

# Mixture of experts meets Prompt-based Continual Learning

Minh Le<sup>3</sup>, An Nguyen<sup>2\*</sup>, Huy Nguyen<sup>1\*</sup>, Trang Nguyen<sup>3\*</sup>,

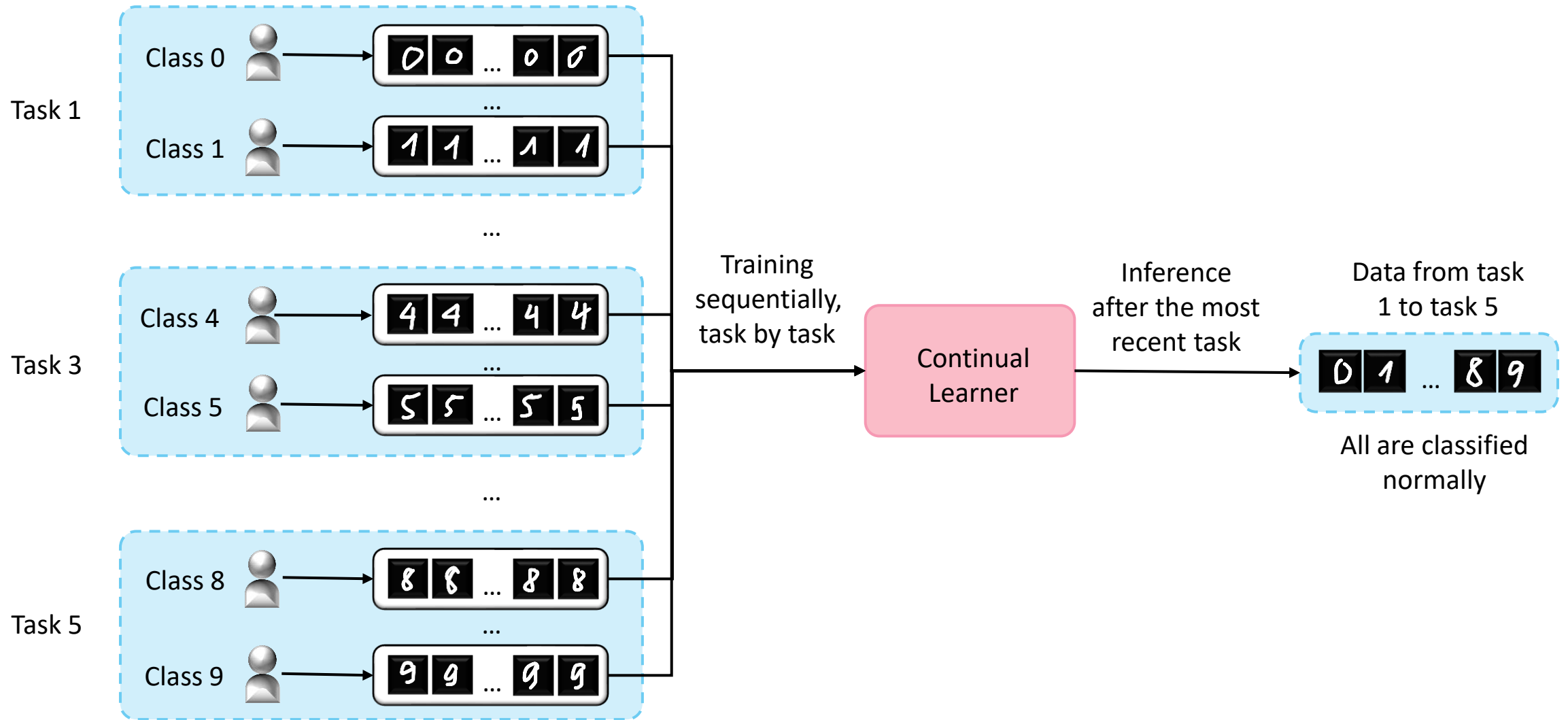
Trang Pham<sup>3\*</sup>, Linh Van Ngo<sup>2</sup>, Nhat Ho<sup>1</sup>

<sup>1</sup> University of Texas at Austin

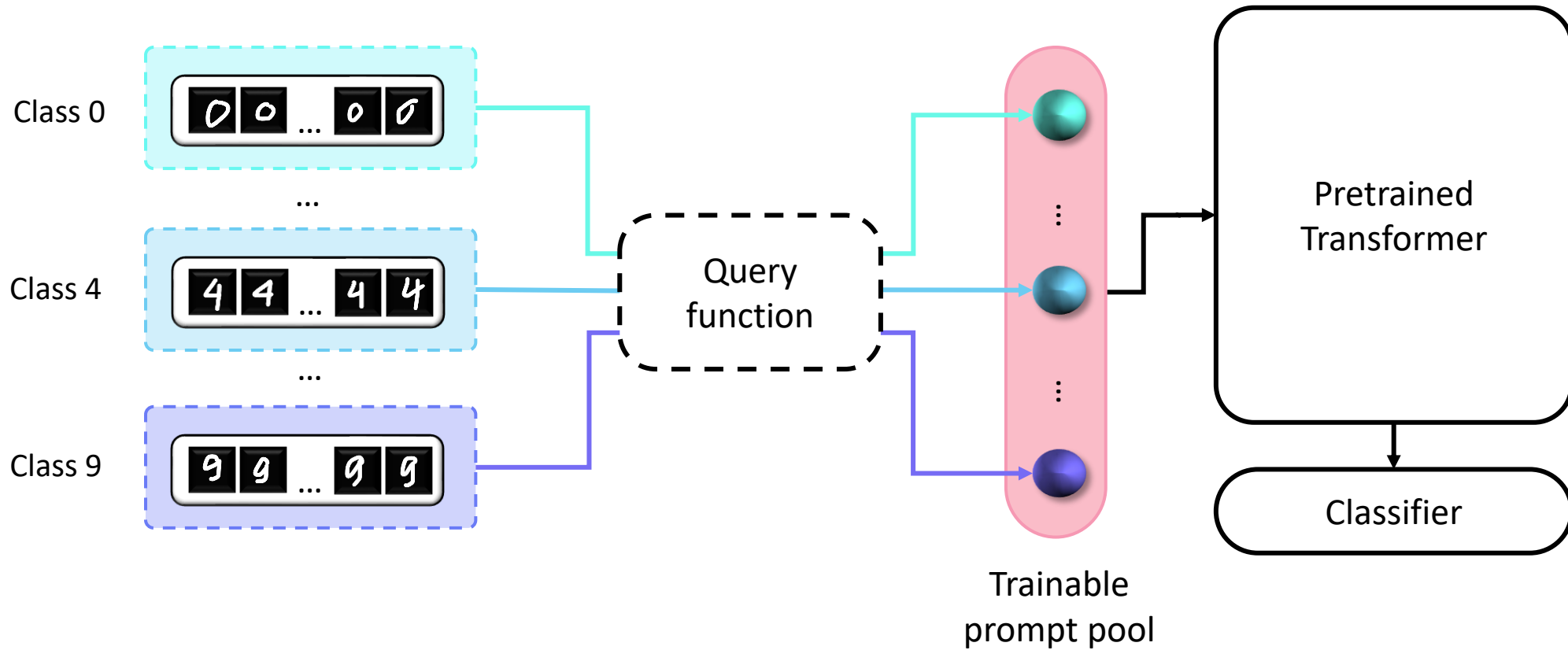
<sup>2</sup> Hanoi University of Science and Technology

<sup>3</sup> VinAI Research

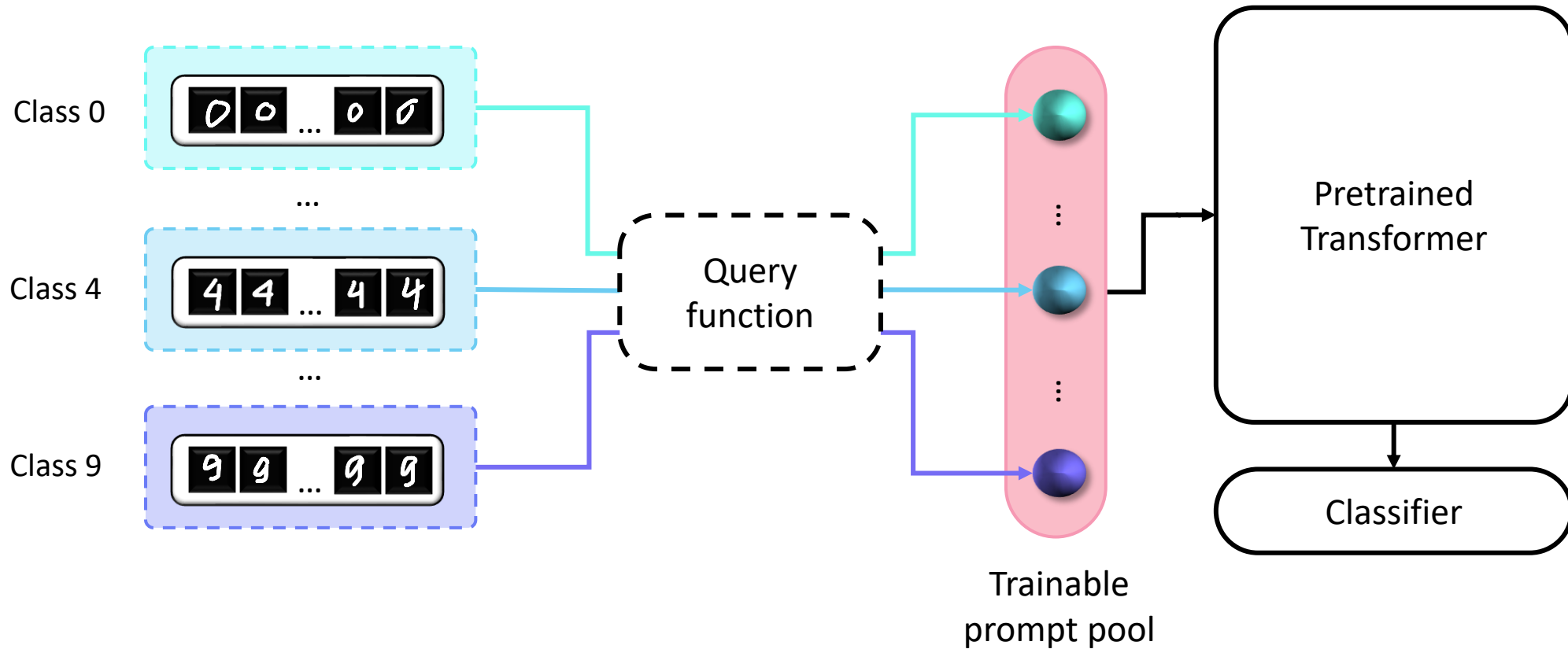
# Continual Learning



# Prompt-based Continual Learning



# Prompt-based Continual Learning



Efficient, astonishing performance  
BUT LACKS A THEORETICAL FOUNDATION!

# Mixture of Experts

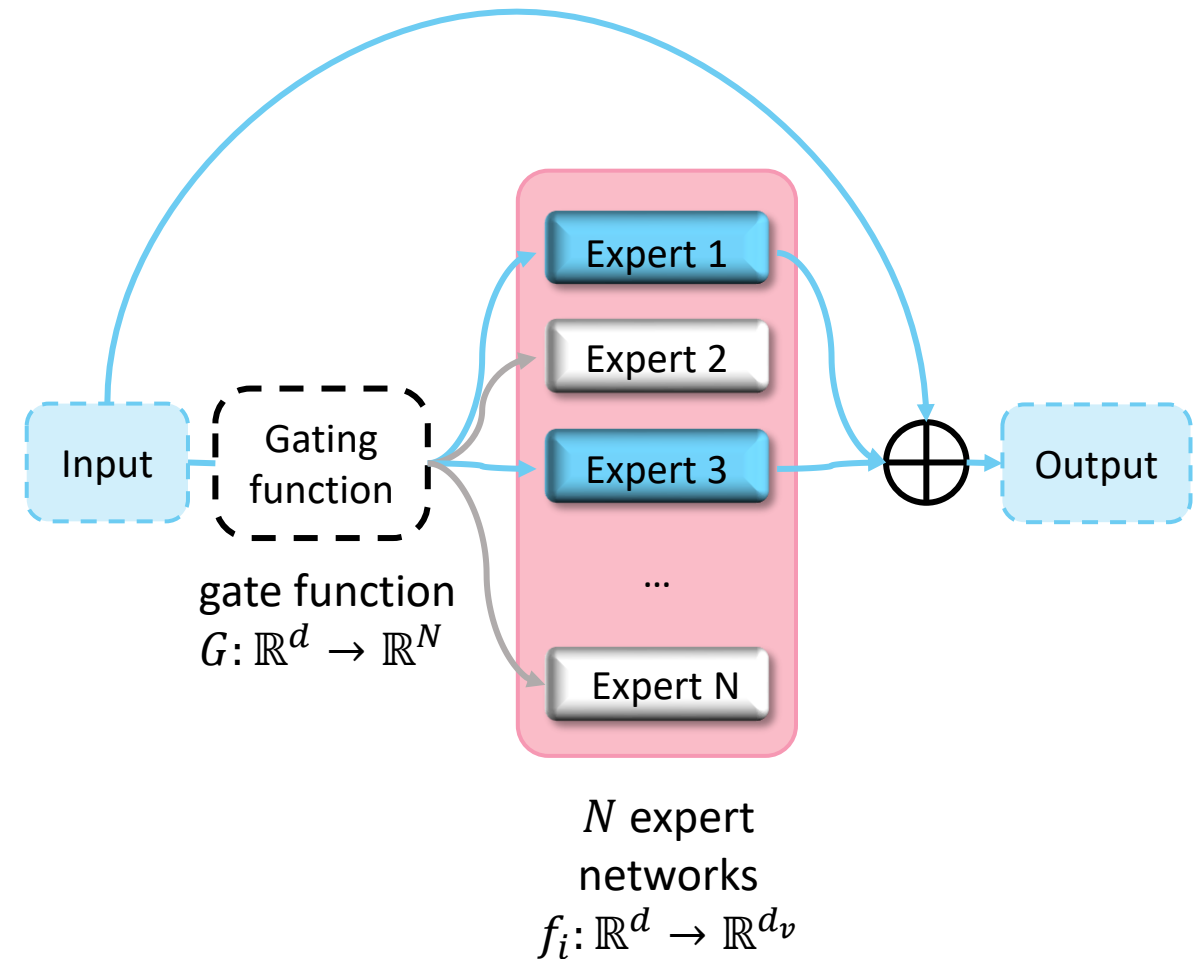
An MoE model consists of:

- a group of  $N$  expert networks  $f_i: \mathbb{R}^d \rightarrow \mathbb{R}^{d_v}$
- a gate function  $G: \mathbb{R}^d \rightarrow \mathbb{R}^N$
- a learned score function  $s_i: \mathbb{R}^d \rightarrow \mathbb{R}$

Given an input  $h \in \mathbb{R}^d$ , its MoE output is computed as:

$$y := \sum_{j=1}^N G(\mathbf{h})_j \cdot f_j(\mathbf{h}) := \sum_{j=1}^N \frac{\exp(s_j(\mathbf{h}))}{\sum_{l=1}^N \exp(s_l(\mathbf{h}))} \cdot f_j(\mathbf{h}),$$

where  $G(\mathbf{h}) := \text{softmax}(s_1(\mathbf{h}), \dots, s_N(\mathbf{h}))$ .



# Mixture of Experts

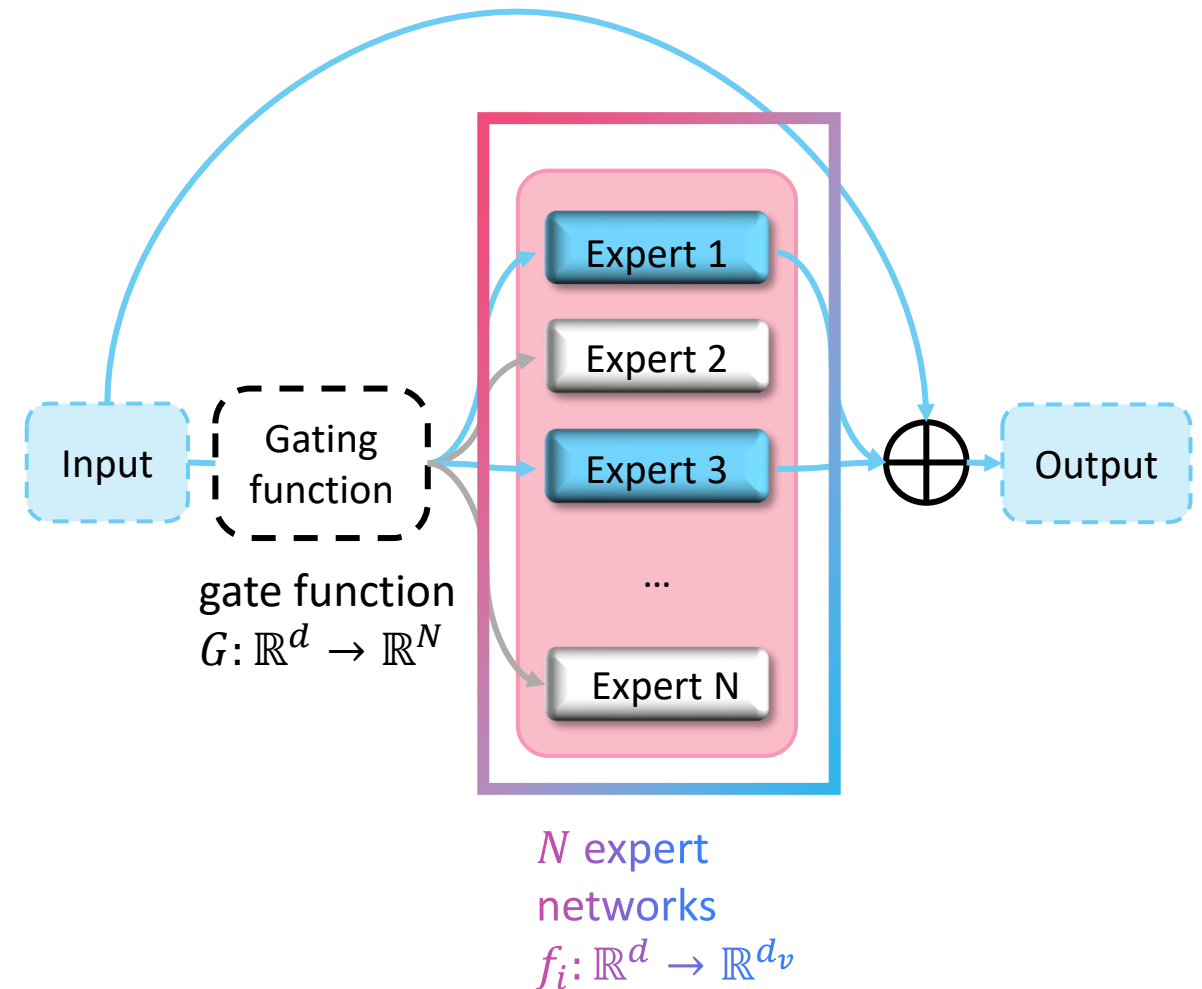
An MoE model consists of:

- a group of  $N$  expert networks  $f_i: \mathbb{R}^d \rightarrow \mathbb{R}^{d_v}$ , for all  $i \in [N]$
- a gate function  $G: \mathbb{R}^d \rightarrow \mathbb{R}^N$
- a learned score function  $s_i: \mathbb{R}^d \rightarrow \mathbb{R}$

Given an input  $h \in \mathbb{R}^d$ , its MoE output is computed as:

$$y := \sum_{j=1}^N G(\mathbf{h})_j \cdot f_j(\mathbf{h}) := \sum_{j=1}^N \frac{\exp(s_j(\mathbf{h}))}{\sum_{l=1}^N \exp(s_l(\mathbf{h}))} \cdot f_j(\mathbf{h}),$$

where  $G(\mathbf{h}) := \text{softmax}(s_1(\mathbf{h}), \dots, s_N(\mathbf{h}))$ .



# Mixture of Experts

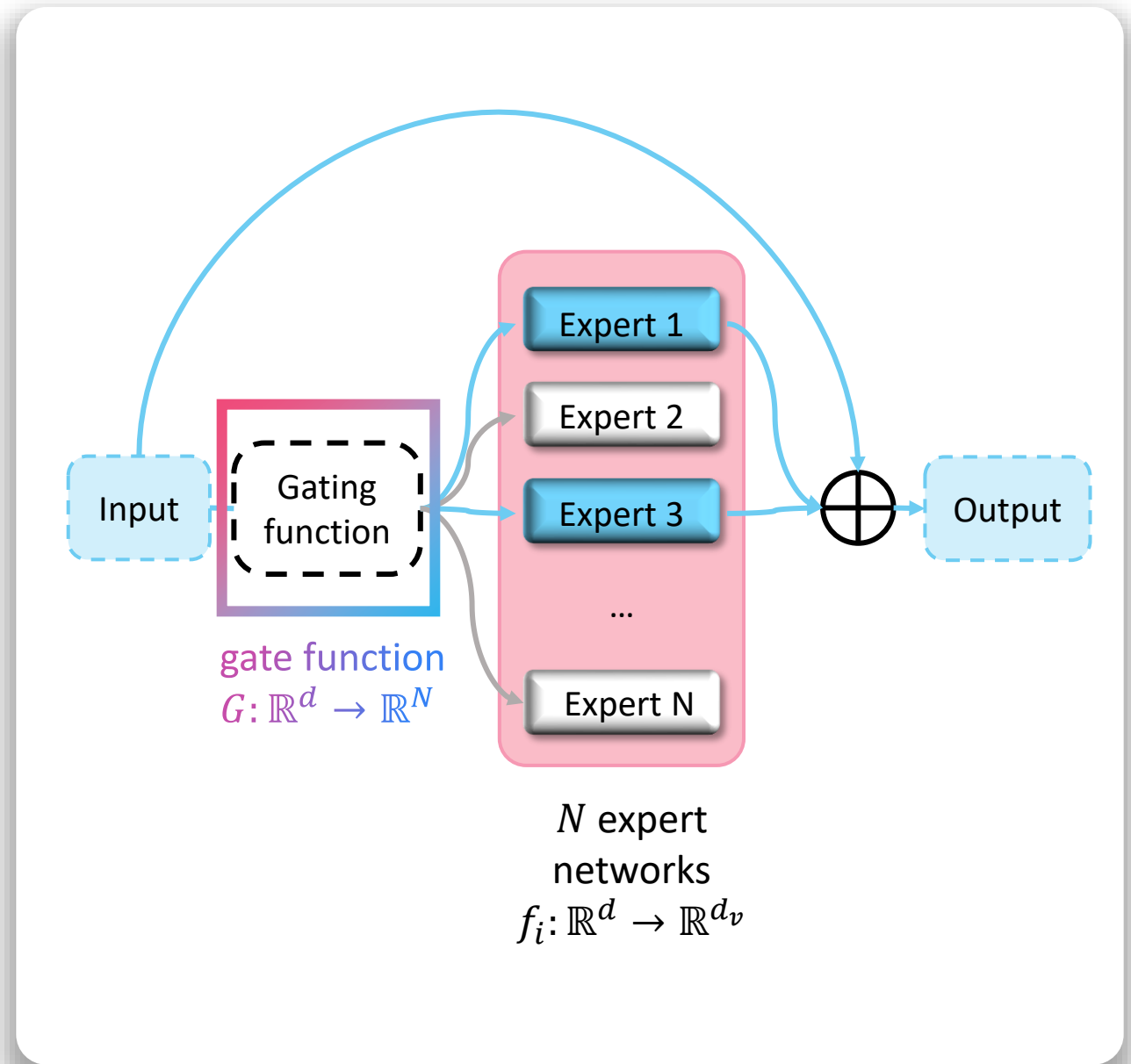
An MoE model consists of:

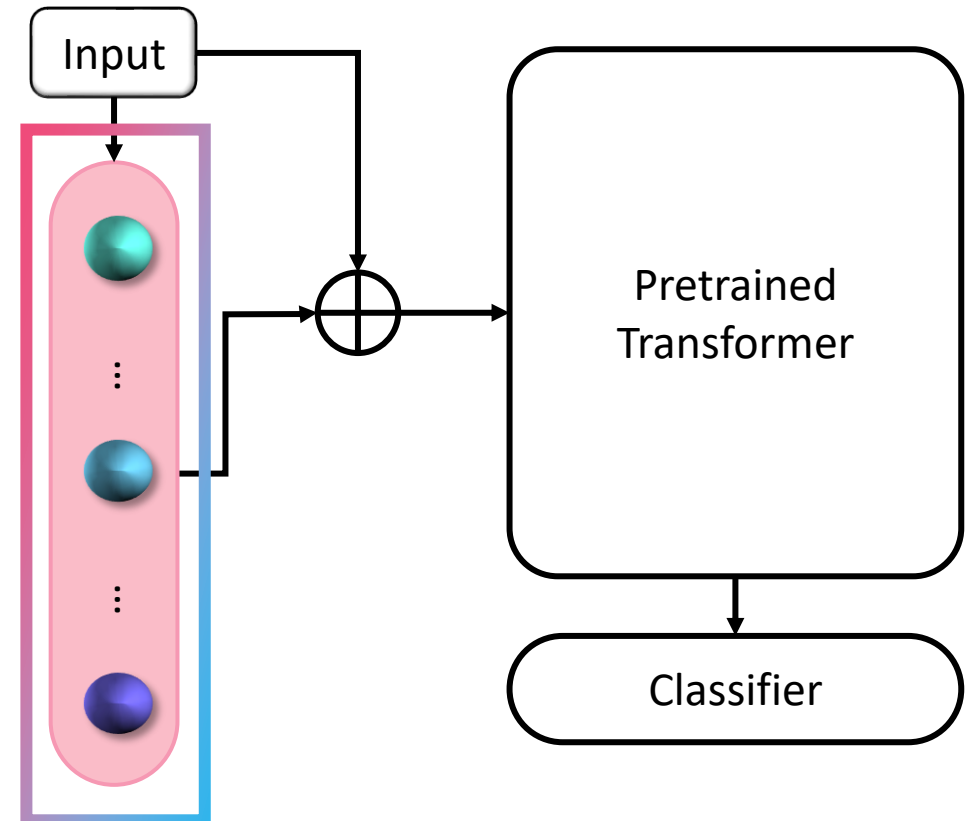
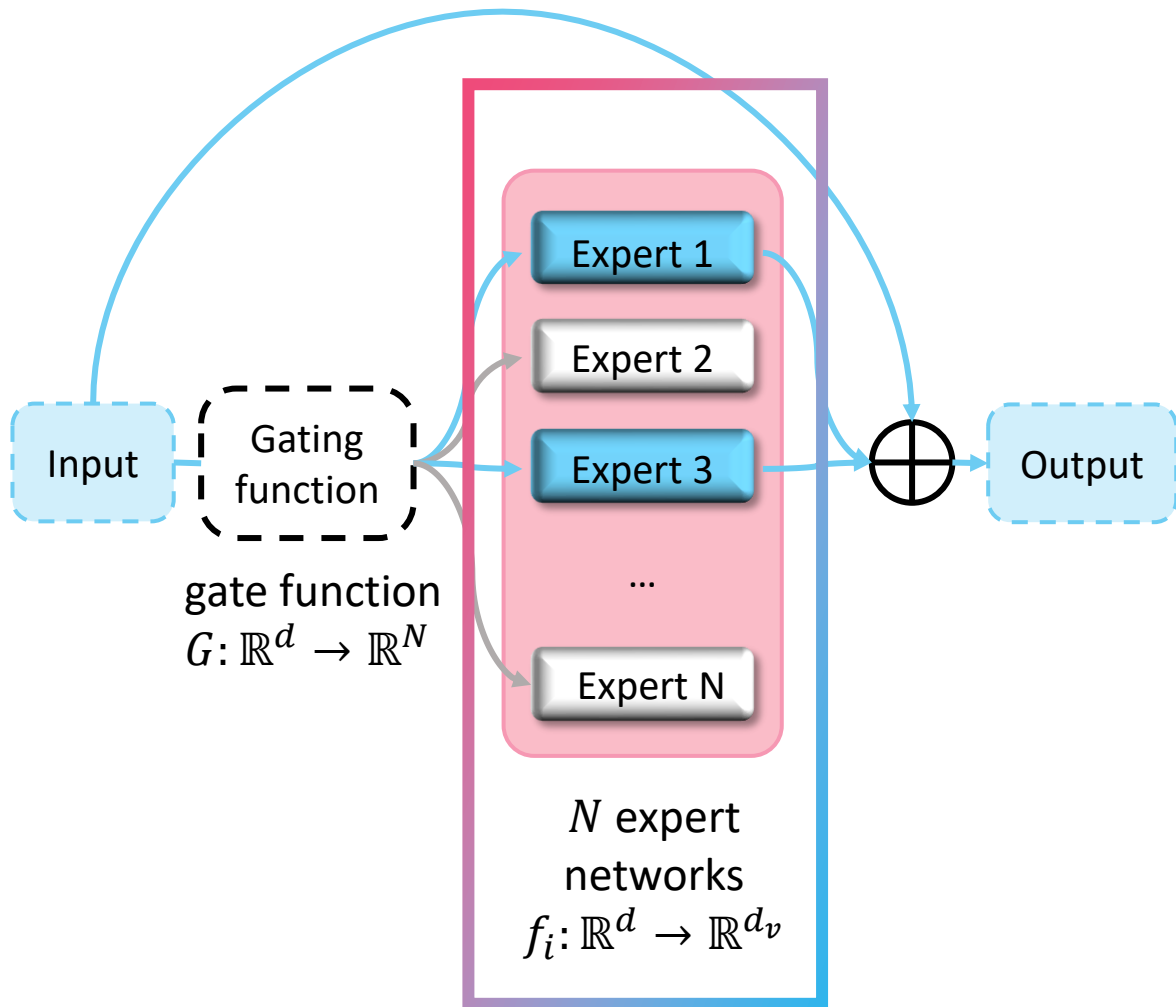
- a group of  $N$  expert networks  $f_i: \mathbb{R}^d \rightarrow \mathbb{R}^{d_v}$ , for all  $i \in [N]$
- a **gate function**  $G: \mathbb{R}^d \rightarrow \mathbb{R}^N$
- learned **score function**  $s_i: \mathbb{R}^d \rightarrow \mathbb{R}$

Given an input  $h \in \mathbb{R}^d$ , its MoE output is computed as:

$$y := \sum_{j=1}^N G(\mathbf{h})_j \cdot f_j(\mathbf{h}) := \sum_{j=1}^N \frac{\exp(s_j(\mathbf{h}))}{\sum_{l=1}^N \exp(s_l(\mathbf{h}))} \cdot f_j(\mathbf{h}),$$

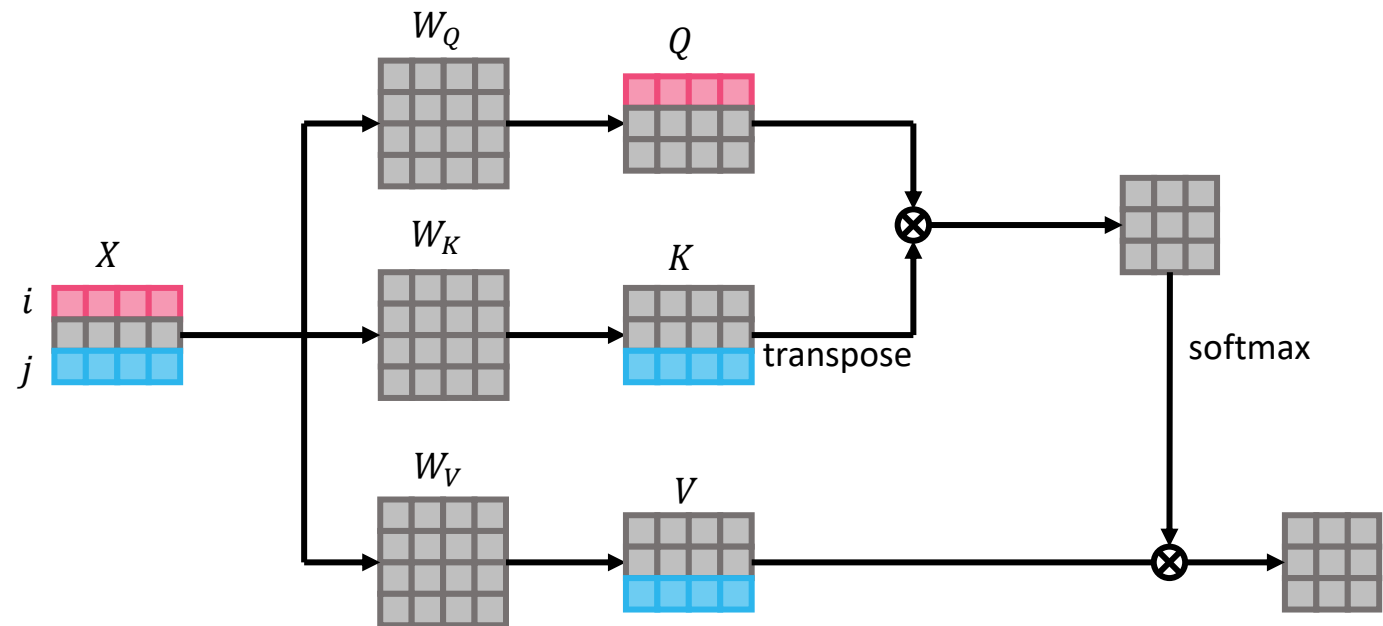
where  $G(\mathbf{h}) := \text{softmax}(s_1(\mathbf{h}), \dots, s_N(\mathbf{h}))$ .







# Mixture of Experts and Self-Attention



Attention Mechanism

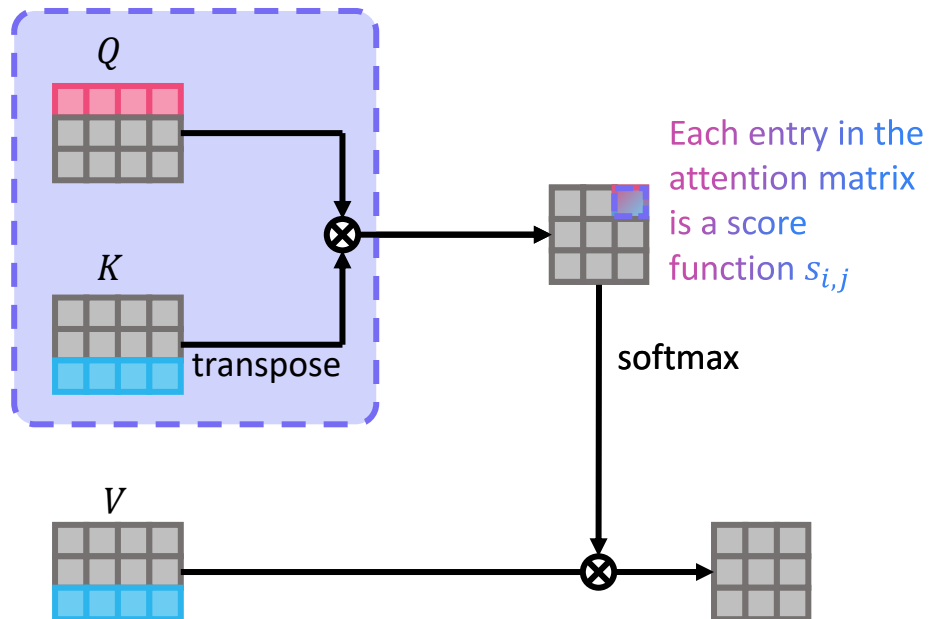
# Mixture of Experts and Self-Attention

- Define  $N$  gating function  $G_i$  with the score function for the  $j^{\text{th}}$  expert of the  $i^{\text{th}}$  gating  $s_{i,j}$ :

$$s_{i,j}(X) := \frac{\mathbf{X}^\top E_i^\top W_l^Q W_l^{K^\top} E_j \mathbf{X}}{\sqrt{d_v}} = \frac{\mathbf{X}_i^\top W_l^Q W_l^{K^\top} \mathbf{X}_j}{\sqrt{d_v}}$$

for  $i, j \in [N]$ .

score function  $s_{i,j}$



# Mixture of Experts and Self-Attention

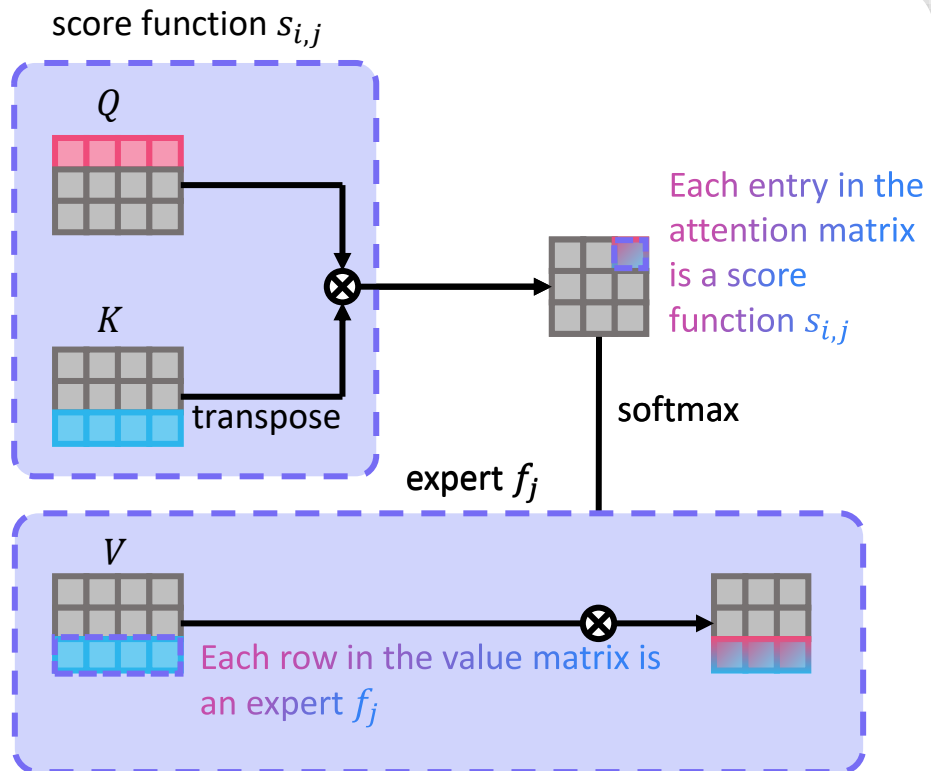
- Define  $N$  gating function  $G_i$  with the score function for the  $j^{\text{th}}$  expert of the  $i^{\text{th}}$  gating  $s_{i,j}$ :

$$s_{i,j}(X) := \frac{\mathbf{X}^\top E_i^\top W_l^Q W_l^{K^\top} E_j \mathbf{X}}{\sqrt{d_v}} = \frac{\mathbf{x}_i^\top W_l^Q W_l^{K^\top} \mathbf{x}_j}{\sqrt{d_v}}$$

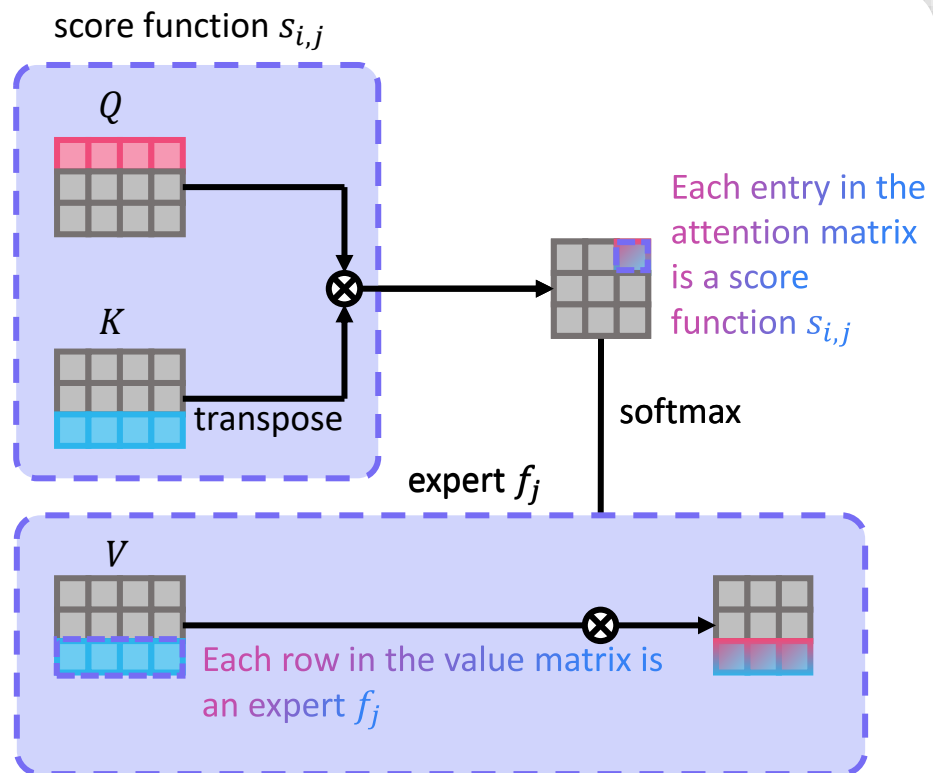
for  $i, j \in [N]$ .

- Define  $N$  experts  $f_j$ :

$$f_j(X) := W_l^{V^\top} E_j \mathbf{X} = W_l^{V^\top} \mathbf{x}_j$$



# Mixture of Experts and Self-Attention



We can express the output of the  $l^{\text{th}}$  head as follows:

$$\mathbf{h}_l = [\mathbf{h}_{l,1}, \dots, \mathbf{h}_{l,N}]^T \in \mathbb{R}^{N \times d_v}$$

$$\mathbf{h}_{l,i} = \sum_{j=1}^N \frac{\exp(s_{i,j}(\mathbf{X}))}{\sum_{k=1}^N \exp(s_{i,k}(\mathbf{X}))} f_j(\mathbf{X})$$

We can interpret each head in a multi-head self-attention layer as a multi-gate mixture of experts architecture.

## Prefix Tuning via the Perspective of Mixture of Experts

- Prefix tuning can be interpreted as the introduction of new experts to customize the pre-trained model for a specific task

$$\mathbf{p}^K = [\mathbf{p}_1^K, \dots, \mathbf{p}_L^K]^\top \in \mathbb{R}^{L \times d}, \mathbf{p}^V = [\mathbf{p}_1^V, \dots, \mathbf{p}_L^V]^\top \in \mathbb{R}^{L \times d}$$

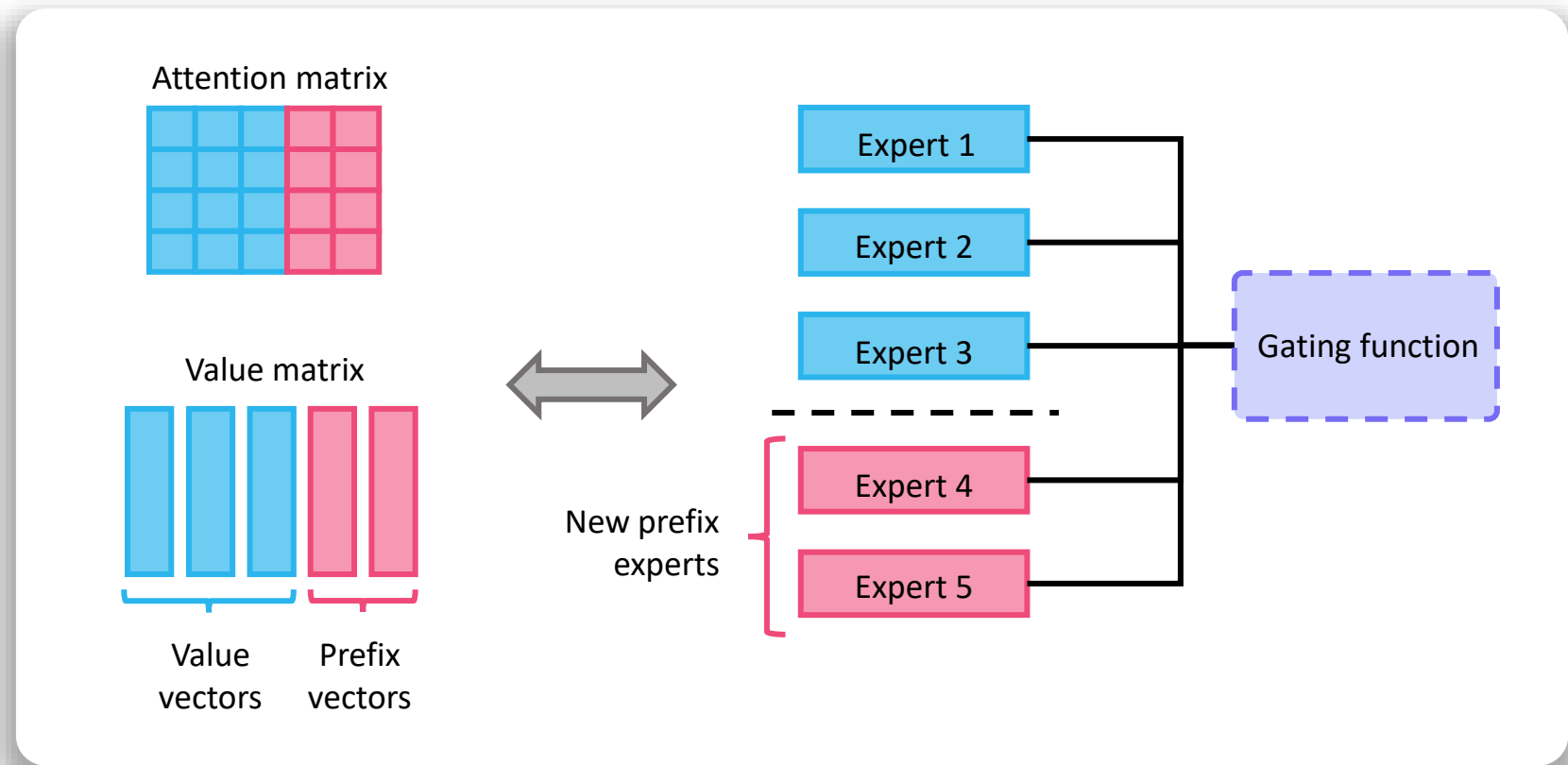
- Define new prefix experts along with their corresponding new score functions:

$$f_{N+j}(\mathbf{x}) := W_l^{V^\top} \mathbf{p}_j^V,$$

$$s_{i,N+j}(x) := \frac{\mathbf{x}^\top E_i^\top W_l^Q W_l^{K^\top} \mathbf{p}_j^K}{\sqrt{d_v}} = \frac{\mathbf{x}_i^\top W_l^Q W_l^{K^\top} \mathbf{p}_j^K}{\sqrt{d_v}}$$

for  $i \in [N]$  and  $j \in [L]$ .

# Linear gating prefix MoE model

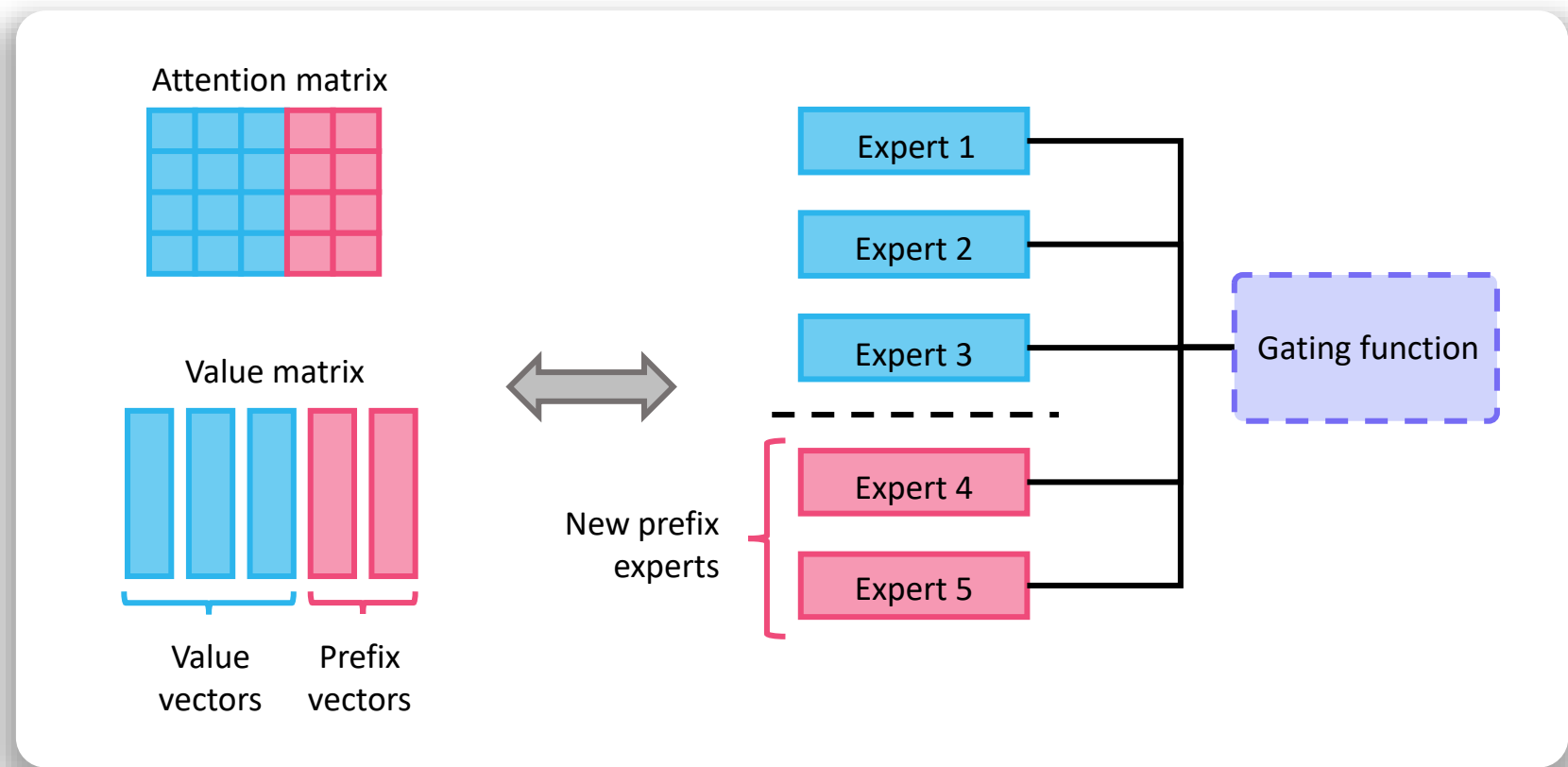


$$\tilde{\mathbf{h}}_l = \text{Attention} \left( \mathbf{X}^Q W_l^Q, \begin{bmatrix} \mathbf{p}^K \\ \mathbf{X}^K \end{bmatrix} W_l^K, \begin{bmatrix} \mathbf{p}^V \\ \mathbf{X}^V \end{bmatrix} W_l^V \right) = [\tilde{\mathbf{h}}_{l,1}, \dots, \tilde{\mathbf{h}}_{l,N}]^\top \in \mathbb{R}^{N \times d_v},$$

$$\tilde{\mathbf{h}}_{l,i} = \sum_{j=1}^N \frac{\exp(s_{i,j}(\mathbf{X}))}{\sum_{k=1}^N \exp(s_{i,k}(\mathbf{X})) + \sum_{k'=1}^L \exp(s_{i,N+k'}(\mathbf{X}))} f_j(\mathbf{X}) \quad \left. \vphantom{\sum_{j=1}^N} \right\} \begin{array}{l} \text{Pretrained} \\ \text{experts} \end{array}$$

$$+ \sum_{j'=1}^L \frac{\exp(s_{i,N+j'}(\mathbf{X}))}{\sum_{k=1}^N \exp(s_{i,k}(\mathbf{X})) + \sum_{k'=1}^L \exp(s_{i,N+k'}(\mathbf{X}))} f_{N+j'}(\mathbf{X}) \quad \left. \vphantom{\sum_{j'=1}^L} \right\} \begin{array}{l} \text{New} \\ \text{experts} \end{array}$$

# Linear gating prefix MoE model



Parameter estimation rate is  $O(1/\log(n)^{\tau})$ .

Requires HUGE amount of data!

$$\tilde{\mathbf{h}}_l = \text{Attention} \left( \mathbf{X}^Q W_l^Q, \begin{bmatrix} \mathbf{p}^K \\ \mathbf{X}^K \end{bmatrix} W_l^K, \begin{bmatrix} \mathbf{p}^V \\ \mathbf{X}^V \end{bmatrix} W_l^V \right) = [\tilde{\mathbf{h}}_{l,1}, \dots, \tilde{\mathbf{h}}_{l,N}]^T \in \mathbb{R}^{N \times d_v},$$

$$\tilde{\mathbf{h}}_{l,i} = \sum_{j=1}^N \frac{\exp(s_{i,j}(\mathbf{X}))}{\sum_{k=1}^N \exp(s_{i,k}(\mathbf{X})) + \sum_{k'=1}^L \exp(s_{i,N+k'}(\mathbf{X}))} f_j(\mathbf{X}) \quad \left. \vphantom{\sum_{j=1}^N} \right\} \text{Pretrained experts}$$

$$+ \sum_{j'=1}^L \frac{\exp(s_{i,N+j'}(\mathbf{X}))}{\sum_{k=1}^N \exp(s_{i,k}(\mathbf{X})) + \sum_{k'=1}^L \exp(s_{i,N+k'}(\mathbf{X}))} f_{N+j'}(\mathbf{X}) \quad \left. \vphantom{\sum_{j'=1}^L} \right\} \text{New experts}$$

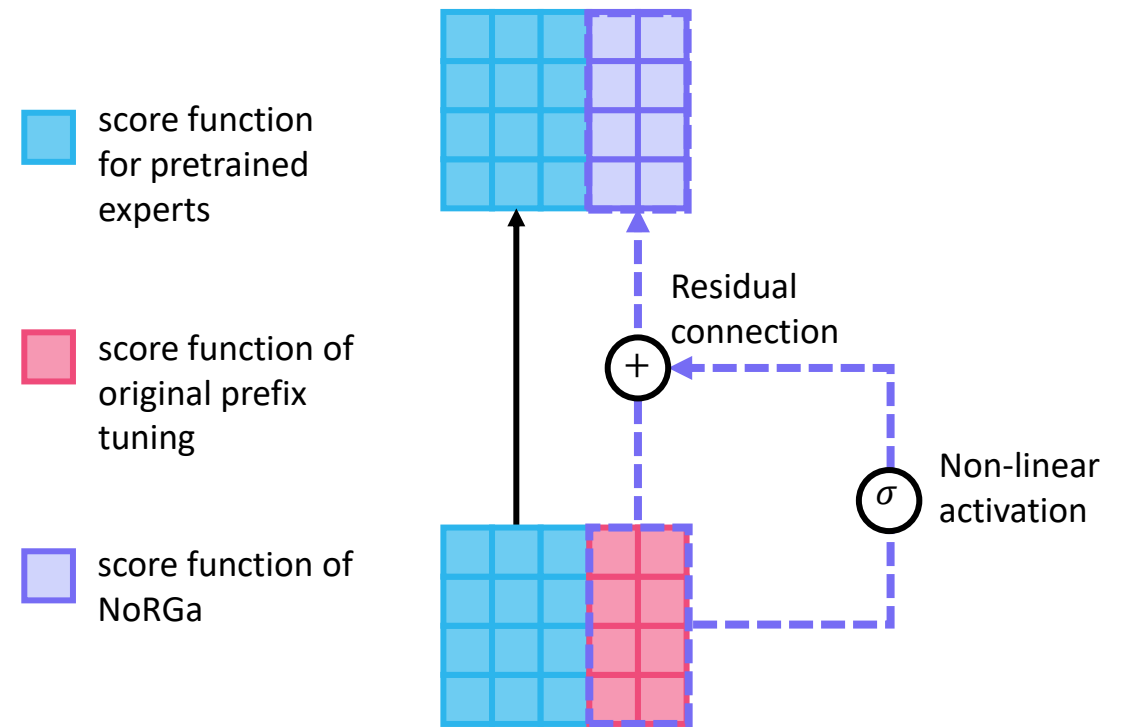
# NoRGa – Non-linear Residual Gate

- Modify the linear gating prefix MoE model:

$$\hat{s}_{i,N+j}(\mathbf{X}) = s_{i,N+j}(\mathbf{X}) + \alpha \cdot \sigma\left(\tau \cdot s_{i,N+j}(\mathbf{X})\right),$$

$i \in [N], j \in [L]$

where  $\alpha, \tau$  are scalar factors,  $\sigma$  is a non-linear activation function





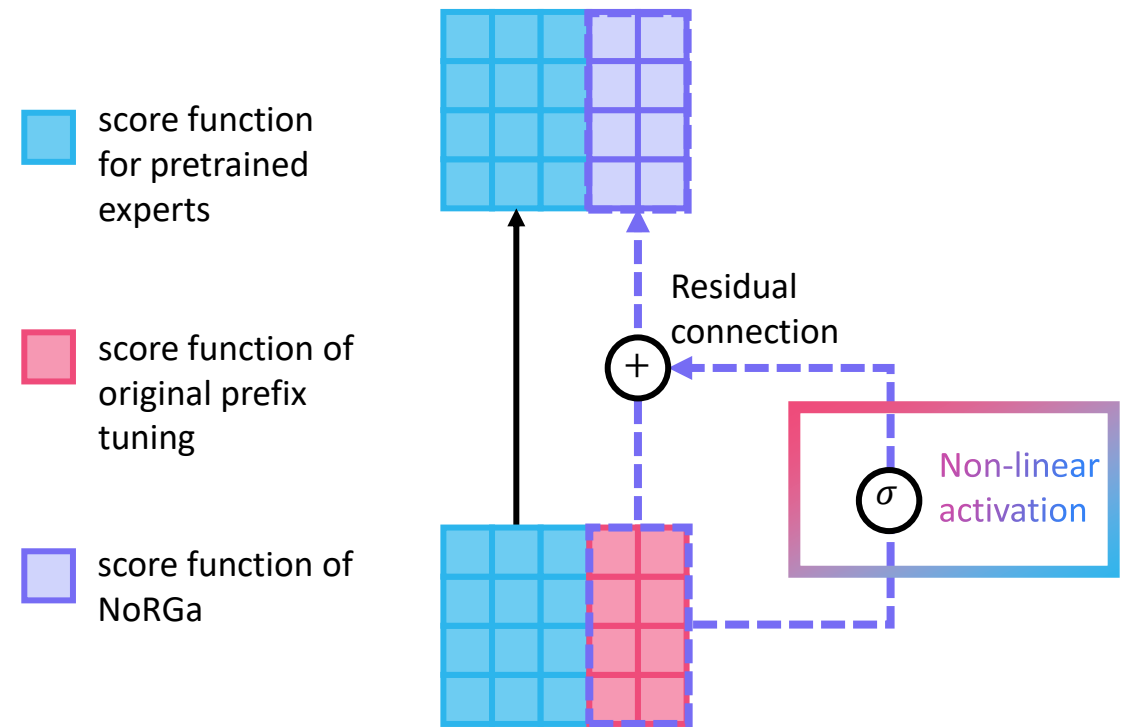
# NoRGa – Non-linear Residual Gate

- Modify the linear gating prefix MoE model:

$$\hat{s}_{i,N+j}(\mathbf{X}) = s_{i,N+j}(\mathbf{X}) + \alpha \cdot \sigma(\tau \cdot s_{i,N+j}(\mathbf{X})),$$

$i \in [N], j \in [L]$

where  $\alpha, \tau$  are scalar factors,  $\sigma$  is a non-linear activation function



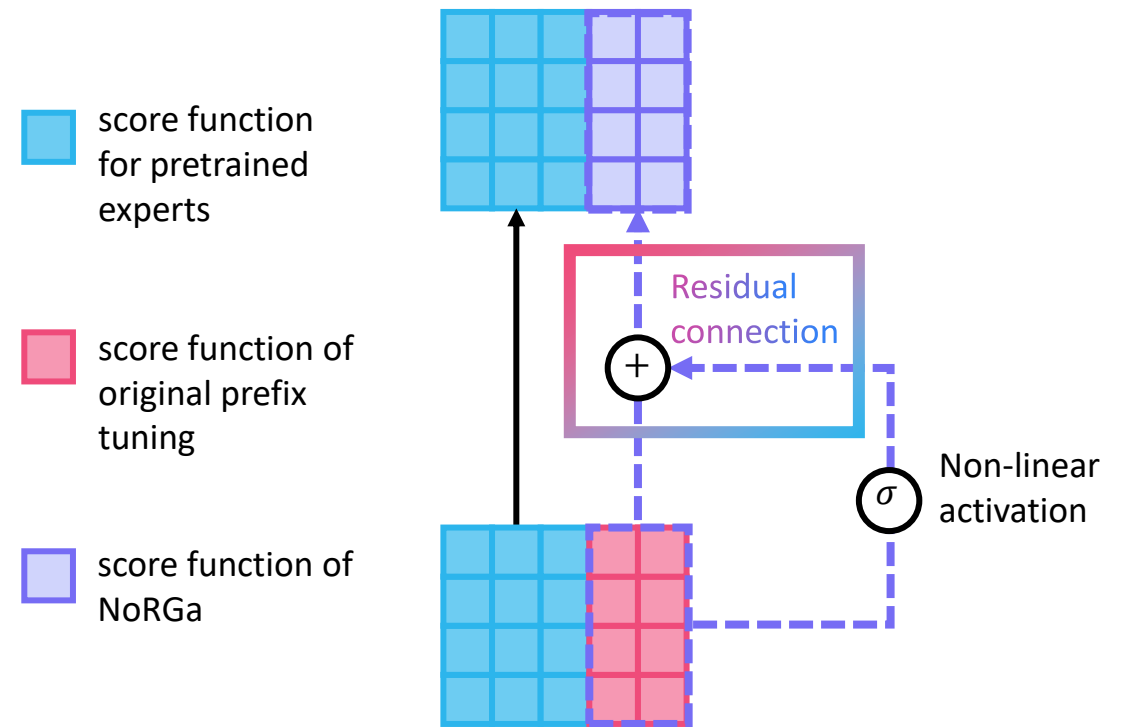
# NoRGa – Non-linear Residual Gate

- Modify the linear gating prefix MoE model:

$$\hat{s}_{i,N+j}(\mathbf{X}) = s_{i,N+j}(\mathbf{X}) + \alpha \cdot \sigma(\tau \cdot s_{i,N+j}(\mathbf{X})),$$

$i \in [N], j \in [L]$

where  $\alpha, \tau$  are scalar factors,  $\sigma$  is a non-linear activation function



# NoRGa – Non-linear Residual Gate

- We prove that estimating parameters in the non-linear residual gating prefix MoE model is statistically efficient in terms of the number of data.

Model	Parameter estimation rate	Number of data
Linear gating prefix MoE	$O(1/\log(n)^\tau)$	Exponential $\exp(\varepsilon^{-\tau})$
Non-linear residual gating prefix MoE	$O\left(\sqrt[4]{\log(n)/n}\right)$	Polynomial $\varepsilon^{-4}$

# Experiments

PTM	Method	Split CIFAR-100			Split ImageNet-R		
		FA (↑)	CA (↑)	FM (↓)	FA (↑)	CA (↑)	FM (↓)
Sup-21K	L2P	83.06 ± 0.17	88.27 ± 0.71	5.61 ± 0.32	67.53 ± 0.44	71.98 ± 0.52	5.84 ± 0.38
	DualPrompt	87.30 ± 0.27	91.23 ± 0.65	3.87 ± 0.43	70.93 ± 0.08	75.67 ± 0.52	5.47 ± 0.19
	S-Prompt	87.57 ± 0.42	91.38 ± 0.69	3.63 ± 0.41	69.88 ± 0.51	74.25 ± 0.55	4.73 ± 0.47
	CODA-Prompt	86.94 ± 0.63	91.57 ± 0.75	4.04 ± 0.18	70.03 ± 0.47	74.26 ± 0.24	5.17 ± 0.22
	HiDe-Prompt	92.61 ± 0.28	94.03 ± 0.01	1.50 ± 0.28	75.06 ± 0.12	76.60 ± 0.01	<b>4.09 ± 0.13</b>
	<b>NoRGa (Ours)</b>	<b>94.48 ± 0.13</b>	<b>95.83 ± 0.37</b>	<b>1.44 ± 0.27</b>	<b>75.40 ± 0.39</b>	<b>79.52 ± 0.07</b>	4.59 ± 0.07
iBOT-21K	L2P	79.13 ± 1.25	85.13 ± 0.05	7.50 ± 1.21	61.31 ± 0.50	68.81 ± 0.52	10.72 ± 0.40
	DualPrompt	78.84 ± 0.47	86.16 ± 0.02	8.84 ± 0.67	58.69 ± 0.61	66.61 ± 0.67	11.75 ± 0.92
	S-Prompt	79.14 ± 0.65	85.85 ± 0.17	8.23 ± 1.15	57.96 ± 1.10	66.42 ± 0.71	11.27 ± 0.72
	CODA-Prompt	80.83 ± 0.27	87.02 ± 0.20	7.50 ± 0.25	61.22 ± 0.35	66.76 ± 0.37	9.66 ± 0.20
	HiDe-Prompt	93.02 ± 0.15	94.56 ± 0.05	<b>1.26 ± 0.13</b>	70.83 ± 0.17	73.23 ± 0.08	<b>6.77 ± 0.23</b>
	<b>NoRGa (Ours)</b>	<b>94.76 ± 0.15</b>	<b>95.86 ± 0.31</b>	1.34 ± 0.14	<b>73.06 ± 0.26</b>	<b>77.46 ± 0.42</b>	6.88 ± 0.49
iBOT-1K	L2P	75.51 ± 0.88	82.53 ± 1.10	6.80 ± 1.70	59.43 ± 0.28	66.83 ± 0.92	11.33 ± 1.25
	DualPrompt	76.21 ± 1.00	83.54 ± 1.23	9.89 ± 1.81	60.41 ± 0.76	66.87 ± 0.41	9.21 ± 0.43
	S-Prompt	76.60 ± 0.61	82.89 ± 0.89	8.60 ± 1.36	59.56 ± 0.60	66.60 ± 0.13	8.83 ± 0.81
	CODA-Prompt	79.11 ± 1.02	86.21 ± 0.49	7.69 ± 1.57	66.56 ± 0.68	73.14 ± 0.57	7.22 ± 0.38
	HiDe-Prompt	93.48 ± 0.11	95.02 ± 0.01	1.63 ± 0.10	71.33 ± 0.21	73.62 ± 0.13	7.11 ± 0.02
	<b>NoRGa (Ours)</b>	<b>94.01 ± 0.04</b>	<b>95.11 ± 0.35</b>	<b>1.61 ± 0.30</b>	<b>72.77 ± 0.20</b>	<b>76.55 ± 0.46</b>	<b>7.10 ± 0.39</b>

# Experiments

Method	Split CIFAR-100		Split CUB-200	
	Sup-21K	iBOT-21K	Sup-21K	iBOT-21K
HiDe-Prompt	92.61	93.02	86.56	78.23
NoRGa tanh	94.36	<b>94.76</b>	90.87	<b>80.69</b>
NoRGa sigmoid	<b>94.48</b>	94.69	<b>90.90</b>	80.18
NoRGa gelu	94.05	94.63	90.74	80.54



## Conclusion

- Reveals a novel connection between **Prefix Tuning**, a popular prompt implementation technique, and **Mixture of Experts**.
- Proposes **Non-linear Residual Gates (NoRGa)**, an innovative gating mechanism.
- Achieves state-of-the-art performance across various **continual learning benchmarks** and pretraining settings.

# THANK YOU FOR YOUR ATTENTION!

---

Mixture of experts meets Prompt-based Continual Learning

**Minh Le<sup>3</sup>, An Nguyen<sup>2\*</sup>, Huy Nguyen<sup>1\*</sup>, Trang Nguyen<sup>3\*</sup>,  
Trang Pham<sup>3\*</sup>, Linh Van Ngo<sup>2</sup>, Nhat Ho<sup>1</sup>**

<sup>1</sup> University of Texas at Austin

<sup>2</sup> Hanoi University of Science and Technology

<sup>3</sup> VinAI Research

