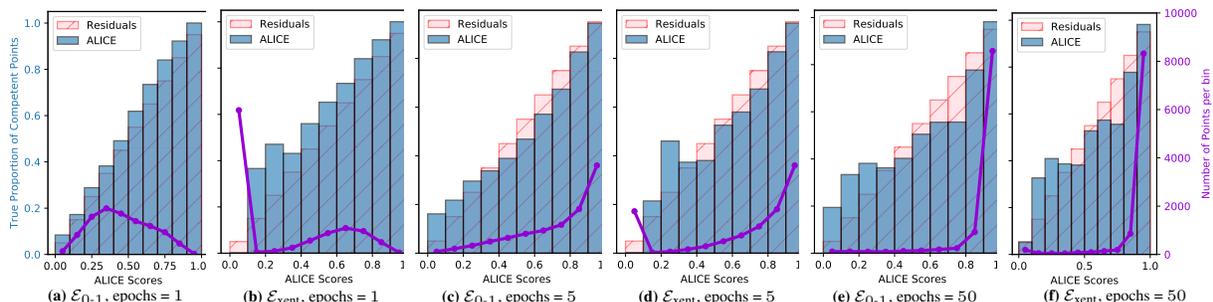


1 The authors would like to thank the reviewers for their thoughtful comments. Our responses are below:  
 2 We have replaced the calibration experiment. In this new experiment, we aim to show that the ALICE score matches  
 3 its semantic meaning: for all points with ALICE score of  $p$ , we expect  $p$  of them to be truly competent. To show this,  
 4 we bin the ALICE scores into tenths ( $[0.0 - 0.1)$ ,  $[0.1 - 0.2)$ , ...,  $[0.9, 1.0)$ ) and plot the true proportion of competent  
 5 points for each bin as a histogram. Note that a perfect competence estimation would result in these histograms roughly  
 6 resembling a  $y = x$  curve. We visualize the difference between our competence estimator and perfect competence  
 7 estimation by showing these residuals as well as the number of points in each bin in Figure 1. Note that ALICE is  
 8 relatively well-calibrated at all stages of training and for all error functions tested. We would like to make clear that all  
 9 mentions of the word *interpretable* refer to this interpretability of the ALICE score — *not* the interpretability of any  
 10 machine learning model’s predictions, as we state in Line 115.



**Figure 1:** ALICE score calibration of ResNet32 trained on CIFAR10. At 50 epochs we reach max validation accuracy. Full experimental details are in the final version.

11 We have replaced the last paragraph of section 4 for further clarity. It now reads: "Note that this metric only evaluates  
 12 how well each estimator *orders* the test points based on competence, and does not consider the actual value of the score.  
 13 We test this since some competence estimators (e.g. TrustScore) only seek to *rank* points based on competence and  
 14 do not care what the magnitude of the final score is. As a technical detail, this means that we cannot parametrize the  
 15 computation of Average Precision by  $\epsilon$  (since some estimators don't output scores in the range  $[0, 1]$ ), and must instead  
 16 parametrize each estimator's AP computation separately by thresholding on that estimator's output."

17 We have redone the distributional uncertainty experiment to follow standard out-of-distribution (OOD) detection  
 18 experiments. We train ResNet32 on CIFAR10 (in-distribution) and estimate the distributional competence (Line 241)  
 19 of images from SVHN (OOD). The results are below in Table 1. We have also revised Table 1 in the paper to have  
 20 underfit, wellfit, and overfit models for each model. The missing results are denoted below in Table 2. We have omitted  
 RF (O) since our random forest did not overfit. Further experimental details will be in the camera-ready version.

**Table 1:** mAP for Distributional Uncertainty ( $\mathcal{E} = \mathcal{E}_D$ ).

**Table 2:** Continuation of Table 1 from the paper. ( $\mathcal{E} = \mathcal{E}_{xent}$ ).

Split	Softmax	TrustScore	Abl. ALICE	ALICE	Model	Softmax	TrustScore	Abl. ALICE	ALICE
10/90	.458 ±0.056	.518 ±0.039	.100 ±0.000	<b>.868 ±0.014</b>	MLP (W)	.989 ±.005	.929 ±.044	.958 ±.042	<b>.998 ±.001</b>
30/70	.693 ±0.034	.721 ±0.026	.300 ±0.000	<b>.946 ±0.007</b>	MLP (O)	.532 ±.062	.768 ±.064	.576 ±.033	<b>.996 ±.003</b>
50/50	.816 ±0.020	.833 ±0.015	.500 ±0.000	<b>.970 ±0.003</b>	RF (W)	.998 ±.002	.898 ±.025	.923 ±.016	<b>.999 ±.000</b>
70/30	.901 ±0.010	.910 ±0.008	.700 ±0.000	<b>.985 ±0.002</b>	SVM (U)	.995 ±.003	.626 ±.046	.496 ±.069	<b>1.00 ±.000</b>
90/10	.970 ±0.003	.972 ±0.002	.900 ±0.000	<b>.997 ±0.001</b>	SVM (W)	<b>1.00 ±.000</b>	.931 ±.048	.963 ±.038	<b>1.00 ±.000</b>

22 We have modified all mentions of PAC Learning to clarify that our method is inspired, not derived, from PAC methods.

23 We have added to lines 39-41 to articulate prior work and motivate our usage of the three types of uncertainty and  
 24 the limitations of these works. It now reads: "Previous attempts to explicitly model these three factors require  
 25 out-of-distribution data, or are not scalable to high dimensional datasets or deep networks <sup>1</sup>".

26 We have edited the notation to distinguish between a finite "label space"  $\mathcal{C}$  and the associated unit simplex, or  
 27 "distributional space"  $\mathcal{Y}$ , and have revised parts of the paper to accommodate (e.g. Eq. 0, the unnumbered before Eq. 1).

28 On Line 70, we remove Eq. 1 and clarify: "The relaxation of the prediction error leads to the generalized notion of  
 29  $\delta$ -**competence**, which we define as  $p(\mathcal{E} < \delta|x, \hat{f})$ . Confidence can be recovered by setting  $\mathcal{E} = \mathcal{E}_{0-1}$  and  $\delta \in (0, 1)$ ."

30 We have revised Line 219 to specifically link model confidence and softmax together so that the role of softmax in line  
 31 246 is clear (another score that does not require ground truth to compute and can be treated as a competence estimator).

32 Several other minor comments (clarifying definitions e.g. "pointwise rankings," "inverse true error," "randomness" etc.,  
 33 clarifying ambiguity between calibration and interpretability) have been assimilated.

<sup>1</sup>Malinin, Andrey, and Mark Gales. "Predictive uncertainty estimation via prior networks." Advances in Neural Information Processing Systems. 2018.; Yarín Gal, "Uncertainty in Deep Learning", Ph.D. thesis, University of Cambridge, 2016; Kimin Lee, Honglak Lee, Kibok Lee, and Jinwoo Shin, "Training confidence-calibrated classifiers for detecting out-of-distribution samples," International Conference on Learning Representations, 2018.