

Combining Ensembles and Data Augmentations can Harm your Calibration

Yeming Wen*, Ghassen Jerfel*, Rafael Muller, Michael W. Dusenberry
Jasper Snoek, Balaji Lakshminarayanan & Dustin Tran
presented by Yeming Wen (work done as an intern at Google brain)

April 2, 2021

- ▶ Ensembles are great to improve model accuracy and calibration.
 - ① Deep Ensembles
 - ② MC-dropout
 - ③ BatchEnsemble (*Wen et al., 2020*)
- ▶ Data augmentation is an orthogonal approach to improve model accuracy and calibration.
 - ① Random Crop and Random Flip
 - ② AugMix (*Hendrycks et al., 2020*)
 - ③ Mixup (*Zhang et al., 2018*)
- ▶ Naturally, we expect ensembles + data augmentation leads to even better calibration.
- ▶ This is not true as our experiments demonstrated.

Motivating examples

- We reported the model accuracy and calibration when various ensemble methods combined with Mixup.

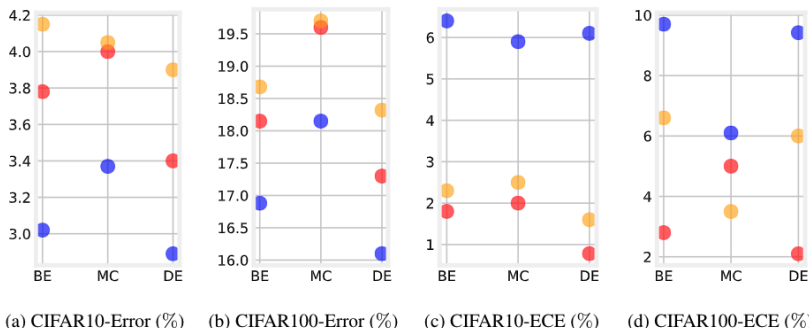


Figure 1: WideResNet 28-10 on CIFAR-10/CIFAR-100. **Red**: Ensembles without Mixup; **Blue**: Ensembles with Mixup; **Orange**: Individual models in ensembles without Mixup. **(a) & (b)**: Applying Mixup to different ensemble methods leads to consistent improvement on test accuracy. **(c) & (d)**: Applying Mixup to different ensemble methods harms calibration. Averaged over 5 random seeds.

Reliability

- ▶ We plot a variant of reliability diagrams (DeGroot and Fienberg, 1983) on BatchEnsemble.
- ▶ Ensembles (underconfidence) + Mixup (underconfidence) leads to too much underconfidence.

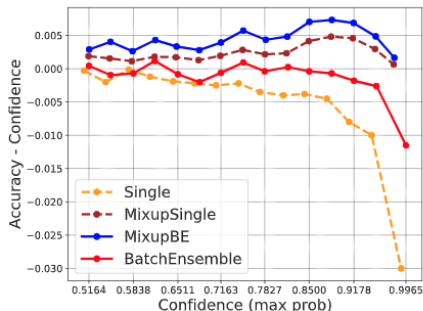


Figure 2: Reliability diagrams on CIFAR-100 with a WideResNet 28-10.

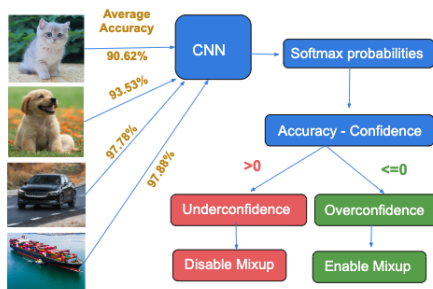
Correction

- ▶ We cannot change the underconfidence of ensembles.
- ▶ Mixup mixes all training examples.
- ▶ Ideally, we only want to mix overconfident examples.
- ▶ In this paper, we proposed to apply mixup on classes which are overconfident in the validation set.
- ▶ Denote the accuracy and confidence of class i as $\text{Acc}(C_i)$ and $\text{Conf}(C_i)$. We adjust Mixup's λ by the sign of $\text{Acc}(C_i) - \text{Conf}(C_i)$,

$$\lambda_i = \begin{cases} 0 & \text{Acc}(C_i) > \text{Conf}(C_i) \\ \lambda & \text{Acc}(C_i) \leq \text{Conf}(C_i). \end{cases} \quad (1)$$

Confidence Adjusted Mixup (CAMixup)

- ▶ We demonstrated class based confidence. One can also use forgetting counts (*Toneva et al., 2018*) during the training.



(a) Proposed CAMixup method.

Results

- CAMixup + Ensembles achieve better trade-off between model accuracy and model calibration.

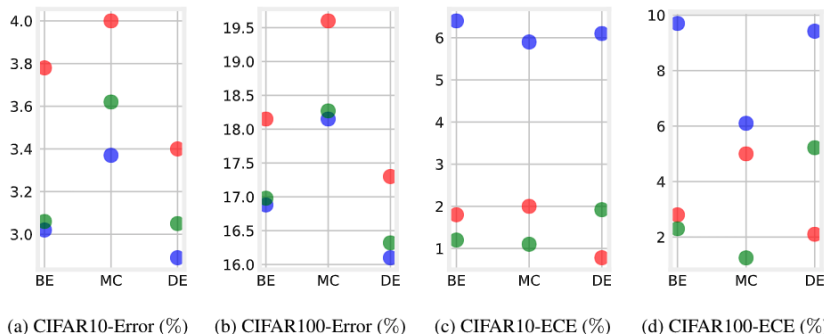


Figure 5: WideResNet 28-10 on CIFAR-10/CIFAR-100. Red: Ensembles without Mixup; Blue: Ensembles with Mixup; Green: Our proposed CAMixup improves both accuracy & ECE of ensembles.

Results

- CAMixup can be combined with AugMix. And this produces the best calibrated model.

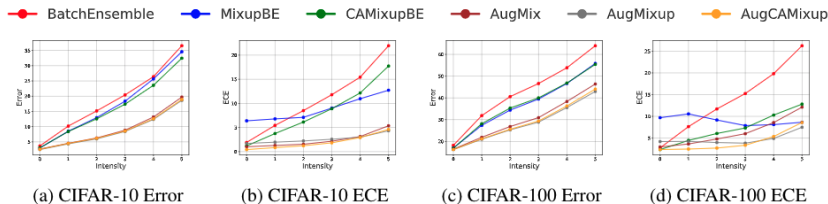


Figure 7: Performance on BatchEnsemble under dataset shift. Mixup and AugMixup improve accuracy and calibration under shift but significantly worsen in-distribution calibration. Our proposed CAMixup and AugCAMixup improve accuracy and calibration.

Method/Metric	CIFAR-10			CIFAR-100		
	Acc(↑)	ECE(↓)	cA/cECE	Acc(↑)	ECE(↓)	cA/cECE
AugMix BE	97.36	1.02%	89.49/2.6%	83.57	2.96%	67.12/7.1%
AugMixup BE	97.52	1.71%	90.05/2.8%	83.77	4.19%	69.26/4.8%
AugCAMixup BE	97.47	0.45%	89.81/ 2.4%	83.74	2.35%	68.71/ 4.4%

Conclusion

- ▶ Not all data augmentation methods can be combined with ensembles to improve calibration.
- ▶ We proposed CAMixup to fix the compounding underconfidence when combining Mixup and ensembles.
- ▶ We showed that combining CAMixup and AugMix generates the best calibrated model.