



南京航空航天大学

Nanjing University of Aeronautics and Astronautics

# Long-Tailed Out-of-Distribution Detection: Prioritizing Attention to Tail

**Yina He<sup>1</sup>, Lei Peng<sup>1</sup>, Yongcun Zhang<sup>1</sup>, Juanjuan Weng<sup>2\*</sup>, Shaozi Li<sup>1,3</sup>, Zhiming Luo<sup>1,3\*</sup>**

<sup>1</sup>Department of Artificial Intelligence, Xiamen University, Xiamen, China

<sup>2</sup>College of Information Science and Technology, Jinan University, Guangzhou, China

<sup>3</sup>Key Laboratory of Multimedia Trusted Perception and Efficient Computing, Ministry of Education of China, Xiamen University

AAAI 2025

---

Current **out-of-distribution (OOD) detection** methods typically assume balanced **in-distribution (ID) data**, while most real-world data follow a **long-tailed distribution**.

Outlier Exposure (OE)  $\mathcal{L}_{OE} = \mathbb{E}_{x,y \sim \mathcal{D}_{in}} [\ell(f(x), y)] + \gamma \mathbb{E}_{x \sim \mathcal{D}_{out}} [\ell(f(x), u)]$ , (1)  
(Hendrycks, Mazeika, and Dietterich 2018)

Supervised Contrastive Learning (SCL)  $\mathcal{L}_{scl}^{in}(z_i, y) = -\log \left\{ \frac{1}{|B_y|} \sum_{p \in B_y} \frac{e^{z_i \cdot z_p / \tau}}{\sum_{j=1}^K \sum_{a \in B_j} e^{z_i \cdot z_a / \tau}} \right\}$ , (2)  
(Khosla et al. 2020)

Logit Adjustment (LA)  $\mathcal{L}_{la}(z_i, y) = -\log \frac{\pi_y e^{\varphi_y(z_i)}}{\sum_{y' \in \mathcal{Y}} \pi_{y'} e^{\varphi_{y'}(z_i)}}$ , (3)  
(Menon et al. 2021)

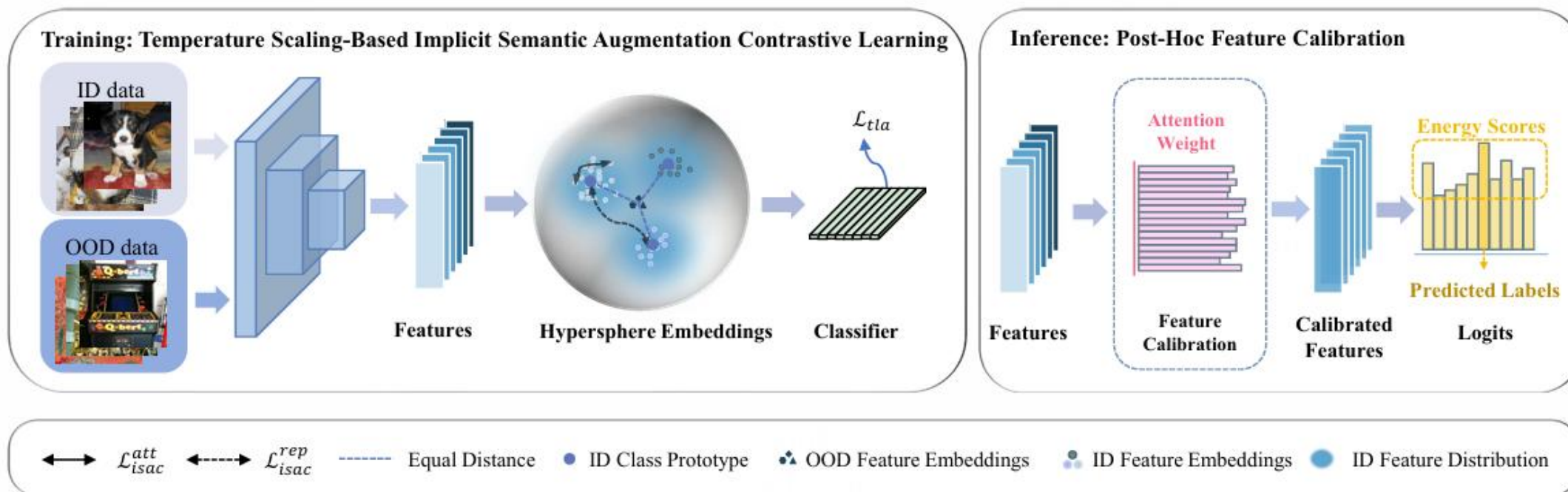
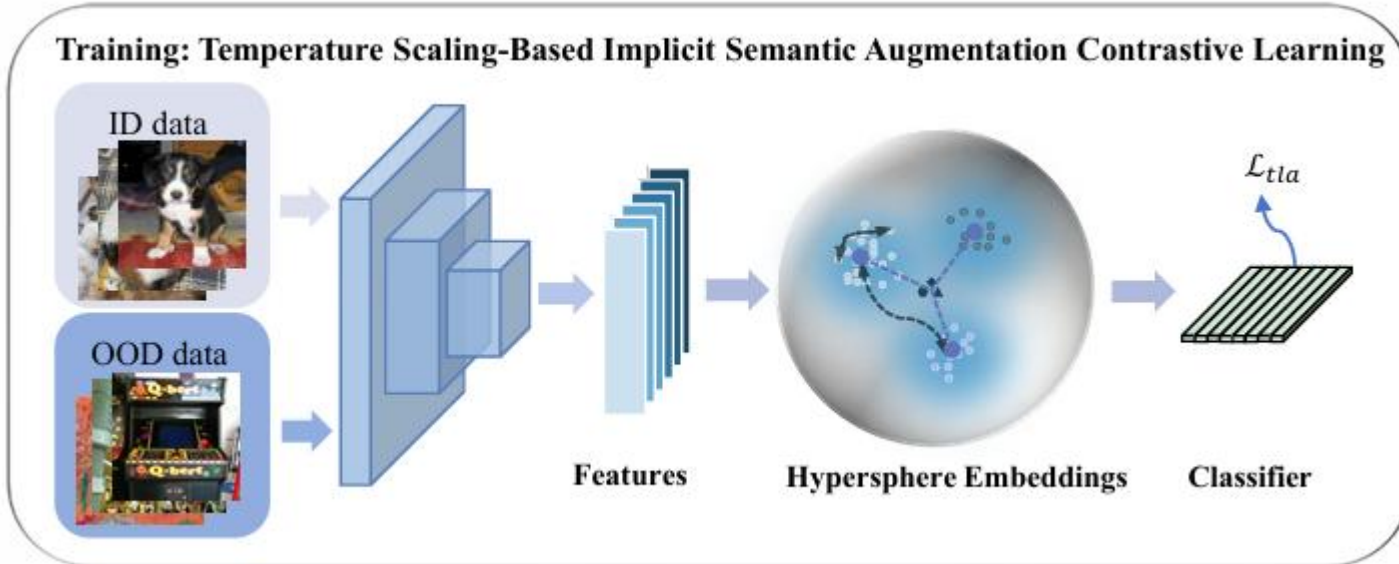


Figure 2: Overview of the proposed framework. The framework consists of a temperature scaling-based implicit semantic augmentation training phase and a feature calibration inference phase. We jointly optimize two complementary terms to encourage desirable hypersphere embeddings: an implicit semantic augmentation contrastive loss to encourage a balanced feature encoder and a temperature scaling-based logit adjustment loss to encourage a balanced high-confidence classifier. Feature calibration fine-tunes features during the inference phase by using an attention weight extracted from the training set, thereby achieving desirable ID classification and OOD detection results.



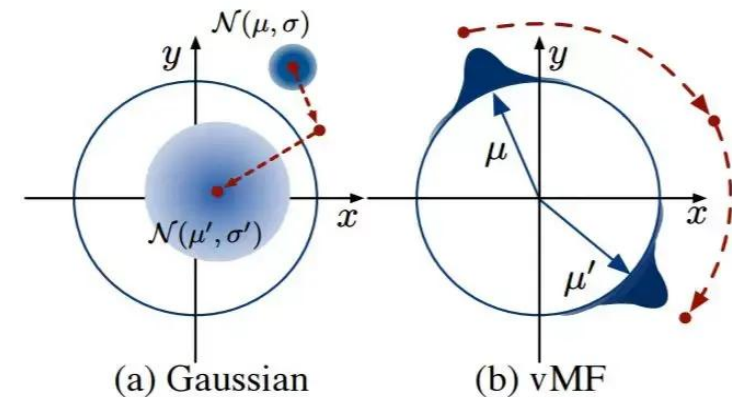
## von Mises-Fisher (vMF) distributions

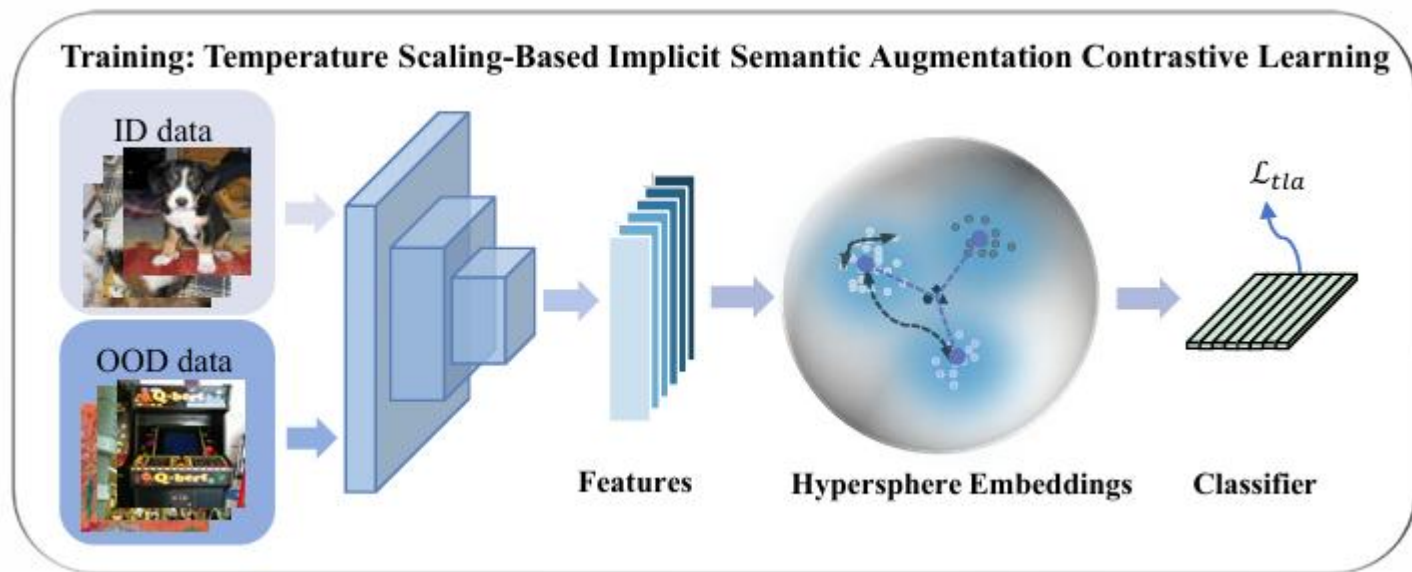
$$P_d(\mathbf{z}; \boldsymbol{\mu}_y, \kappa_y) = Z_d(\kappa_y) e^{\kappa_y \boldsymbol{\mu}_y^\top \mathbf{z}},$$

- $\mathbf{z} \in \mathbb{R}^d$  is a unit vector
- $\boldsymbol{\mu}_y$  is the class prototype
- $\kappa_y \geq 0$  indicates the concentration of the distribution
- $Z_d(\kappa_y)$  serving as the normalization factor.

a mixture of vMF distributions

$$P_d(\mathbf{z}) = \sum_{y=1}^K P_d(y) P_d(\mathbf{z}|y) = \sum_{y=1}^K \pi_y Z_d(\kappa) e^{\kappa \boldsymbol{\mu}_y^\top \mathbf{z}}, \quad (5)$$





## Temperature Scaling-Based Logit Adjustment

ICLR2021

$$\mathcal{L}_{tla}(z_i, y) = -\log \frac{\pi_y e^{\varphi_y(z_i)/\varepsilon}}{\sum_{y' \in \mathcal{Y}} \pi_{y'} e^{\varphi_{y'}(z_i)/\varepsilon}} \quad (7)$$

- hyperparameter  $\varepsilon$

通过数学推导直接计算无限采样下的期望损失:

$$\mathcal{L}_{isac}(z_i, y) = \log \left\{ \sum_{j=1}^K \frac{\pi_j Z_d(\tilde{\kappa}_y) Z_d(\kappa_j)}{\pi_y Z_d(\kappa_y) Z_d(\tilde{\kappa}_j)} \right\}, \quad (6)$$

- $\tilde{\kappa}_j = \|\kappa_j \boldsymbol{\mu}_j + \mathbf{z}_i / \tau\|_2$ , 类别  $j$  的原分布参数,  $i$  样本的方向向量
- $\mathbf{z}_j \sim \text{vMF}(\boldsymbol{\mu}_j, \kappa_j)$
- $\tau$  is a temperature parameter

Finally, the overall loss function is as follows:

$$\mathcal{L} = \mathcal{L}_{isac} + \alpha \mathcal{L}_{tla} + \beta \mathcal{L}_{out},$$

- $\alpha$  and  $\beta$  are hyperparameters

The **penultimate feature layer** is more correlated with the final classification.

- The attention weight extracted from a **class-balanced ID dataset** and an **OOD dataset**

$$\mathcal{X}^{cb} = \mathcal{D}_{in}^{cb} \cup \mathcal{D}_{out}. \quad x_i \in \mathcal{X}^{cb}$$

- d-dimensional feature embedding

$$z_i = f(x_i, \theta) = [z_i^1, z_i^2, \dots, z_i^d]^T$$

- The score of  $z_i$  being predicted as class  $y_i$  is

$$S_{y_i}(z_i) = \varphi(z_i)$$

- The importance of k-th dimension of  $z_i$  is defined as  $I^k(z_i) = \frac{\partial S_{y_i}(z_i)}{\partial z_i^k} \cdot z_i^k$

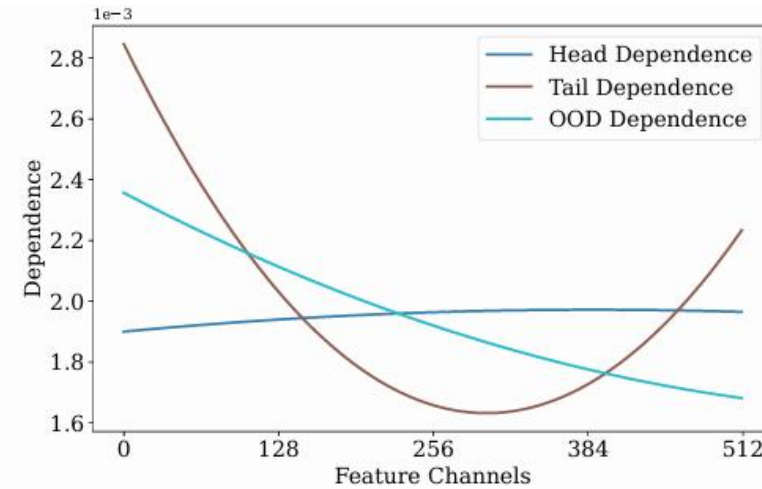


Figure 3: It visualizes the dependence of OOD, head class, and tail class samples on feature channels, showing that these three types of samples rely on different feature channels.

$$A = \frac{1}{|N|} \sum_{j=1}^K \left\{ \sum_{i \in N_j^{in}} \frac{I(z_i^{in})}{\pi_j} - \sum_{i \in N_j^{out}} \frac{I(z_i^{out})}{\pi_j} \right\}, \quad (9)$$

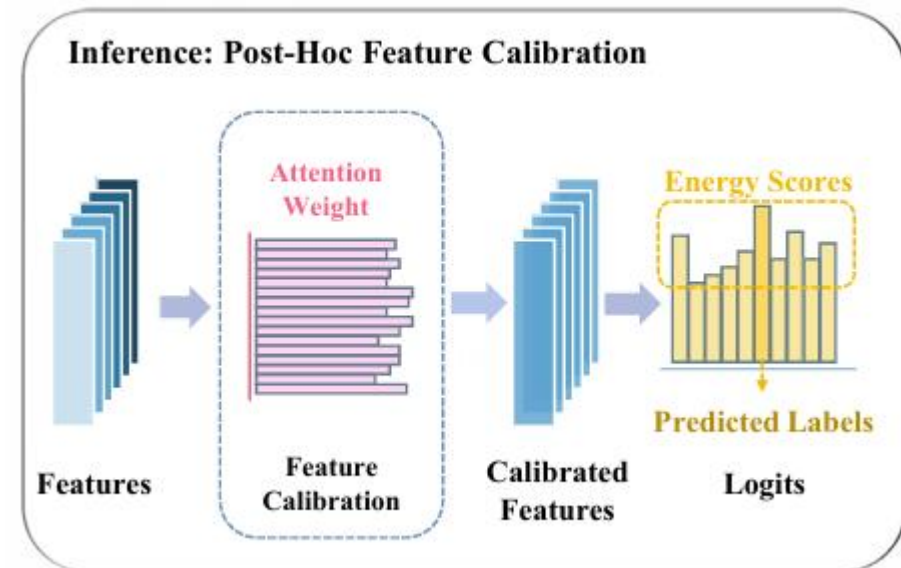
They scale A to the range [0, 2], where values between [0,1] are attenuated and values between [1,2] are enhanced.

$$A = \frac{1}{|N|} \sum_{j=1}^K \left\{ \sum_{i \in N_j^{in}} \frac{I(z_i^{in})}{\pi_j} - \sum_{i \in N_j^{out}} \frac{I(z_i^{out})}{\pi_j} \right\}, \quad (9)$$

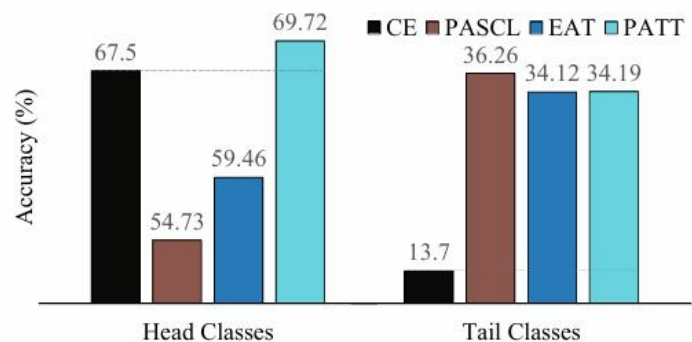
$$z^{cal} = z \odot A^{scale}$$

- Instead of using maximum softmax probability(MSP) as OOD scores, we are inspired by energyOE(Liu et al.2020) and use energy scores as OOD scores in the inference phase, which is defined as follows:

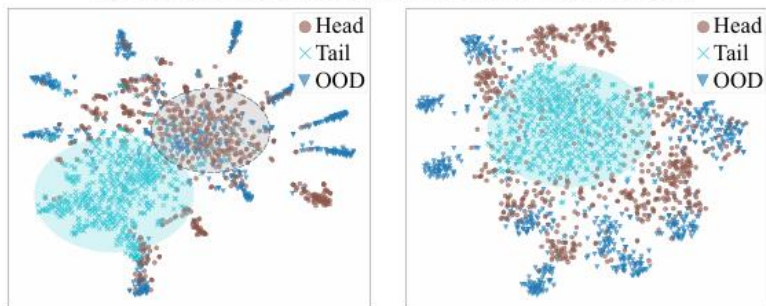
$$S_{ood}(z_i) = \log \sum_{j=1}^K e^{\varphi_j(z_i^{cal})}.$$



- They conduct experiments on widely used datasets, i.e., CIFAR10-LT, CIFAR100-LT (Cao et al. 2019), and ImageNet-LT (Liu et al. 2019) as ID training sets ( $D_{in}$ ).



(a) Separate ACC for head and tail on ImageNet-LT.



(b) PASCL's feature distribution

(c) PATT's feature distribution

$D_{out}^{test}$	Method	AUROC $\uparrow$	AUPR-in $\uparrow$	AUPR-out $\uparrow$	FPR95 $\downarrow$
Texture	OE	92.59 $\pm$ 0.4	96.01 $\pm$ 1.4	83.32 $\pm$ 1.7	25.10 $\pm$ 1.1
	PASCL	93.16 $\pm$ 0.4	96.57 $\pm$ 1.2	84.80 $\pm$ 1.5	<b>23.26</b> $\pm$ 0.9
	<b>Ours</b>	<b>93.96</b> $\pm$ 0.6	<b>97.69</b> $\pm$ 0.8	<b>86.49</b> $\pm$ 0.8	26.65 $\pm$ 1.6
SVHN	OE	95.10 $\pm$ 1.0	91.59 $\pm$ 0.5	97.14 $\pm$ 0.8	16.15 $\pm$ 1.5
	PASCL	96.63 $\pm$ 0.9	92.89 $\pm$ 0.5	98.06 $\pm$ 0.6	12.18 $\pm$ 3.3
	<b>Ours</b>	<b>98.21</b> $\pm$ 0.7	<b>97.19</b> $\pm$ 0.6	<b>98.50</b> $\pm$ 0.5	<b>5.73</b> $\pm$ 2.0
CIFAR100	OE	83.40 $\pm$ 0.3	84.06 $\pm$ 0.3	80.93 $\pm$ 0.6	56.96 $\pm$ 0.9
	PASCL	84.43 $\pm$ 0.2	85.32 $\pm$ 0.5	82.99 $\pm$ 0.5	57.27 $\pm$ 0.9
	<b>Ours</b>	<b>85.36</b> $\pm$ 0.2	<b>86.01</b> $\pm$ 0.3	<b>83.29</b> $\pm$ 0.5	<b>51.12</b> $\pm$ 0.9
Tiny ImageNet	OE	86.14 $\pm$ 0.3	89.88 $\pm$ 0.7	79.33 $\pm$ 0.7	47.78 $\pm$ 0.7
	PASCL	87.14 $\pm$ 0.2	90.22 $\pm$ 0.5	81.54 $\pm$ 0.4	47.69 $\pm$ 0.6
	<b>Ours</b>	<b>88.62</b> $\pm$ 0.2	<b>90.82</b> $\pm$ 0.7	<b>84.54</b> $\pm$ 0.1	<b>41.30</b> $\pm$ 1.7
LSUN	OE	91.35 $\pm$ 0.2	93.06 $\pm$ 0.3	87.62 $\pm$ 0.8	27.86 $\pm$ 0.7
	PASCL	<b>93.17</b> $\pm$ 0.15	82.59 $\pm$ 0.3	91.76 $\pm$ 0.5	26.40 $\pm$ 1.0
	<b>Ours</b>	91.64 $\pm$ 0.3	<b>93.16</b> $\pm$ 0.6	<b>92.27</b> $\pm$ 0.1	<b>24.41</b> $\pm$ 1.5
Place365	OE	90.07 $\pm$ 0.3	82.09 $\pm$ 0.4	95.15 $\pm$ 0.2	34.04 $\pm$ 0.9
	PASCL	91.43 $\pm$ 0.2	82.59 $\pm$ 0.2	96.28 $\pm$ 0.1	33.40 $\pm$ 0.9
	<b>Ours</b>	<b>91.95</b> $\pm$ 0.6	<b>82.63</b> $\pm$ 0.3	<b>97.81</b> $\pm$ 0.2	<b>30.15</b> $\pm$ 1.9
Average	OE	89.77 $\pm$ 0.3	89.45 $\pm$ 0.4	87.25 $\pm$ 0.6	34.65 $\pm$ 0.5
	PASCL	90.99 $\pm$ 0.2	90.18 $\pm$ 0.4	89.24 $\pm$ 0.3	33.36 $\pm$ 0.8
	<b>Ours</b>	<b>91.62</b> $\pm$ 0.5	<b>91.25</b> $\pm$ 0.5	<b>90.48</b> $\pm$ 0.3	<b>29.89</b> $\pm$ 1.4

(a) Comparison of PATT to PASCL and OE on six OOD datasets.

Table 1: Comparison results on CIFAR10-LT. The best results are shown in bold, and the second-best results are underlined.

ID Dataset	ISAC	TLA	FC	AUROC $\uparrow$	AUPR-in $\uparrow$	AUPR-out $\uparrow$	FPR95 $\downarrow$	ACC $\uparrow$	ACC-t $\uparrow$
CIFAR10-LT	$\times$	$\times$	$\times$	71.21	74.84	64.37	58.09	78.10	63.00
	$\checkmark$	$\times$	$\times$	74.10	78.70	66.43	49.09	79.55	66.47
	$\times$	$\checkmark$	$\times$	84.70	83.68	81.12	48.18	81.11	70.77
	$\times$	$\times$	$\checkmark$	79.75	80.48	74.38	52.17	79.55	67.50
	$\checkmark$	$\checkmark$	$\times$	91.06	90.90	88.92	32.35	82.09	70.00
	<b>PATT</b>			<b>91.62</b>	<b>91.25</b>	<b>90.48</b>	<b>29.89</b>	<b>84.77</b>	<b>79.67</b>
CIFAR100-LT	$\times$	$\times$	$\times$	71.70	72.77	65.95	72.22	44.68	10.97
	$\checkmark$	$\times$	$\times$	74.43	76.80	67.50	64.91	45.49	17.00
	$\times$	$\checkmark$	$\times$	72.10	73.88	66.35	69.80	49.01	27.12
	$\times$	$\times$	$\checkmark$	72.95	74.43	66.51	70.79	44.90	12.72
	$\checkmark$	$\checkmark$	$\times$	75.48	77.19	68.87	64.25	49.56	27.12
	<b>PATT</b>			<b>76.25</b>	<b>77.84</b>	<b>70.37</b>	<b>62.94</b>	<b>50.07</b>	<b>31.03</b>
ImageNet-LT	$\times$	$\times$	$\times$	67.87	44.28	69.74	88.76	45.13	9.24
	$\checkmark$	$\times$	$\times$	68.39	45.31	79.39	83.39	48.77	13.79
	$\times$	$\checkmark$	$\times$	72.74	51.17	86.25	82.08	50.27	31.58
	$\times$	$\times$	$\checkmark$	72.07	47.84	82.07	82.36	46.17	14.01
	$\checkmark$	$\checkmark$	$\times$	73.77	50.10	87.32	81.92	55.06	33.17
	<b>PATT</b>			<b>74.13</b>	<b>51.41</b>	<b>87.43</b>	<b>80.57</b>	<b>55.14</b>	<b>34.19</b>

Table 4: Ablation results of three key modules for PATT on CIFAR10-LT, CIFAR100-LT and ImageNet-LT.

The **penultimate feature layer** is more correlated with the final classification.

- The attention weight extracted from a **class-balanced ID dataset** and an **OOD dataset**

$$\mathcal{X}^{cb} = \mathcal{D}_{in}^{cb} \cup \mathcal{D}_{out}. \quad x_i \in \mathcal{X}^{cb}$$

- d-dimensional feature embedding

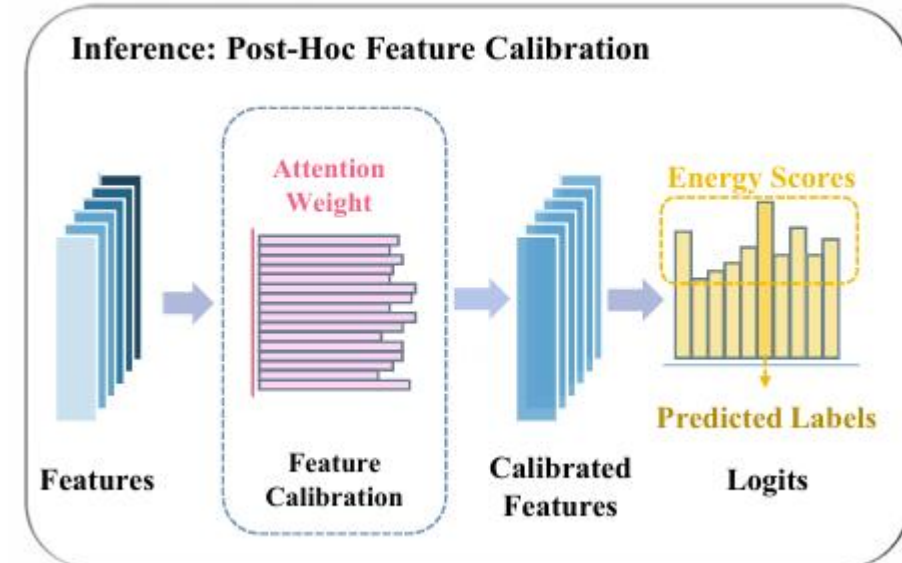
$$z_i = f(x_i, \theta) = [z_i^1, z_i^2, \dots, z_i^d]^T$$

- The score of  $z_i$  being predicted as class  $y_i$  is

$$S_{y_i}(z_i) = \varphi(z_i)$$

- The importance of k-th dimension of  $z_i$  is defined as  $I^k(z_i) = \frac{\partial S_{y_i}(z_i)}{\partial z_i^k} \cdot z_i^k$

Class-Balanced?





南京航空航天大学  
Nanjing University of Aeronautics and Astronautics

Thanks

---