

## Take-Home Message

### Proposed method:

- We show that the underlying cause of miscalibration in adaptation is with the increase of logit ranges and demonstrated that the zero-shot baselines are better calibrated.
- We provide two solutions (normalization, penalty) during training and an unsupervised scaling during inference time to constrain the logit range based on the zero-shot logits.

### Results:

- Our solutions reduce miscalibration error in popular OOD classification benchmarks for both adapters and prompt learning while keeping the discriminative performance.
- Incorporating our approaches decreases the logit range with typical increase in logit norm.

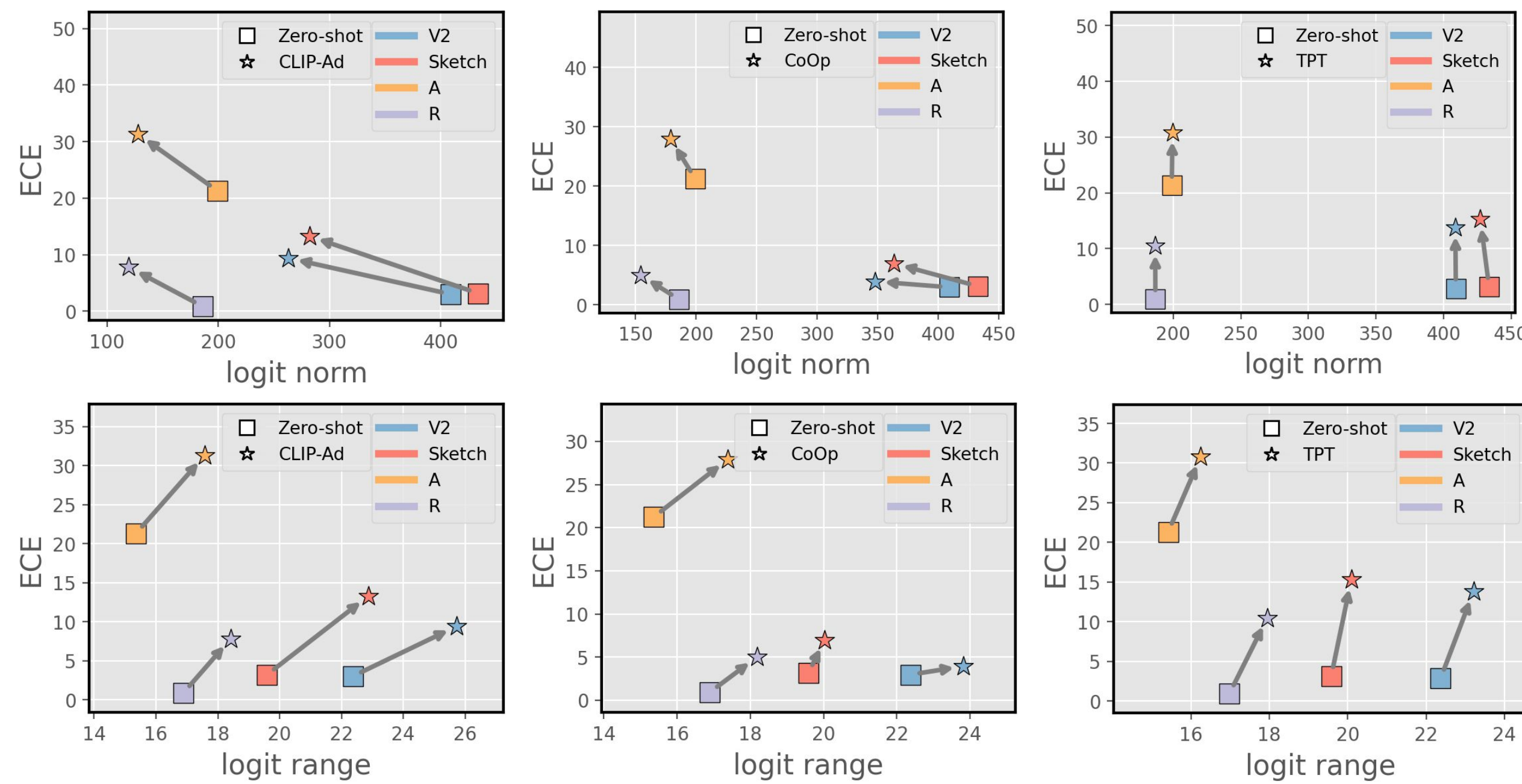
## Introduction

### Motivation:

- Deep learning is undergoing a paradigm shift with pre-trained large-scale language-vision models, such as CLIP [1].
- Adapters [2], Prompt Learning [3], and TPT [4] methods have been developed to generalize for unseen related-domains.
- These methods have improved the discriminative performance of a zero-shot baseline, but calibration is significantly degraded.

### Background and observations:

- Recent literature [5] suggests that the miscalibration is caused by increasing the logit norm during training.
- We expose that the underlying cause of miscalibration is, in fact, the increase of the logit ranges instead of norm.



### Contributions:

- We empirically demonstrate that popular CLIP adaptation strategies, substantially degrade the calibration capabilities of the zero-shot baseline in the presence of distributional drift.
- We present a simple, and model-agnostic solution, scaling the logit range of each sample based on the zero-shot logits.
- Comprehensive experiments on popular OOD classification benchmarks demonstrate the effectiveness of our approaches.

## Method

### Formulation

The logits used in training the main objective  $\mathcal{H}(\mathbf{Y}, \mathbf{P})$  are constrained to the range of its zero-shot prediction by the following constrained problem:

$$\begin{aligned} &\text{minimize} && \mathcal{H}(\mathbf{Y}, \mathbf{P}) \\ &\text{subject to} && l_i^{\text{ZS-min}} \mathbf{1} \leq \mathbf{l}_i \leq l_i^{\text{ZS-max}} \mathbf{1} \quad \forall i \in \mathcal{D} \end{aligned}$$

where  $\mathbf{l}_i$  is the logit magnitude of sample  $\mathbf{x}_i$ , and  $l_i^{\text{ZS-min}}$  and  $l_i^{\text{ZS-max}}$  are the min and max logit magnitudes of its zero-shot prediction.

### Sample-adaptive logit scaling (SaLS)

The logit normalization of sample  $\mathbf{x}_i$  at inference time is given by:

$$\mathbf{l}'_i = \frac{(l_i^{\text{ZS-max}} - l_i^{\text{ZS-min}})}{(l_i^{\text{max}} - l_i^{\text{min}})} (\mathbf{l}_i - l_i^{\text{min}} \mathbf{1}) + l_i^{\text{ZS-min}} \mathbf{1}$$

where  $l_i^{\text{max}} = \max_j(l_{ij})$  and  $l_i^{\text{min}} = \min_j(l_{ij})$

### Zero-shot logit normalization during training (ZS-Norm)

The learning objective with normalized logit ( $\mathbf{l}'_i$ ) is given by:

$$\mathcal{H}(\mathbf{Y}, \mathbf{P}) = - \sum_{i \in \mathcal{S}} \sum_{k=1}^K y_{ik} \log \frac{\exp(l'_{ik})}{\sum_{j=1}^K \exp(l'_{ij})}$$

where  $\mathbf{l}'_i$  denotes the zero-shot normalized logit vector of  $\mathbf{x}_i$

### Integrating explicit constraints in the objective (Penalty)

The objective function with ReLU penalties is given by:

$$\min_{\theta} \mathcal{H}(\mathbf{Y}, \mathbf{P}) + \lambda \sum_{i \in \mathcal{S}} \sum_{k=1}^K (\text{ReLU}(l_{ik} - l_i^{\text{ZS-max}}) + \text{ReLU}(l_i^{\text{ZS-min}} - l_{ik}))$$

where  $\lambda$  controls the trade-off between the main loss and penalties.

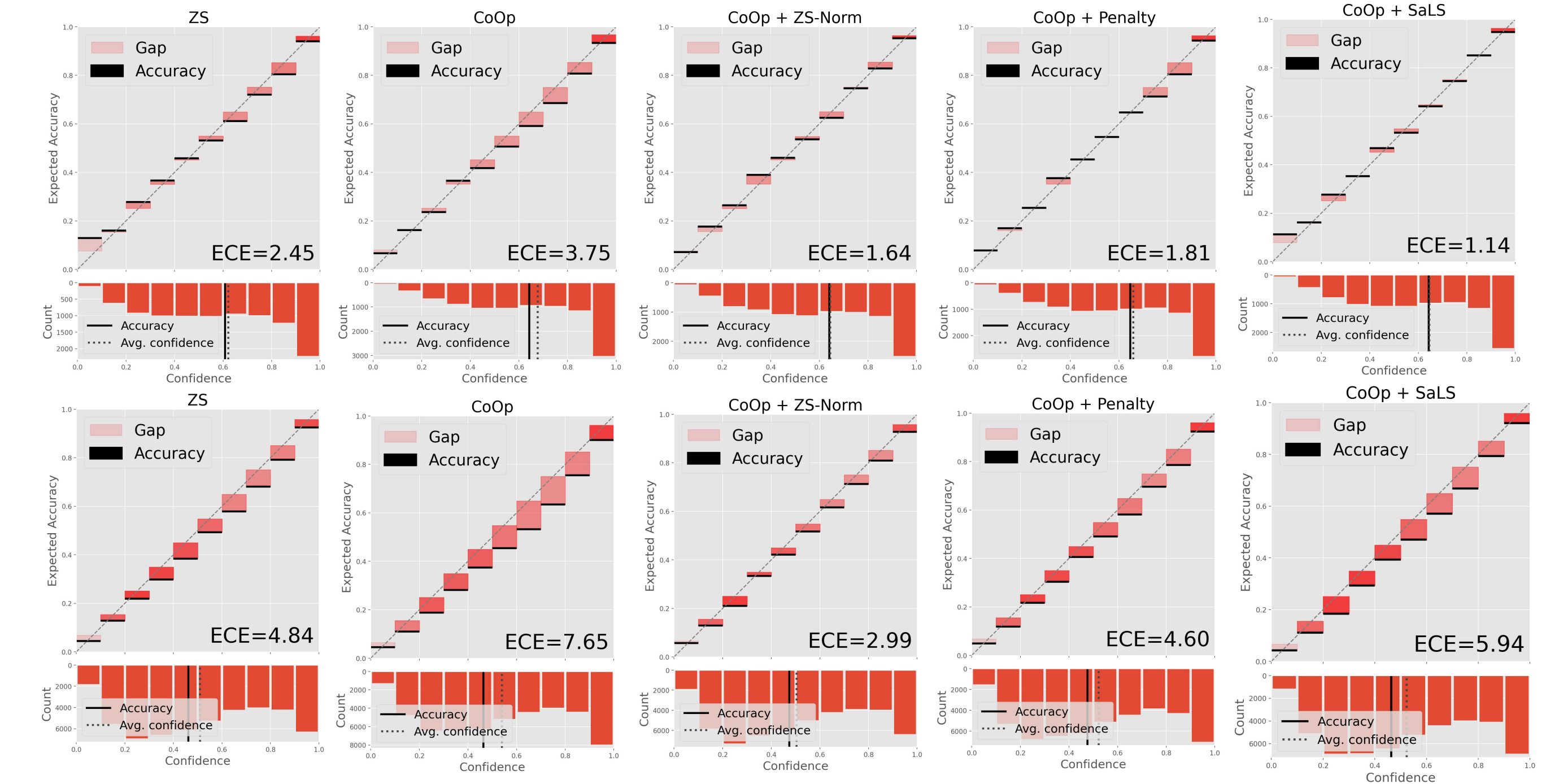
**Repository** : <https://github.com/Bala93/CLIPCalib>

## Results

**(1) Quantitative performance:** Logit range scaling provided improved calibration on different prompt learning methods.

Method	Avg. ACC	OOD ECE
Zero-Shot	57.15	4.78
CoOp	58.41	6.61
w/ <b>ZS-Norm</b>	58.75 <sub>(+0.34)</sub> ↑	<b>4.35</b> <sub>(−2.26)</sub> ↓
w/ <b>Penalty</b>	<b>59.18</b> <sub>(+0.77)</sub> ↑	4.91 <sub>(−1.70)</sub> ↓
w/ <b>SaLS</b>	58.41	4.90 <sub>(−1.71)</sub> ↓
CoCoOp	59.74	4.83
w/ <b>ZS-Norm</b>	59.90 <sub>(+0.16)</sub> ↑	3.94 <sub>(−0.89)</sub> ↓
w/ <b>Penalty</b>	<b>60.20</b> <sub>(+0.46)</sub> ↑	<b>3.89</b> <sub>(−0.94)</sub> ↓
w/ <b>SaLS</b>	59.74	4.81 <sub>(−0.00)</sub> ~
MaPLe	60.07	4.13
w/ <b>ZS-Norm</b>	60.09 <sub>(+0.02)</sub> ↑	<b>3.59</b> <sub>(−0.14)</sub> ↓
w/ <b>Penalty</b>	<b>60.62</b> <sub>(+0.55)</sub> ↑	3.78 <sub>(−0.35)</sub> ↓
w/ <b>SaLS</b>	60.07	4.38 <sub>(+0.25)</sub> ↑

**(2) Qualitative results:** Reliability plot for Prompt learning method CoOp with ImageNetV2 (Top), and ImageNetSketch (Bottom).



## References

- [1] Radford et al. Learning visual models from natural language supervision. In *ICML*, 2021.
- [2] Gao et al. Clip-adapter: Better vision-language models with feature adapters. *IJCV*, 2024.
- [3] Zhou et al. Learning to prompt for vision-language models. *IJCV*, 2022.
- [4] Shu et al. Test-time prompt tuning for vision-language models. *NeurIPS*, 2022.
- [5] Wei et al. Mitigating neural network overconfidence with logit normalization. In *ICML*, 2022.