



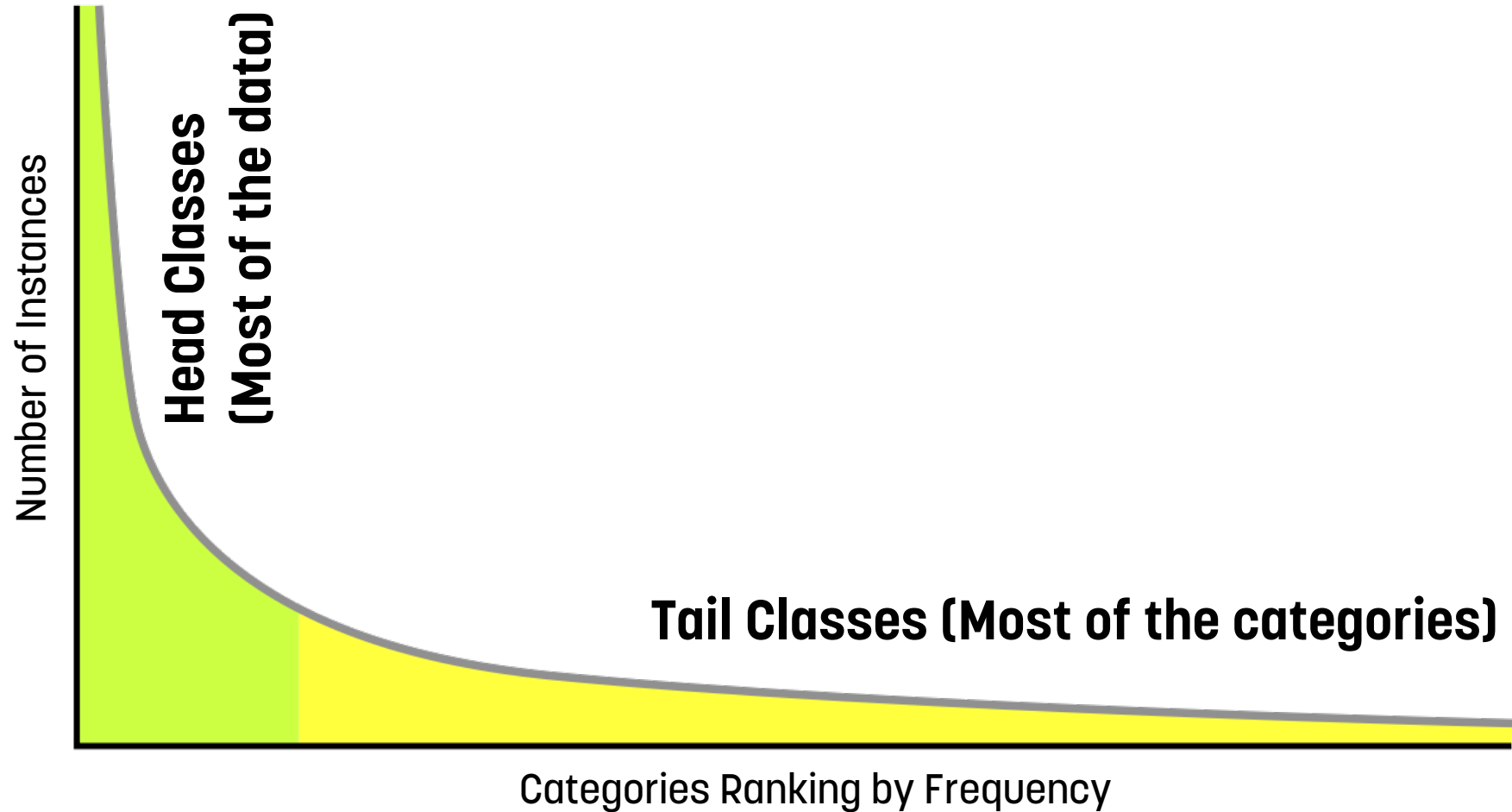
# Balanced Product of Calibrated Experts for Long-Tailed Recognition

Emanuel Sanchez Aimar<sup>1</sup>, Arvi Jonnarth<sup>1,2</sup>,  
Michael Felsberg<sup>1</sup>, Marco Kuhlmann<sup>1</sup>

<sup>1</sup>Linköping University, Sweden    <sup>2</sup>Husqvarna Group, Sweden

CVPR 2023 THU-AM-331

# What is long-tailed recognition?



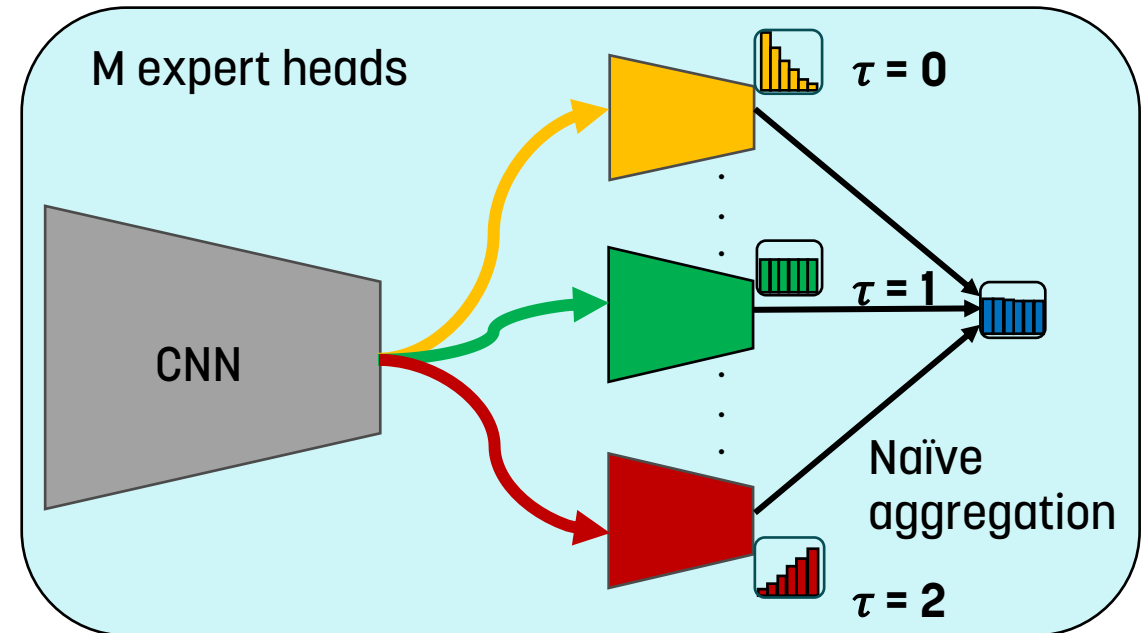
## Why is it challenging?

- **Label distribution shift:** training and test distributions tend to be different (long-tailed vs balanced)
- Leads to **uncalibrated models**

# Balanced Product of Calibrated Experts (BalPoE)

We propose a simple ensemble-based framework for LT recognition

- Diverse
- Unbiased
- Well-calibrated



- + . How to learn multiple
- experts while ensuring the  
desired properties?

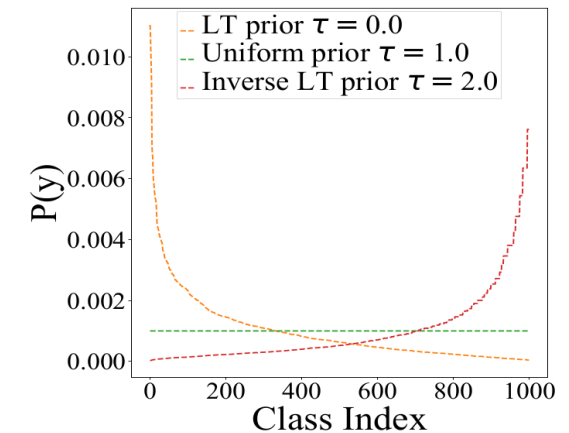
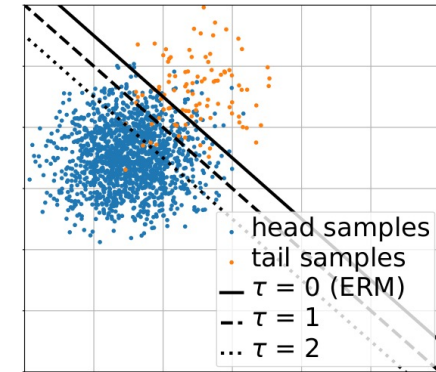
# Expert learning

- Logit-adjusted loss for  $\lambda$ -expert with  $\lambda_y \equiv 1 - \tau_y$

$$\ell_{\tau}(y, f(x)) = -\log \frac{e^{f_y(x) + \tau_y \log \mathbb{P}^{\text{train}}(y)}}{\sum_{j \in \mathcal{Y}} e^{f_j(x) + \tau_j \log \mathbb{P}^{\text{train}}(j)}}$$

- Overall ensemble loss:

$$\ell^{\text{total}}(y, \bar{f}(x)) = \frac{1}{|S_{\lambda}|} \sum_{\lambda \in S_{\lambda}} \ell_{1-\lambda}(y, f^{\lambda}(x))$$



# Expert aggregation

- Naïve fusion with uniform weights:

$$\bar{p}(x, y) \equiv \exp [\bar{f}_y(x)] \equiv \exp \left[ \frac{1}{|S_\lambda|} \sum_{\lambda \in S_\lambda} f_y^\lambda(x) \right]$$

- BalPoE attains the average bias of all its experts,
  - Controlled by  $\bar{\lambda} \equiv \frac{1}{|S_\lambda|} \sum_{\lambda \in S_\lambda} \lambda$
- **Simple constraint:**  $\bar{\lambda}$  is zero => unbiased ensemble

# Revisiting our assumptions

- Logit adjustment assumes a **bayes-optimal classifier** for the training distribution (Menon'20):

$$e^{f_y(x)} \propto \mathbb{P}^{\text{train}}(y|x)$$

- A weaker but necessary requirement is **perfect calibration** (Bröcker'09)
- We leverage Mixup (Zhang'18) to meet the calibration assumption

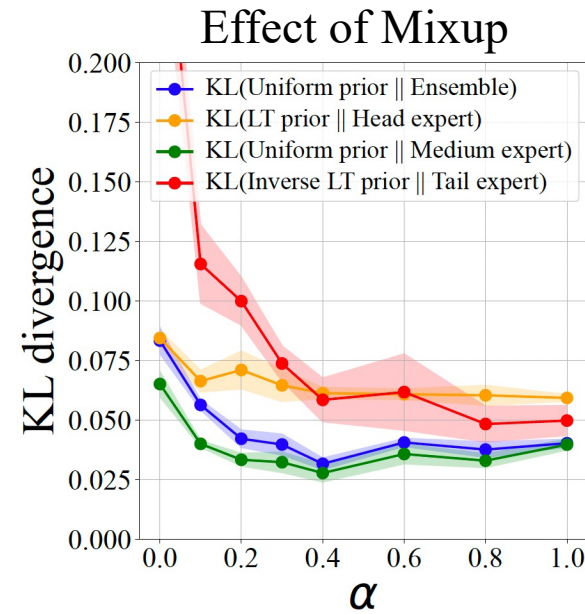
Menon et al. "Long-tail learning via logit adjustment." (ICLR 2020)

Bröcker et al. Reliability, sufficiency, and the decomposition of proper scores. (Quarterly Journal of the Royal Meteorological Society 2009).

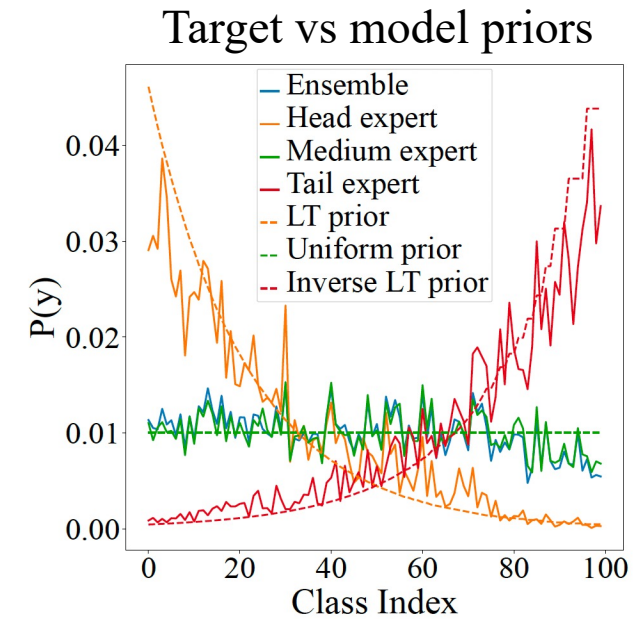
Zhang et al. Mixup: Beyond empirical risk minimization. (ICLR 2018)



Mixup  
 promotes  
 expert  
 calibration



(a)



(b)

Figure 3. (a) KL divergence of target priors vs expected marginal. (b) Target priors vs expected marginals. Marginals estimated by averaging predictions over CIFAR-100-LT-100 test set.

# Experiments



# SOTA accuracy under the balanced test set

		CIFAR-100-LT			ImageNet-LT		iNat	
		IR $\rightarrow$	10	50	100	256	256	500
Method $\downarrow$	BB $\rightarrow$	R32	R32	R32	R50	RX50	R50	R50
<i>Stronger DA</i>								
BCL [71]	(CVPR'22)	64.9	56.6	51.9	56.0	57.1	71.8	
<b>BalPoE (ours)</b>		<b>66.3</b>	<b>58.7</b>	<b>54.7</b>	<b>59.7</b>	<b>61.6</b>	<b>73.5</b>	
<i>Longer training</i>								
PaCo [12]	(ICCV'21)	64.2	56.0	52.0	57.0	58.2	73.2	
CMO+BS [48]	(CVPR'22)	<u>65.3</u>	56.7	51.7	58.0	-	74.0	
BCL [71]	(CVPR'22)	-	-	53.9	-	-	-	
NCL [38]	(CVPR'22)	-	<u>58.2</u>	<u>54.2</u>	<u>59.5</u>	60.5	<u>74.9</u>	
SADE [69]	(NeurIPS'22)	<u>65.3</u>	57.3	52.2	-	<u>61.2</u>	74.5	
<b>BalPoE (ours)</b>		<b>68.1</b>	<b>60.1</b>	<b>55.9</b>	<b>60.8</b>	<b>62.0</b>	<b>76.9</b>	

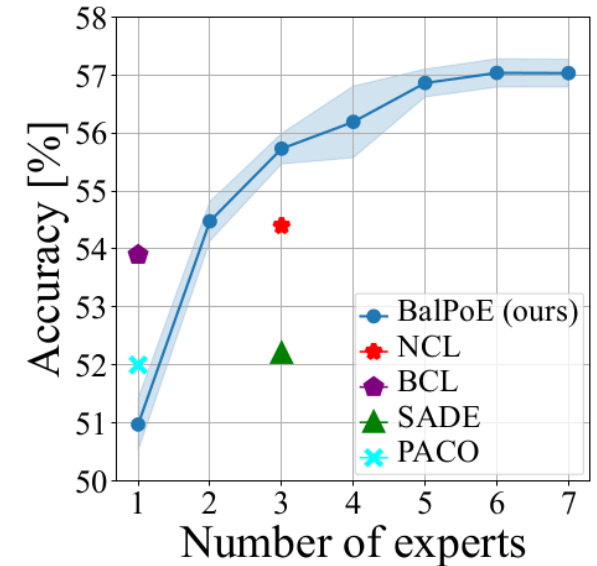


Table 1. Test accuracy (%) on CIFAR-100-LT, ImageNet-LT, and iNaturalist 2018 for different imbalance ratios (IR) and backbones (BB). Notation: R32=ResNet32, R50=ResNet50, RX50=ResNeXt50. DA denotes data augmentation. \*: reproduced results. \*\*: reproduced with mixup. †: From [66]. ‡: From [70].

# Improvements across the board

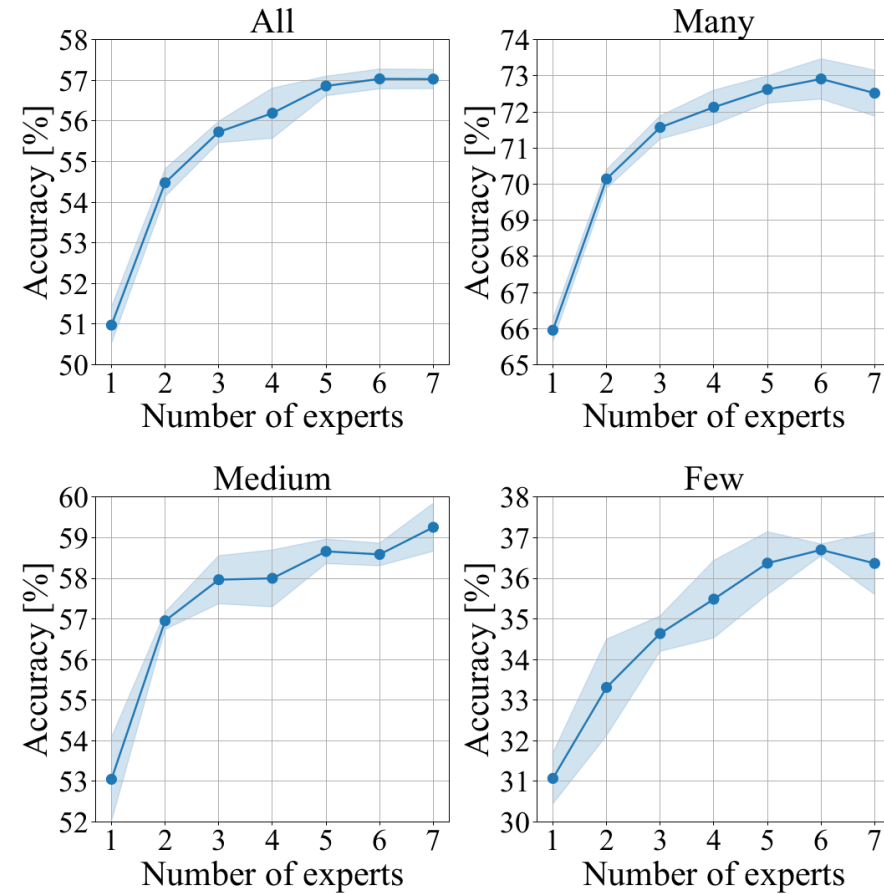


Figure 5. Test accuracy vs number of experts on CIFAR-100-LT-100 for all, many-shot, medium-shot and few-shot classes.

## Robustness under varying degrees of distribution shift

Table 3. Test accuracy (%) on multiple target distributions for CIFAR-100-LT-100 and Imagenet-LT (ResNeXt50). *Prior*: test class distribution is used. \*: Prior implicitly estimated from test data by self-supervised learning. †: results from [70].

Method	prior	CIFAR-100-LT-100					Imagenet-LT				
		Fwd-LT		Uni	Bwd-LT		Fwd-LT		Uni	Bwd-LT	
		50	5	1	5	50	50	5	1	5	50
Softmax <sup>†</sup>	✗	63.3	52.5	41.4	30.5	17.5	66.1	56.6	48.0	38.6	27.6
MiSLAS <sup>†</sup>	✗	58.8	53.0	46.8	40.1	32.1	61.6	56.3	51.4	46.1	39.5
LADE <sup>†</sup>	✗	56.0	51.0	45.6	40.0	34.0	63.4	57.4	52.3	46.8	40.7
RIDE <sup>†</sup>	✗	63.0	53.6	48.0	38.1	29.2	<b>67.6</b>	61.7	56.3	51.0	44.0
SADE	✗	58.4	53.1	49.4	42.6	35.0	65.5	62.0	58.8	54.7	49.8
<b>BalPoE</b>	✗	<b>65.1</b>	<b>54.8</b>	<b>52.0</b>	<b>44.6</b>	<b>36.1</b>	<b>67.6</b>	<b>63.3</b>	<b>59.8</b>	<b>55.7</b>	<b>50.8</b>
LADE <sup>†</sup>	✓	62.6	52.7	45.6	41.1	41.6	65.8	57.5	52.3	48.8	49.2
SADE	*	65.9	54.8	49.8	44.7	42.4	69.4	63.0	58.8	55.5	53.1
<b>BalPoE</b>	✓	<b>70.3</b>	<b>59.3</b>	<b>52.0</b>	<b>46.9</b>	<b>46.1</b>	<b>72.5</b>	<b>64.6</b>	<b>59.8</b>	<b>57.2</b>	<b>56.9</b>

# Confidence Calibration

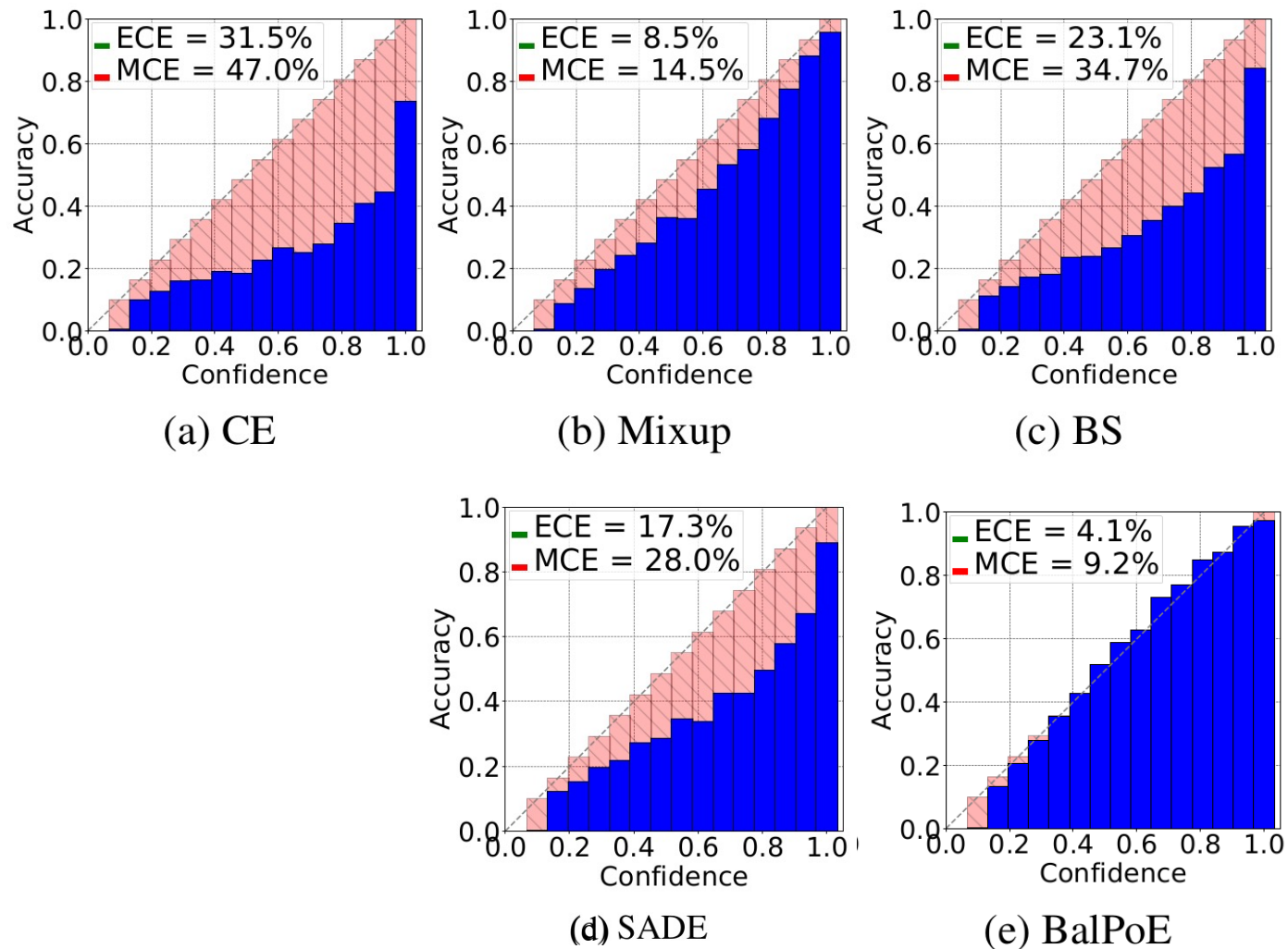


Figure 2. Reliability plots for (a) CE, (b) mixup, (c) BS, (d) SADE, and (e) BalPoE. Computed over **CIFAR-100-LT-100** test set. We also report the Expected Calibration Error (ECE) and Maximum Calibration Error (MCE).



SOTA  
 calibration  
 under the  
 balanced  
 test distribution

Table 2. Expected calibration error (ECE), maximum calibration error (MCE), and test accuracy (ACC) on CIFAR-100-LT-100.  $\star$ : Our reproduced results.  $\dagger$ : from [65].  $\ddagger$ : trained with ERM.

Method $\downarrow$	CIFAR-100-LT-100		
	ECE $\downarrow$	MCE $\downarrow$	ACC $\uparrow$
CE $\star$	32.0 $\pm$ 0.4	47.3 $\pm$ 1.8	38.8 $\pm$ 0.6
Bayias [65]	24.3	39.7	43.5
TLC [37]	22.8	-	49.0
<b>BalPoE<math>\ddagger</math></b> (ours)	17.6 $\pm$ 0.4	28.9 $\pm$ 0.9	49.2 $\pm$ 0.5
Mixup $\star$ [68]	9.6 $\pm$ 0.8	15.9 $\pm$ 1.5	40.8 $\pm$ 0.7
Remix $\dagger$ [9]	33.6	51.0	41.9
UniMix+Bayias [65]	23.0	37.4	45.5
MiSLAS [70]	<b>4.8</b>	-	47.0
<b>BalPoE</b> (ours)	<b>4.9<math>\pm</math>1.0</b>	<b>11.3<math>\pm</math>1.6</b>	<b>52.0<math>\pm</math>0.5</b>

# Conclusions

- We extend the logit adjustment formulation for training a balanced product of experts
- We show that the ensemble is Fisher-consistent for the balanced error, given that a simple constraint is fulfilled
- We find that expert calibration is important for achieving an unbiased ensemble using naïve expert aggregation



+  
•  
◦

# Thank you for listening!

