will it be to use your eventually released code to do that?"* We hope these experiments will  
32 generalize, but we believe it should be quite easy to add new datasets and experiments to our code (which has now been  
33 released). We will add details of how to do this, along with recommendations, to the paper.

34 *"... Does more capacity mean better OOD calibration?"* We discuss the tradeoffs in terms of memory and computational  
35 complexity in the appendix section A.8. However, the question of whether the additional capacity of ensembling gives  
36 an advantage is an interesting one. We will add an experiment to the paper where we train and evaluate higher capacity  
37 networks for the dropout and vanilla methods to explore this.

38 *"details about ensemble":* we used just a plain ensemble (no adversarial training, no heteroscedastic loss).

39 **R3:** *"Have the authors investigated the performance on different OOD datasets that are more / less similar to the*  
40 *source dataset...?"* We used the corruptions benchmark (Hendrycks and Dietterich, 2019) as it allows us to control the  
41 amount of similarity to source dataset, and see if consistent orderings are observed.

42 *"analysis / actionable takeaways."* This was limited due to space constraints, but we will add this to the supplement.

43 *"I found the experimental section somewhat disorganized ... would have made for a significantly better read."* Thanks  
44 for pointing this out. We will address the layout and presentation, particularly w.r.t. these figures, for the camera ready.

45 *"... While SVI does seem to have lower Brier scores with shift, accuracies don't appear to be any better – how was this*  
46 *observation made?"* Valid point. This wasn't well worded, but we meant specifically with regard to uncertainty as it's  
47 significantly better in Brier score but not accuracy - i.e. much better calibrated under significant shift. We will reword.

48 *"... What are these additional properties that ECE and entropy capture?"* Each proper scoring rule induces a calibration  
49 measure, see [Brocker, "Reliability, sufficiency, and the decomposition of proper scores", Quarterly Journal of the Royal  
50 Meteorological Society, 2009]. However, ECE is not the result of such decomposition and has no corresponding proper  
51 scoring rule; we instead chose to include ECE because it is popularly used. Each proper scoring rule is also associated  
52 with a corresponding entropy function and Shannon entropy is that for log probability [see Gneiting & Raftery, Journal  
53 of the American Statistical Association, 2007]. We will reword this and clarify our reasoning for the camera ready.