Queen Mary University of London

# Semi-Supervised Learning under Class Distribution Mismatch

Yanbei Chen[1]   Xiatian Zhu[2]   Wei Li[1]   Shaogang Gong[1]

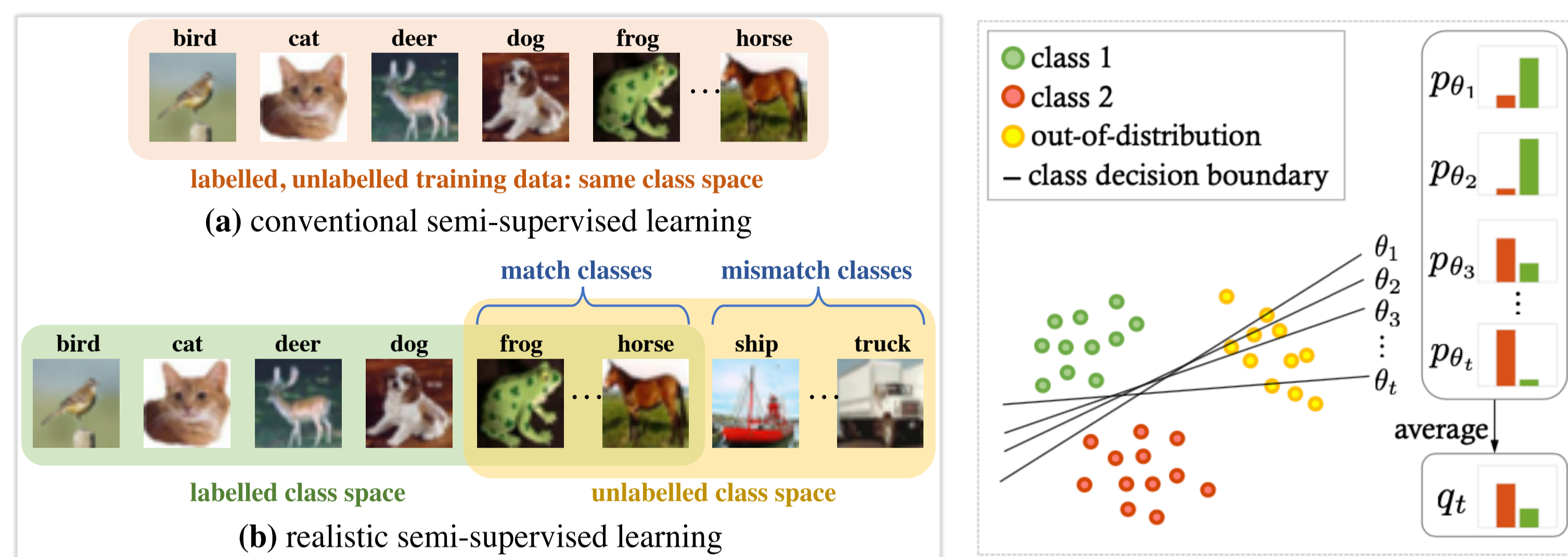yanbei.chen@qmul.ac.uk   eddy.zhuxt@gmail.com   w.li@qmul.ac.uk   s.gong@qmul.ac.uk

[1]Queen Mary University of London   [2]Vision Semantics Ltd., London, UK

## Introduction

### Problem

- Semi-supervised learning (SSL) aims for model optimisation with limited labelled data and abundant unlabelled data.
- In conventional SSL, the labelled and unlabelled data sets are assumed to come from an identical class distribution.
- **In realistic SSL, class distribution mismatch often exists** between two sets. We consider this realistic SSL challenge.



(a) conventional semi-supervised learning
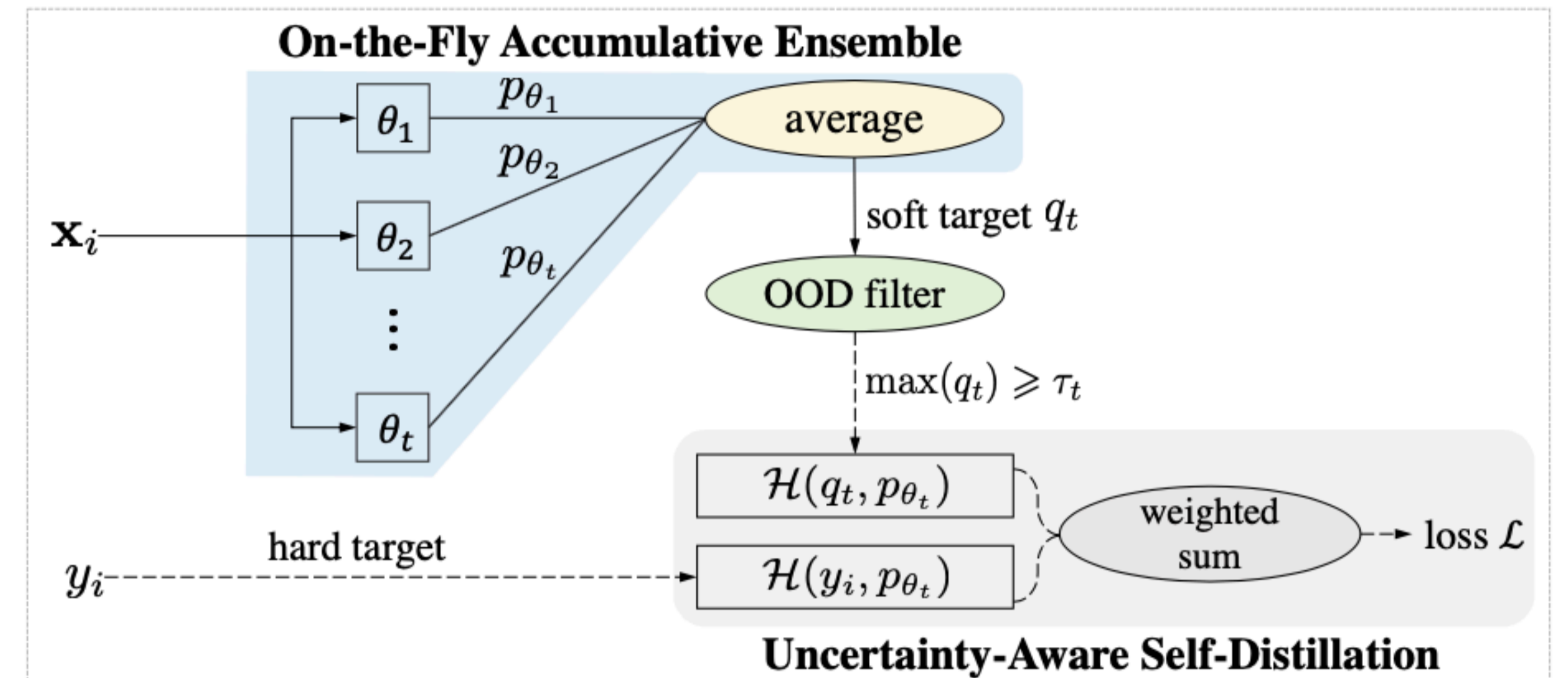
(b) realistic semi-supervised learning

### Key Contribution

- A novel **Uncertainty-Aware Self-Distillation (UASD)** formulation, which accumulatively aggregates network predictions on-the-fly for joint *Self-Distillation* and *Out-of-Distribution (OOD) Filtering*. Our formulation is aware of the *uncertainty* of whether an unlabelled sample likely lies in- or out-of-distribution, and selectively learns from the unconstrained unlabelled data.

## Method Overview

### Method



- On-the-Fly Accumulative Ensemble

$$q_t(y|\mathbf{x}_i) = \frac{1}{t}\sum_{j=0}^{t-1} p(y|\mathbf{x}_i;\theta_j)$$

- Unlabelled Training Data Filtering
  - derive a predictive confidence score on each sample
  $$c_t(\mathbf{x}_i) = \max(q_t(y|\mathbf{x}_i))$$
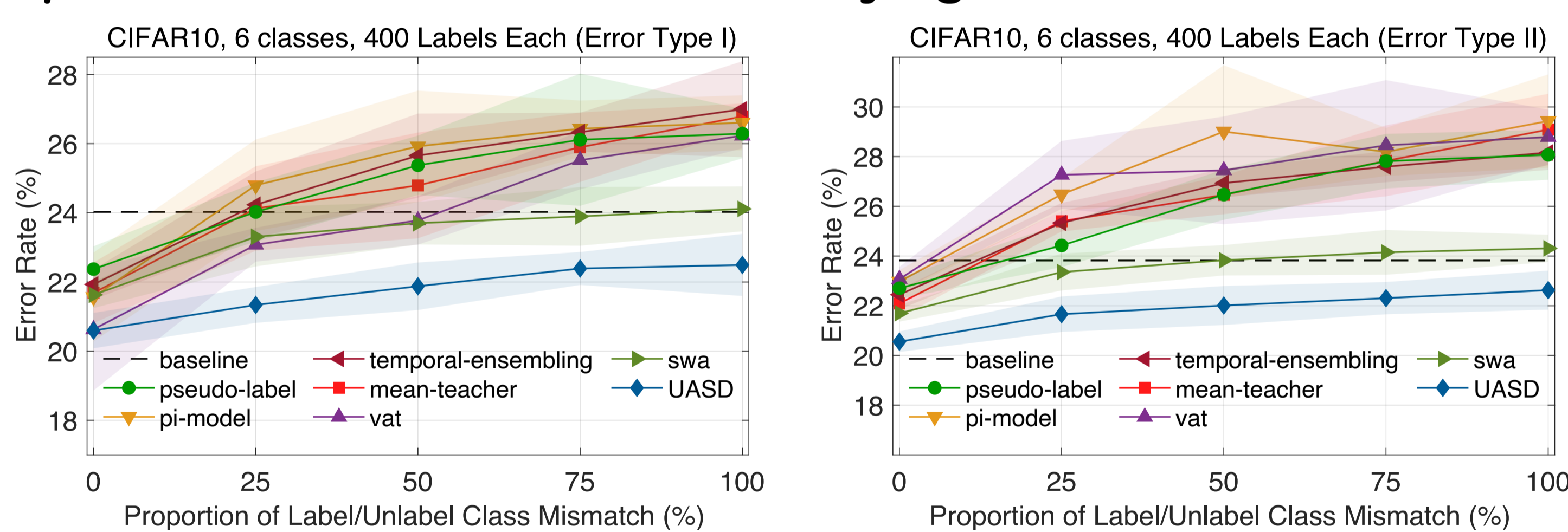  - define an OOD filter to discard samples with low confidence
  $$f(\mathbf{x}_i;\tau_t) = \begin{cases} 1, & \text{if } c_t(\mathbf{x}_i) \geqslant \tau_t, \text{ selected} \\ 0, & \text{if } c_t(\mathbf{x}_i) < \tau_t, \text{ rejected} \end{cases}$$
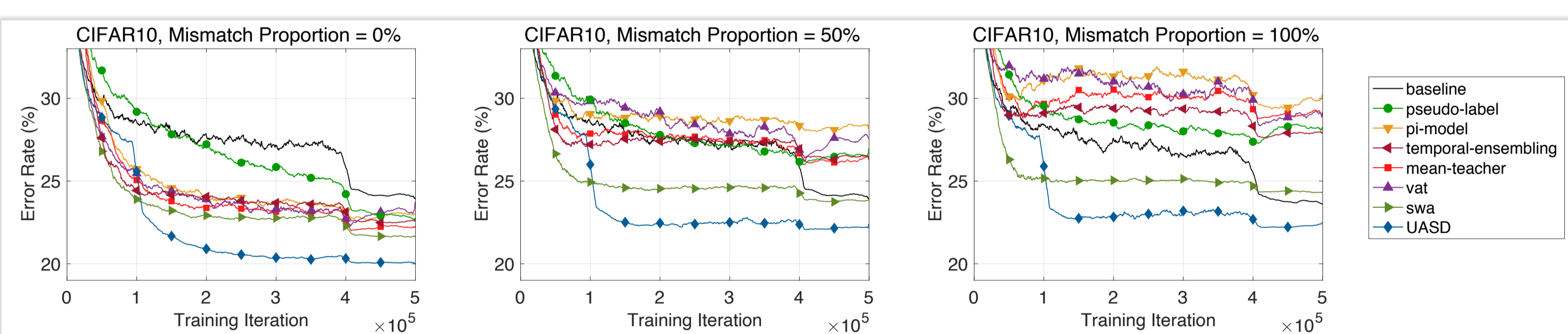
- Uncertainty-Aware Self-Distillation
$$\mathcal{L} = \mathcal{H}(y_{\text{true}}, p_\theta) + w(t)f(\cdot;\tau_t) \cdot \mathcal{H}(q_t, p_\theta)$$

## Experiments

### Experiments on CIFAR10 under varying class mismatch rate



*Left*: test errors with lowest validation errors. *Right*: the median of test errors in last 20 epochs.



*Smoothed learning curves* averaged over five runs under different class mismatch rate (0/50/100%).
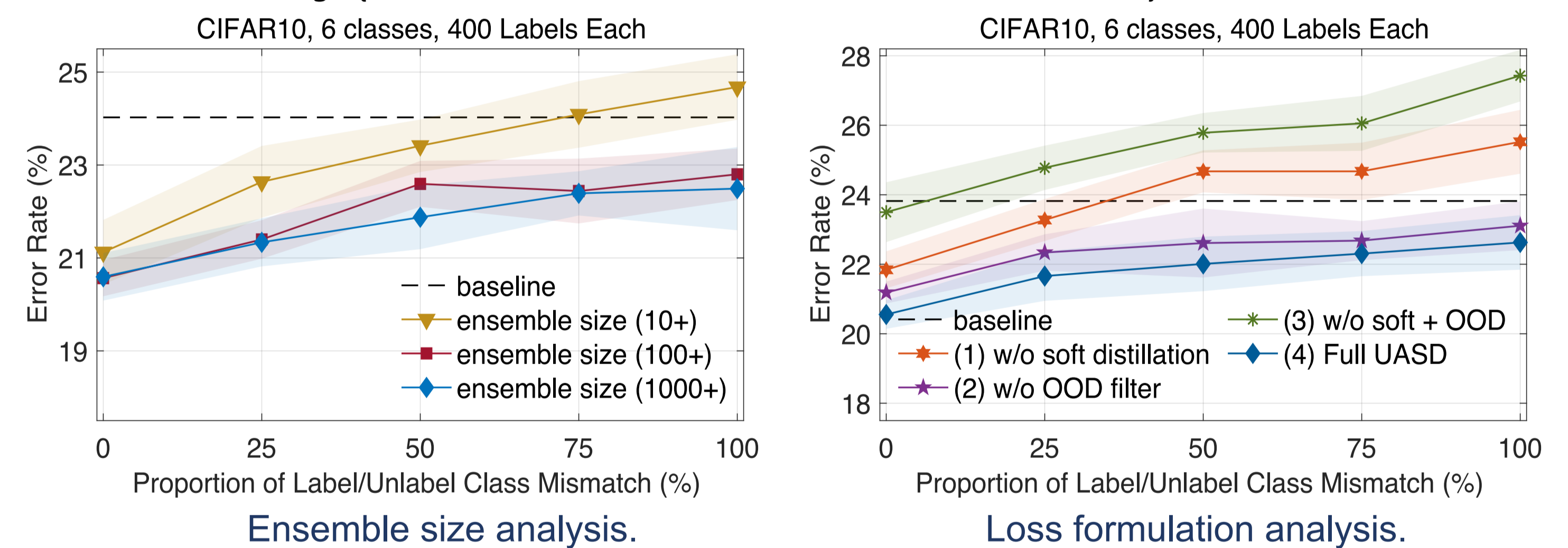
### Experiments on CIFAR100, TinyImageNet and Cross-Dataset

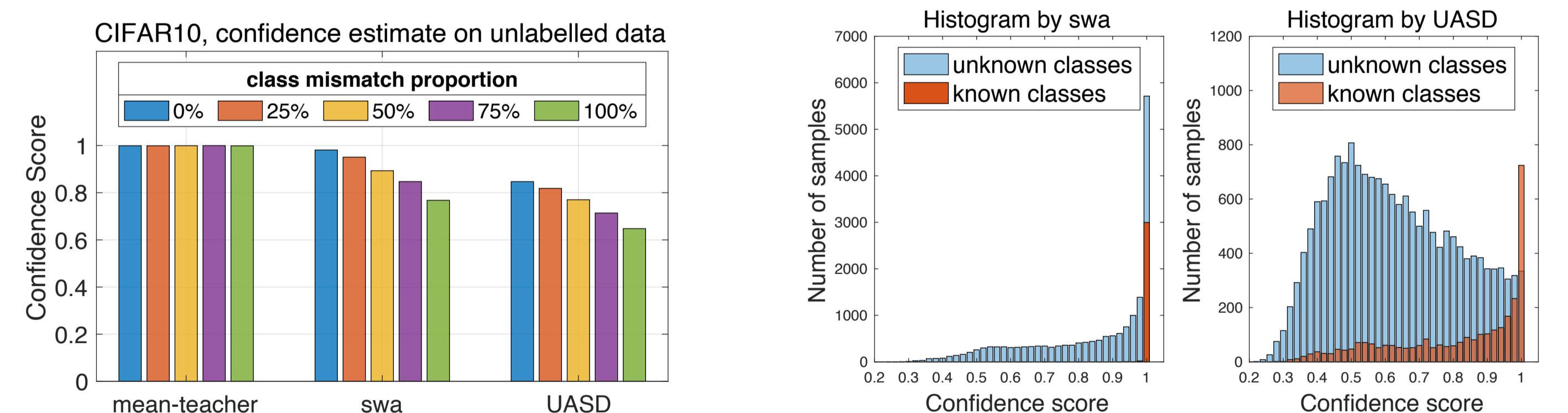| Method | CIFAR100 | TinyImageNet | CIFAR100 + TinyImageNet |
|---|---|---|---|
| baseline | $39.79 \pm 1.19$ | $61.64 \pm 0.59$ | $48.31 \pm 0.63$ |
| pseudo-label | $43.30 \pm 0.57$ | $62.41 \pm 0.57$ | $53.3 \pm 0.73$ |
| VAT | $43.78 \pm 1.15$ | $63.75 \pm 0.69$ | $50.55 \pm 0.55$ |
| Π-Model | $42.96 \pm 0.46$ | $61.79 \pm 0.67$ | $53.05 \pm 2.21$ |
| Temporal Ensembling | $41.27 \pm 0.76$ | $60.69 \pm 0.31$ | $47.88 \pm 0.64$ |
| Mean-Teacher | $40.98 \pm 0.98$ | $60.54 \pm 0.31$ | $49.67 \pm 1.95$ |
| SWA | $37.66 \pm 0.48$ | $57.97 \pm 0.42$ | $44.61 \pm 0.52$ |
| Ours | $35.93 \pm 0.60$ | $57.15 \pm 0.76$ | $42.83 \pm 0.25$ |

*CIFAR100, TinyImageNet* mismatch rate: 50%; *CIFAR100+TinyImageNet* mismatch rate: 86.5%.

## Further Analysis

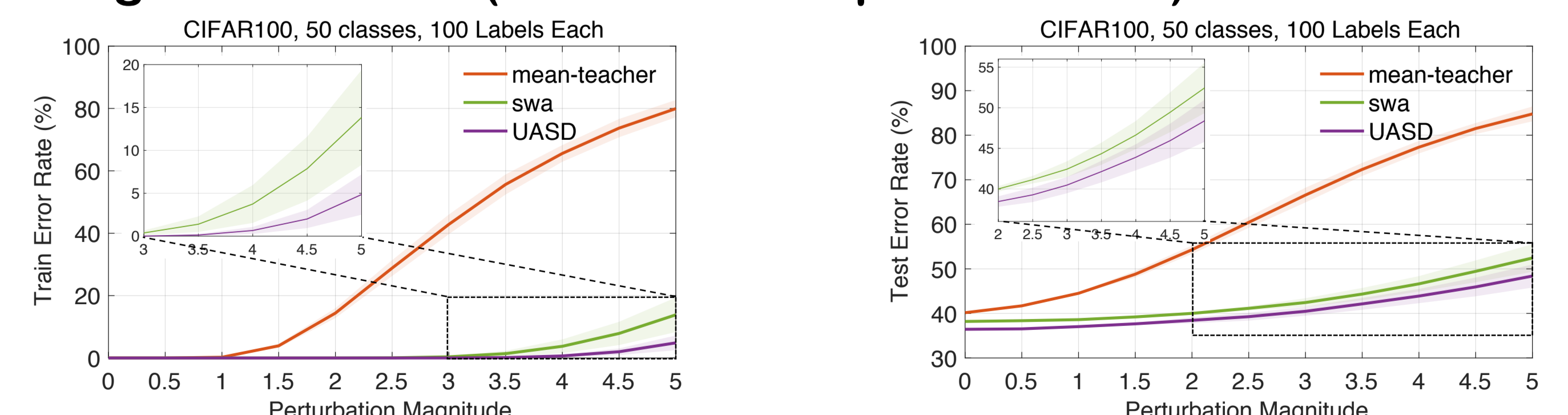### Ablation Study (Ensemble size & Loss formulation)



Ensemble size analysis.

Loss formulation analysis.

### Confidence Calibration



**Conclusion:** Confidence scores estimated by UASD can better delimit known and unknown.

### Model generalization (robustness to perturbation)



## Reference

[1] Realistic evaluation of deep semi-supervised learning algorithms. Oliver, A.; Odena, A.; Raffel, C. A.; Cubuk, E. D.; and Goodfellow, I. NeurIPS2018.

[2] Simple and scalable predictive uncertainty estimation using deep ensembles. Lakshminarayanan, B.; Pritzel, A.; and Blundell, C. NeurIPS2017.

[3] There are many consistent explanations of unlabeled data: Why you should average. Athiwaratkun, B.; Finzi, M.; Izmailov, P.; and Wilson, A. G.  ICLR2019

[4] Temporal ensembling for semi-supervised learning. Laine, S., Aila, T. ICLR2017

[5] Mean teachers are better role models. Tarvainen, A., and Valpola, H. NeurIPS2017.