



NEURAL INFORMATION  
PROCESSING SYSTEMS

# Visual Prompt Tuning in Null Space for Continual Learning

---

Yue Lu<sup>1</sup>, Shizhou Zhang<sup>1\*</sup>, De Cheng<sup>2\*</sup>, Yinghui Xing<sup>1</sup>,  
Nannan Wang<sup>2</sup>, Peng Wang<sup>1</sup>, Yanning Zhang<sup>1</sup>

<sup>1</sup> School of Computer Science, Northwestern Polytechnical University, China

<sup>2</sup> School of Telecommunications Engineering, Xidian University, China



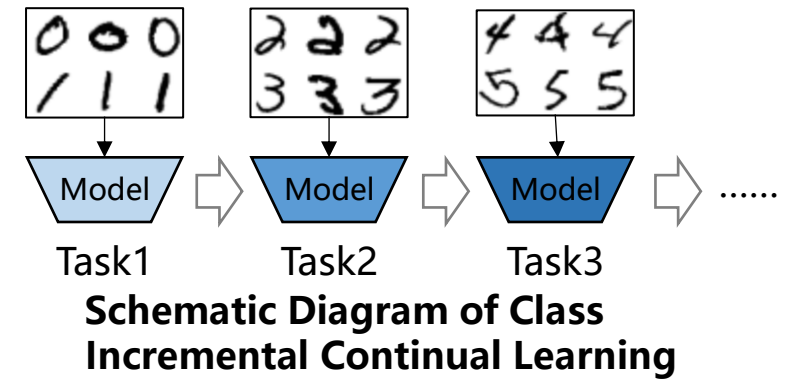
西北工業大學  
NORTHWESTERN POLYTECHNICAL UNIVERSITY



西安電子科技大學  
XIDIAN UNIVERSITY

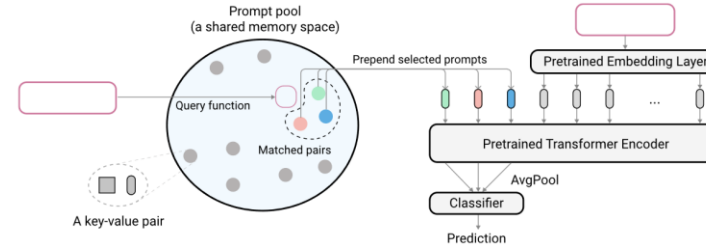
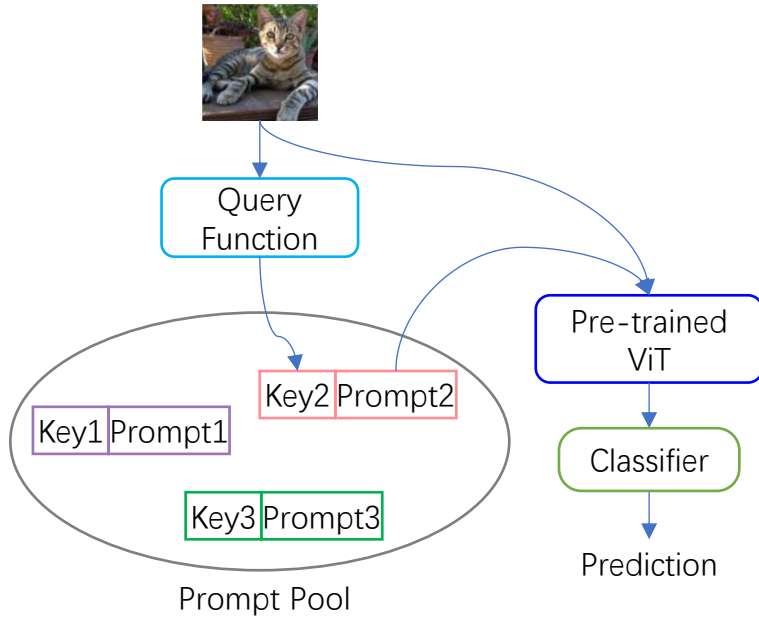
**Problem definition:** The model is trained on learning tasks that come one after another. **The data of the learned tasks are no longer visible.** In the inference stage, the model can have good discrimination ability for **all the data of the learned tasks.**

**Class-incremental continual learning:** Each task contains several categories that **have not been learned before**, and the model is trained task by task. In the inference stage, the model has good classification ability for all known category data when **the task to which the input image belongs is unknown.**

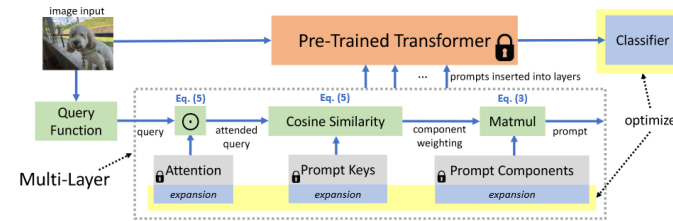


## Continual Learning Based on Visual Prompt Tuning

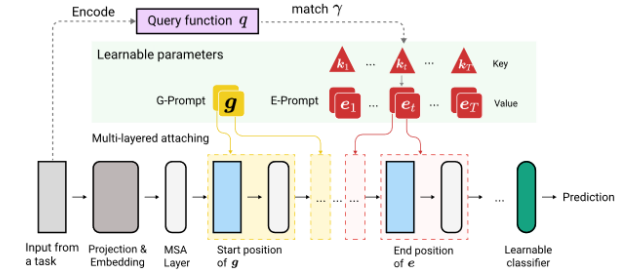
Encoding task knowledge into a prompt pool, selecting and updating task-relevant prompts within the pool based on the training data.



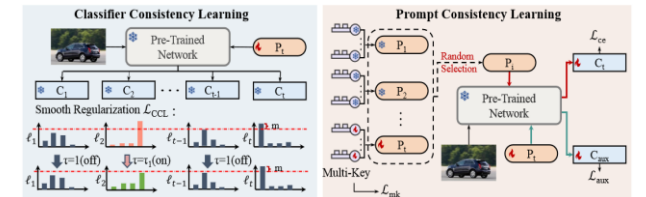
L2P(CVPR22)



CODA-Prompt(CVPR23)



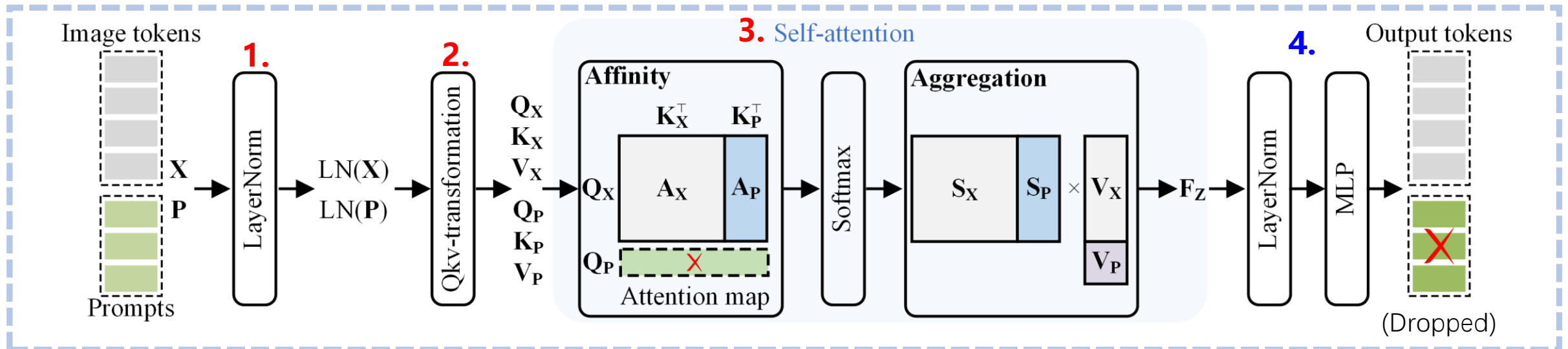
DualPrompt(ECCV22)



CPrompt(CVPR24)

**Training interference:** The features of the data from previous tasks change as related prompts are updated, leading to catastrophic forgetting.

## Analysis of the ViT Layer Forward Process $f_{\text{ViT}}(\mathbf{X}|\mathbf{P})$



**Input:**  $\mathbf{X} \in \mathbb{R}^{N \times D}$ ,  $\mathbf{P} \in \mathbb{R}^{M \times D}$ , where  $N, M$  represent the number of image tokens and prompts, respectively, and  $D$  denotes the token dimension.

$\mathbf{Z} = [\mathbf{X}; \mathbf{P}] \in \mathbb{R}^{(N+M) \times D}$ , representing the concatenated tokens along channel dimension.

**1. LayerNorm  $\text{LN}(\cdot)$ :**

$$\text{LN}(\mathbf{Z}) = \frac{\mathbf{Z} - \mu_{\mathbf{Z}}}{\sigma_{\mathbf{Z}}} \odot \alpha + \beta$$

**2. QKV- Transformation:**

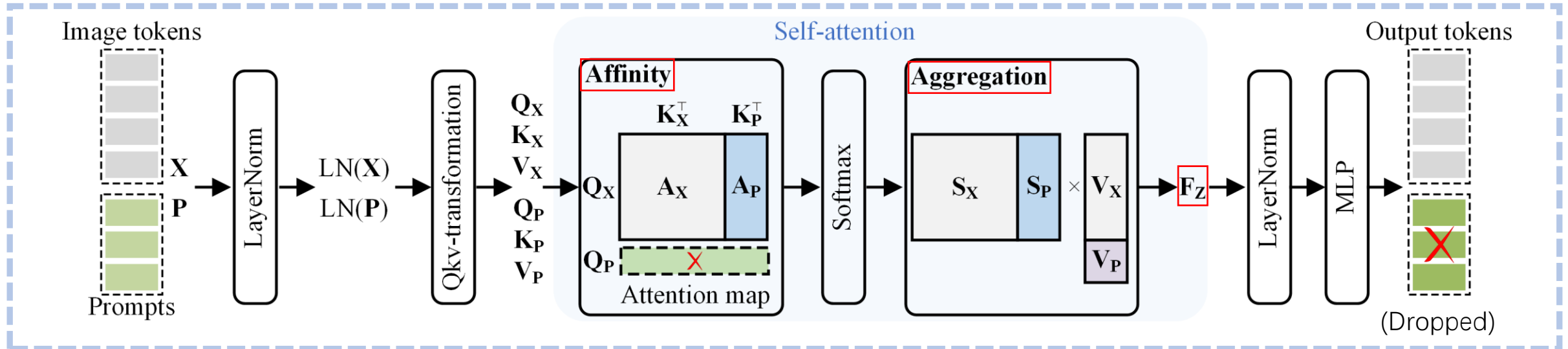
$$\begin{aligned} \mathbf{Q}_Z &= \text{LN}(\mathbf{Z})\mathbf{W}_q + \mathbf{b}_q \\ \mathbf{K}_Z &= \text{LN}(\mathbf{Z})\mathbf{W}_k + \mathbf{b}_k \\ \mathbf{V}_Z &= \text{LN}(\mathbf{Z})\mathbf{W}_v + \mathbf{b}_v \end{aligned}$$

**3. Self-attention  $f_{\text{SA}}$ :**

$$\mathbf{F}_Z = \text{softmax}\left(\frac{\mathbf{Q}_X \mathbf{K}_Z^\top}{\sqrt{D}}\right) \mathbf{V}_Z$$

**4. Feed-Forward Network:**  
LayerNorm and MLP

## Analyzing $f_{ViT}(X_t|P_t) = f_{ViT}(X_t|P_{t+1})$



To ensure that the output tokens of  $X_t$  remain consistent between tasks  $t$  and  $t + 1$ , i.e.,

$$f_{ViT}(X_t|P_t) = f_{ViT}(X_t|P_{t+1})$$

We define  $Z_t = [X_t; P_t], Z_{t+1} = [X_t; P_{t+1}]$ , This requires satisfying the condition:

$$F_{Z_t} = F_{Z_{t+1}}$$

Simplification of  $F_{Z_t} = F_{Z_{t+1}}$ :

$$\begin{cases} A_Z = f_{aff}(Q_X, K_Z) = \frac{Q_X K_Z^T}{\sqrt{D}} = \frac{Q_X [K_X^T \ K_P^T]}{\sqrt{D}} \\ S_Z = \text{softmax}(A_Z) = \text{softmax}(A_X \ A_P) = [S_X \ S_P] \\ F_Z = f_{agg}(S_Z, V_Z) = S_Z V_Z = [S_X \ S_P] \begin{bmatrix} V_X \\ V_P \end{bmatrix} \end{cases}$$

## Two Sufficient Conditions

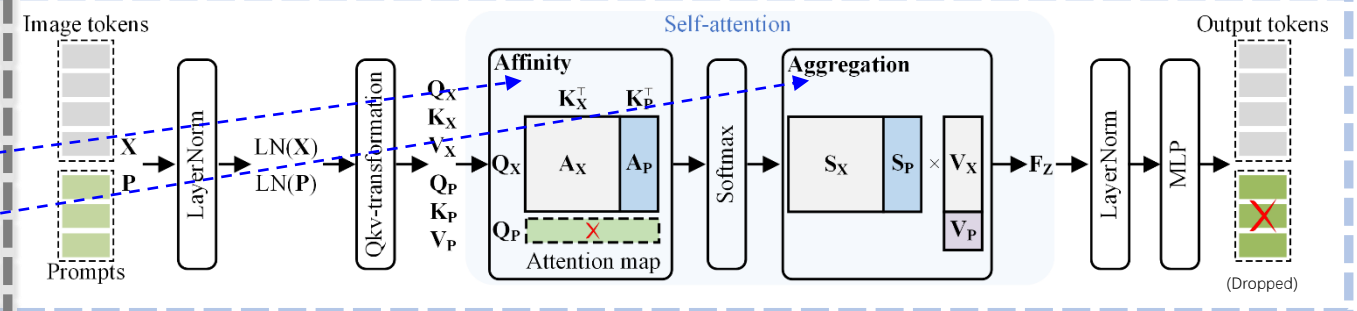
Two sufficient conditions to satisfy

$$\mathbf{F}_{Z_t} = \mathbf{F}_{Z_{t+1}}$$

$$\begin{cases} f_{\text{aff}}(\mathbf{Q}_{X_t}, \mathbf{K}_{Z_t}) = f_{\text{aff}}(\mathbf{Q}_{X_t}, \mathbf{K}_{Z_{t+1}}) \\ f_{\text{agg}}(\mathbf{S}_{Z_t}, \mathbf{V}_{Z_t}) = f_{\text{agg}}(\mathbf{S}_{Z_{t+1}}, \mathbf{V}_{Z_{t+1}}) \end{cases}$$

①

②



①  $f_{\text{aff}}(\mathbf{Q}_{X_t}, \mathbf{K}_{Z_t}) = f_{\text{aff}}(\mathbf{Q}_{X_t}, \mathbf{K}_{Z_{t+1}})$ , each side of the equation is:

$$\begin{cases} f_{\text{aff}}(\mathbf{Q}_{X_t}, \mathbf{K}_{Z_t}) = \mathbf{Q}_{X_t} [\mathbf{K}_{X_t}^T \ \mathbf{K}_{P_t}^T] = [\mathbf{Q}_{X_t} \mathbf{K}_{X_t}^T \ \mathbf{Q}_{X_t} [\text{LN}(\mathbf{P}_t) \mathbf{W}_k + \mathbf{b}_k]^T] \\ f_{\text{aff}}(\mathbf{Q}_{X_t}, \mathbf{K}_{Z_{t+1}}) = \mathbf{Q}_{X_t} [\mathbf{K}_{X_t}^T \ \mathbf{K}_{P_{t+1}}^T] = [\mathbf{Q}_{X_t} \mathbf{K}_{X_t}^T \ \mathbf{Q}_{X_t} [\text{LN}(\mathbf{P}_{t+1}) \mathbf{W}_k + \mathbf{b}_k]^T] \end{cases}$$

$$\mathbf{A}_Z = f_{\text{aff}}(\mathbf{Q}_X, \mathbf{K}_Z) = \frac{\mathbf{Q}_X [\mathbf{K}_X^T \ \mathbf{K}_P^T]}{\sqrt{D}}$$

From ①  $f_{\text{aff}}(\mathbf{Q}_{X_t}, \mathbf{K}_{Z_t}) = f_{\text{aff}}(\mathbf{Q}_{X_t}, \mathbf{K}_{Z_{t+1}})$ , we derive  $\mathbf{A}_{Z_t} = \mathbf{A}_{Z_{t+1}}$  and  $\mathbf{S}_{Z_t} = \mathbf{S}_{Z_{t+1}}$

②  $f_{\text{agg}}(\mathbf{S}_{Z_t}, \mathbf{V}_{Z_t}) = f_{\text{agg}}(\mathbf{S}_{Z_{t+1}}, \mathbf{V}_{Z_{t+1}})$ , each side of the equation is:

$$\begin{cases} f_{\text{agg}}(\mathbf{S}_{Z_t}, \mathbf{V}_{Z_t}) = \mathbf{S}_{X_t} \mathbf{V}_{X_t} + \mathbf{S}_{P_t} \mathbf{V}_{P_t} = \mathbf{S}_{X_t} \mathbf{V}_{X_t} + \mathbf{S}_{P_t} [\text{LN}(\mathbf{P}_t) \mathbf{W}_v + \mathbf{b}_v] \\ f_{\text{agg}}(\mathbf{S}_{Z_{t+1}}, \mathbf{V}_{Z_{t+1}}) = f_{\text{agg}}(\mathbf{S}_{Z_t}, \mathbf{V}_{Z_{t+1}}) = \mathbf{S}_{X_t} \mathbf{V}_{X_t} + \mathbf{S}_{P_t} [\text{LN}(\mathbf{P}_{t+1}) \mathbf{W}_v + \mathbf{b}_v] \end{cases}$$

$$\begin{aligned} \mathbf{S}_Z &= \text{softmax}(\mathbf{A}_Z) = [\mathbf{S}_X \ \mathbf{S}_P] \\ \mathbf{F}_Z &= f_{\text{agg}}(\mathbf{S}_Z, \mathbf{V}_Z) = [\mathbf{S}_X \ \mathbf{S}_P] \begin{bmatrix} \mathbf{V}_X \\ \mathbf{V}_P \end{bmatrix} \end{aligned}$$

Simplifying conditions ①②, we obtain: 
$$\begin{cases} \mathbf{Q}_{X_t} \mathbf{W}_k^T \text{LN}(\mathbf{P}_t)^T = \mathbf{Q}_{X_t} \mathbf{W}_k^T \text{LN}(\mathbf{P}_{t+1})^T = \mathbf{Q}_{X_t} \mathbf{W}_k^T \text{LN}(\mathbf{P}_t + \Delta \mathbf{P})^T \\ \mathbf{S}_{P_t} \text{LN}(\mathbf{P}_t) \mathbf{W}_v = \mathbf{S}_{P_t} \text{LN}(\mathbf{P}_{t+1}) \mathbf{W}_v = \mathbf{S}_{P_t} \text{LN}(\mathbf{P}_t + \Delta \mathbf{P}) \mathbf{W}_v \end{cases}$$

## Simplification of the LayerNorm Term

$$\begin{cases} \mathbf{Q}_{X_t} \mathbf{W}_k^\top \text{LN}(\mathbf{P}_t)^\top = \mathbf{Q}_{X_t} \mathbf{W}_k^\top \text{LN}(\mathbf{P}_t + \Delta\mathbf{P})^\top & \textcircled{3} \\ \mathbf{S}_{P_t} \text{LN}(\mathbf{P}_t) \mathbf{W}_v = \mathbf{S}_{P_t} \text{LN}(\mathbf{P}_t + \Delta\mathbf{P}) \mathbf{W}_v & \textcircled{4} \end{cases}$$

$$\text{LN}(\mathbf{P}) = \frac{\mathbf{P} - \mu_{\mathbf{P}}}{\sigma_{\mathbf{P}}} \odot \alpha + \beta \quad \text{LN}(\mathbf{P}) = \frac{\mathbf{P}_t + \Delta\mathbf{P} - \mu_{\mathbf{P}_t + \Delta\mathbf{P}}}{\sigma_{\mathbf{P}_t + \Delta\mathbf{P}}} \odot \alpha + \beta$$

For the term  $\text{LN}(\mathbf{P}_t + \Delta\mathbf{P})$ , it cannot be expressed directly in terms of  $\text{LN}(\mathbf{P}_t)$  and  $\Delta\mathbf{P}$ .

### Assumption of Distribution Invariance:

Assuming the distribution of prompts remains unchanged during training:  $\begin{cases} \mu_{\mathbf{P}_t + \Delta\mathbf{P}} = \mu_{\mathbf{P}_t} \\ \sigma_{\mathbf{P}_t + \Delta\mathbf{P}} = \sigma_{\mathbf{P}_t} \end{cases}$

Relationship between  $\text{LN}(\mathbf{P}_t + \Delta\mathbf{P})$  and  $\text{LN}(\mathbf{P}_t)$  using  $\Delta\mathbf{P}$ :

$$\text{LN}(\mathbf{P}_t + \Delta\mathbf{P}) = \frac{\mathbf{P}_t + \Delta\mathbf{P} - \mu_{\mathbf{P}_t + \Delta\mathbf{P}}}{\sigma_{\mathbf{P}_t + \Delta\mathbf{P}}} \odot \alpha + \beta = \frac{\mathbf{P}_t + \Delta\mathbf{P} - \mu_{\mathbf{P}_t}}{\sigma_{\mathbf{P}_t}} \odot \alpha + \beta = \left( \frac{\mathbf{P}_t - \mu_{\mathbf{P}_t}}{\sigma_{\mathbf{P}_t}} \odot \alpha + \beta \right) + \frac{\Delta\mathbf{P}}{\sigma_{\mathbf{P}_t}} \odot \alpha = \text{LN}(\mathbf{P}_t) + \frac{\Delta\mathbf{P}}{\sigma_{\mathbf{P}_t}} \odot \alpha$$

From  $\text{LN}(\mathbf{P}_t + \Delta\mathbf{P}) = \text{LN}(\mathbf{P}_t) + \frac{\Delta\mathbf{P}}{\sigma_{\mathbf{P}_t}} \odot \alpha$ , simplify  $\textcircled{3}\textcircled{4}$ :

$$\begin{cases} \mathbf{Q}_{X_t} \mathbf{W}_k^\top \text{LN}(\mathbf{P}_t)^\top = \mathbf{Q}_{X_t} \mathbf{W}_k^\top \text{LN}(\mathbf{P}_t)^\top + \frac{\mathbf{Q}_{X_t} \mathbf{W}_k^\top \Delta\mathbf{P}^\top}{\sigma_{\mathbf{P}_t}} \odot \alpha^\top & \textcircled{5} \\ \mathbf{S}_{P_t} \text{LN}(\mathbf{P}_t) \mathbf{W}_v = \mathbf{S}_{P_t} \text{LN}(\mathbf{P}_t) \mathbf{W}_v + \frac{\mathbf{S}_{P_t} \Delta\mathbf{P} \mathbf{W}_v}{\sigma_{\mathbf{P}_t}} \odot \alpha & \textcircled{6} \end{cases}$$

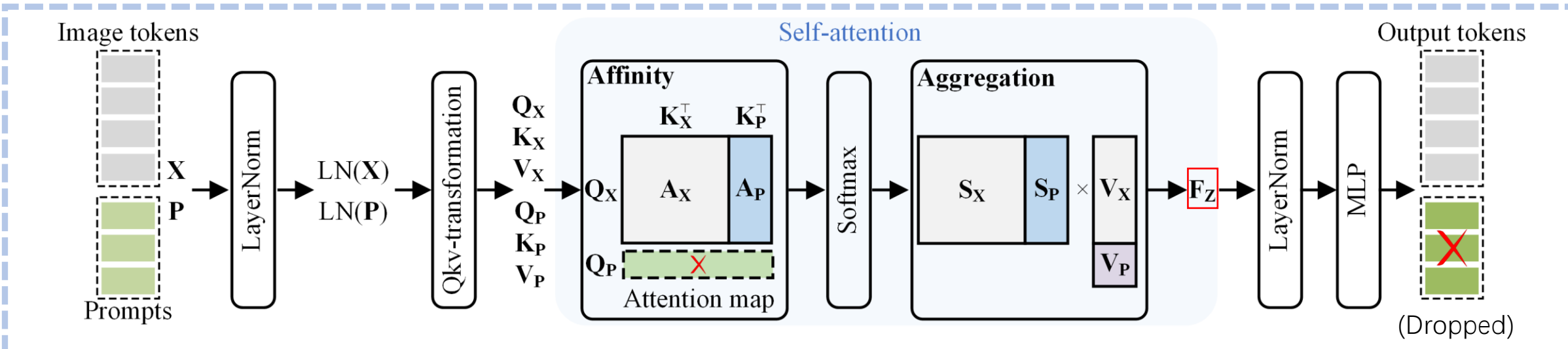
From  $\textcircled{5}\textcircled{6}$ :  $\begin{cases} (\mathbf{Q}_{X_t} \mathbf{W}_k^\top) \Delta\mathbf{P}^\top = 0 \\ \mathbf{S}_{P_t} \Delta\mathbf{P} \mathbf{W}_v = 0 \end{cases}$

It is sufficient to satisfy conditions  $\textcircled{7}\textcircled{8}$ :

$$\begin{cases} (\mathbf{Q}_{X_t} \mathbf{W}_k^\top) \Delta\mathbf{P}^\top = 0 & \textcircled{7} \\ \mathbf{S}_{P_t} \Delta\mathbf{P} = 0 & \textcircled{8} \end{cases}$$



## Conclusion



To achieve the following objective:

$$f_{ViT}(\mathbf{X}_t | \mathbf{P}_t) = f_{ViT}(\mathbf{X}_t | \mathbf{P}_{t+1})$$

Specifically, this requires:

$$\mathbf{F}_{Z_t} = \mathbf{F}_{Z_{t+1}}$$

Converting into two sufficient conditions:

$$\begin{cases} f_{\text{aff}}(\mathbf{Q}_{X_t}, \mathbf{K}_{Z_t}) = f_{\text{aff}}(\mathbf{Q}_{X_t}, \mathbf{K}_{Z_{t+1}}) \\ f_{\text{agg}}(\mathbf{S}_{Z_t}, \mathbf{V}_{Z_t}) = f_{\text{agg}}(\mathbf{S}_{Z_{t+1}}, \mathbf{V}_{Z_{t+1}}) \end{cases}$$

Introducing a constraint on the prompt distribution's variation

$$\begin{cases} \mu_{\mathbf{P}_{t+\Delta\mathbf{P}}} = \mu_{\mathbf{P}_t} \\ \sigma_{\mathbf{P}_{t+\Delta\mathbf{P}}} = \sigma_{\mathbf{P}_t} \end{cases}$$

This ultimately leads to the following two conditions called **consistency conditions**.

$$\begin{cases} (\mathbf{Q}_{X_t} \mathbf{W}_k^T) \Delta \mathbf{P}^T = \mathbf{0} \quad \textcircled{7} \\ \mathbf{S}_{\mathbf{P}_t} \Delta \mathbf{P} = \mathbf{0} \quad \textcircled{8} \end{cases}$$



## Optimization of Consistency Conditions

**Two-step optimization scheme:**

- 1) Use the projection matrix  $\mathbf{B}_1$  ( $\Delta \mathbf{P}^\top = \mathbf{B}_1 \mathbf{P}_G^\top$ ) to make  $\Delta \mathbf{P}^\top$  orthogonal to the subspace spanned by  $\mathbf{Q}_{\mathbf{X}_t} \mathbf{W}_k^\top$ ;
- 2) Use the projection matrix  $\mathbf{B}_2$  ( $\Delta \mathbf{P} = \mathbf{B}_2 \mathbf{P}_G$ ) to make  $\Delta \mathbf{P}$  orthogonal to the subspace spanned by  $\mathbf{S}_{\mathbf{P}_t}$ .

**Computing  $\mathbf{B}_1$  and  $\mathbf{B}_2$  via null space:**

For  $\mathbf{B}_1$ : SVD  $\left( (\mathbf{Q}_{\mathbf{X}_t} \mathbf{W}_k^\top)^\top \mathbf{Q}_{\mathbf{X}_t} \mathbf{W}_k^\top \right)$ , let the matrix of right singular vectors corresponding to the singular values that are (or close to) zero be  $\mathbf{U}_{1,0}$ , then  $\mathbf{B}_1 = \mathbf{U}_{1,0} \mathbf{U}_{1,0}^\top$

For  $\mathbf{B}_2$ : SVD  $(\mathbf{S}_{\mathbf{P}_t}^\top \mathbf{S}_{\mathbf{P}_t})$ , let the matrix of right singular vectors corresponding to the singular values that are (or close to) zero be  $\mathbf{U}_{2,0}$ , then  $\mathbf{B}_2 = \mathbf{U}_{2,0} \mathbf{U}_{2,0}^\top$

Finally, the update rule from gradient  $\mathbf{P}_G$  to update  $\Delta \mathbf{P}$  is given by:

$$\Delta \mathbf{P} = \mathbf{B}_2 \mathbf{P}_G \mathbf{B}_1 = (\mathbf{U}_{2,0} \mathbf{U}_{2,0}^\top) \mathbf{P}_G (\mathbf{U}_{1,0} \mathbf{U}_{1,0}^\top) \quad \textcircled{9}$$

**The constraint on the distributional variation of the prompt**

$$\begin{cases} \mu_{\mathbf{P}_{t+\Delta \mathbf{P}}} = \mu_{\mathbf{P}_t} \\ \sigma_{\mathbf{P}_{t+\Delta \mathbf{P}}} = \sigma_{\mathbf{P}_t} \end{cases}$$

is achieved by introducing a loss function  $\mathcal{L}_{\text{LN}}$ :

$$\mathcal{L}_{\text{LN}} = \|\mu_{\mathbf{P}_{t+1}} - \mu_{\mathbf{P}_t}\|_1 + \|\sigma_{\mathbf{P}_{t+1}} - \sigma_{\mathbf{P}_t}\|_1 \quad \textcircled{10}$$

- 1)  $\textcircled{9}$ : Applying null space projection to the gradient
- 2)  $\textcircled{10}$ : Constraining the distribution of the prompt

## Extension to Multi-Head Attention

The Transformer architecture generally employs multi-head attention, so the derived consistency conditions ⑦⑧ need to be extended to multi-head attention.

$$\begin{cases} (\mathbf{Q}_{X_t} \mathbf{W}_k^\top) \Delta \mathbf{P}^\top = \mathbf{0} & \text{⑦} \\ \mathbf{S}_{P_t} \Delta \mathbf{P} = \mathbf{0} & \text{⑧} \end{cases}$$

Assuming there are  $H$  attention heads, let the  $h$ -th attention head for  $h \in [1, 2, \dots, H]$  have its respective components  $\mathbf{Q}_{X_t.h}, \mathbf{W}_{k.h}, \mathbf{S}_{P_t.h}$ . To **satisfy the consistency conditions across all attention heads**, we define:

$$\forall h \in [1, 2, \dots, H], \begin{cases} (\mathbf{Q}_{X_t.h} \mathbf{W}_{k.h}^\top) \Delta \mathbf{P}^\top = \mathbf{0} \\ \mathbf{S}_{P_t.h} \Delta \mathbf{P} = \mathbf{0} \end{cases}$$

$\mathbf{\Omega}_{1,t} = [\mathbf{Q}_{X_t.1} \mathbf{W}_{k.1}^\top; \dots; \mathbf{Q}_{X_t.H} \mathbf{W}_{k.H}^\top]$ , representing the concatenated matrices  $\mathbf{Q}_{X_t.h} \mathbf{W}_{k.h}^\top$  from each attention head.

$\mathbf{\Omega}_{2,t} = [\mathbf{S}_{P_t.1}; \dots; \mathbf{S}_{P_t.H}]$  representing the concatenated  $\mathbf{S}_{P_t.h}$  matrices from each attention head.

Using block matrix operations, the consistency conditions for all attention heads can be expressed as follows:

$$\begin{cases} \mathbf{\Omega}_{1,t} \Delta \mathbf{P}^\top = \mathbf{0} \\ \mathbf{\Omega}_{2,t} \Delta \mathbf{P} = \mathbf{0} \end{cases}$$

The prompt distribution constraint loss function  $\mathcal{L}_{LN} = \|\mu_{P_{t+1}} - \mu_{P_t}\|_1 + \|\sigma_{P_{t+1}} - \sigma_{P_t}\|_1$  remains independent of the number of attention heads, so no further extension is required.

## Key Settings

### **Model used:**

- a) ViT-Base16 model pre-trained on ImageNet-21K, with prompts of length 4 inserted into all ViT layers (referred to as VPT);**
- b) CLIP model, with prompts added to the image encoder, while the text encoder is frozen.**

### **Evaluation benchmarks:**

**Four class-incremental test benchmarks, including 10-task CIFAR-100, 20-task CIFAR-100, 10-task ImageNet-R, and 10-task DomainNet-200.**

### **Evaluation metrics:**

**Accuracy (Acc., the higher, the better) and Forgetting (the lower, the better), with the average results of three runs reported. Accuracy is the primary focus metric.**

## Comparing with Baseline

"-Seq" : baseline

"-NSP<sup>2</sup>" : our method

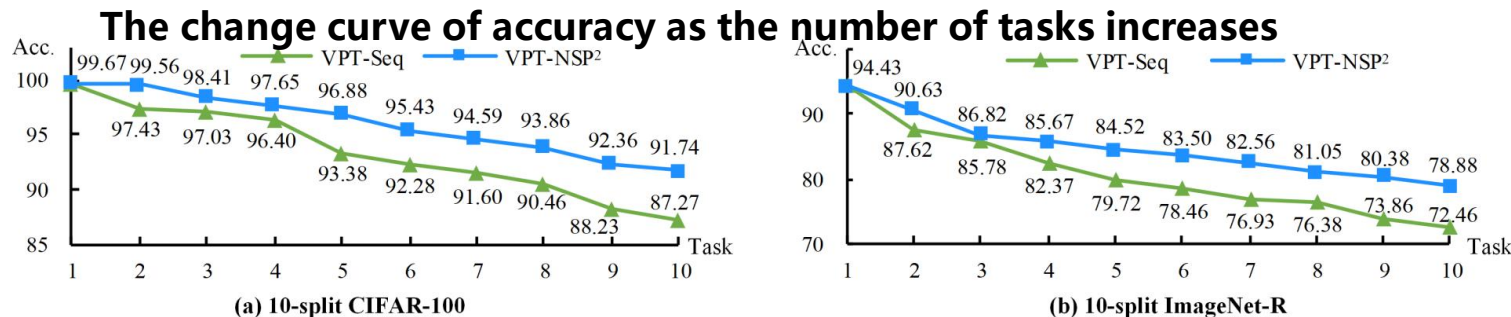
"Upper-bound": Jointly training all tasks at once, can be considered the upper bound of performance in continual learning.

Comparison with the baseline using VPT and CLIP models

Method	10S-CIFAR-100		20S-CIFAR-100		10S-ImageNet-R		10S-DomainNet	
	Acc. ↑	Forgetting ↓	Acc. ↑	Forgetting ↓	Acc. ↑	Forgetting ↓	Acc. ↑	Forgetting ↓
VPT-Seq	87.27	12.33	82.36	17.36	72.46	19.41	73.28	25.65
VPT-NSP <sup>2</sup>	<b>91.74</b>	<b>3.28</b>	<b>89.89</b>	<b>4.91</b>	<b>78.88</b>	<b>5.06</b>	<b>83.54</b>	<b>8.54</b>
Upper-bound	93.87	-	93.87	-	84.60	-	89.25	-
CLIP-Seq	72.91	15.13	71.37	17.89	75.69	19.21	67.73	35.60
CLIP-NSP <sup>2</sup>	<b>80.96</b>	<b>12.45</b>	<b>79.83</b>	<b>13.77</b>	<b>82.17</b>	<b>6.42</b>	<b>77.04</b>	<b>18.33</b>
Upper-bound	84.52	-	84.52	-	84.86	-	81.65	-

Accuracy increased by **4%~10%**  
Forgetting decreased by **9%~17%**.

Accuracy increased by **6%~9%**  
Forgetting decreased by **3%~17%**.



Our method consistently outperforms the baseline, with its advantages becoming more pronounced over time.

## Comparing with Existing Methods

Comparison with existing methods based on ImageNet-21K pre-trained VPT, numbers following  $\pm$  representing the standard deviation

Method	Venue	10S-CIFAR-100		20S-CIFAR-100		10S-ImageNet-R		10S-DomainNet	
		Acc.	Forgetting	Acc.	Forgetting	Acc.	Forgetting	Acc.	Forgetting
L2P [40]	CVPR'22	83.83 $\pm$ 0.04	7.63 $\pm$ 0.30	80.10 $\pm$ 0.72 $\ddagger$	-	61.57 $\pm$ 0.66	9.73 $\pm$ 0.47	81.17 $\pm$ 0.83 $\dagger$	8.98 $\pm$ 1.25
DualPrompt [39]	ECCV'22	86.51 $\pm$ 0.33	5.16 $\pm$ 0.09	82.02 $\pm$ 0.32 $\ddagger$	-	68.13 $\pm$ 0.49	4.68 $\pm$ 0.20	81.70 $\pm$ 0.78 $\dagger$	8.04 $\pm$ 0.31
CODA-P [32]	CVPR'23	86.25 $\pm$ 0.74	1.67 $\pm$ 0.26	-	-	75.45 $\pm$ 0.56	1.64 $\pm$ 0.10	80.04 $\pm$ 0.79 $\dagger$	10.16 $\pm$ 0.35
ESN [38]	AAAI'23	86.34 $\pm$ 0.52	4.76 $\pm$ 0.14	80.56 $\pm$ 0.94 $\ddagger$	-	62.61 $\pm$ 0.96 $\ddagger$	-	79.22 $\pm$ 2.04 $\dagger$	10.62 $\pm$ 2.12
APG [33]	ICCV'23	89.35	6.01	88.64	6.51	73.27	8.59	-	-
LAE [10]	ICCV'23	85.59 $\pm$ 0.46	-	83.93 $\pm$ 0.28	-	72.66 $\pm$ 0.63	-	-	-
DualP-LGCL [15]	ICCV'23	87.23 $\pm$ 0.21	5.10 $\pm$ 0.15	-	-	69.46 $\pm$ 0.04	4.20 $\pm$ 0.06	-	-
C-LN [23]	ICCVW'23	86.95 $\pm$ 0.37	6.98 $\pm$ 0.43	-	-	76.36 $\pm$ 0.51	8.31 $\pm$ 1.28	-	-
EvoPrompt [18]	AAAI'24	87.97 $\pm$ 0.30	2.60 $\pm$ 0.42	84.64 $\pm$ 0.14	3.98 $\pm$ 0.24	76.83 $\pm$ 0.08	2.78 $\pm$ 0.06	-	-
OVOR-Deep [12]	ICLR'24	85.99 $\pm$ 0.89	6.42 $\pm$ 2.03	-	-	76.11 $\pm$ 0.21	7.16 $\pm$ 0.34	-	-
DualP-PGP [26]	ICLR'24	86.92 $\pm$ 0.05	5.35 $\pm$ 0.19	83.74 $\pm$ 0.01	7.91 $\pm$ 0.15	69.34 $\pm$ 0.05	4.53 $\pm$ 0.04	-	-
InfLoRA [20]	CVPR'24	87.06 $\pm$ 0.25	-	-	-	75.65 $\pm$ 0.14	-	-	-
EASE [45]	CVPR'24	87.76	-	85.80	-	76.17	-	-	-
CPrompt [11]	CVPR'24	87.82 $\pm$ 0.21	5.06 $\pm$ 0.50	-	-	77.14 $\pm$ 0.11	5.97 $\pm$ 0.68	82.97 $\pm$ 0.34	7.45 $\pm$ 0.93
VPT-NSP <sup>2</sup>	This work	<b>91.74</b> $\pm$ 0.63	3.28 $\pm$ 0.45	<b>89.89</b> $\pm$ 0.72	4.91 $\pm$ 0.59	<b>78.88</b> $\pm$ 0.50	5.06 $\pm$ 0.26	<b>83.54</b> $\pm$ 0.77	8.54 $\pm$ 0.48

**Our method Achieves state-of-the-art performance on all evaluation benchmarks.**



## Ablation Study

Ablation study on  $B_1$ ,  $B_2$  and  $\mathcal{L}_{LN}$

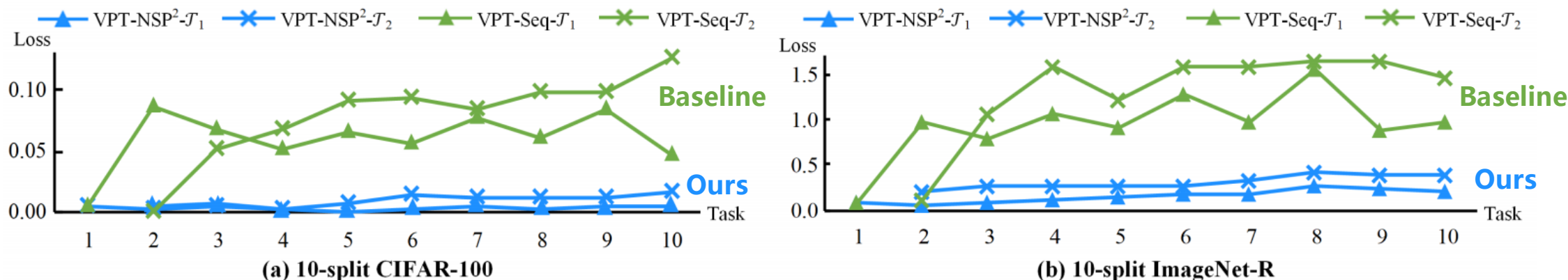
$B_1$	$B_2$	$\mathcal{L}_{LN}$	10S-CIFAR-100		20S-CIFAR-100		10S-ImageNet-R		10S-DomainNet	
			Acc. $\uparrow$	Forgetting $\downarrow$	Acc. $\uparrow$	Forgetting $\downarrow$	Acc. $\uparrow$	Forgetting $\downarrow$	Acc. $\uparrow$	Forgetting $\downarrow$
			87.27	12.33	82.36	17.36	72.46	19.41	73.28	25.65
✓			90.58	6.91	88.13	10.27	78.05	8.14	82.31	10.89
	✓		88.74	10.85	83.32	16.48	74.71	14.69	78.87	17.81
✓	✓		91.33	4.22	88.96	6.42	78.37	6.25	83.17	8.95
✓		✓	91.42	3.94	88.46	8.64	78.30	6.31	83.13	9.32
	✓	✓	89.36	9.32	86.67	11.59	75.27	13.35	79.45	16.50
✓	✓	✓	<b>91.74</b>	<b>3.28</b>	<b>89.89</b>	<b>4.91</b>	<b>78.88</b>	<b>5.06</b>	<b>83.54</b>	<b>8.54</b>

The first consistency condition, the orthogonal subspace projection matrix  $B_1$ , has the most significant impact on performance. However,  $B_2$  and  $\mathcal{L}_{LN}$  are also indispensable. The highest accuracy and lowest forgetting rate are achieved only when all three components are used together.

## Analysis of Reducing Training Interventions

As the number of learning tasks increases, if the issue of training interventions becomes more severe, the model's loss on the old task data will be higher. Conversely, if training interventions can be minimized or eliminated, the model's loss on the old task data will not increase.

The loss curves on the training data of the 1st and 2nd tasks as the number of learning tasks increases, comparing the **proposed method** (VPT-NSP2) with the **baseline method** (VPT-Seq).



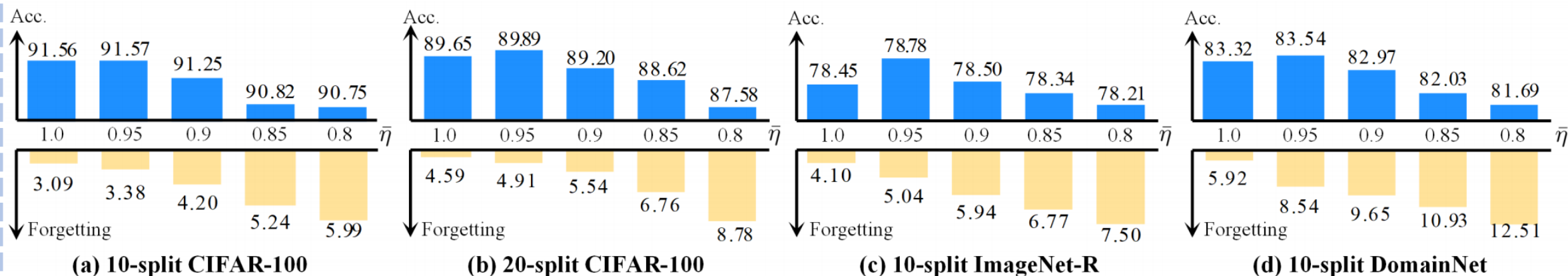
The proposed method maintains **almost no change in the loss** on old task data, demonstrating its ability to eliminate training interventions, thereby preventing forgetting.



## Stability-Plasticity Trade-off

The stability-plasticity trade-off in the proposed method is represented by  $\Delta P = [\eta_2 B_2 + (1 - \eta_2)I]P_G[\eta_1 B_1 + (1 - \eta_1)I]$ , where the orthogonal subspace projection matrices  $B_1, B_2$  are weighted and fused with the identity matrix  $I$ . The parameters  $\eta_1, \eta_2 \in [0, 1]$  control the weights, allowing for the adjustment of the trade-off between plasticity and stability.

When the weights  $\eta_1, \eta_2$  are set to the same value, denoted as  $\bar{\eta}$ , the changes in **accuracy** and **forgetting rate** as  $\bar{\eta}$  decreases from 1.0 to 0.8 are shown in the following figure.

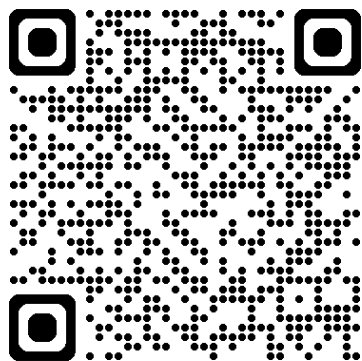


Accuracy is influenced by both stability and plasticity, and the best performance is achieved when a good balance between the two is attained. An increase in plasticity means a weakened ability to retain knowledge from old tasks, which results in a gradual increase in the forgetting rate.

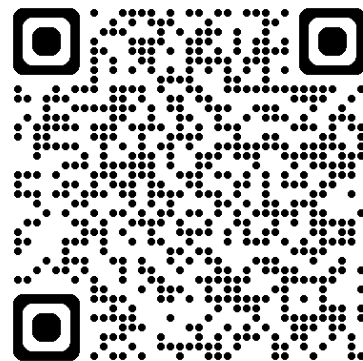


NEURAL INFORMATION  
PROCESSING SYSTEMS

# Thank You for Watching!



**Paper (ArXiv)**



**Code (Github)**



**Author: Yue Lu (WeChat)**



西北工業大學  
NORTHWESTERN POLYTECHNICAL UNIVERSITY



西安電子科技大學  
XIDIAN UNIVERSITY