



西安交通大学
XI'AN JIAOTONG UNIVERSITY



中国科学院深圳先进技术研究院
SHENZHEN INSTITUTE OF ADVANCED TECHNOLOGY
CHINESE ACADEMY OF SCIENCES



The 38th Annual AAAI
Conference on Artificial
Intelligence

FEBRUARY 20-27, 2024 | VANCOUVER, CANADA

Evolving Parameterized Prompt Memory for Continual Learning

Muhammad Rifki Kurniawan¹, Xiang Song¹, Zhiheng
Ma³, Yuhang He², Yihong Gong^{1,2}, Qi Yang⁴, Xing Wei¹

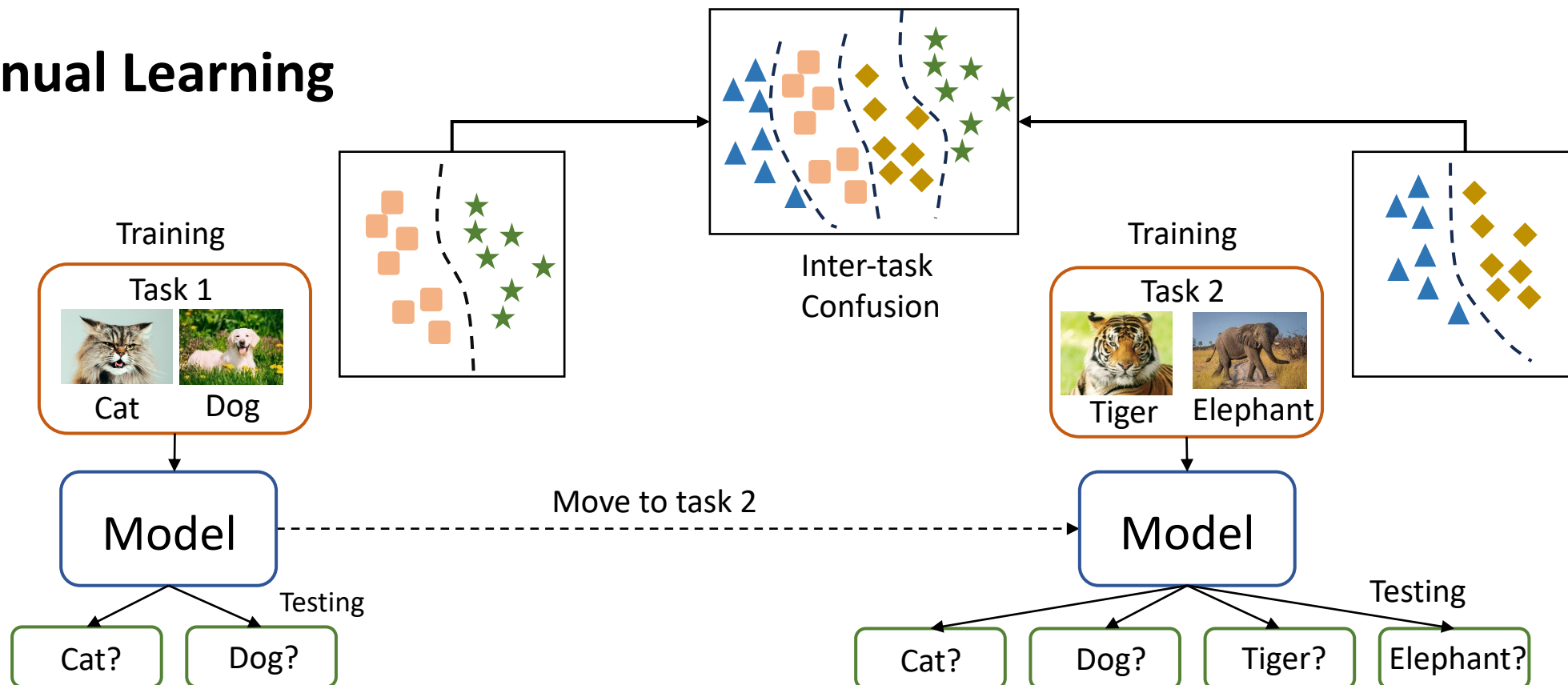
¹School of Software Engineering, Xi'an Jiaotong University

²College of Artificial Intelligence, Xi'an Jiaotong University

³Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences

⁴School of Computer Science and Technology, Xi'an Jiaotong University

Continual Learning



Training 1: These are {cat, dog}

Test-time 1: What is this [img] among {cat, dog}?

Training 2: These are {tiger, elephant}

Test-time 2: What is this [img] among {cat, dog, tiger, elephant}?

Continual Learning from Pre-trained

Existing CL From Pre-trained

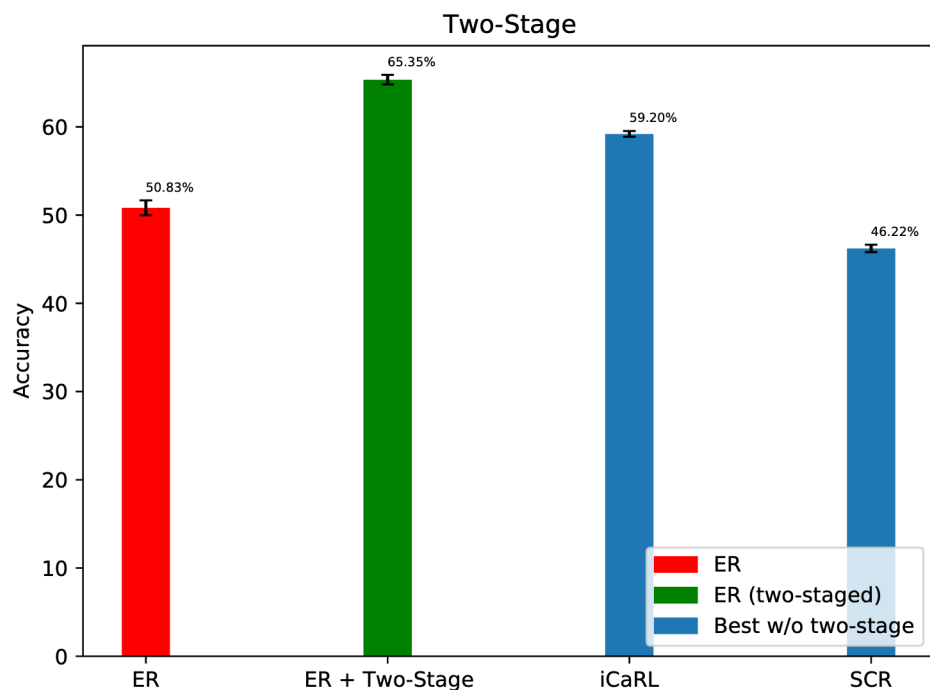
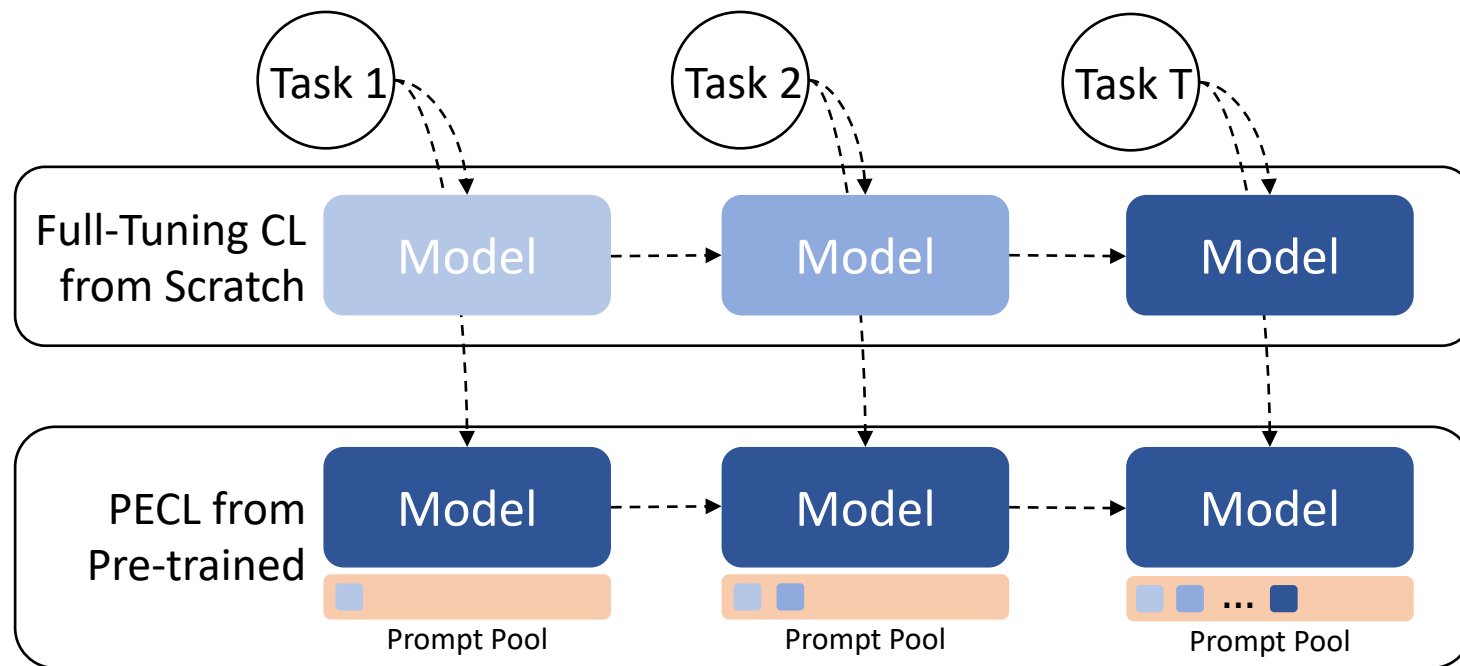


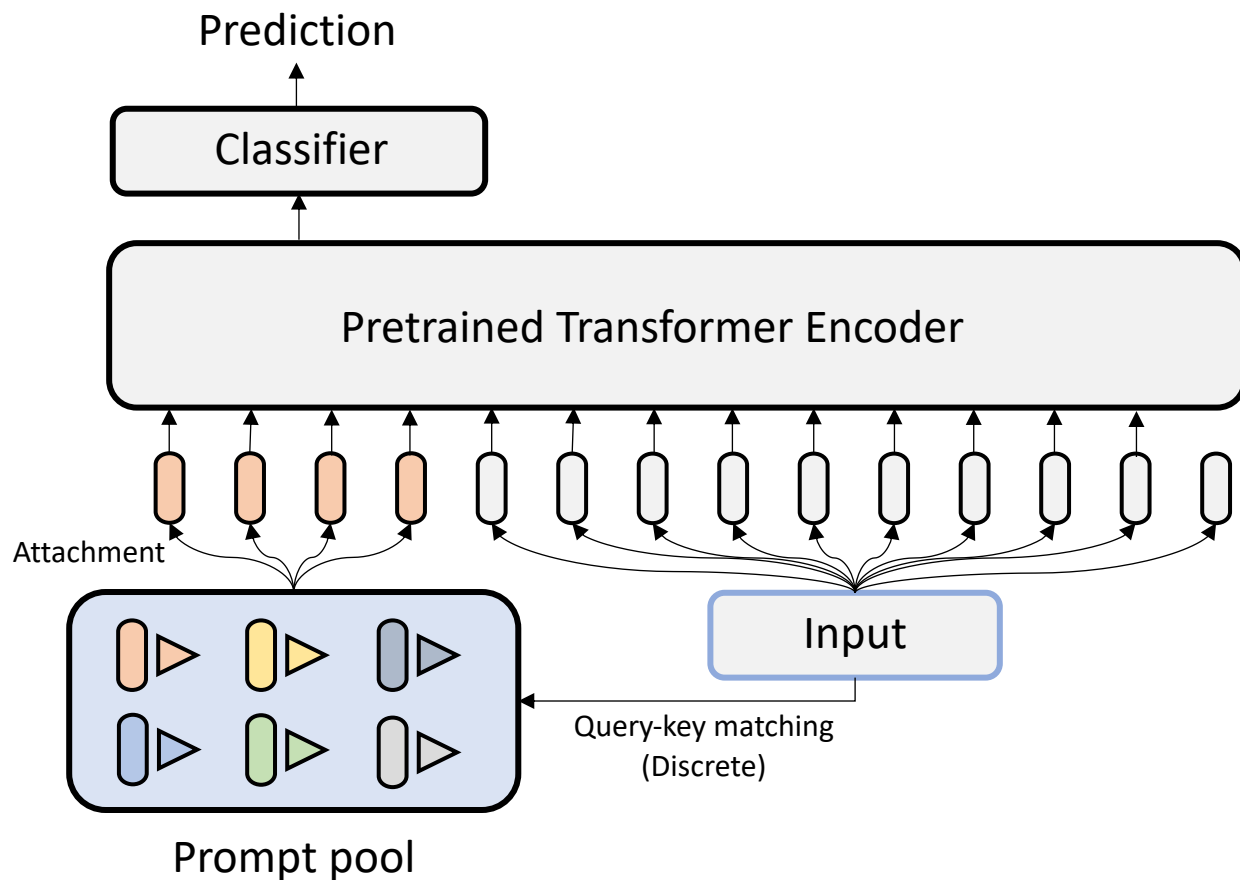
Image from [1].

From Full-tuning CL to Parameter Efficient CL

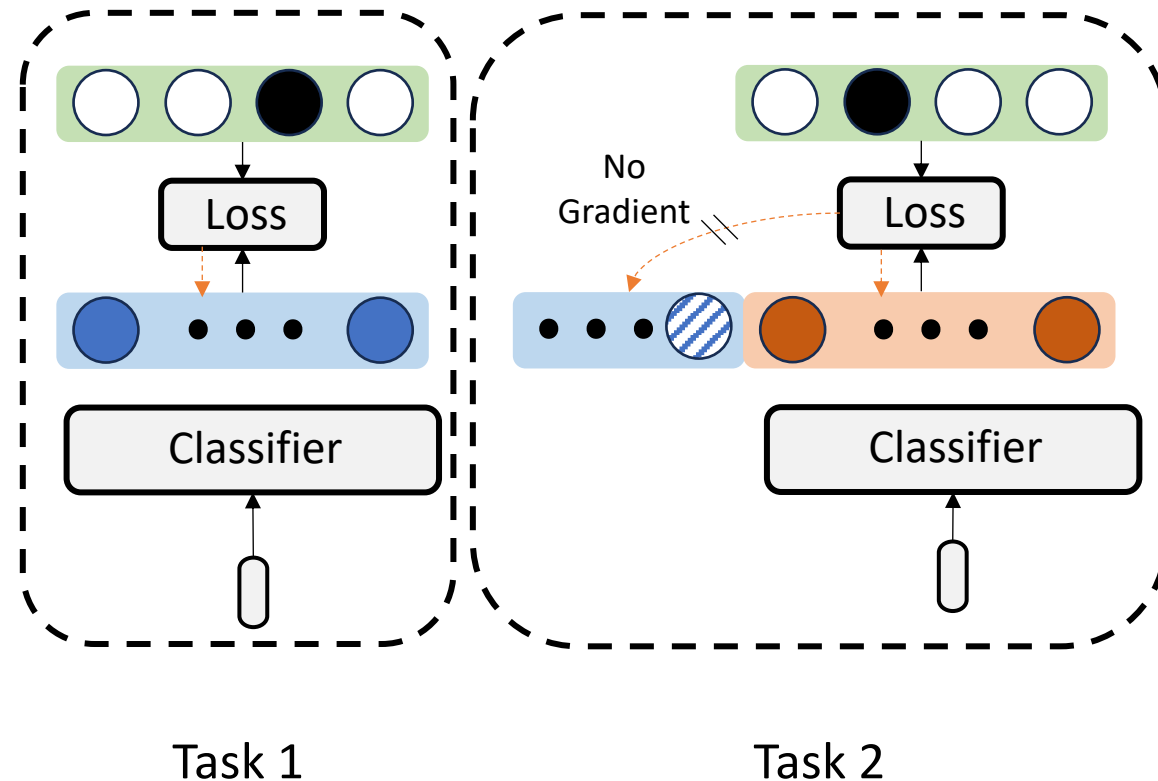


[1] K. Y. Lee, Y. Zhong, and Y. X. Wang, "Do Pre-Trained Models Benefit Equally in Continual Learning?" in Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), January 2023, pp. 6485-6493.

Learning To Prompt In Continual Learning



Architecture of Common
Prompt-based CL
(L2P, DualPrompt, S-Prompt)



Learning Classifier in Isolation
(Not backward Compatible)

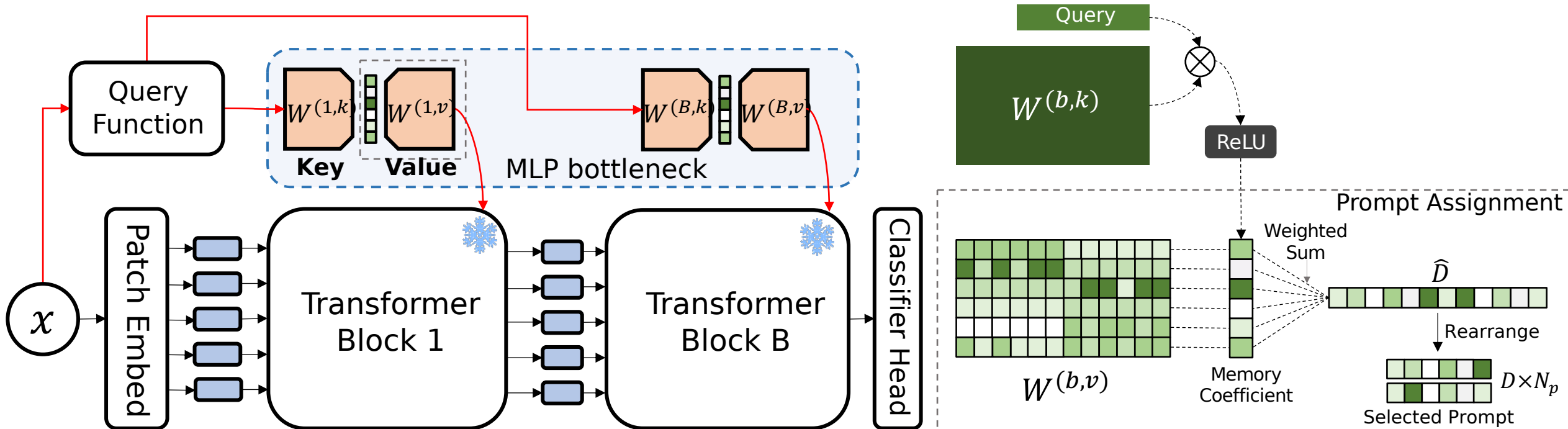
Continual Learning from Pre-trained

Method	Prompt Selection	Dynamically Expand?	Backward Compatible?	Additional Params M	Additional Params %
L2P	Discrete	✓	✗	0.89	1.21
DualPrompt	Discrete	✓	✗	0.95	1.28
CODA-P	Continuous	✓	✗	3.84	1.25
EvoPrompt	Continuous	✗	✓	0.29	0.69

Ours, Evolving Prompt (EvoPrompt):

- **Continuous** prompting
- **Non-expanded** prompt memory
- **Backward-compatible** classifier

Parameterized Prompt Memory



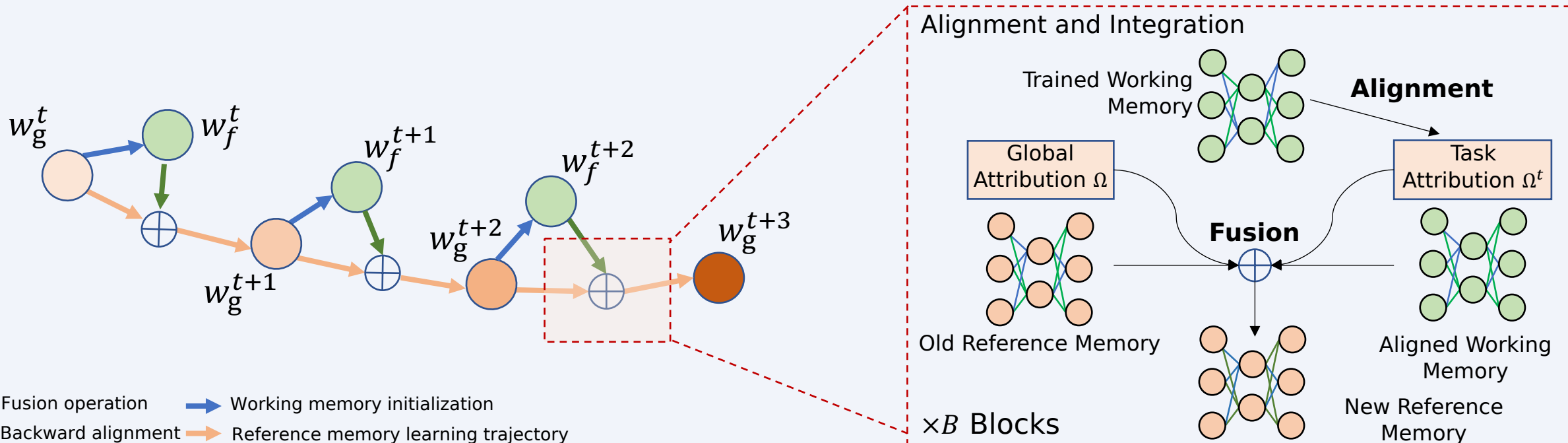
$$p_b = f_{\mathbf{W}^{(b)}} \left(q(x); \mathbf{W}^{(b)} \right),$$

$$= \underbrace{\text{ReLU} \left(q(x) \cdot \mathbf{W}^{(b,k)} \right)}_{\text{Query-key matching} \rightarrow \text{Prompt Coefficient}} \cdot \mathbf{W}^{(b,v)},$$

Query-key matching \rightarrow Prompt Coefficient

How to learn this memory without **catastrophic forgetting**?

Evolving Prompt Memory via Incremental Fusion



Alignment:

$$\hat{\mathbf{W}}_{f,\ell}^t \leftarrow \text{diag} \left(\frac{1}{\beta_\ell} \right) \mathbf{P}_\ell \mathbf{W}_{f,\ell}^t \mathbf{P}_{\ell-1}^T \text{diag} \left(\frac{1}{\beta_{\ell-1}} \right)$$

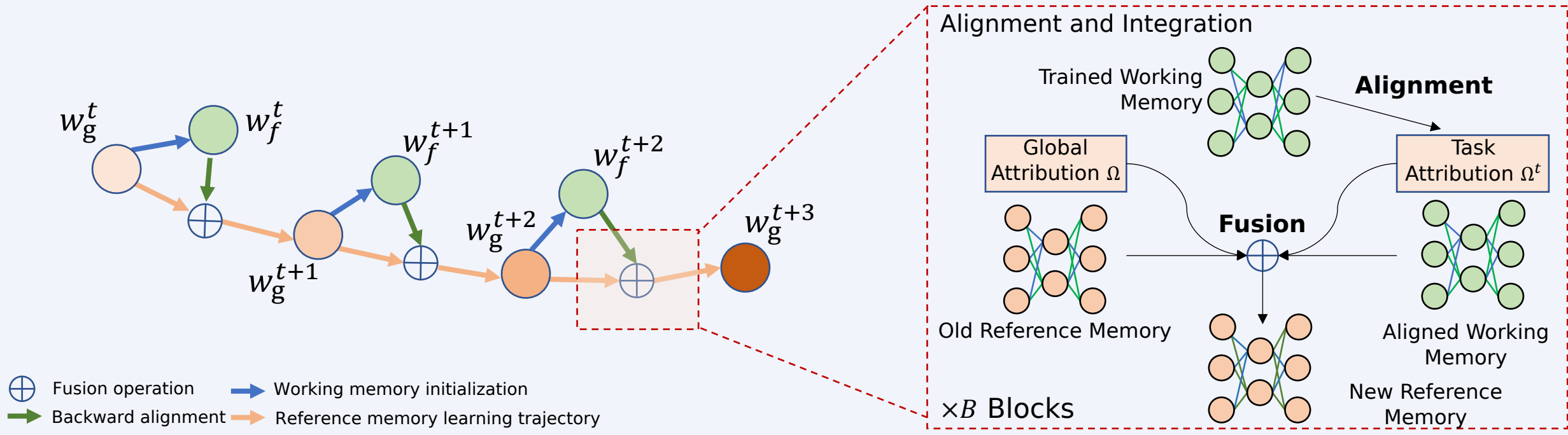
Fusion:

$$\mathbf{W}_{g,\ell}^{t+1} \leftarrow \lambda_\ell \hat{\mathbf{W}}_{f,\ell}^t + (1 - \lambda_\ell) \mathbf{W}_{g,\ell}^t,$$

Solving optimal transport (OT) problem:

$$\begin{aligned} \mathbf{P}_\ell = \min_{\mathbf{P}_\ell \in \mathbb{R}_+^{N_\ell \times N_\ell}} \text{tr} \left(\mathbf{P}_\ell^T \mathbf{D}_\ell \right) &= \text{OT}(\alpha_\ell, \beta_\ell, \mathbf{D}_\ell), \\ \text{s.t. } \mathbf{P}_\ell \mathbf{1}_\ell &= \alpha_\ell, \mathbf{P}_\ell^T \mathbf{1}_n = \beta_\ell, \end{aligned}$$

Evolving Prompt Memory via Incremental Fusion



Attribution-aware aggregation momentum:

$$\lambda_\ell = \lambda_{\text{base},\ell} + \Delta \hat{\Omega}_\ell \lambda_{\text{adj},\ell}, \quad \lambda \in [0, 1],$$

$$\lambda_{\text{adj},\ell} = \begin{cases} \lambda_{\text{base},\ell}, & \text{if } \hat{\Omega}_\ell^t - \hat{\Omega}_\ell < 0 \\ 1 - \lambda_{\text{base},\ell}, & \text{if } \hat{\Omega}_\ell^t - \hat{\Omega}_\ell \geq 0. \end{cases}$$

Set to 0.5 $\rightarrow \Omega_{\ell(j)}^t - \Omega_{\ell(j)}$

Momentum offset

Computing local task-specific attribution and global attribution:

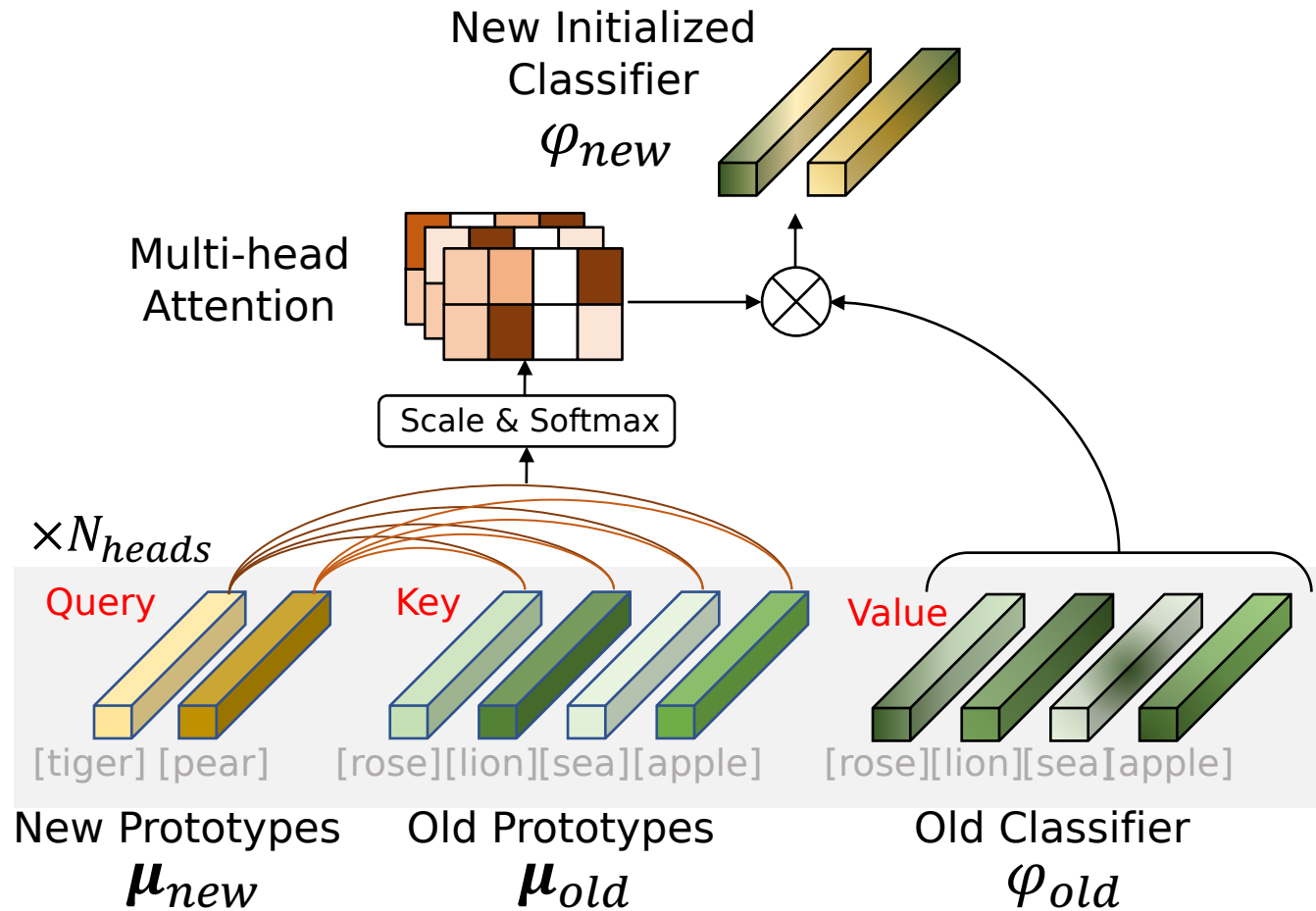
$$\Omega_{\ell(j)}^t = \frac{1}{|\mathcal{X}^t|} \sum_{x_i^t \in \mathcal{X}^t} \text{RELU}(f_{n_{\ell(j)}}(x_i^t)), \forall i, \hat{y}_i = y_i,$$

Task-specific attribution

$$\Omega_{\ell(j)} = \max(\Omega_{\ell(j)}, \Omega_{\ell(j)}^t),$$

Global attribution

Compositional Classifier Initialization (CCI)



Prototypical attention:

$$\varphi_{new} = \text{PA} = \text{Softmax} \left(\frac{d(\mu_{new}, \mu_{old})}{\tau} \right) \varphi_{old},$$

Extending PA into multi-head:

$$\begin{aligned} \varphi_{new} &= \text{MPA}(\mu_{new}, \mu_{old}, \varphi_{old}) \\ &= \text{Concat} \left(\text{PA}_1(\mu_{new,1}, \mu_{old,1}, \varphi_{old,1}), \right. \\ &\quad \left. \dots, \text{PA}_k(\mu_{new,k}, \mu_{old,k}, \varphi_{old,k}) \right). \end{aligned}$$

Empirical Results on Class Incremental Learning

- Split CIFAR-100: 50,000 training images, 1,000 testing images, 100 classes.

Method	Split CIFAR-100							
	5 Steps		10 Steps		20 Steps		Avg	Avg
	Acc.(↑)	Forget.(↓)	Acc.(↑)	Forget.(↓)	Acc.(↑)	Forget.(↓)	Acc.(↑)	Forget.(↓)
FT-seq	73.17 \pm 0.75	2.95 \pm 0.56	62.77 \pm 2.30	20.73 \pm 2.05	55.97 \pm 2.95	32.74 \pm 2.97	63.97 (+0.00)	18.81 (−0.00)
LP-seq	71.69 \pm 0.61	1.36 \pm0.27	66.90 \pm 0.53	13.08 \pm 0.32	60.98 \pm 0.74	21.27 \pm 1.20	66.52 (+2.55)	11.90 (−6.91)
NME-seq	78.30	7.70	78.33	1.14	78.33	2.68	78.32 (+14.35)	3.84 (−14.97)
L2P	86.53 \pm 0.14	7.67 \pm 0.20	84.97 \pm 8.21	8.21 \pm 0.22	83.39 \pm 0.41	10.18 \pm 0.24	84.96 (+20.99)	8.69 (−10.12)
DualPrompt	88.26 \pm 0.33	5.72 \pm 0.43	86.83 \pm 0.37	6.21 \pm 0.35	84.11 \pm 0.45	8.75 \pm 0.38	86.40 (+22.43)	6.89 (−11.92)
ESN	88.09 \pm 0.21	<u>5.18 \pm0.13</u>	85.96 \pm 0.14	4.54 \pm 0.35	82.71 \pm 0.51	6.44 \pm 0.31	85.59 (+21.62)	5.39 (−13.42)
CODA-P-S	88.90 \pm 0.26	<u>6.29 \pm0.27</u>	86.33 \pm 0.25	6.29 \pm 0.52	81.71 \pm 0.47	9.41 \pm 0.22	85.65 (+21.68)	7.33 (−11.48)
CODA-P	89.16 \pm0.26	6.08 \pm 0.33	87.31 \pm 0.14	5.95 \pm 0.41	81.69 \pm 0.38	9.85 \pm 0.58	86.05 (+22.08)	7.29 (−11.52)
EvoPrompt-S	88.69 \pm 0.16	9.93 \pm 0.22	87.95 \pm 0.13	2.38 \pm 0.14	84.98 \pm0.36	3.42 \pm 0.39	87.20 (+23.23)	5.24 (−13.57)
EvoPrompt	<u>88.97 \pm0.41</u>	10.12 \pm 0.35	87.97 \pm0.30	<u>2.60 \pm0.42</u>	<u>84.64 \pm0.14</u>	3.98 \pm 0.24	<u>87.19 (+23.22)</u>	5.57 (−13.24)
Upper-bound [†]	90.85 \pm 0.12	-	90.85 \pm 0.12	-	90.85 \pm 0.12	-	90.85	-

+1.15% Acc
-2.05% Forget

- Split ImageNet-R: 24,000 training images, 6,000 testing images, 200 classes.

Method	Split ImageNet-R							
	5 Steps		10 Steps		20 Steps		Avg	Avg
	Acc.(↑)	Forget.(↓)	Acc.(↑)	Forget.(↓)	Acc.(↑)	Forget.(↓)	Acc.(↑)	Forget.(↓)
FT-seq	61.41 \pm 0.38	5.76 \pm 0.48	50.28 \pm 2.29	24.28 \pm 1.73	39.25 \pm 0.90	40.38 \pm 0.77	50.31 (+0.00)	23.48 (−0.00)
LP-seq	59.83 \pm 0.33	1.50 \pm0.41	55.30 \pm 0.12	7.85 \pm 0.10	51.97 \pm 0.34	13.87 \pm 0.21	53.64 (+3.33)	7.74 (−15.74)
NME-seq	61.06	6.64	61.40	0.76	61.76	2.89	61.41 (+11.10)	3.43 (−20.05)
L2P	66.63 \pm 0.33	6.65 \pm 0.38	64.05 \pm 0.39	10.05 \pm 0.26	60.34 \pm 0.17	14.44 \pm 0.61	63.67 (+13.36)	10.38 (−13.10)
DualPrompt	71.06 \pm 0.35	4.19 \pm 0.25	69.71 \pm 0.25	5.44 \pm 0.12	66.26 \pm 0.46	8.74 \pm 0.33	69.01 (+18.70)	6.12 (−17.36)
ESN	73.42 \pm 0.40	<u>3.79 \pm0.55</u>	71.07 \pm 0.29	4.99 \pm 0.49	64.77 \pm 0.71	6.65 \pm 1.24	69.75 (+19.44)	5.14 (−18.34)
CODA-P-S	73.80 \pm 0.40	5.56 \pm 0.64	71.95 \pm 0.41	5.92 \pm 0.35	69.67 \pm 0.35	6.23 \pm 0.40	71.81 (+21.50)	5.90 (−17.58)
CODA-P	73.77 \pm 0.48	6.60 \pm 0.52	72.42 \pm 0.40	6.26 \pm 0.61	70.18 \pm 0.43	5.53 \pm 0.21	72.12 (+21.81)	6.13 (−17.35)
EvoPrompt-S	76.79 \pm 0.23	9.84 \pm 0.15	76.22 \pm 0.16	2.33 \pm 0.24	74.68 \pm0.51	2.70 \pm 0.19	75.90 (+25.59)	4.96 (−18.52)
EvoPrompt	77.16 \pm0.18	9.89 \pm 0.30	76.83 \pm0.08	2.78 \pm 0.06	74.41 \pm 0.23	2.56 \pm0.22	76.13 (+25.82)	5.08 (−18.40)
Upper-bound [†]	79.13 \pm 0.18	-	79.13 \pm 0.18	-	79.13 \pm 0.18	-	79.13	-

+4.01% Acc
-1.05% Forget

Empirical Results on Domain Incremental Learning

- Total:
 - fixed 50 classes
 - 11 domains
 - 120,000 images
- Training: 8 domains
- Testing: 3 unseen domains
- Metrics using final accuracy.



Method	Test Acc. (\uparrow)	Δ Acc. (\uparrow)
NME-seq	78.20	+00.00
EWC [†]	74.82 \pm 0.60	-3.38
LwF [†]	75.45 \pm 0.40	-2.75
L2P [†]	78.33 \pm 0.06	+0.13
S-iPrompts [‡]	83.13 \pm 0.51	+4.93
S-liPrompts [‡]	89.06 \pm 0.86	+10.86
ESN [‡]	91.80 \pm 0.31	+13.60
EvoPrompt-S	94.77 \pm 0.50	+16.57
EvoPrompt	95.27 \pm 0.15	+17.07
Upper-bound	91.32 \pm 0.23	-

+3.47% Acc

Empirical Results on Online Continual Incremental

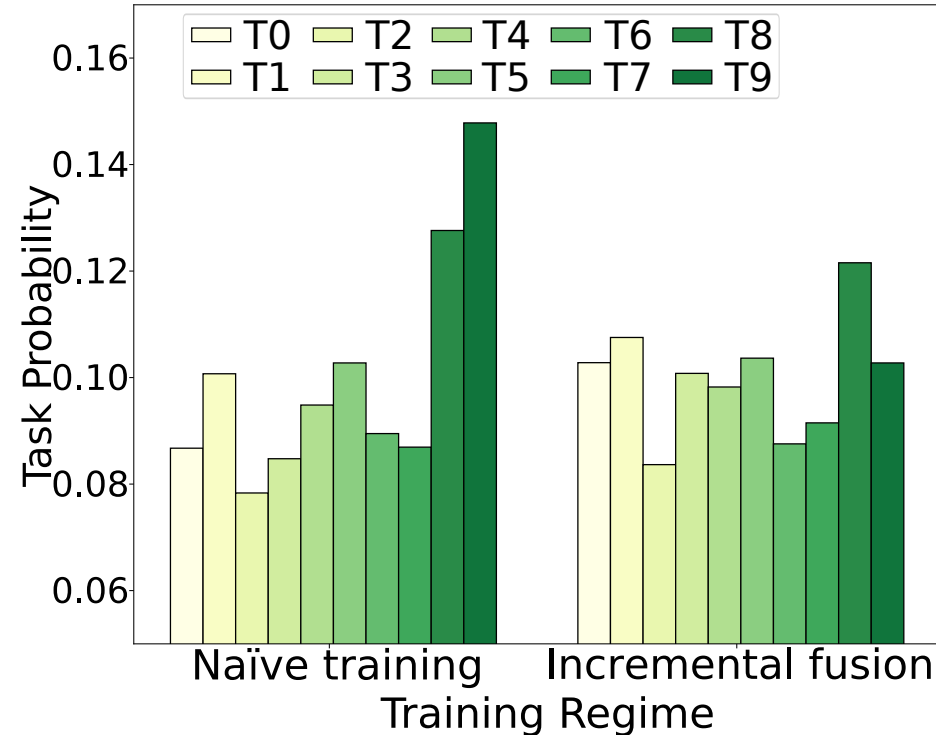
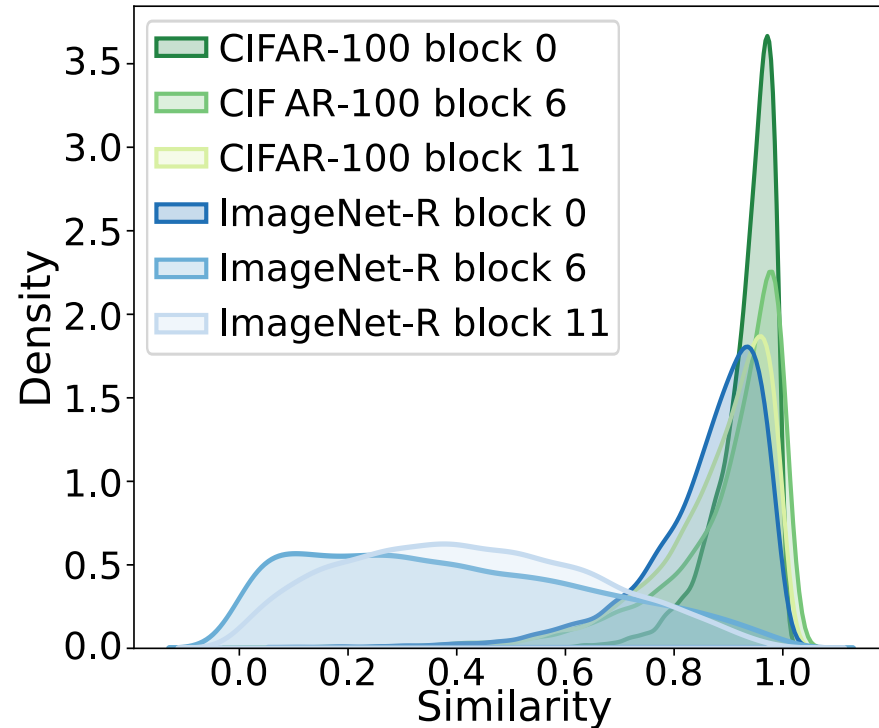
- The model encounters the samples in a single pass, technically with epoch set to 1.

Method	Split CIFAR-100		Split ImageNet-R	
	Acc.(↑)	Forget.(↓)	Acc.(↑)	Forget.(↓)
FT-seq	35.39±1.00	32.98±1.53	7.51±3.73	11.22±0.71
L2P	80.49±0.28	8.74±0.44	57.52±0.18	6.54±0.34
DualPrompt	82.17±0.34	7.52±0.21	61.09 ±0.18	4.40 ±0.62
ESN	74.17 ±1.14	10.59 ±1.39	-	-
CODA-P-S	79.46 ±0.06	11.92 ±1.27	64.60 ±0.99	6.09 ±1.18
CODA-P	81.07 ±0.38	10.10 ±0.84	66.47 ±0.33	5.42 ±0.87
EvoPrompt-S	84.23 ±0.57	1.64 ±0.29	73.56 ±0.42	3.82 ±0.24
EvoPrompt	84.72 ±0.94	0.89 ±0.72	74.05 ±0.48	3.66 ±0.36

+3.65% Acc +7.58% Acc
-9.21% Forget -1.76% Forget

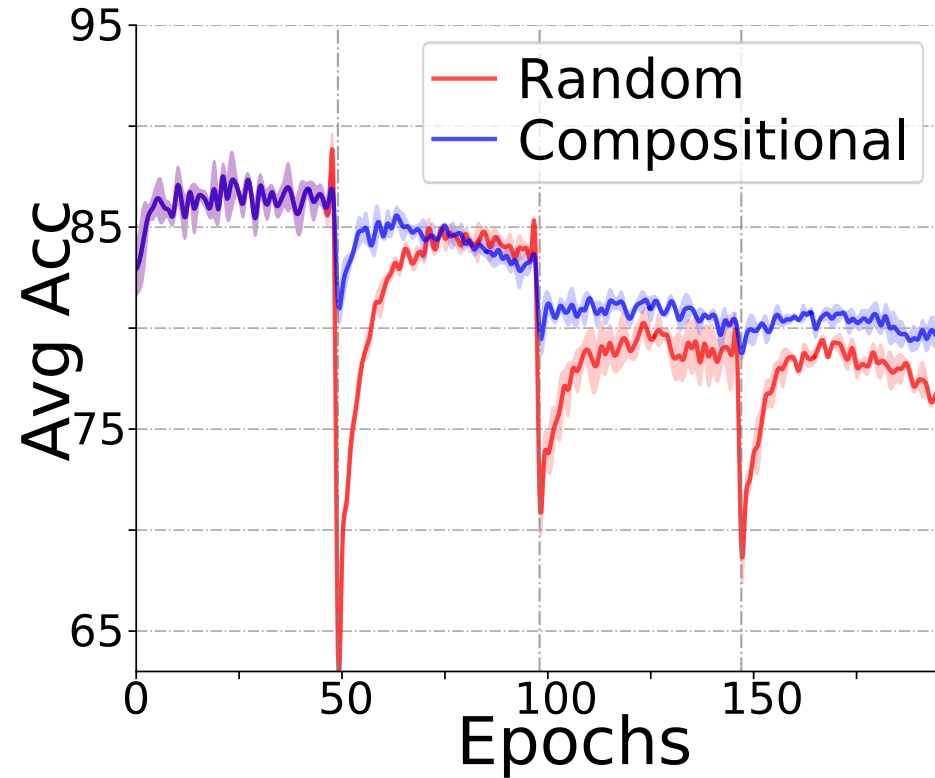
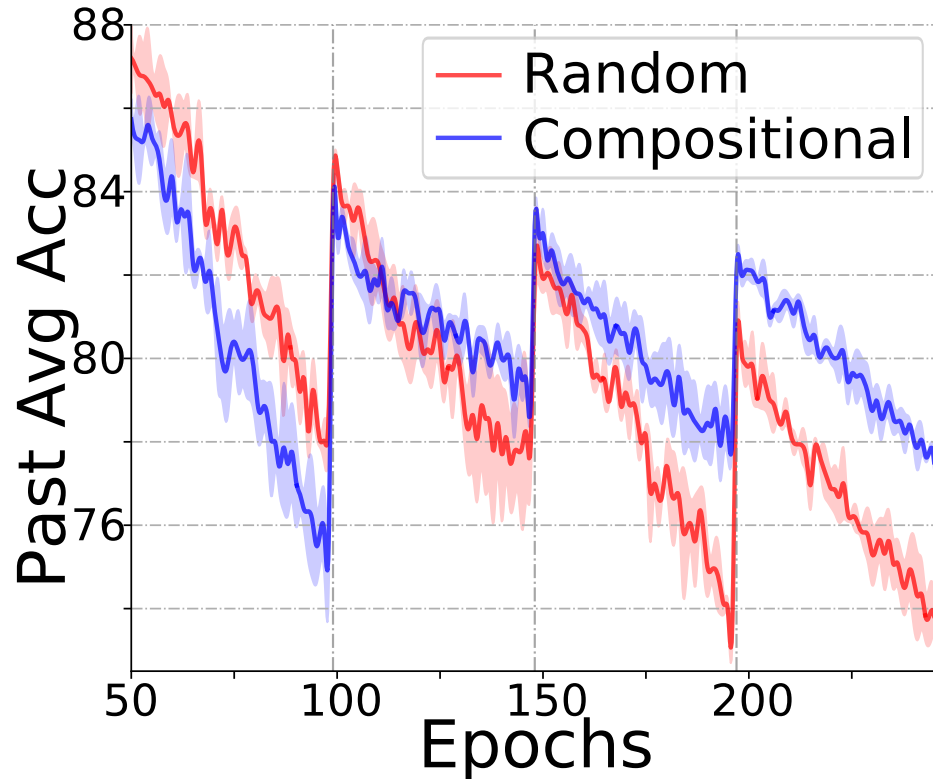
- Ours is better on **knowledge transfer** thus **improve consolidation**
- Key components for knowledge transfer: 1) WPM init from RPM 2) Compositional classifier initialization.

Further Analysis – Assignment Diversity and Recency Bias



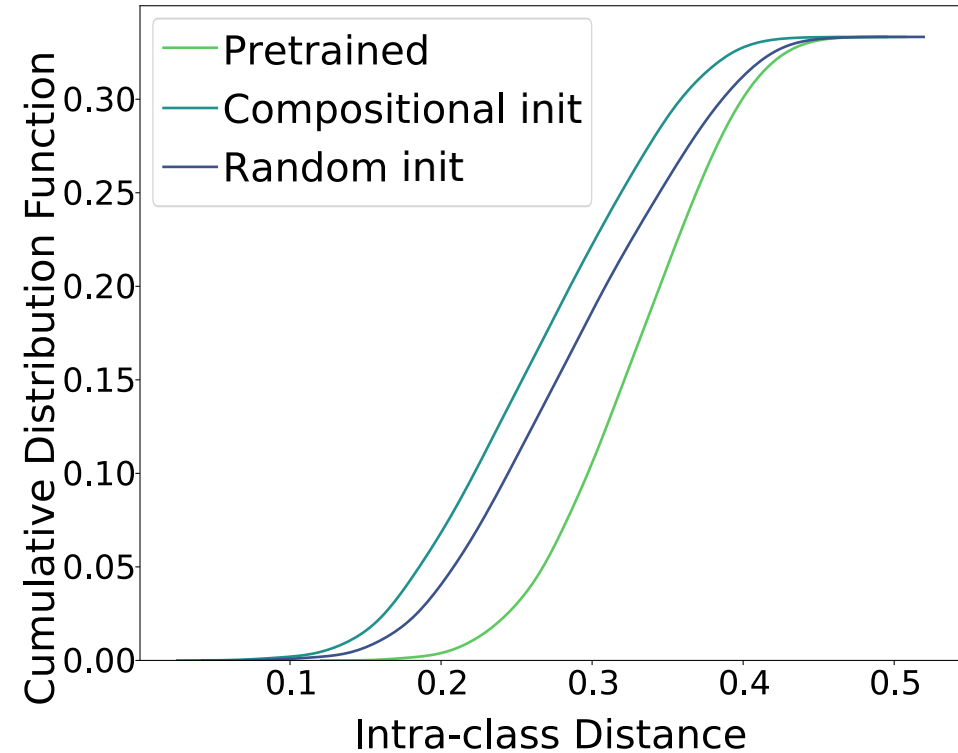
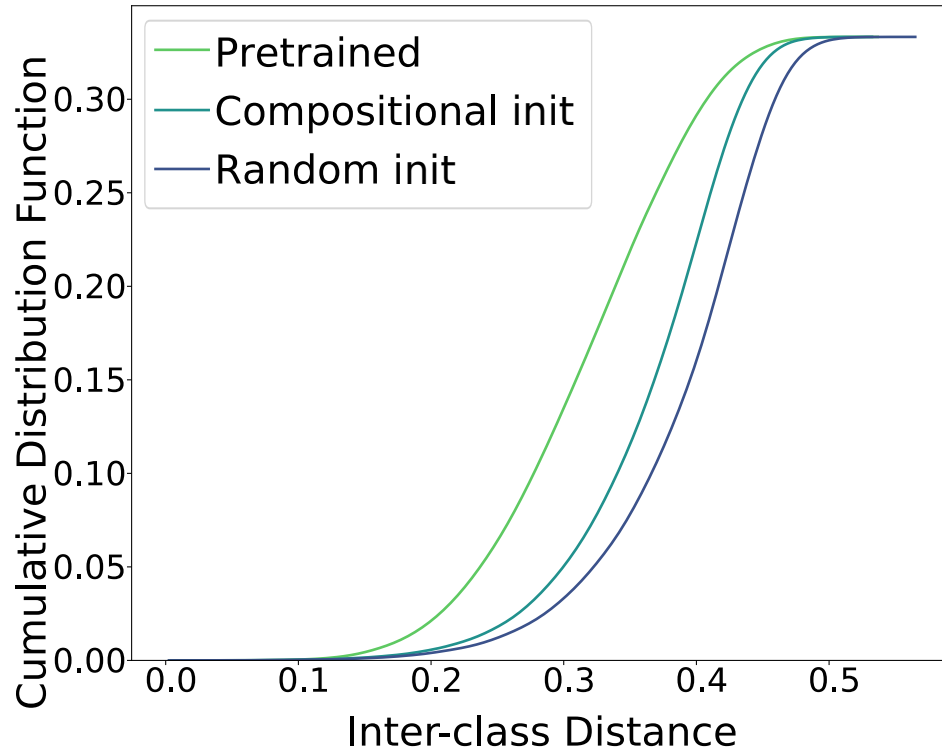
- Coefficient variability is **dataset-dependent**.
- Incremental fusion maintains **parameter balance**, ensuring **general solutions**, and **mitigating recency bias**.

Further Analysis – Stability Gap



- **No stability gap is observed** in either random or CCI.
- CCI exhibits stable performance, **smooth transitions** between tasks, and an **accelerated assimilation** of current knowledge

Further Analysis – Representation Compactness and Discriminativeness



- CCI producing more **compact intra-class structure**.
- CCI maintains **equilibrium** in inter-class margins, resulting in **smaller inter-class distances**, thereby improving **backward compatibility**.



The 38th Annual AAAI
Conference on Artificial
Intelligence

FEBRUARY 20-27, 2024 | VANCOUVER, CANADA

Thank You

github.com/MIV-XJTU/EvoPrompt



Code



Blog



Paper



西安交通大学

XI'AN JIAOTONG UNIVERSITY



中国科学院深圳先进技术研究院

SHENZHEN INSTITUTE OF ADVANCED TECHNOLOGY
CHINESE ACADEMY OF SCIENCES