



# Learning Multi-Source and Robust Representations for Continual Learning

Fei Ye<sup>1</sup>, Yongcheng Zhong<sup>1</sup>, Qihe Liu<sup>1</sup>, Adrian G. Bors<sup>2</sup>, Jingling Sun<sup>1</sup>, Rongyao Hu<sup>1</sup>, Shijie Zhou<sup>1</sup>

<sup>1</sup>University of Electronic Science and Technology of China

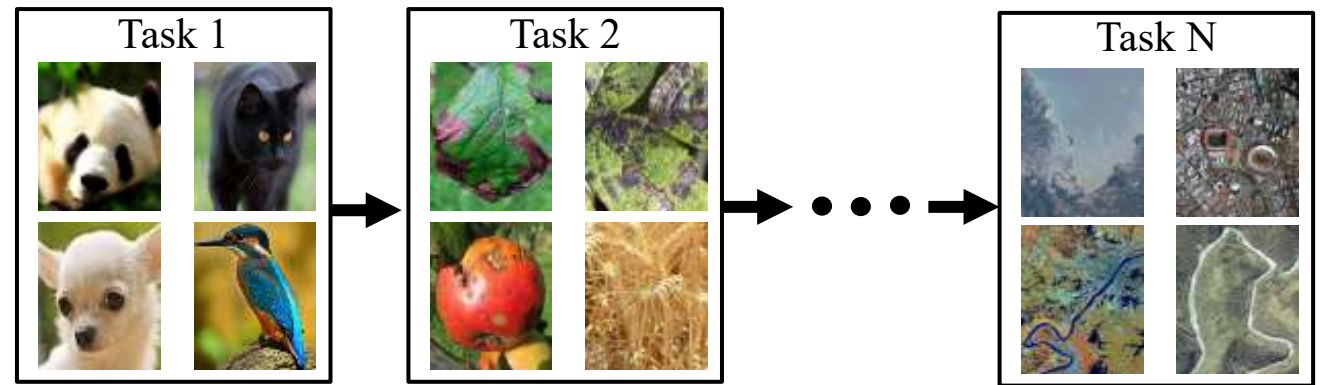
<sup>2</sup>University of York, UK

# Background

## Continual Learning (CL)

- Goal: Enable models to learn continuously from a sequence of tasks.
- Key challenge: Catastrophic forgetting — new knowledge overwrites old ones.

- Two essential abilities:  
**Plasticity** – adapt quickly to new tasks.  
**Stability** – preserve past knowledge.



## Existing Solutions

Category	Key Idea	Limitation
Rehearsal-based	Store a small memory buffer of old samples	Limited memory → poor scalability
Dynamic expansion	Add new sub-networks for new tasks	Leads to parameter growth
Regularization-based	Constrain parameter updates	Over-regularization → reduced plasticity

# Motivation

## Limitation of Pretrained Backbones(CL)

- Recent CL methods use pretrained ViTs or CNNs to improve stability.
- However:
  - Relying on a **single pretrained backbone** restricts adaptability.
  - Fixed representations often fail to generalize to new domains.
  - Updating too many layers → instability; freezing too many → rigidity.

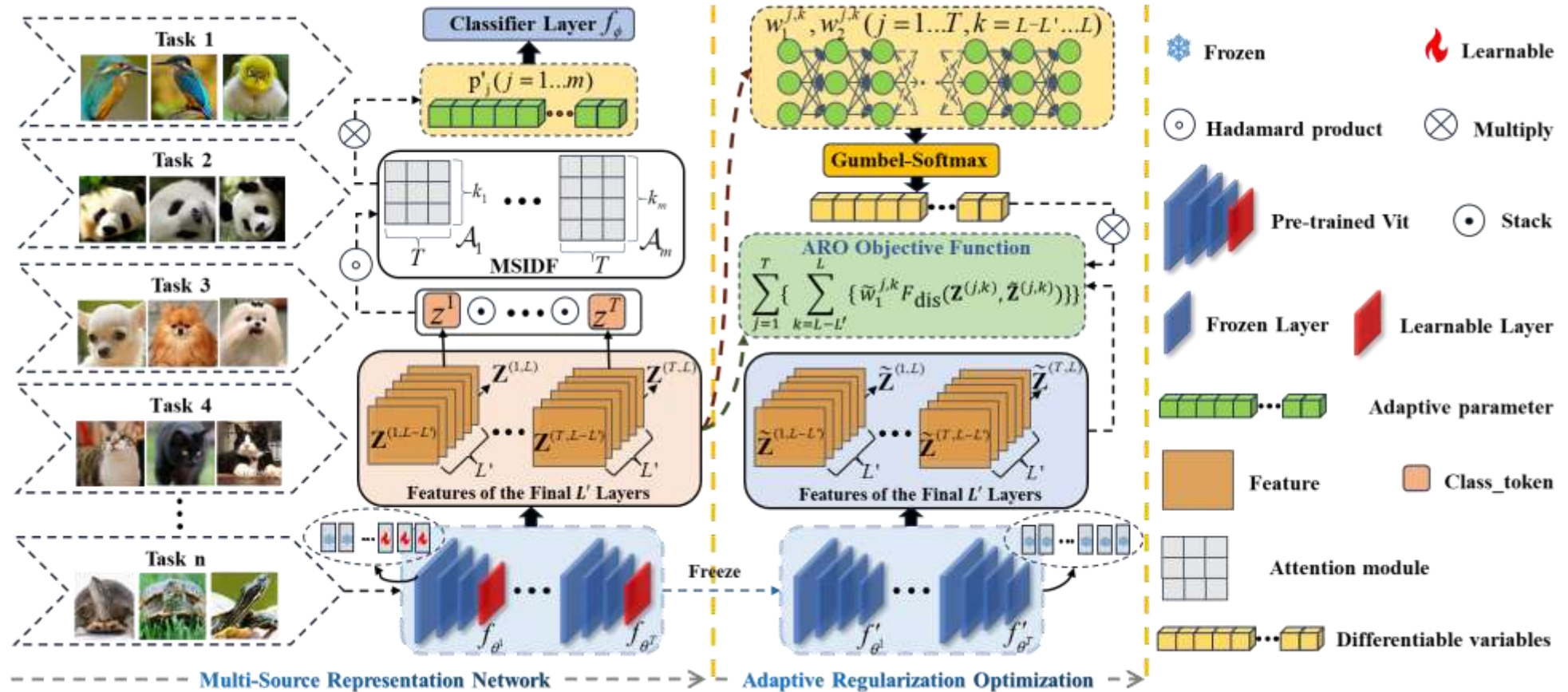
## Our Motivation

Can we **leverage multiple pretrained models** to build a **robust and adaptive representation space**, achieving both **stability** and **plasticity** without architectural expansion?

## LMSRR Framework

- A novel framework for learning multi-source and robust representations in continual learning.
- Core idea:
  - Fuse features from **multiple ViT backbones**
  - Adapt representations through **dynamic and adaptive optimization**

# Framework: LMSRR



LMSRR dynamically fuses multi-source features and adaptively regulates representation learning, achieving an optimal **balance between stability and plasticity** in continual learning.

# Method I: Multi-Scale Interaction and Dynamic Fusion

(1) Stack features:

$$\tilde{\mathbf{z}}_s = f_{\theta^1}(\mathbf{x}_s) \bullet \cdots \bullet f_{\theta^T}(\mathbf{x}_s)$$

(3) Adaptive fusion:

- Learn weights  $p_j \rightarrow \text{softmax} \rightarrow p'_j = \frac{\exp(p_j)}{\sum_{c=1}^m \exp(p_c)}$

- Final representation:

$$\mathbf{z}_s = \sum_{j=1}^m p'_j \cdot \mathbf{z}_s^j$$

(2) Multi-scale attention:

- Apply  $m$  learnable attention modules  $\mathcal{A}_j$  with window size  $k_j$ .
- Each module highlights relevant cross-backbone patterns via  $W^j$ .

Advantages:

- Output dimension fixed (independent of  $T$ ).
- Automatically focuses on most informative backbone combinations per input.

**Goal of MSIDF:** Fuse multi-source ViT features without redundancy, with adaptivity.

# Method II : Multi-Level Representation Optimization (MLRO)

Mechanism:

- Keep a frozen copy of each ViT from previous task ( $f'_{\theta_j}$ )
- For current task, extract features from last  $L'$  layers of both:
  - Active network:  $\mathbf{Z}^{(j,k)}$
  - Frozen network:  $\tilde{\mathbf{Z}}^{(j,k)}$
- Minimize L2 distance between them:

$$\mathcal{L}_{\text{MLRO}} = \sum_{j=1}^T \sum_{k=L-L'}^L \|\mathbf{Z}^{(j,k)} - \tilde{\mathbf{Z}}^{(j,k)}\|_2$$

**Goal of MLRO:** Prevent representation drift during fine-tuning → preserve stability.

# Method III: Adaptive Regularization Optimization (ARO)

**Problem:** MLRO applies uniform regularization  $\rightarrow$  may over-constrain some layers.

**Solution:** Introduce learnable gate per layer to control regularization strength.

**Key Idea:**

For each layer  $k$  in backbone  $j$ , learn switch  $(w_1^{j,k}, w_2^{j,k})$

Use Gumbel-Softmax to get differentiable weight  $\tilde{w}_1^{\{j,k\}}$

$$w_1^{j,k} = \frac{\exp((\log(w_1^{j,k}) + g_1) / \tau)}{\sum_{t=1}^2 \{\exp((\log(w_t^{j,k}) + g_t) / \tau)\}}$$

Modify loss:

$$L_{ARO} = \sum_{j=1}^T \left\{ \sum_{k=L-L'}^L \left\{ w_1^{j,k} F_{\text{dis}}(\mathbf{Z}^{(j,k)}, \mathbf{Z}^{(j,k)}) \right\} \right\}$$

**Why it works:**

If a layer is critical for new task :

$$\tilde{w}_1^{j,k} \approx 0 \rightarrow \text{less}$$

constraint.

If a layer encodes old knowledge :

$$\tilde{w}_1^{j,k} \approx 1 \rightarrow \text{strong}$$

constraint.

# Experiments Results

## Standard datasets results

Buffer	Method	CIFAR-10		Tiny ImageNet		R-MNIST
		Average	Last	Average	Last	Domain-IL
-	EWC <a href="#">[51]</a>	68.29±3.92	<b>97.07±0.74</b>	19.20±0.31	75.15±3.18	77.35±5.77
	SI <a href="#">[63]</a>	68.05±5.91	94.18±0.88	36.32±0.13	65.80±3.25	71.91±5.83
	LwF <a href="#">[37]</a>	63.29±2.35	96.75±0.35	15.85±0.58	77.95±3.60	-
	PNN <a href="#">[50]</a>	95.13±0.72	96.63±0.10	67.84±0.29	69.03±1.01	-
	DAP <a href="#">[27]</a>	<b>97.13±2.06</b>	96.05±3.39	<b>92.49±0.60</b>	<b>94.95±1.20</b>	<b>88.58±2.53</b>
200	ER <a href="#">[49]</a>	91.19±0.94	97.50±0.35	38.17±2.00	79.40±0.28	85.01±1.90
	GEM <a href="#">[39]</a>	90.44±0.94	96.60±0.35	-	-	80.80±1.15
	A-GEM <a href="#">[12]</a>	83.88±1.49	97.90±0.07	22.77±0.03	78.65±3.32	81.91±0.76
	iCaRL <a href="#">[48]</a>	88.99±2.13	97.07±0.32	28.19±1.47	47.45±0.78	-
	FDR <a href="#">[7]</a>	91.01±0.68	97.78±0.24	40.36±0.68	81.40±0.70	85.22±3.35
	GSS <a href="#">[3]</a>	88.80±2.89	97.42±0.24	-	-	79.50±0.41
	HAL <a href="#">[11]</a>	82.51±3.20	94.60±0.14	-	-	84.02±0.98
	DER <a href="#">[8]</a>	91.40±0.92	97.80±0.28	40.22±0.67	79.15±0.21	90.04±2.61
	DER++ <a href="#">[8]</a>	91.92±0.60	97.72±0.38	40.87±1.16	78.35±0.49	90.43±1.87
	DER++(re) <a href="#">[56]</a>	92.01±3.03	97.65±3.03	47.61±8.87	81.40±1.41	91.64±2.26
	<b>Ours</b>	<b>98.85±0.05</b>	<b>99.35±0.21</b>	<b>92.08±0.31</b>	<b>96.00±0.01</b>	<b>94.20±1.24</b>
500	ER <a href="#">[49]</a>	93.61±0.27	97.15±0.28	48.64±0.46	80.80±1.69	88.91±1.44
	GEM <a href="#">[39]</a>	92.16±0.69	96.63±0.17	-	-	81.15±1.98
	A-GEM <a href="#">[12]</a>	89.48±1.45	97.40±0.78	25.33±0.49	81.00±0.42	80.31±6.29
	iCaRL <a href="#">[48]</a>	88.22±2.62	96.57±0.10	31.55±3.27	50.65±1.20	-
	FDR <a href="#">[7]</a>	93.29±0.59	97.32±0.24	49.88±0.71	81.10±0.56	89.67±1.63
	GSS <a href="#">[3]</a>	91.02±1.57	96.97±0.24	-	-	81.58±0.58
	HAL <a href="#">[11]</a>	84.54±2.36	94.22±0.60	-	-	85.00±0.96
	DER <a href="#">[8]</a>	93.40±0.39	97.90±0.28	51.78±0.88	79.30±1.13	92.24±1.12
	DER++ <a href="#">[8]</a>	93.88±0.50	98.10±0.01	51.91±0.68	76.20±5.23	92.77±1.05
	DER++(re) <a href="#">[56]</a>	93.06±0.38	97.75±0.38	54.06±0.79	79.65±1.34	93.28±0.75
	<b>Ours</b>	<b>99.15±0.05</b>	<b>99.48±0.04</b>	<b>92.75±0.32</b>	<b>96.23±0.40</b>	<b>96.97±1.58</b>
1000	ER <a href="#">[49]</a>	95.34±0.16	97.67±0.67	55.92±0.90	80.30±0.82	90.42±1.07
	GEM <a href="#">[39]</a>	93.67±0.32	97.37±0.17	-	-	81.15±1.98
	A-GEM <a href="#">[12]</a>	85.61±2.01	97.45±0.42	24.29±1.28	79.65±2.19	81.30±5.33
	iCaRL <a href="#">[48]</a>	91.40±1.06	96.85±0.35	63.87±0.25	54.00±2.82	-
	FDR <a href="#">[7]</a>	94.02±0.64	97.60±0.56	56.05±0.71	80.25±0.49	91.68±1.01
	GSS <a href="#">[3]</a>	91.79±2.16	96.10±1.70	-	-	82.25±2.42
	HAL <a href="#">[11]</a>	87.33±1.46	92.27±3.21	-	-	89.33±2.01
	DER <a href="#">[8]</a>	92.33±0.61	97.72±0.07	56.62±1.13	78.50±0.42	93.13±0.28
	DER++ <a href="#">[8]</a>	94.99±0.26	97.94±0.08	58.05±0.52	79.95±0.35	93.82±0.39
	DER++(re) <a href="#">[56]</a>	93.66±1.00	97.40±0.01	61.91±1.15	80.45±3.18	93.37±0.58
	<b>Ours</b>	<b>99.21±0.06</b>	<b>99.43±0.03</b>	<b>93.24±0.24</b>	<b>96.10±0.57</b>	<b>97.05±0.04</b>

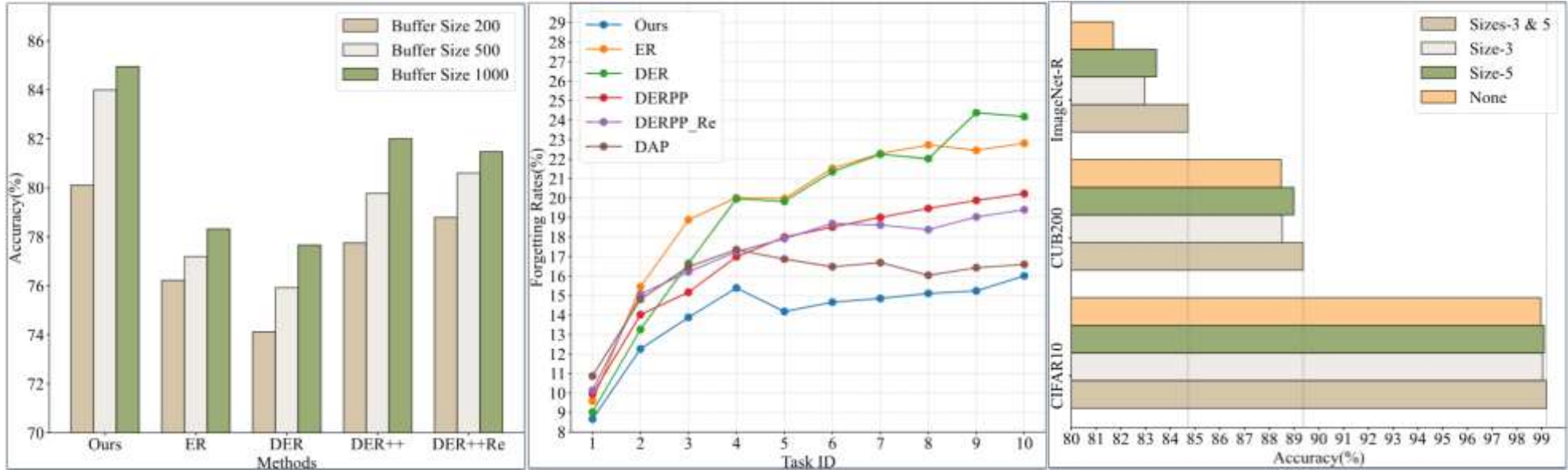


# Experiments Results

## Complex datasets results

Method	CIFAR-100		CUB-200		Imagenet-R		Cars196	
	Average	Last	Average	Last	Average	Last	Average	Last
ER [49]	73.37 $\pm$ 0.43	93.35 $\pm$ 1.34	30.57 $\pm$ 4.81	35.57 $\pm$ 14.86	24.85 $\pm$ 4.06	45.85 $\pm$ 0.01	30.52 $\pm$ 4.4	54.32 $\pm$ 5.07
A-GEM [12]	48.06 $\pm$ 0.57	92.80 $\pm$ 0.32	13.22 $\pm$ 0.31	42.18 $\pm$ 0.01	16.87 $\pm$ 2.65	47.56 $\pm$ 12.31	8.07 $\pm$ 0.15	16.45 $\pm$ 7.41
FDR [7]	76.29 $\pm$ 1.44	93.60 $\pm$ 1.34	23.94 $\pm$ 0.07	45.58 $\pm$ 0.19	15.74 $\pm$ 3.69	42.14 $\pm$ 10.75	31.41 $\pm$ 1.30	58.36 $\pm$ 1.17
GSS [3]	57.50 $\pm$ 1.93	92.80 $\pm$ 2.98	27.04 $\pm$ 0.28	42.01 $\pm$ 0.08	17.83 $\pm$ 0.88	33.44 $\pm$ 6.75	34.67 $\pm$ 2.27	56.80 $\pm$ 4.15
DER [8]	74.93 $\pm$ 1.06	93.25 $\pm$ 0.35	26.19 $\pm$ 2.07	51.79 $\pm$ 1.08	18.26 $\pm$ 1.67	25.26 $\pm$ 0.47	39.75 $\pm$ 0.36	68.02 $\pm$ 5.20
DER++ [8]	75.64 $\pm$ 0.60	92.60 $\pm$ 0.14	33.40 $\pm$ 1.48	49.83 $\pm$ 1.63	22.87 $\pm$ 5.83	43.10 $\pm$ 10.51	35.39 $\pm$ 3.38	60.56 $\pm$ 8.45
DER++refresh [56]	77.71 $\pm$ 0.85	93.40 $\pm$ 1.13	35.77 $\pm$ 3.20	50.85 $\pm$ 0.47	23.74 $\pm$ 3.03	31.00 $\pm$ 0.01	33.94 $\pm$ 2.46	60.29 $\pm$ 4.73
CoFiMA [41]	94.21 $\pm$ 0.47	96.13 $\pm$ 0.59	<b>90.66<math>\pm</math>0.76</b>	92.54 $\pm$ 0.28	83.76 $\pm$ 0.53	85.86 $\pm$ 0.58	87.28 $\pm$ 0.54	90.33 $\pm$ 0.45
DAP [27]	90.11 $\pm$ 0.33	92.30 $\pm$ 2.12	71.83 $\pm$ 1.44	72.23 $\pm$ 2.85	83.22 $\pm$ 1.25	84.61 $\pm$ 2.85	39.79 $\pm$ 1.85	65.35 $\pm$ 2.21
L2P [57]	95.36 $\pm$ 0.12	96.80 $\pm$ 0.14	86.30 $\pm$ 0.21	90.81 $\pm$ 0.24	<b>86.01<math>\pm</math>0.30</b>	87.50 $\pm$ 0.90	79.55 $\pm$ 0.86	84.45 $\pm$ 0.12
<b>Ours</b>	<b>95.76<math>\pm</math>0.08</b>	<b>98.70<math>\pm</math>0.37</b>	88.91 $\pm$ 0.64	<b>94.31<math>\pm</math>0.12</b>	84.35 $\pm$ 0.52	<b>88.43<math>\pm</math>0.15</b>	<b>90.14<math>\pm</math>0.06</b>	<b>95.32<math>\pm</math>0.39</b>

# Ablation Study



(a) Same backbone comparison.

(b) The forgetting curve.

(c) Configuration comparison.

(a) Comparison of performance of various models with varying buffer sizes on ImageNet-R, where each model uses the same backbone. (b) Comparison of forgetting curves of the proposed approach with other benchmark methods on ImageNet-R. (c) Performance variations of the proposed MSIDF method under different configurations.

# Conclusion

✓ We propose LMSRR, a novel continual learning framework that orchestrates multiple pre-trained ViTs to learn robust, adaptive representations.

✓ Our method introduces three key innovations:

- MSIDF: Dynamically fuses multi-source features via learnable attention and adaptive weighting.
- MLRO: Preserves stability by minimizing representation drift across tasks.
- ARO: Relieves over-regularization through layer-wise adaptive gating.

✓ LMSRR achieves state-of-the-art performance without expanding the model or using task-specific parameters.