# Learning De-Biased Representations for Remote-Sensing Imagery

Conference on Neural Information Processing Systems 2024 (NeurIPS'24)

# Background & Motivation

Challenges in RS domain • Current Solutions & Limits • Our Key Observations

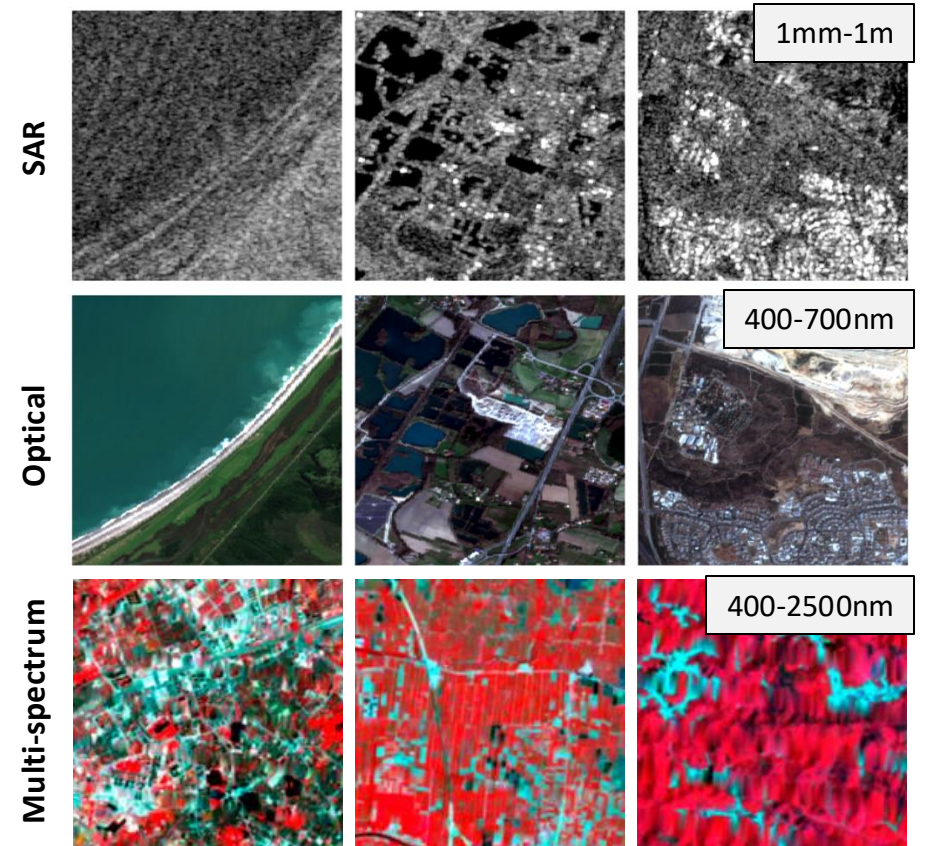# Remote Sensing Domain

- **Definition**

  Remote sensing images are captured from an overhead perspective by spaceborne or airborne sensors, which present unique viewpoints compared to natural images.

- **Multiple Spectrums**

  - Optical RS (ORS): 400-700nm

  - Multi-spectral RS (MSRS): 400-2500nm

  - Synthetic Aperture Radar (SAR): 1mm-1m

- **Key Applications**

  - Environmental monitoring
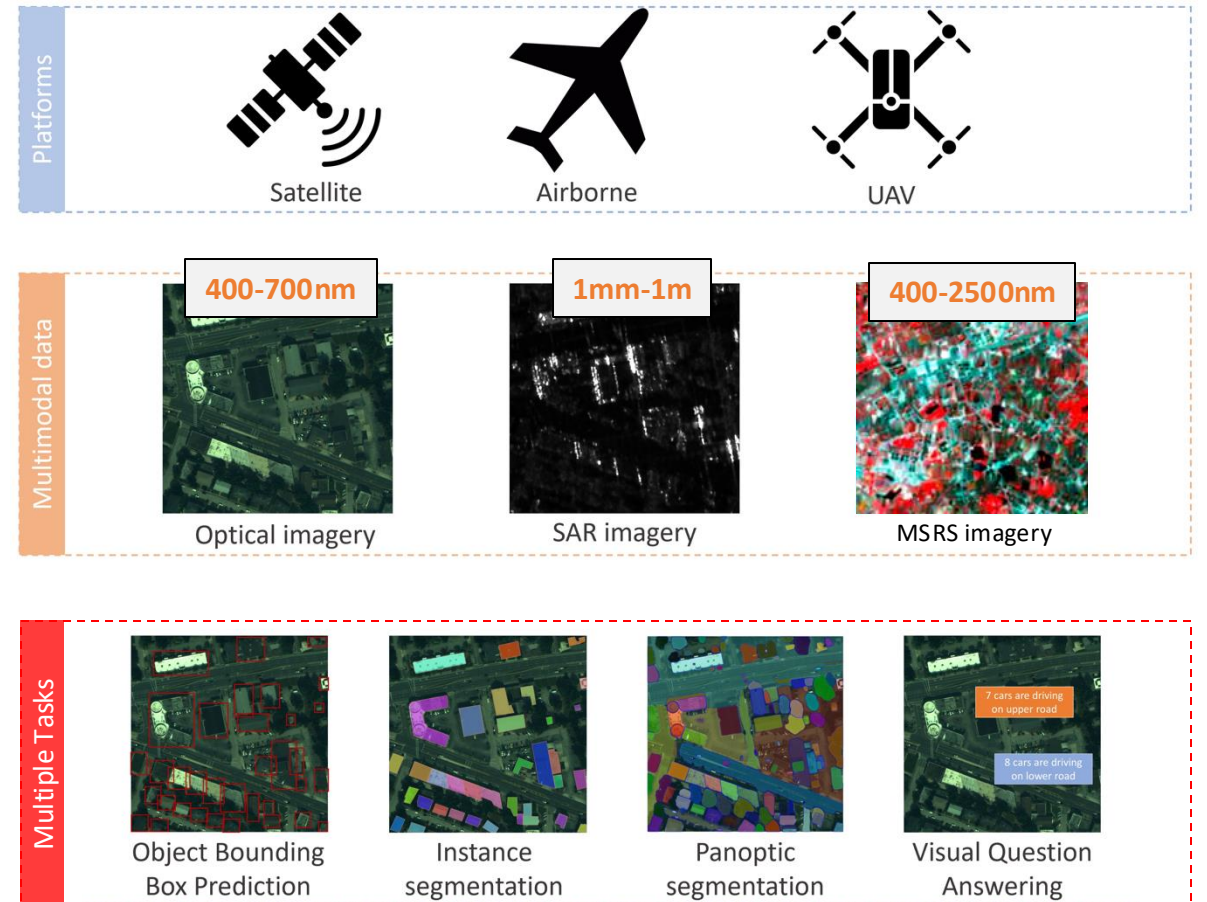
  - Resource management

  - Disaster response



SAR — 1mm-1m

Optical — 400-700nm

Multi-spectrum — 400-2500nm

Source: EUSI Database

*Remote Sensing data are diverse and complex, requiring heavy processing costs.*

# Challenges in RS Data

- **RS Data Diversity and Complexity**
  - Various data **source & processing tech**
  - Various **spectrums**
  - Various downstream **tasks**

**Platforms**

Satellite          Airborne          UAV

**Multimodal data**

400-700nm          1mm-1m          400-2500nm

Optical imagery          SAR imagery          MSRS imagery

**Multiple Tasks**

7 cars are driving on upper road

8 cars are driving on lower road

Object Bounding Box Prediction          Instance segmentation          Panoptic segmentation          Visual Question Answering

# Challenges in RS Data

Learning **robust and generic representations** is desirable!

# Parameter Efficient **Transfer Learning**

- **Self-supervised Training from Scratch**
  - Data scarcity in certain spectrums (*e.g.*, **SAR** imagery)
  - Constraints in **model scale** and **data scale**
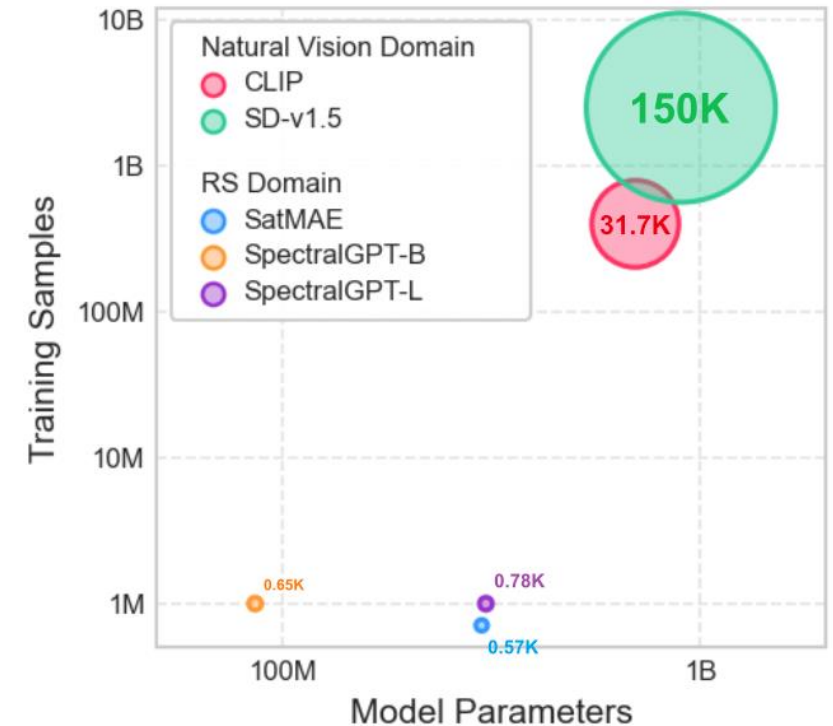  - Constraints in **training GPU time**



Figure: **Compare foundation models.** The bubble figure shows model scale, data scale and training time of five representative foundation models. Numbers near to bubbles are training GPU-hour. Models from RS domain uses less training GPU-hours compared with natural vision domain.
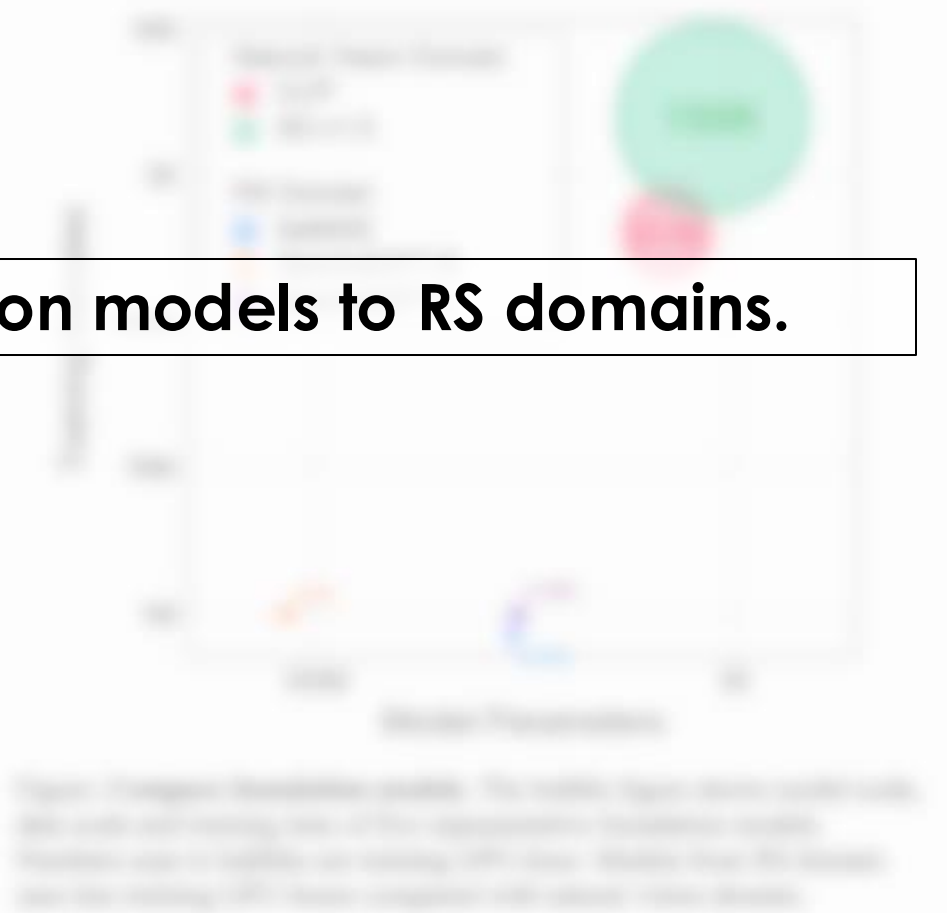
# Parameter Efficient **Transfer Learning**

- Self-supervised Training from Scratch
  - Data scarcity in certain spectrums (e.g., SAR imagery)
  - Constraints in model scale and data scale
  - Constraints in training GPU flow

We propose to **transfer existing foundation models to RS domains.**

*Why do we need parameter efficient?*

# **Parameter Efficient** Transfer Learning

- **Transfer Learning Setups**
  - Adaptation **from natural vision domain to RS domain**
  - Adaptation **between RS spectrums**

- **Zero-Shot and Fine-tuning**
  - Fine-tuning suffers from 1) catastrophic forgetting, 2) long training time, and 3) high VRAM usage.
  - Even zero-shot outperforms fine-tuning.

- **Parameter Efficient Transfer Learning (PEFT)**
  - **LoRA - Lo**w **R**ank **A**daptation
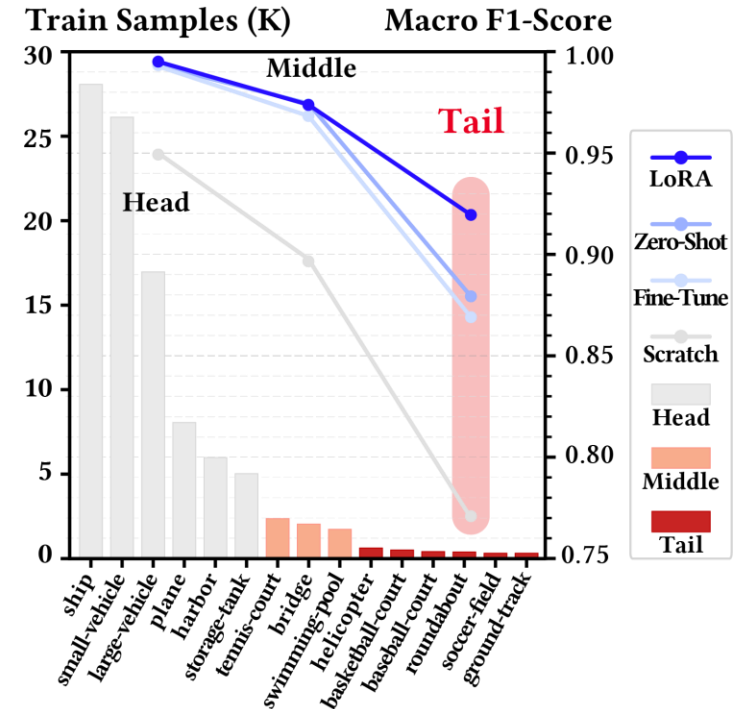  - Both fine-tuning, zero-shot and PEFT suffers from **long-tailed** distribution issue.



Figure: **Performance of Natural to ORS adaptation setting.** The debLoRA achieves highest performance, especially for tail class.

# Insights & Design

Key Observations • Framework • Core Components • Algorithm Explanation

# Key Observation – Biased Representation Space

- **Biased Representation Space**
  - When learnt on long-tailed data, LoRA's adapted **feature space** of LoRA **is biased**[2].
  - Validation samples of <span style="color:blue">**head class**</span> are mostly **correctly** classified.
  - Validation samples of <span style="color:red">**tail class**</span> are **wrongly** classified as head class.
  - Key Challenge: **Train/Val distribution mismatch** for tail classes.
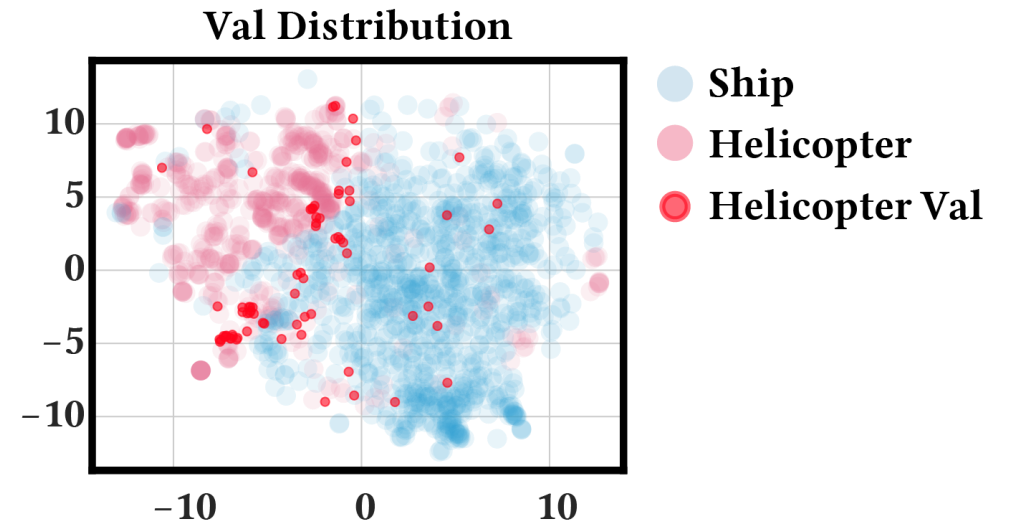
**Val Distribution**

- Ship
- Helicopter
- Helicopter Val

Figure: **Feature distribution of training samples.** For clearer visualization, we pick representative head class "Helicopter" and tail class "Ship" from DOTA v1 dataset as an example.
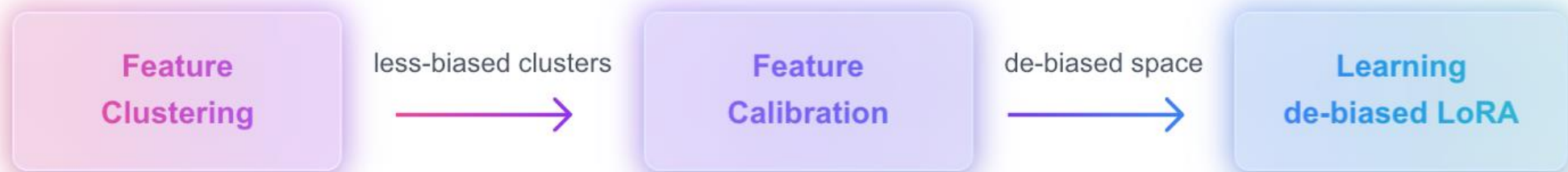
[2] We define feature space $Z$ as biased if $Vol(Z_h) \gg Vol(Z_t)$, and $\exists\, z_t \in Z_t : P(z_t \in Z_h) > P(z_t \in Z_t)$, where $Z_h$ and $Z_t$ denotes the feature spaces of head and tail classes respectively, $Vol(\cdot)$ denotes feature space volume, and $P(\cdot)$ denotes the probability distribution predicted by the model.

# Framework of Our Approach

- **Three key components**

  - **Feature clustering** – Unsupervised clustering to find less biased prototypes.

  - **Feature calibration** – Use less-biased prototypes to calibrate tail class features.

  - **debLoRA learning** – Learn a LoRA module to capture this de-bias mapping.
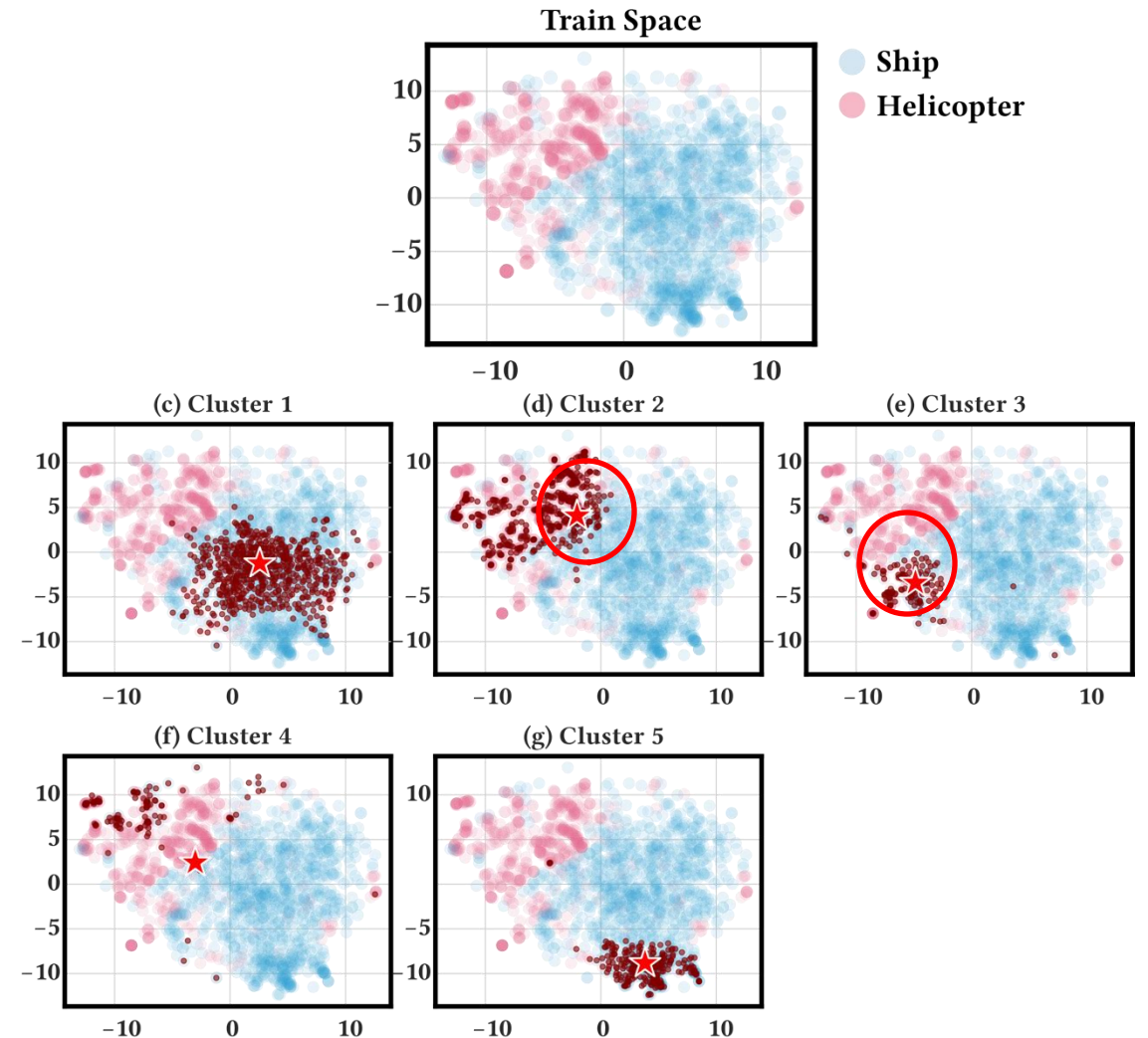
# Feature Clustering

- **Feature clustering**

  o We conduct K-Means clustering over training samples' feature space.

  $$\min_{\mu_k} \sum_{i=1}^{N} \min_{k} \|z_i - \mu_k\|^2, s.t. \forall k, n_k \geq \frac{N}{K \cdot \rho},$$

  where $\mu_k$ and $n_k$ denote the center and size of the $k$-th cluster, respectively.

  o Some cluster centers are contributed by both head and tail classes, and hence is less biased (*e.g.*, clusters 2 and 3).



**Train Space**

- Ship
- Helicopter

(c) Cluster 1 (d) Cluster 2 (e) Cluster 3

(f) Cluster 4 (g) Cluster 5
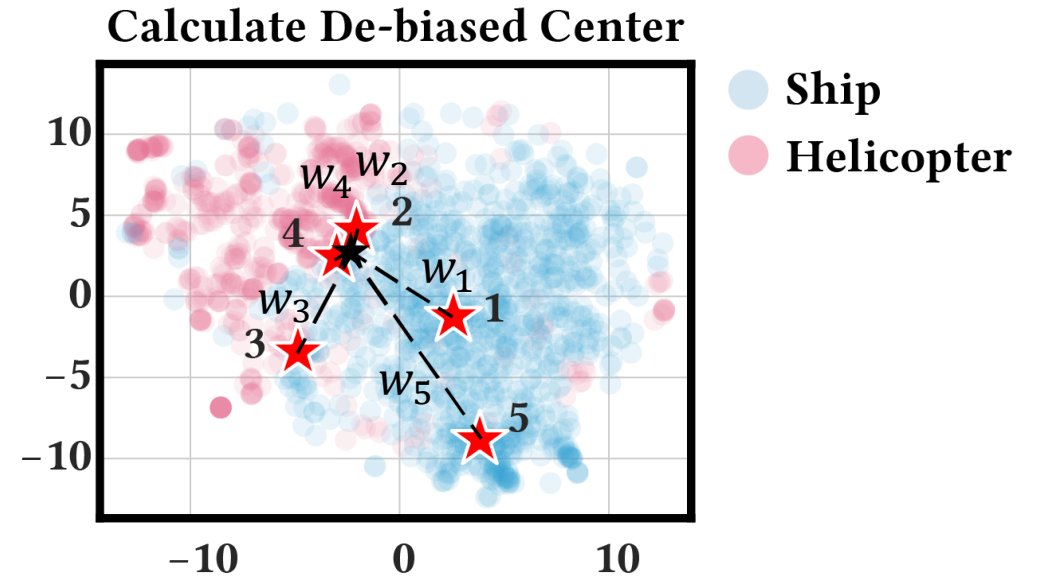
# Construct De-Biased Center

- **De-Biased Center**
  - ○ We calculate de-biased representation center for each tail class:

$$\hat{\mu}_c = \sum_k w_k \cdot \mu_k \, , w_k = \frac{n_k}{n_c} \, ,$$

  here weight $w_k$ proportion to the fraction of class $c$ samples in $k$-th cluster.

  - ○ This ensures that the de-biased center $\hat{\mu}_c$ is not dominated by head classes



**Calculate De-biased Center**

# Feature Calibration

- **Tail Class Calibration**

  ○ De-Biased Center are **closer to validation** samples.

  ○ We calibrate tail class features $z$ by moving them close to de-biased center $\hat{\mu}$ :

  $$\tilde{z} = \alpha z + (1 - \alpha)\hat{\mu},$$

  where $\alpha = \min(1, \frac{10}{ir})$ empirically.

- **Learning debLoRA**

  ○ We learn an LoRA module with training objective

  $$\min_{\phi} \frac{1}{D_t} \sum_{x \in D_t} \left\| g_{\phi}(f_{\theta}(x)) - \tilde{z} \right\|^2$$

**Feature Calibration & Learn**



Legend:
- Ship
- Helicopter
- ★ Original Biased Center
- ★ De-Biased Center
- debLoRA