

# Fine-Tuning is Fine, if Calibrated

Zheda Mai<sup>\*1</sup>, Arpita Chowdhury<sup>\*1</sup>, Ping Zhang<sup>\*1</sup>, Cheng-Hao Tu<sup>1</sup>, Hong-You Chen<sup>1</sup>, Vardaan Pahuja<sup>1</sup>, Tanya Berger-Wolf<sup>1</sup>, Song Gao<sup>2</sup>, Charles Stewart<sup>3</sup>, Yu Su<sup>1</sup>, Wei-Lun Chao<sup>1</sup>

<sup>1</sup>The Ohio State University, <sup>2</sup>University of Wisconsin Madison, <sup>3</sup>Rensselaer Polytechnic Institute

## Highlights

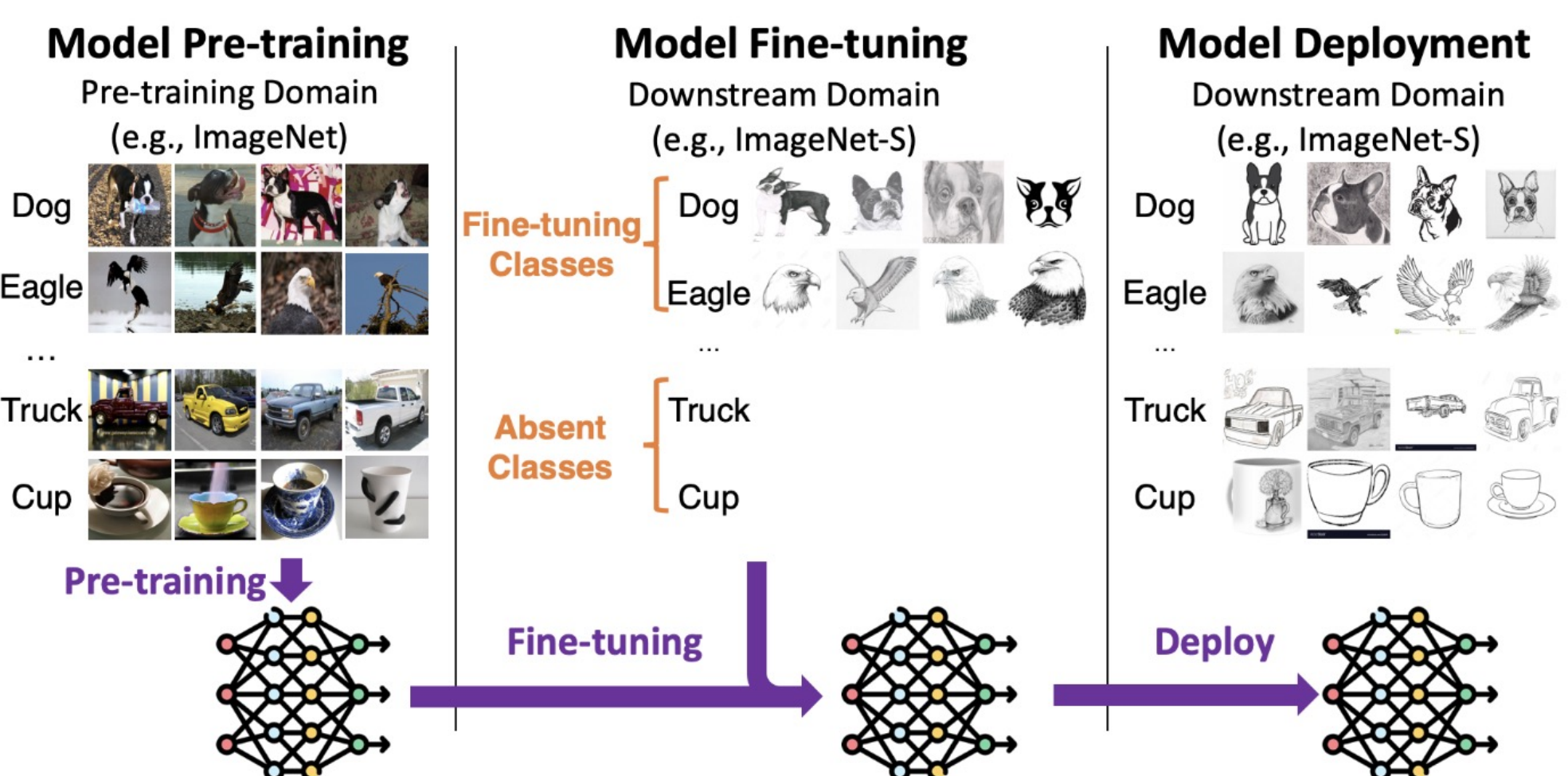
- Fine-tuning (FT) a **pre-trained** model that recognized many classes with only a **subset** of classes from a **new domain** significantly **degrades** absent class accuracies.
- Accuracy degradation does **NOT** come from **feature deterioration**, but stems from **biased logits** toward fine-tuned classes.
- Simple **post-processing calibration** restores model's capabilities and reveals improved features.

## Motivation

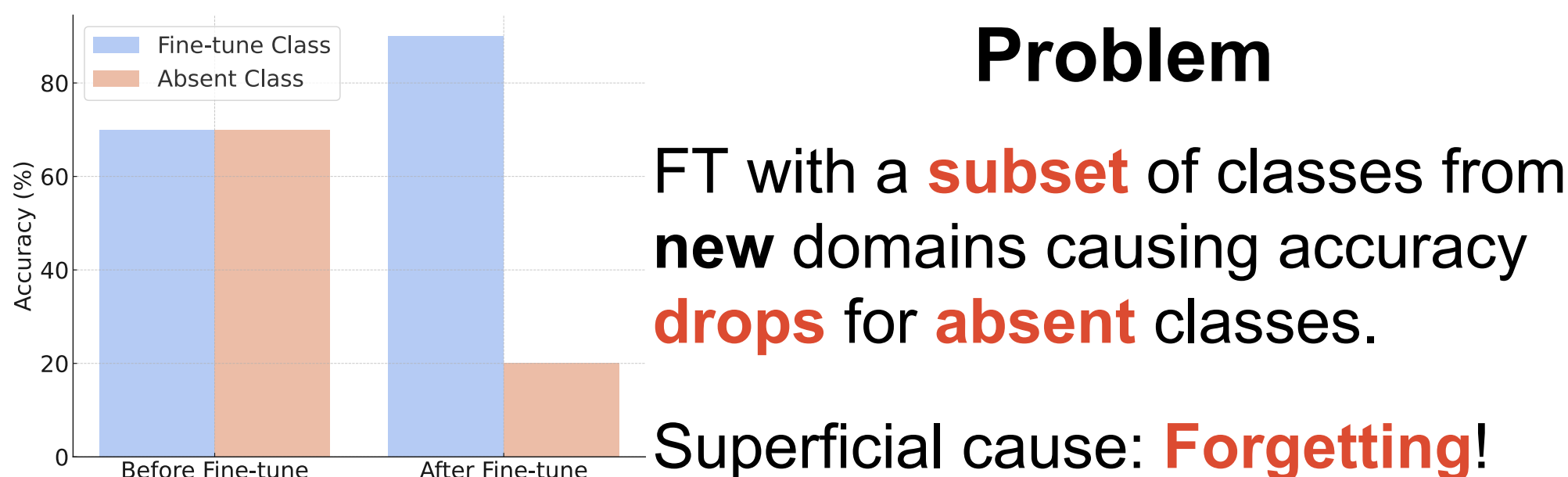
Fine-tuning pre-trained models on a **subset** of target classes is often unavoidable, e.g., camera trap, some species may **NOT** appear within the collection time.



## Setting

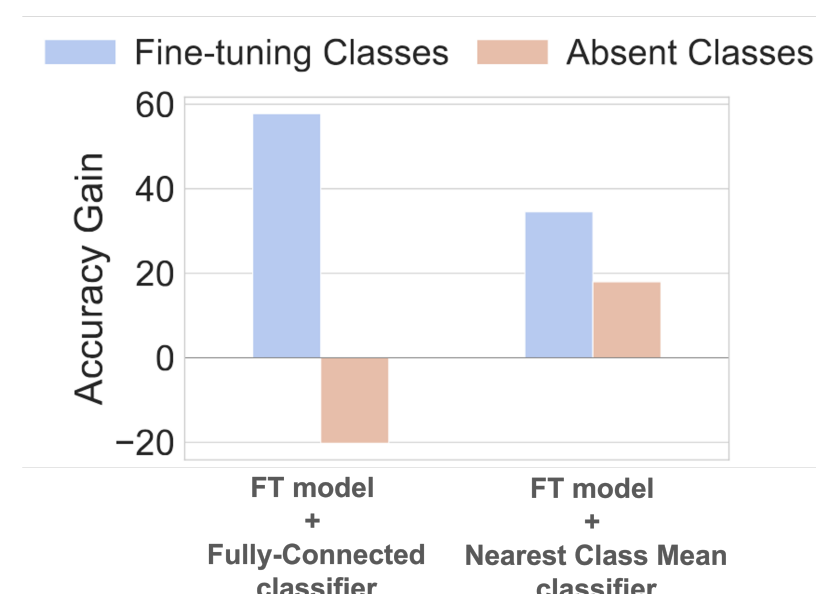


## Problem



## What are damaged in the fine-tuned (FT) model?

### Is the FT feature extractor damaged?



**NO!** FT with **subset** of classes can **adapt** the **feature extractor** to the new domain and **improve absent** classes.

### Is the FT classifier damaged?

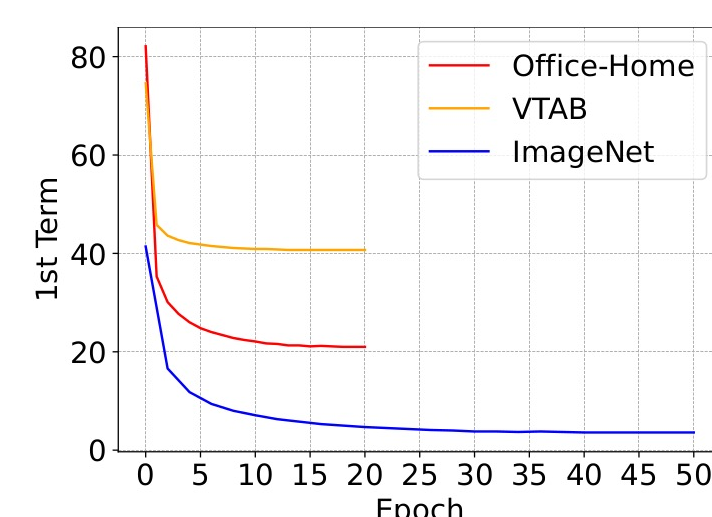
$$p(c|x) = \frac{\exp(w_c^\top f_\theta(x))}{\sum_{c' \in \mathcal{Y}} \exp(w_{c'}^\top f_\theta(x))} = \frac{z_c(x)}{\sum_{c' \in (S \cup U)} z_{c'}(x)}$$

$$= \frac{\sum_{c' \in U} z_{c'}(x)}{\sum_{c' \in S} z_{c'}(x) + \sum_{c' \in U} z_{c'}(x)} \times \frac{z_c(x)}{\sum_{c' \in U} z_{c'}(x)}$$

**SoftMax** decomposition: probability that  $x$  belongs to an **absent** class  $c$

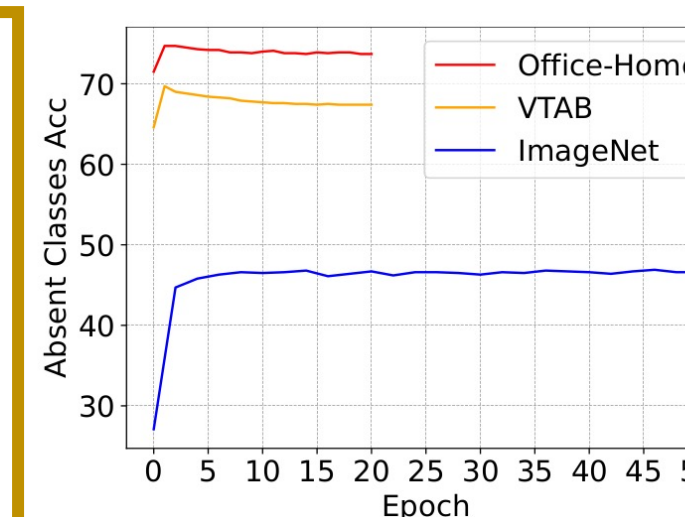
probability that  $x$  belongs to any **absent** classes  $U$

probability that within the **absent** classes  $U$ ,  $x$  belongs to class  $c$ .



FT model tends to classify **absent** class as **fine-tuning** classes

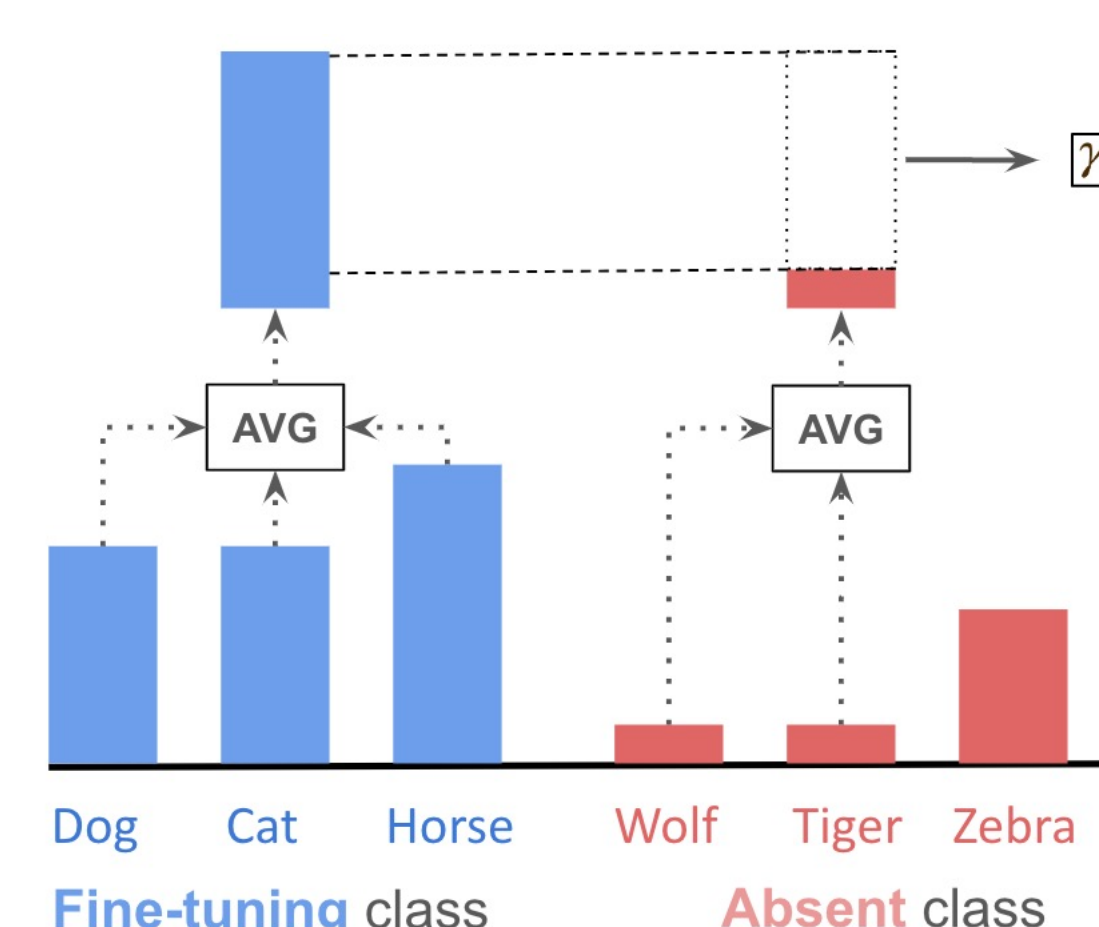
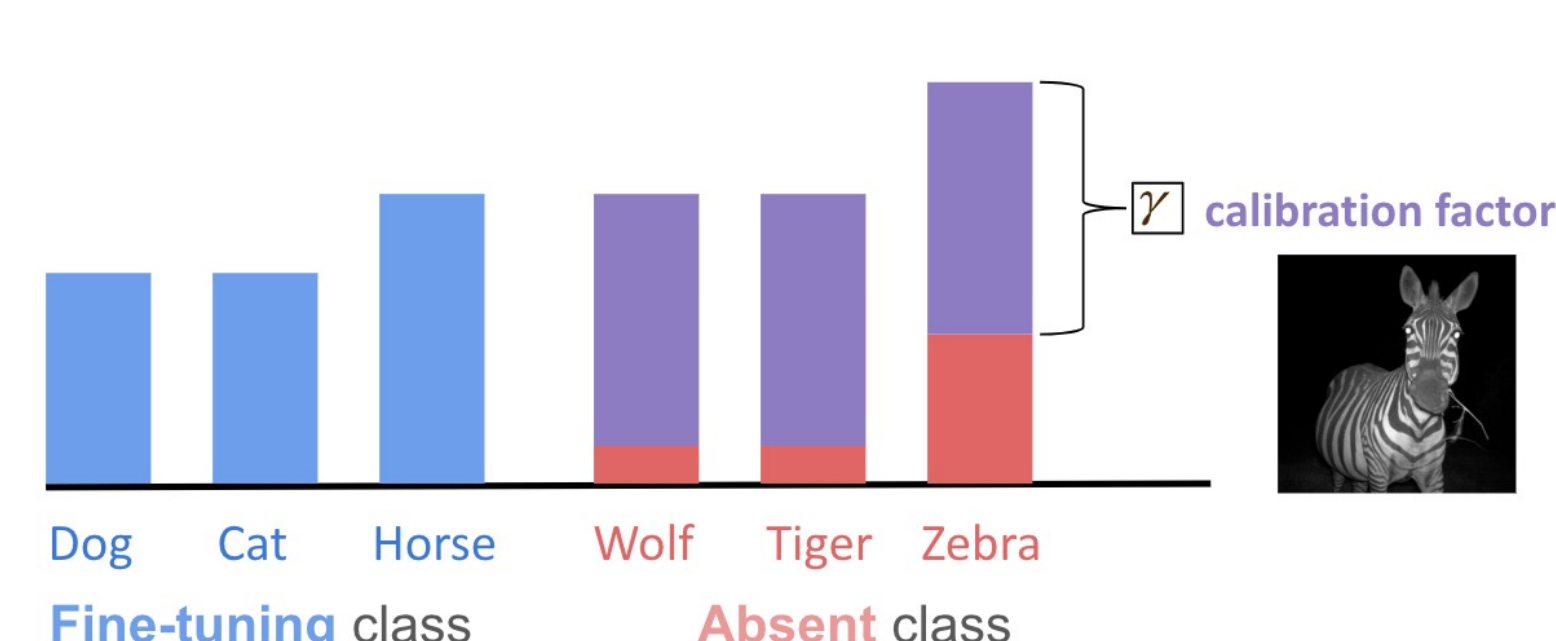
FT model's ability to discriminate absent class sample **within absent** classes is **improved**



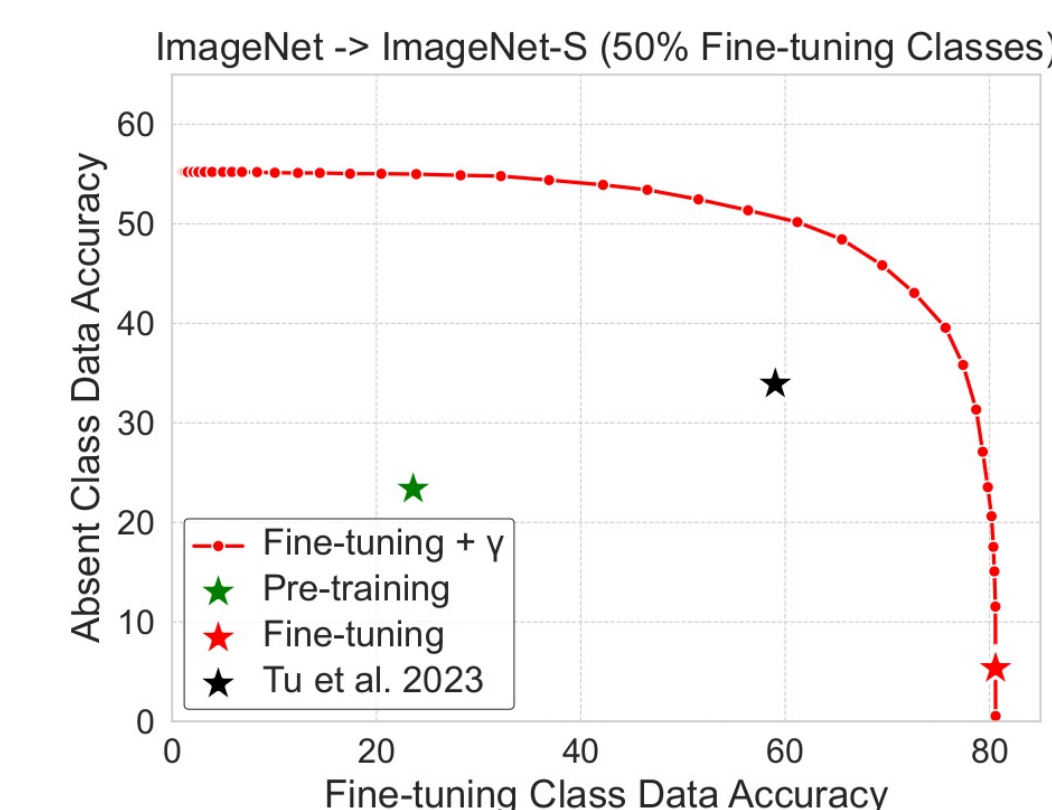
**Root cause:** FT model's **biased logits** towards fine-tuning classes

## Post-Processing Calibration for the Rescue

$$\hat{y} = \arg \max_{c \in \mathcal{Y}} w_c^\top f_\theta(x) + \gamma \mathbb{1}[c \in U].$$

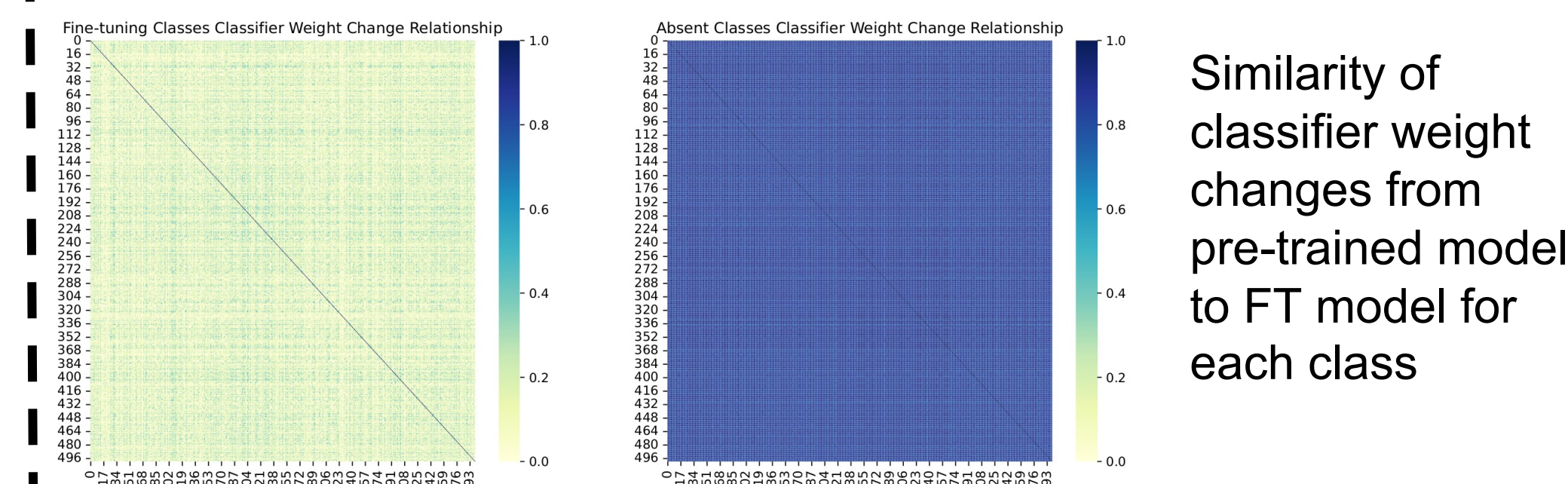


## Results



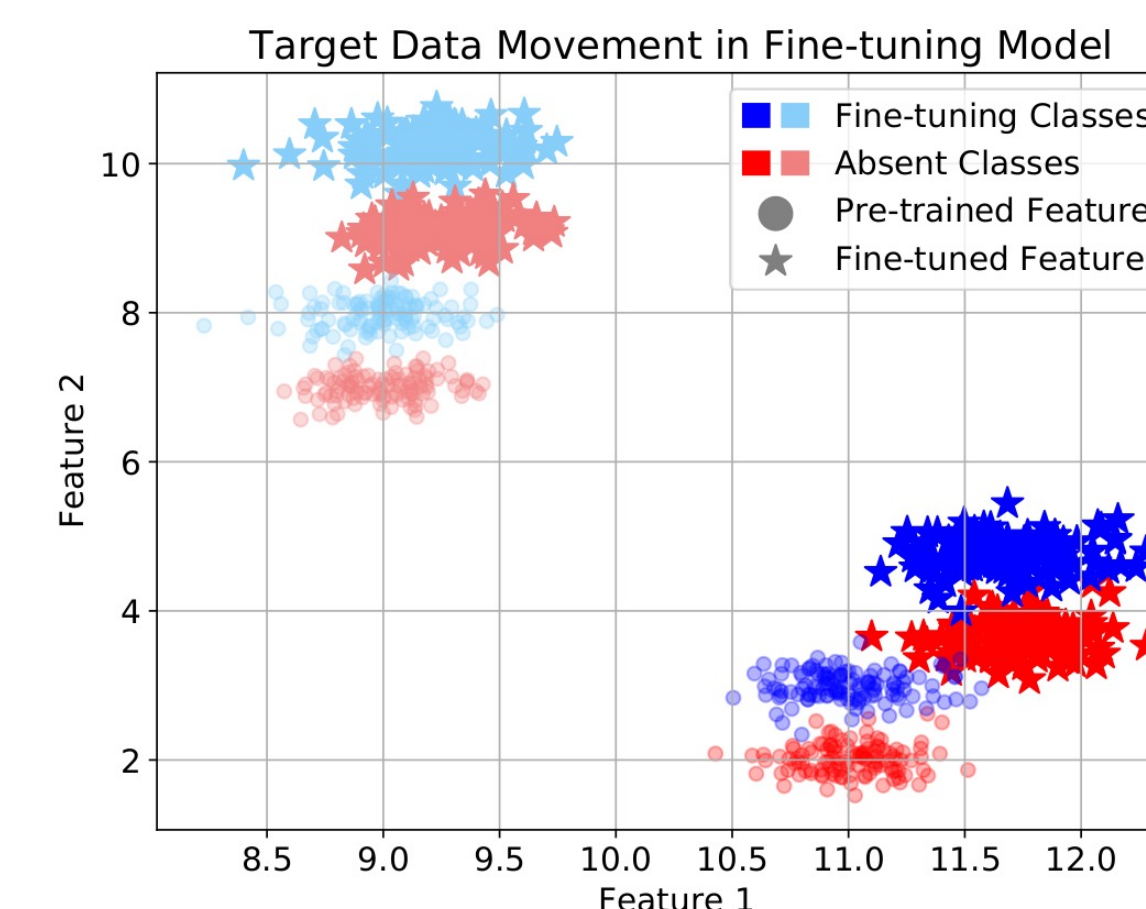
- FT (\*) + **post processing calibration** outperforms SOTA (\*)
- Observations are **robust** to data split, fine-tuning class size and optimizer.

## Why is absent class relationship preserved?



Compared with **fine-tuning** classes, the **weight changes** for **absent** classes are highly similar

## Why do absent class features improve after FT?



Toy example:

- absent** class features moving in a **similar direction** as nearby **fine-tuning** class features.

**Acknowledgement:** This research is supported by grants from the NSF (ICICLE: OAC2112606).