

1 We thank the reviewers for their thoughtful feedback, which will help us to improve the manuscript. We are encouraged that they
 2 found our formulation and idea to be novel(**R1,R3,R4**), intuitive (**R1,R2,R4**), impactful(**R2**), clever(**R4**), well-motivated(**R4**), and
 3 model agnostic(**R1**). We are glad they found empirical experiments to be detailed and comprehensive throughout with strong and
 4 promising results(**R1,R2,R3,R4**). All reviewers found the paper is well-written and results are reproducible. We are gratified that **R4**
 5 recognized the importance of AvUC loss "*differentiable approximation to the AvUC metric is pretty exciting*" and found that the
 6 paper resoundingly addresses the problem to get reliable predictive uncertainties. We thank the reviewers for their kind words and
 7 constructive feedback. We address reviewer comments below and will incorporate all feedback in the manuscript.

8 First, we would like to clarify that AvUC (accuracy vs uncertainty calibration) loss is devised to improve uncertainty calibration that
 9 can be **combined** with existing losses without modifying the underlying principles (e.g. ELBO for Bayesian NN, cross-entropy for
 10 non-Bayesian NN classifier). AvUC enables uncertainty calibration overcoming the challenge of unavailability of ground truth for
 11 uncertainty. AvUC accounts for quality of principled aleatoric and epistemic uncertainties, which are important for many applications.

12 **R1:Strong baseline:** Though ensemble provides higher accuracy and performs well under lower data shift intensities, we observe
 13 SVI-AvUC is more robust under high data shift intensities and out-of-distribution settings (which is common in real-world) based
 14 on the results evident from Table 1 and Fig 1,2,3. To demonstrate superiority of AvUC loss theoretically, for the camera ready
 15 version we will include a discussion section as suggested and justify how AvUC serves as loss-calibrated inference to obtain reliable
 16 uncertainties from approximate Bayesian inference (SVI) and discuss why AvUC is a proper uncertainty calibration loss in general.

17 *Application to regression problems:* Thank you for pointing us conformal prediction direction. We strongly believe our algorithm
 18 can be extended to regression problems with modifications in the formulation of loss function leveraging the relationship between
 19 prediction accuracy and uncertainty bounds. We are excited about this suggestion and we will seriously consider this for future work.

20 **R2:Comparison to PAVPU[35]:** Though our work is motivated from PAVPU metric, PAVPU is not a differentiable function to be used
 21 as a training loss and was originally proposed for uncertainty evaluation. We propose a trainable uncertainty calibration loss with
 22 differentiable approximation to AvU metric as defined in Eqns 4.To the best of our knowledge, this is the first work on uncertainty
 23 calibration that leverages the relationship between accuracy and uncertainty to train the model towards well-calibrated uncertainties.

24 **R2, R3: Theoretical justification:** Theoretically AvUC loss will be perfect 0 only when the model’s uncertainty is perfectly calibrated
 25 (AvU=1). In appendix D.1, we show how AvU score and AvUC loss are related to each other during training. As noted in Eqns
 26 3 and 4, AvUC loss attempts to maximize AvU which will indirectly push the values of uncertainties up or down based on the
 27 accuracy of predictions. When classification accuracy do not match uncertainty, $AvU \rightarrow 0$ and $\mathcal{L}_{AvUC} \rightarrow \infty$ forcing the gradient
 28 computation exert towards $\mathcal{L}_{AvUC} \rightarrow 0$, which will happen when AvU score is pushed higher ($AvU \rightarrow 1$), enabling the model to
 29 provide well-calibrated uncertainties. We will include a detailed section discussing why AvUC is a proper uncertainty calibration
 30 loss and theoretically justifying how it serves as loss-calibrated inference method (as suggested by **R4**), in the manuscript.

31 **R3:"Comparison with temperature scaling on top of SVI (SVI-TS):"** valid point; we have run experiments to compare with SVI-TS
 32 and will include the results for this method in manuscript. Results from the experiment for ResNet50/ImageNet in below table.

| method | ECE↓ [data shift intensity (test / 1 / 2 / 3 / 4 / 5)] | UCE↓ [data shift intensity (test / 1 / 2 / 3 / 4 / 5)] |
|-----------|--|--|
| SVI-TS | 0.024 / 0.030 / 0.048 / 0.072 / 0.098 / 0.123 | 0.117 / 0.164 / 0.199 / 0.235 / 0.269 / 0.294 |
| SVI-AvUTS | 0.019 / 0.027 / 0.029 / 0.041 / 0.059 / 0.080 | 0.088 / 0.144 / 0.142 / 0.173 / 0.203 / 0.226 |

33 *"Comparisons have been made only on top of SVI":* In appendix D.6, comparison of AvUTS on top of vanilla is presented. Also
 34 based on request from **R4**, we provide results of AvUC results applied to non-Bayesian NN here below (line 49-50). In the final
 35 version, we will include the evaluations on top of both SVI and Vanilla baselines (Bayesian and non-Bayesian method). In L194-198
 36 of main manuscript, we mention why we chose SVI as a baseline to evaluate our methods.

37 *Whether AvUC is a proper loss; Comparison to existing losses including MMCE[11] and cross entropy within temperature scaling:*
 38 In appendix D.1, we show how AvUC loss converges based on uncertainty calibration (AvU score). When compared to existing
 39 losses including MMCE, AvUC loss accounts for quality of principled aleatoric and epistemic uncertainties while training the model,
 40 improves uncertainty calibration, can be combined with existing losses for training or used solely for post-hoc calibration and
 41 overcomes the challenge of unavailability of ground truth uncertainty. Specifically, MMCE enables confidence calibration i.e. it
 42 accounts for the probability of the predicted class, but not overall model’s predictive uncertainty.

43 **R4:** Thanks for pointing the references on loss-calibrated inference methods, we agree including these citations will help to make
 44 the work more theoretically grounded and extendable as suggested. We will incorporate the feedback in final version. Regarding
 45 *training the NN solely through the AvUC loss*, we were not able to achieve training convergence in this case as AvUC is devised as an
 46 uncertainty calibration loss. However, we are able to perform temperature scaling solely with AvUC loss (AvUTS method).

47 *Would it be possible to train a non-Bayesian NN:* We trained vanilla baseline (ResNet20/CIFAR10) with AvUC loss in addition to
 48 cross-entropy as requested, the results are better (lower ECE and UCE) than solely training with cross-entropy. Also, in appendix
 49 D.6 we have presented the results of AvUTS applied to non-Bayesian NN (ResNet50/ImageNet).

| method | ECE↓ [data shift intensity (test / 1 / 2 / 3 / 4 / 5)] | UCE↓ [data shift intensity (test / 1 / 2 / 3 / 4 / 5)] |
|--------------|--|--|
| Vanilla | 0.046 / 0.098 / 0.139 / 0.183 / 0.236 / 0.315 | 0.038 / 0.085 / 0.122 / 0.162 / 0.212 / 0.285 |
| Vanilla-AvUC | 0.034 / 0.074 / 0.102 / 0.138 / 0.185 / 0.250 | 0.017 / 0.044 / 0.064 / 0.092 / 0.134 / 0.189 |

50 *Are the results in the paper performed over different trials?:* Yes, we used five different trails for CIFAR10 experiments (In Fig 2 and
 51 F9, shading in the plots shows the standard errors from different trails). For ImageNet experiments, results are reported over one trial.

52 *Results presentation:* Thank you for suggesting the presentation improvements of reported results, we will address these in camera
 53 ready version. We observed SVI-AvUC results show better spearman’s ρ compared to other methods indicating least correlation
 54 of calibration errors w.r.t dataset shift [ImageNet/ECE - Vanilla:1.0, Ensemble: 0.82, SVI-AvUTS: 0.94, SVI-AvUC: **0.2**]. We
 55 will report relative ranks and include spearman rank correlation coefficient covering all the baselines, metrics and datasets in the
 56 manuscript as suggested. We have quantified the distribution separation of uncertainties in Fig 3 with Wasserstein distance [Vanilla:
 57 2.73, Temp scaling: 2.90, Ensemble: 3.02, Dropout: 3.38, SVI: 3.73, SVI-AvUTS: 3.94, SVI-AvUC: **4.29**].