

# Enhancing Few-Shot Class-Incremental Learning via Training-Free Bi-Level Modality Calibration

Yiyang Chen<sup>1</sup>, Tianyu Ding<sup>2</sup>, Lei Wang<sup>3</sup>, Jing Huo<sup>1</sup>, Yang Gao<sup>1</sup>, Wenbin Li<sup>1,4\*</sup>

<sup>1</sup>State Key Laboratory for Novel Software Technology, Nanjing University, China

<sup>2</sup>Applied Sciences Group, Microsoft, USA      <sup>3</sup>University of Wollongong, Australia

<sup>4</sup>Shenzhen Research Institute of Nanjing University, Shenzhen, China

CVPR 2025

## ▶ Class-Incremental Learning:

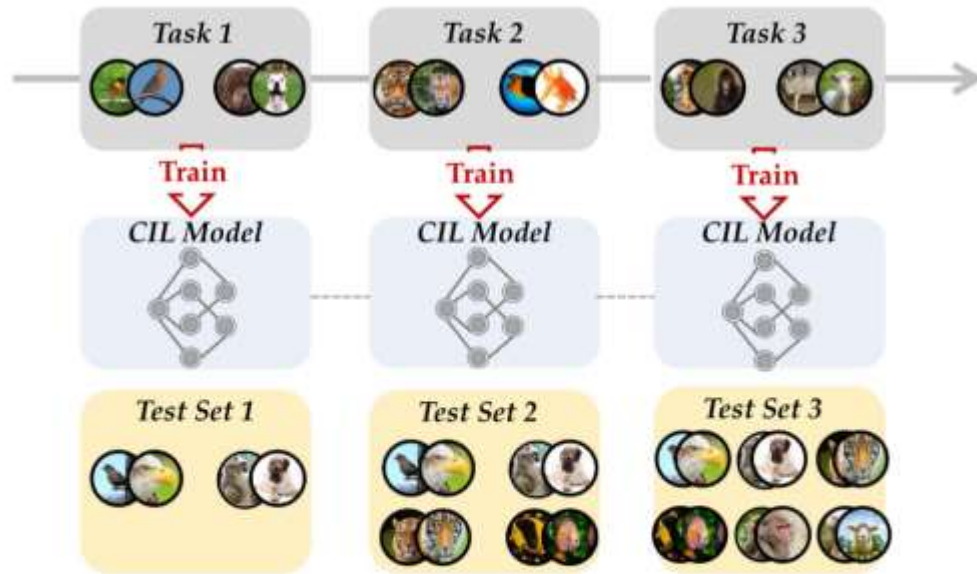


Fig. 1. The setting of CIL. Non-overlapping classes arrive sequentially, and the model needs to learn to classify all the classes incrementally. After learning each task, the model is evaluated among all seen classes. An ideal model should perform well in the newly learned classes and remember the former without forgetting.

## Major Challenges:

- catastrophic forgetting
- overfitting

## ▶ Few-Shot Class-Incremental Learning (FSCIL):

Base phase(task 0): sufficient samples

Incremental phase(task 1~n): a few samples

# Motivation

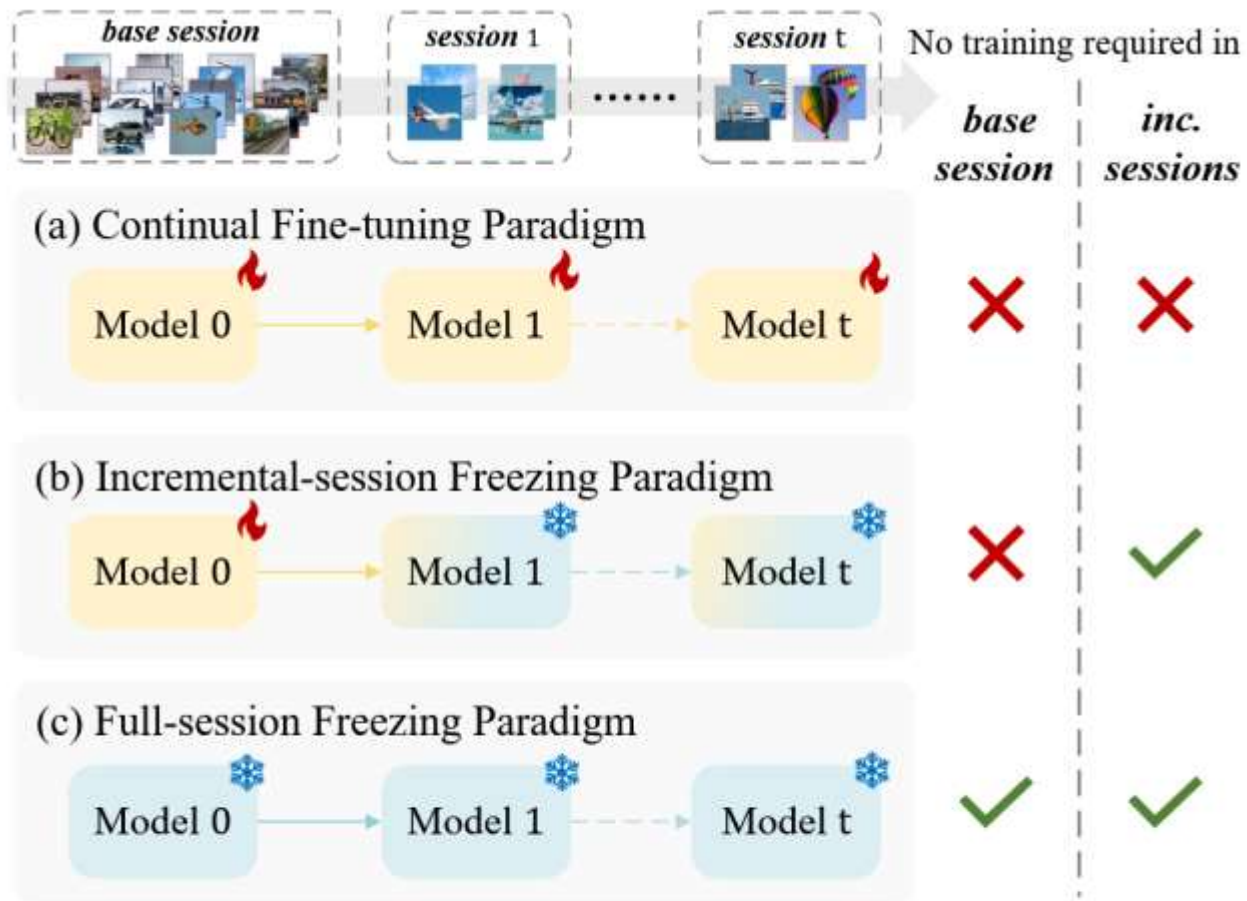


Figure 1. Model updating paradigms in FSCIL. (a) Continuous model updates during incremental phases. (b) Leveraging the base task to learn a generalizable model. (c) Our proposed framework: training-free in both base and incremental sessions.

How can we effectively leverage pre-trained knowledge and domain-specific visual priors to adapt the large-scale model using few samples *without additional training*?



**Bi-level Modality Calibration (BiMC)**

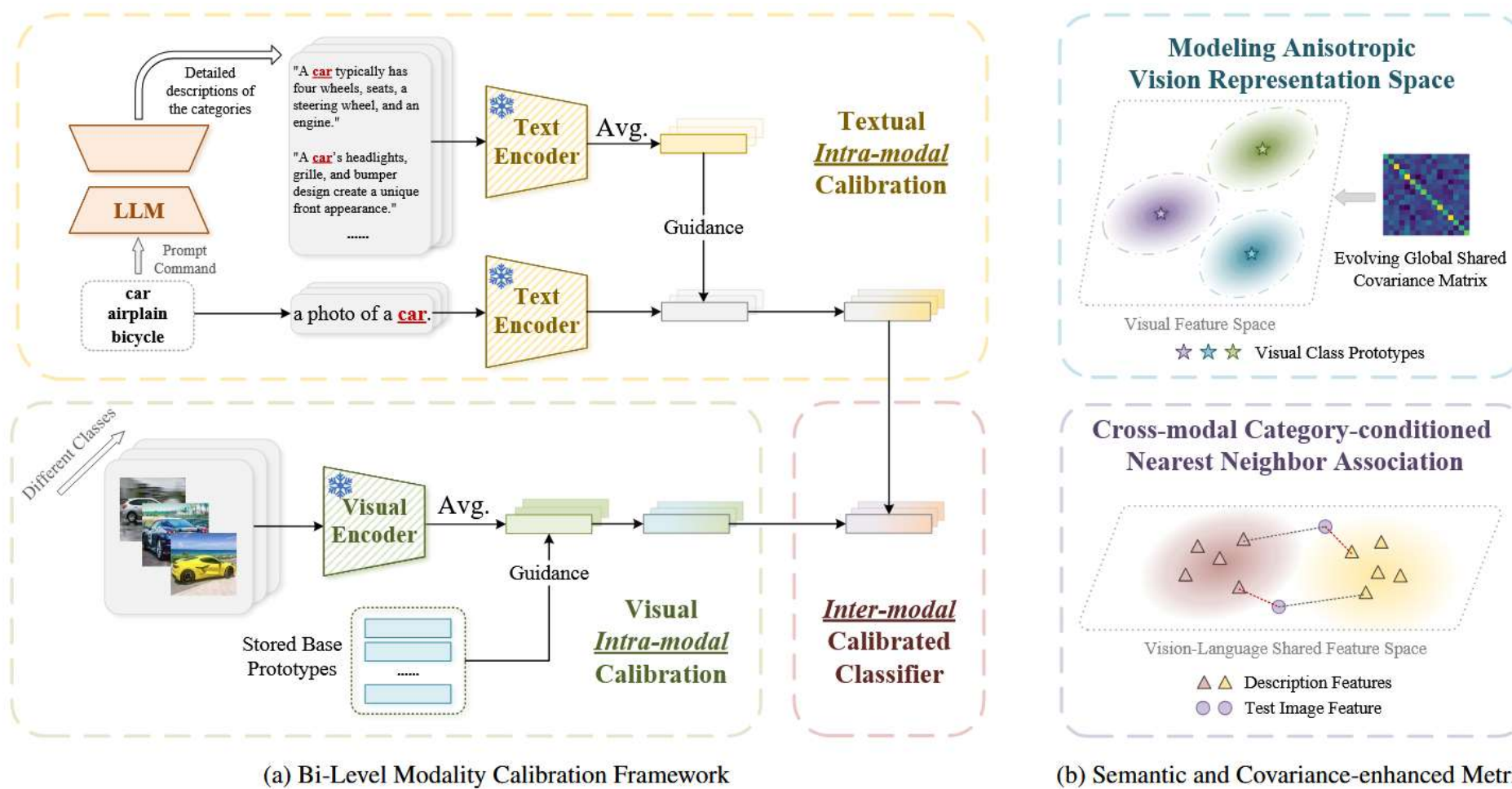
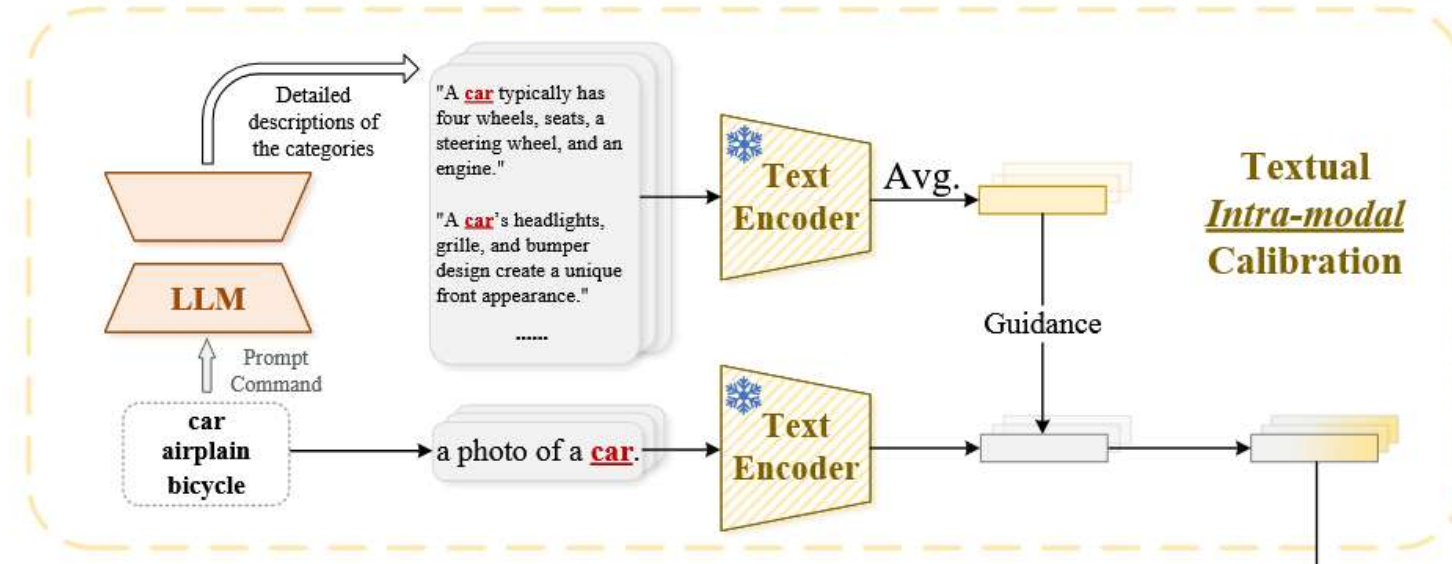


Figure 2. Overview of our framework. In each task adaptation phase, we construct a calibrated classifier through our proposed bi-level calibration framework, which involves *intra-modal calibration* and *inter-modal calibration*. To enhance performance, we introduce a globally shared covariance metric for visual feature modeling, complemented by a category-conditioned nearest-neighbor scoring strategy.

## Intra-modal classifier calibration

### ► Textual modality



$$\tilde{\mu}_c^T = (1 - \lambda_T)w_c + \lambda_T \left( \frac{1}{n_c} \sum_{j=1}^{n_c} \frac{g(t_{c,j})}{\|g(t_{c,j})\|_2} \right)$$

"a photo of a [CLS]."

$\lambda_T$  : the intensity of intra-modal calibration within the text modality

$t_{c,j}$  : j-th LLM-generated description for class c

$n_c$  : the number of descriptions for class c

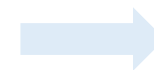
$g(\cdot)$ : the text encoder

These templates fails to capture category-specific nuances, making them particularly inadequate for fine-grained classification tasks

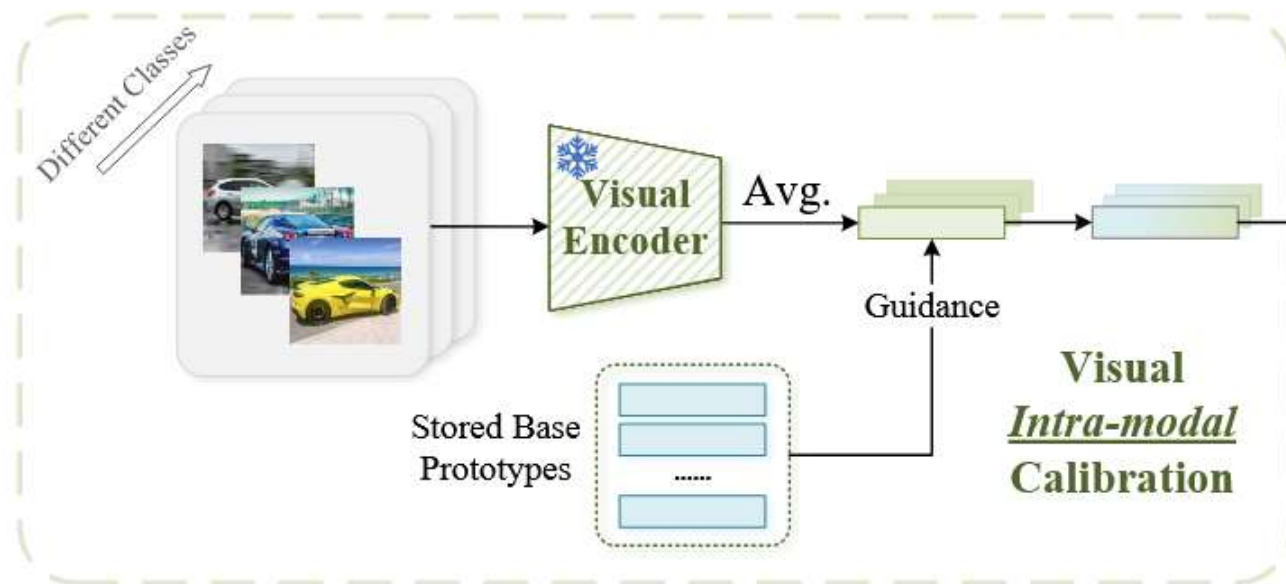
## Intra-modal classifier calibration

► Visual modality

the limited number of samples in incremental phase inevitably leads to biased estimates of new class prototypes.



using base class to calibrate



$$\tilde{\mu}_c^I = \begin{cases} \mu_c^I & , t = 0 \\ (1 - \lambda_I)\mu_c^I + \lambda_I \sum_{b=1}^{|\mathcal{Y}_0|} s_{b,c} \mu_b^I & , t > 0, \end{cases}$$

$\lambda_I$  : the strength of the visual intra-modal calibration

$s_{b,c}$  : the normalized cosine similarity between the visual prototypes of class b (**from the base classes**) and class c

$\langle \cdot, \cdot \rangle$ : cosine similarity

where

$$\mu_c^I = \frac{1}{m_c} \sum_{j=1}^{m_c} \frac{f(\mathbf{x}_{c,j})}{\|f(\mathbf{x}_{c,j})\|_2}, \quad s_{b,c} = \frac{e^{\tau \cdot \langle \mu_b, \mu_c \rangle}}{\sum_{i=1}^{|\mathcal{Y}_0|} e^{\tau \cdot \langle \mu_i, \mu_c \rangle}}$$

$\mu_c^I$ : naive visual prototype of class c.

$m_c$  : the number of training images in class c

$f(\cdot)$ : the visual encoder

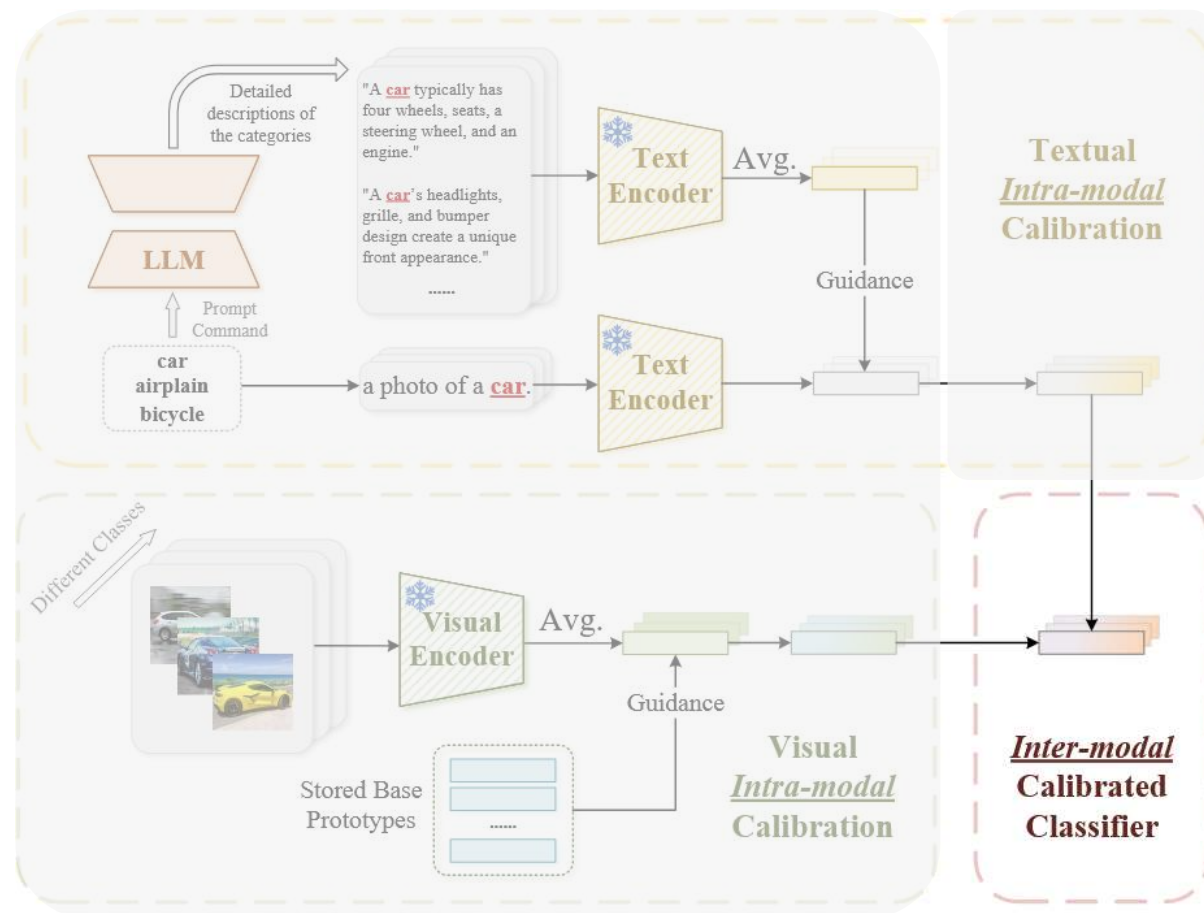
$x_{c,j}$  : j-th image belonging to class c

## Inter-modal classifier calibration

$$\mu_c = \beta \tilde{\mu}_c^T + (1 - \beta) \tilde{\mu}_c^I.$$

### ► Inference

$$s_c^{\text{calib}} = \frac{f(\mathbf{x})^\top \mu_c}{\|f(\mathbf{x})\|_2 \cdot \|\mu_c\|_2}.$$



## Semantic and covariance-enhanced metric

previous metrics effectively mitigate modal bias, but they primarily measure sample distances from distribution centers. This isotropic approach struggles to capture higher-order information and complex data relationships.

leveraging statistical information from existing data for more comprehensive and efficient measurements

avoid singularity

► For task  $t$ :

$$\tilde{\Sigma}^t = \Sigma^t + \frac{\gamma}{d} \text{tr}(\Sigma^t) \mathbf{I}_d,$$

$\gamma$ : the strength of regularization

$d$ : feature dimension

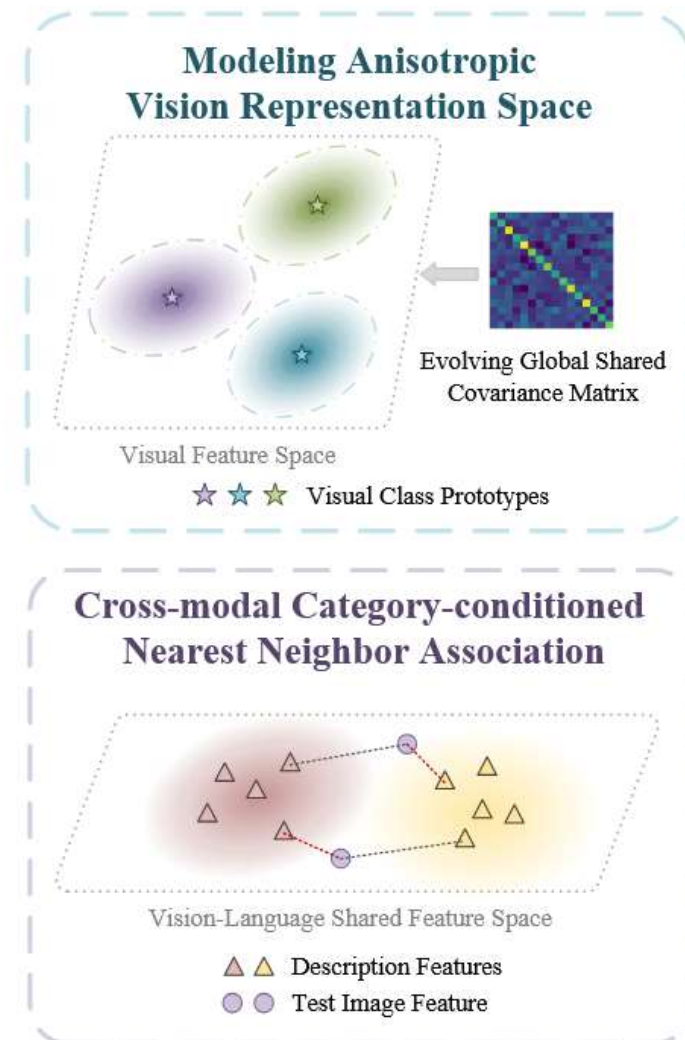
$\text{tr}(\cdot)$ : trace operator of a matrix

continuously evolving shared covariance matrix

$$\tilde{\Sigma}_G^t = \frac{|\mathcal{Y}_{t-1}|}{|\mathcal{Y}_t|} \tilde{\Sigma}_G^{t-1} + \left(1 - \frac{|\mathcal{Y}_{t-1}|}{|\mathcal{Y}_t|}\right) \tilde{\Sigma}^t.$$

Given a test sample  $\mathbf{x}$ , calculate its score as:

$$\mathbf{s}_c^{\text{cov}}(\mathbf{x}) = -\frac{1}{d} (f(\mathbf{x}) - \tilde{\mu}_c^I)^T \tilde{\Sigma}_G^{-1} (f(\mathbf{x}) - \tilde{\mu}_c^I).$$



(b) Semantic and Covariance-enhanced Metric

## Cross-modal category nearest neighbor metric

- ✓ To better leverage LLM-generated category descriptions

$$s_c^{\text{nn}} = \max_j \{ z_{c,j}^\top \cdot \mathbf{v} \},$$

where 
$$z_{c,j} = \frac{g(\mathbf{t}_{c,j})}{\|g(\mathbf{t}_{c,j})\|_2} \quad \mathbf{v} = \frac{f(\mathbf{x})}{\|f(\mathbf{x})\|_2}.$$

## Inference score reorganization strategy

base-derived covariance matrices perform poorly in distinguishing new class data

- covariance from extensive base data fails to align with new classes
- the limited data available for new classes prevents accurate covariance estimation

$$\mathbf{p}_c = \begin{cases} \alpha \mathbf{p}_c^{\text{calib}} + (1 - \alpha) \mathbf{p}_c^{\text{cov}} & , c \in \mathcal{Y}_0 \\ \alpha \mathbf{p}_c^{\text{calib}} + (1 - \alpha) \mathbf{p}_c^{\text{nn}} & , c \notin \mathcal{Y}_0. \end{cases}$$

for new classes

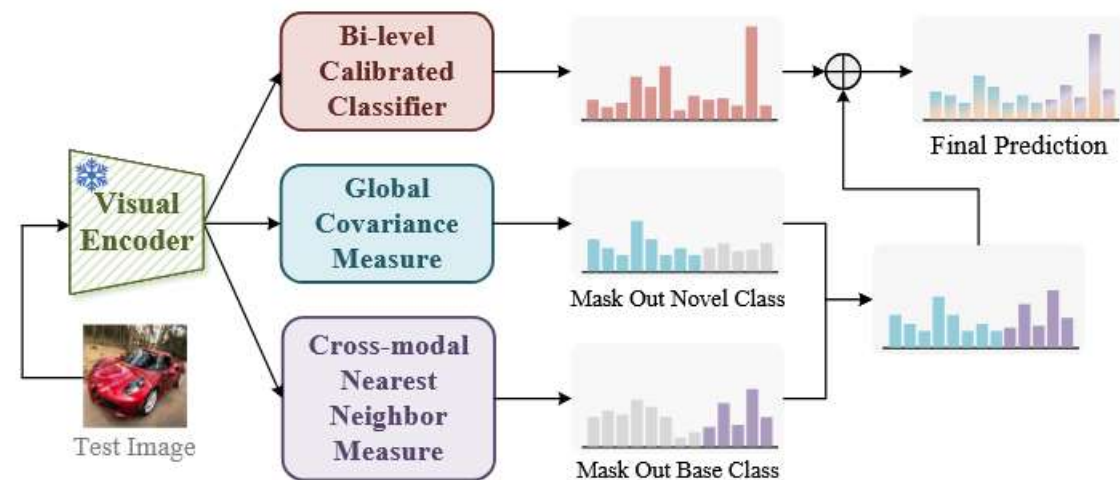


Figure 3. In the inference phase, an ensemble strategy with category masking is utilized.

# Experiments

Table 1. Detailed session-wise accuracy, average accuracy (Avg) and performance degradation (PD) comparison on *miniImageNet* dataset. **V** and **L** represent the visual and language modalities, respectively. BiMC refers to the results obtained solely through the bi-level calibration framework, whereas BiMC<sup>†</sup> incorporates the ensemble classifier strategy. The best results on each each sessions are indicated in **bold**, while the second-best results are underlined.  $\uparrow$  means higher is better, while  $\downarrow$  means lower is better.

Method	Modality	Accuracy in each session(%) $\uparrow$										Avg $\uparrow$	PD $\downarrow$
		0	1	2	3	4	5	6	7	8			
CLIP Zero-Shot [24]	<b>L</b>	91.27	91.25	89.74	89.43	88.98	88.42	86.96	86.64	86.15	88.76	5.12	
Visual Prototype [29]	<b>V</b>	92.58	91.88	89.37	88.24	88.42	87.86	86.33	86.16	86.09	88.55	6.49	
TEEN [35]	<b>V</b>	92.58	92.06	89.69	88.60	88.80	88.31	86.84	86.71	86.66	88.92	5.92	
FeCAM [8]	<b>V</b>	94.42	93.71	91.61	90.59	90.62	90.22	88.43	88.17	88.20	90.66	6.22	
BiMC	<b>V-L</b>	<u>94.90</u>	<u>94.80</u>	<u>93.21</u>	<u>92.91</u>	<u>92.88</u>	<u>92.60</u>	<u>91.78</u>	<u>91.88</u>	<u>91.81</u>	<u>92.97</u>	3.09	
BiMC <sup>†</sup>	<b>V-L</b>	<b>95.47</b>	<b>95.34</b>	<b>93.80</b>	<b>93.63</b>	<b>93.51</b>	<b>93.18</b>	<b>92.51</b>	<b>92.56</b>	<b>92.40</b>	<b>93.60</b>	<b>3.07</b>	

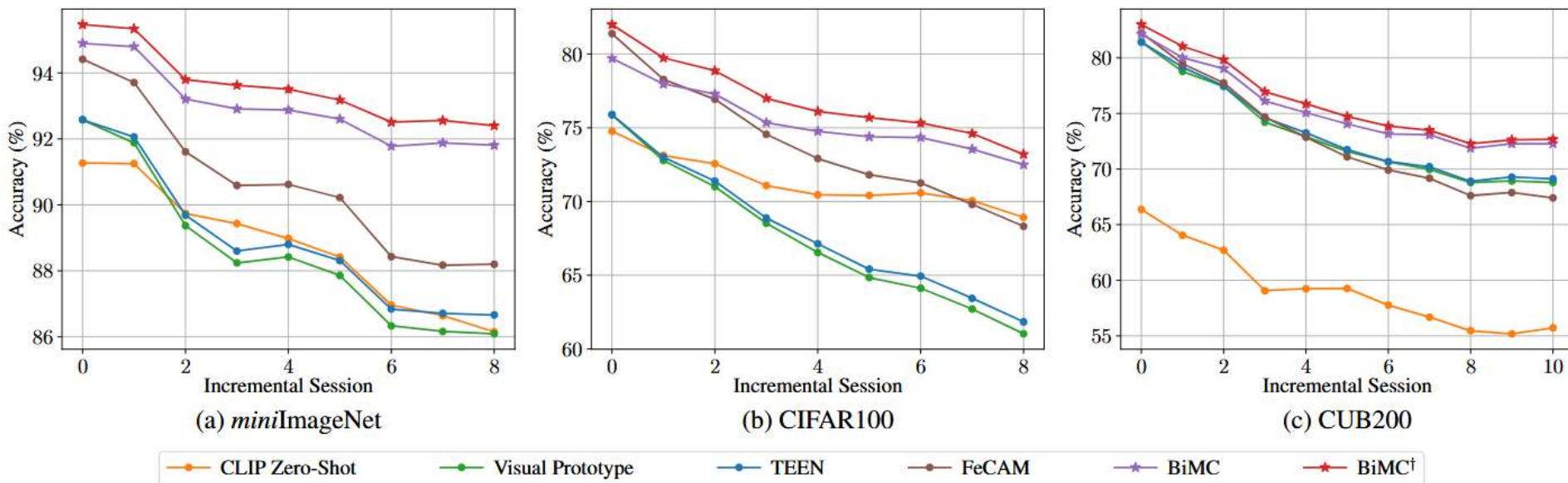


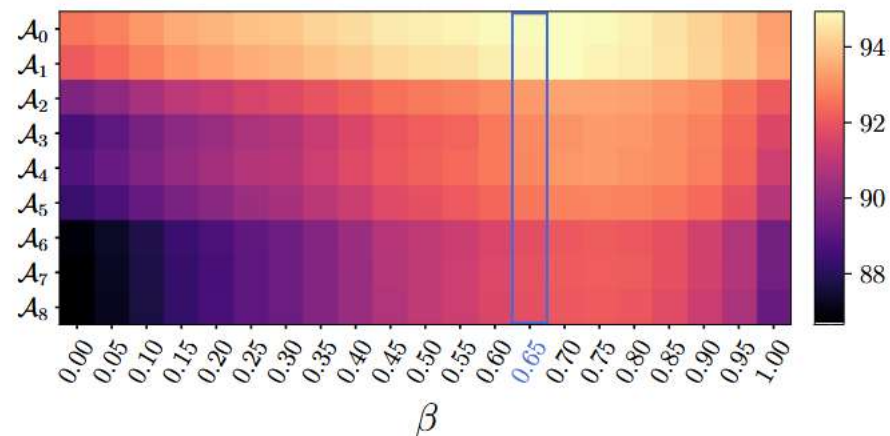
Figure 4. Performance curve of each incremental sessions on (a) *miniImageNet*, (b) CIFAR100 and (c) CUB200 datasets.

## Comparison with trainable methods

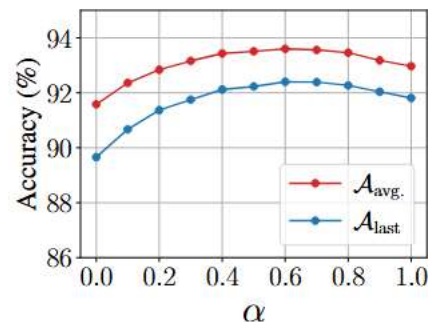
Table 5. Comparison with three trainable methods on the *miniImageNet* dataset.  $N_p$  is the number of parameters which require training.  $\Delta\mathcal{A}_{\text{last}}$  reflects the last task’s performance gap between our method and the comparative one.

Method	$N_p$	$\mathcal{A}_{\text{base}}$	$\mathcal{A}_{\text{last}}$	$\mathcal{A}_{\text{avg.}}$	$\Delta\mathcal{A}_{\text{last}}$
CPE-CLIP [5]	400k	90.23	82.77	86.13	+9.63
CLIP-M <sup>3</sup> [6]	46k	96.00	92.50	94.10	-0.10
LP-DiF [10]	8.1k	96.34	91.68	93.76	+0.72
BiMC <sup>†</sup>	<b>0</b>	95.47	92.40	93.60	0.00

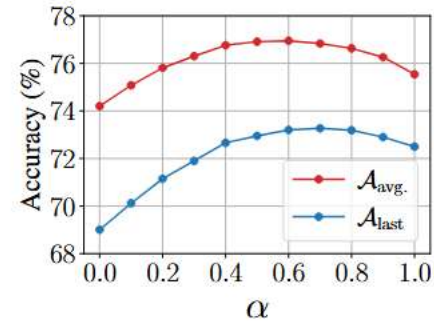
## Analysis on hyper-parameters



(a) Influence of  $\beta$  and session-wise accuracy on *miniImageNet*



(b) Influence of  $\alpha$  on *miniImageNet*



(c) Influence of  $\alpha$  on CIFAR100

Figure 7. Analysis on hyper-parameters  $\beta$  and  $\alpha$ . For  $\beta$ , we present heatmaps illustrating accuracy under varying  $\beta$  across different sessions. For  $\alpha$ , we report the average task accuracy  $\mathcal{A}_{\text{avg.}}$  and the final task accuracy  $\mathcal{A}_{\text{last}}$  under different settings.

Table 2. Ablation studies of bi-level calibration framework on *miniImageNet*. **Vis.** and **Lang.** individually represent classifiers derived from the visual and language modalities. **Intra-C** and **Inter-C** refer to the strategies of intra-modal and inter-modal calibration, respectively.

Vis	Lang	Intra-C	Inter-C	Accuracy in each session(%) $\uparrow$								Avg $\uparrow$	PD $\downarrow$	
				0	1	2	3	4	5	6	7			8
✓				92.58	91.88	89.37	88.24	88.42	87.86	86.33	86.16	86.09	88.55	6.49
✓		✓		92.58	92.06	89.69	88.60	88.80	88.31	86.84	86.71	86.66	88.92	5.92
	✓			91.27	91.25	89.74	89.43	88.98	88.42	86.96	86.64	86.15	88.76	5.12
	✓	✓		93.27	93.32	91.91	91.51	91.21	90.73	89.28	89.23	89.01	91.05	4.26
✓	✓		✓	94.68	94.58	92.87	92.49	92.56	92.31	91.59	91.58	91.54	92.69	3.14
✓	✓	✓	✓	<b>94.90</b>	<b>94.80</b>	<b>93.21</b>	<b>92.91</b>	<b>92.88</b>	<b>92.60</b>	<b>91.78</b>	<b>91.88</b>	<b>91.81</b>	<b>92.97</b>	<b>3.09</b>

Table 4. Ablation study of Semantic and Covariance-enhanced Metric on the *miniImageNet* dataset. To verify their effectiveness, we measured the performance of the base task  $\mathcal{A}_{\text{base}}$ , the performance of last task  $\mathcal{A}_{\text{last}}$  and the average performance  $\mathcal{A}_{\text{avg}}$  across all tasks. Additionally, we report the accuracy of base  $\mathcal{A}_{\text{last}}^b$  and novel  $\mathcal{A}_{\text{last}}^n$  class in last task.

Ablation	$\mathcal{A}_{\text{base}}$	$\mathcal{A}_{\text{last}}$	$\mathcal{A}_{\text{last}}^b$	$\mathcal{A}_{\text{last}}^n$	$\mathcal{A}_{\text{avg}}$
BiMC	94.90	91.81	93.43	89.38	92.97
+ MGC	<b>95.47</b>	91.51	<b>94.37</b>	87.22	92.91
+ CMNN	94.85	92.15	93.30	90.42	93.18
BiMC <sup>†</sup> (w/o mask.)	<b>95.47</b>	92.20	94.12	89.32	93.37
BiMC <sup>†</sup>	<b>95.47</b>	<b>92.40</b>	93.20	<b>91.20</b>	<b>93.60</b>

Thanks!