

1 We thank the reviewers for their comments and actionable suggestions on improving the paper. Below we address the  
 2 most pressing concerns. We paraphrase some of the comments for brevity.

Method	Test Accuracy	ECE
Baseline	95.12	0.023
Mixup ( $\alpha=0.4$ )	96.16	0.019
Mixup ( $\alpha=1.0$ )	96.04	0.1
Label Smoothing ( $\epsilon=0.1$ )	95.51	0.089
ERL ( $\kappa=0.1$ )	95.55	0.046

(a) CIFAR-10/ResNet-18

Method	Test Accuracy	ECE
Baseline	78.28	0.049
Mixup ( $\alpha=0.5$ )	79.57	0.035
Mixup ( $\alpha=1.0$ )	79.54	0.091
Label Smoothing ( $\epsilon=0.1$ )	79.08	0.066
ERL ( $\kappa=1.0$ )	78.47	0.6

(b) CIFAR-100/ResNet-18

3 **Comment:** "Justify that the baseline models are well trained, and compare with existing baselines that use ResNet-18  
 4 for CIFAR-10/100" (R2). We provide additional results on ResNet-18 for both CIFAR-10 and 100. Our baselines (w.  
 5 no mixup) match the baseline accuracies reported in related work. We also provide the expected calibration error (ECE)  
 6 for the best performing model as well as the mixup model that used  $\alpha = 1.0$  as suggested by reviewer 2. We find  
 7 that lower  $\alpha$  gives slightly better classification and significantly better ECE. Note that ECE can be high both due to  
 8 the model being overconfident as well as under-confident, the latter being the case for  $\alpha = 1.0$  since this causes the  
 9 resulting training signal to have higher entropies than with smaller  $\alpha$ .

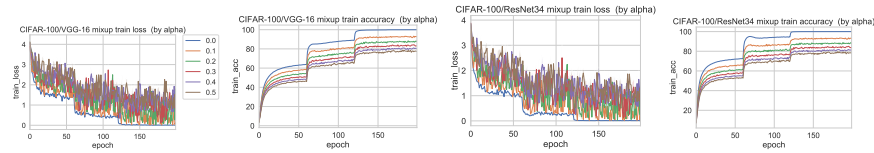


Figure 1: Train loss and accuracy for mixup for various alphas. Baseline corresponds to  $\alpha = 0$

10 **Comment:** "Provide training train loss curves for Figures 1 and 2. " (R2) We show training curves for some of the experiments  
 11 in the paper in above Figure. Your intuition is correct: for the baseline (i.e when  $\alpha = 0.$ ), over-fitting on the training  
 12 set is indeed correlated with transitioning to overconfidence. The baseline train loss and accuracy approach 0 and  
 13 100% respectively (i.e., over-fitting), while in the mixup case (non-zero  $\alpha$ 's), the strong data augmentation prevents  
 14 over-fitting and thus restricts the model from making overconfident predictions. This behavior is sustained even if one  
 15 trains for much longer (see next section)

16 **Comment:** "Will the mixup models become overconfident if trained for longer?" (R2) Below we provide the ECE vs  
 17 epoch for both CIFAR-10 and CIFAR-100 for the mixup models trained for 1000 epochs (original experiments only  
 18 used 200 epochs). We see that the mixup model, even when trained for much longer, continues to have a low calibration  
 19 error, suggesting that the mixing of data has a sustained inhibitive effect on over-fitting the training data (the training  
 20 loss for mixup continues to be significantly higher than baseline even after extended training) and preventing the model  
 21 from becoming overconfident.

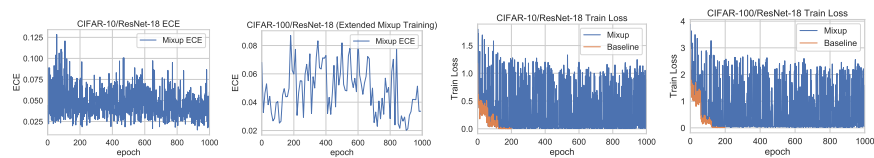


Figure 2

22 As for why mixup improves calibration (R2), please see discussion in Section 4: the strong data augmentation and label  
 23 smoothening are both contributive factors: one can view mixup training as training with infinite data (since the model  
 24 never sees the same data point twice) in which case true posteriors are learnt according to statistical learning theory,  
 25 but in addition the label softening (which prevents the winning logits from becoming arbitrarily large) also prevents  
 26 overconfidence. Note that mixup models can turn out to be *underconfident* if  $\alpha$  is large. In fact, this is also related  
 27 to manifold intrusion: a mixed-up sample is more likely to lie away from the original manifold and thus be affected  
 28 by manifold intrusion if  $\alpha$  is large. In our experiments, we see the resulting models are prone to under-fitting and  
 29 under-confidence. We will include a discussion on ROC and AUC curves for mixup in the final version (R1). As for  
 30 comparing calibration of mixup with temperature scaling (R3), this produces almost perfectly calibrated scores since it  
 31 is a post-training calibration approach. We will incorporate comparisons with model ensembles and manifold mixup  
 32 (R3) in the final version; we expect the latter to also produce well-calibrated scores since it is a generalization of mixup.