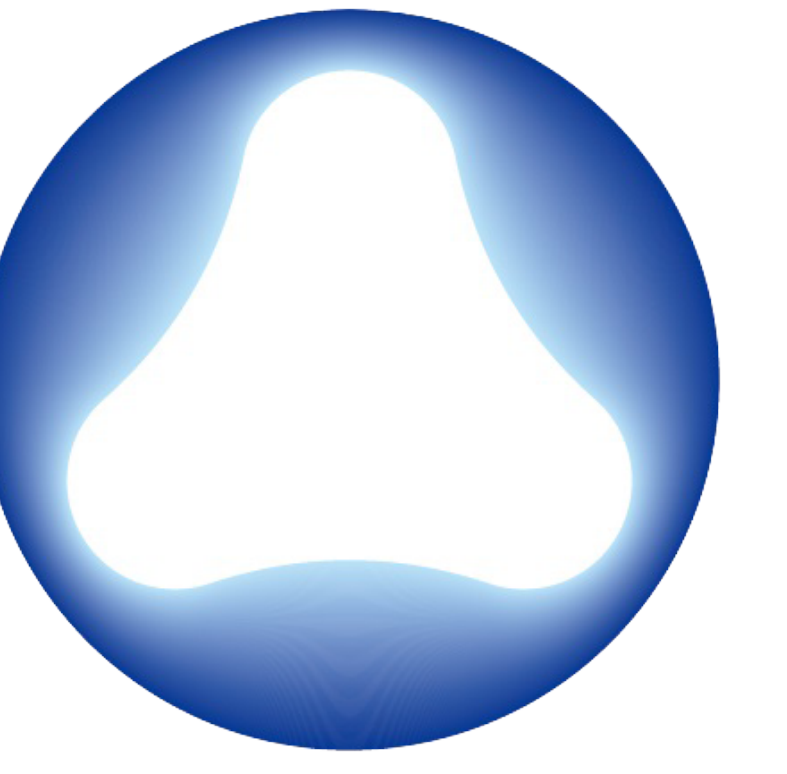


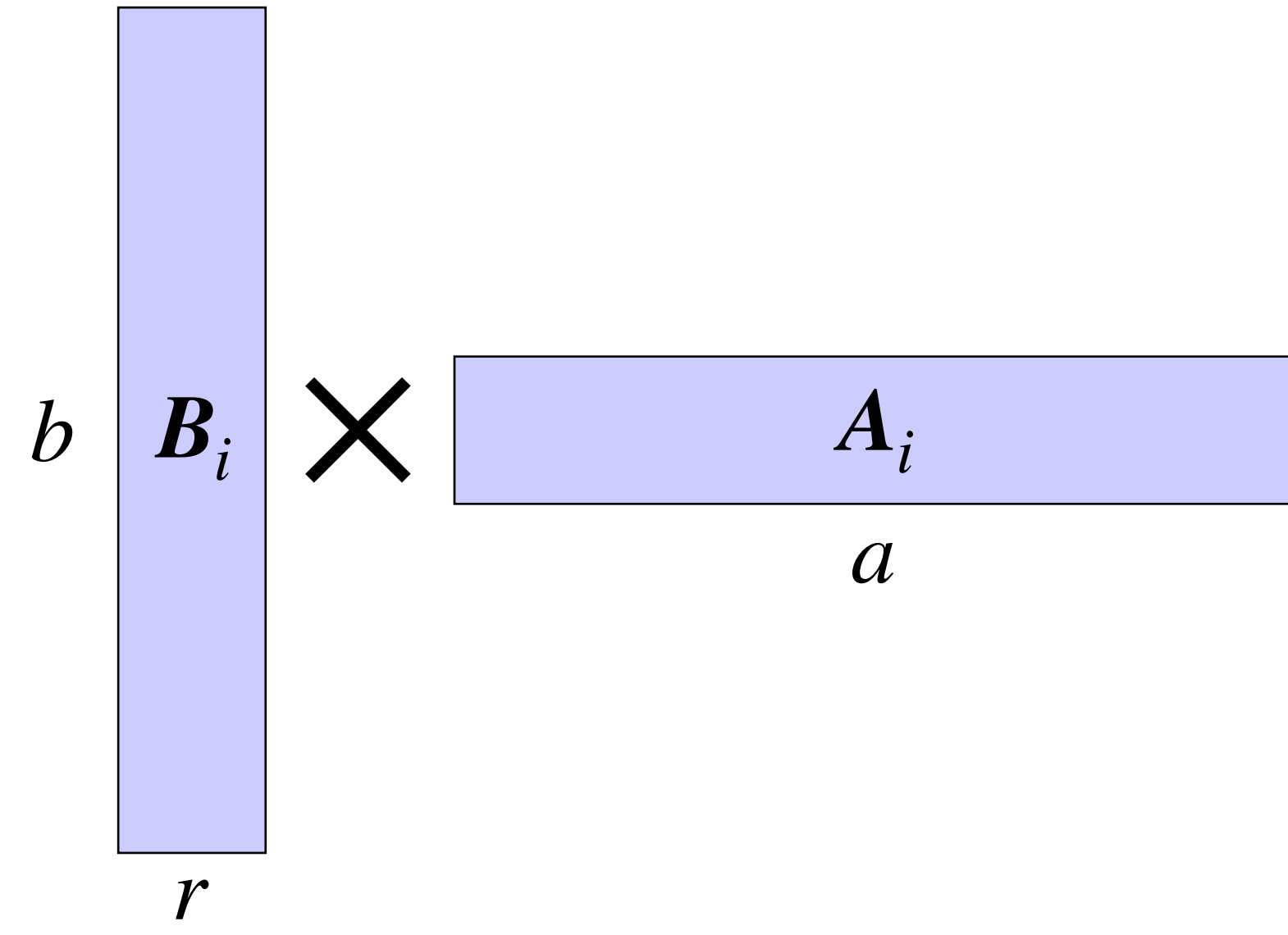


RaSA: Rank-Sharing Low-Rank Adaptation

Zhiwei He Zhaopeng Tu Xing Wang Xingyu Chen Zhijie Wang Jiahao Xu Tian Liang
Wenxiang Jiao Zhuosheng Zhang Rui Wang



Background & Motivation



$$\mathbf{W} + \Delta\mathbf{W} = \mathbf{W} + \frac{\alpha}{r} \mathbf{B}\mathbf{A} \quad (\mathbf{B} \in \mathbb{R}^{b \times r}, \mathbf{A} \in \mathbb{R}^{r \times a})$$

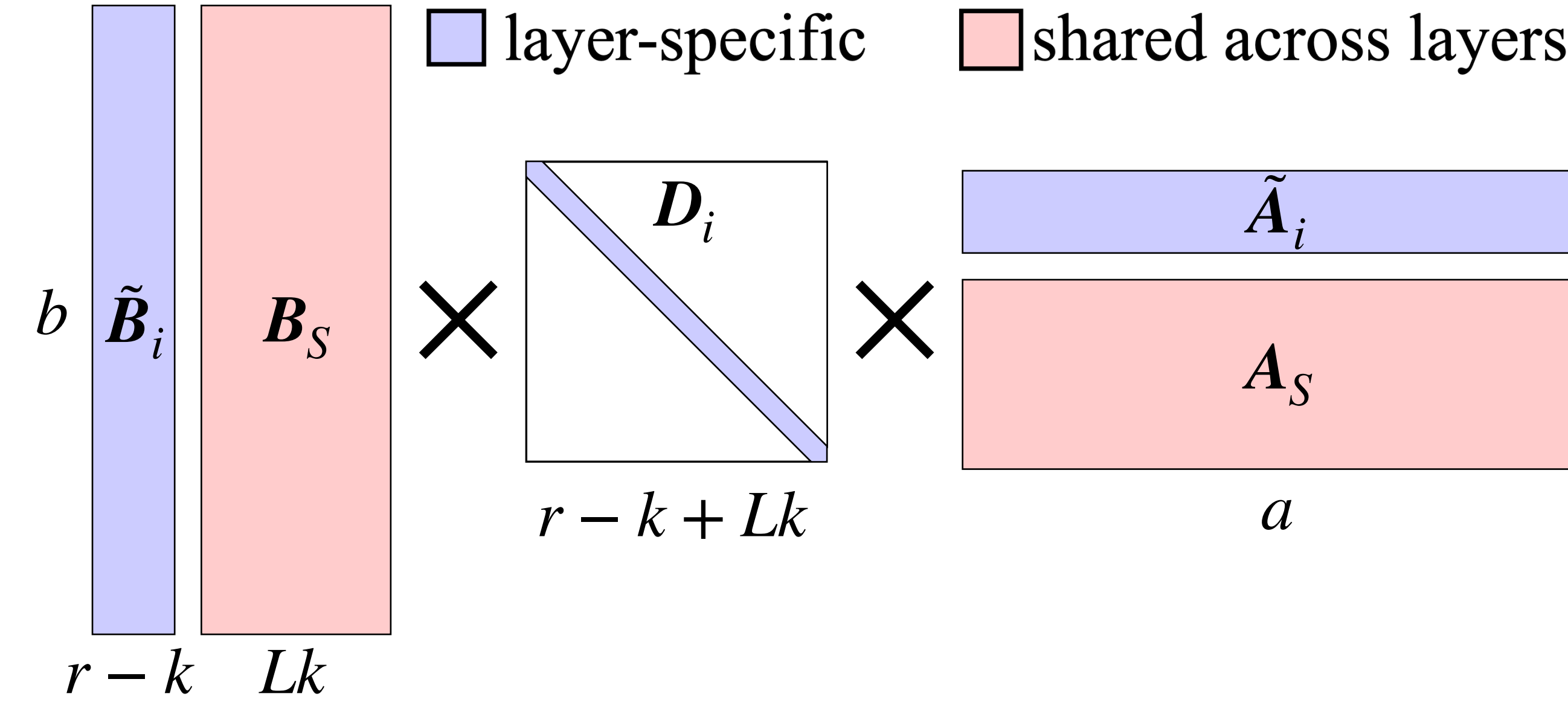
LoRA

✓ LoRA still lags behind full fine-tuning (FFT), particularly in scenarios involving large training datasets and complex tasks such as mathematical reasoning and code generation. A plausible explanation for this performance gap is that the low-rank constraint limits the expressive capacity of LoRA [1-2].

✓ Recent studies still indicate redundancy in LoRA's parameters, reducing LoRA by 1000 times without performance loss [3-6].

! This contradiction suggests that LoRA's parameters are still not being fully utilized.

Method



$$\mathbf{W}_i + \Delta\mathbf{W}_i = \mathbf{W}_i + \underbrace{\begin{bmatrix} \tilde{\mathbf{B}}_i & \mathbf{B}_S \end{bmatrix}}_{\mathbb{R}^{b \times (r-k+Lk)}} \underbrace{\mathbf{D}_i}_{\mathbb{R}^{(r-k+Lk) \times (r-k+Lk)}} \underbrace{\begin{bmatrix} \tilde{\mathbf{A}}_i \\ \mathbf{A}_S \end{bmatrix}}_{\mathbb{R}^{(r-k+Lk) \times a}}$$

RaSA

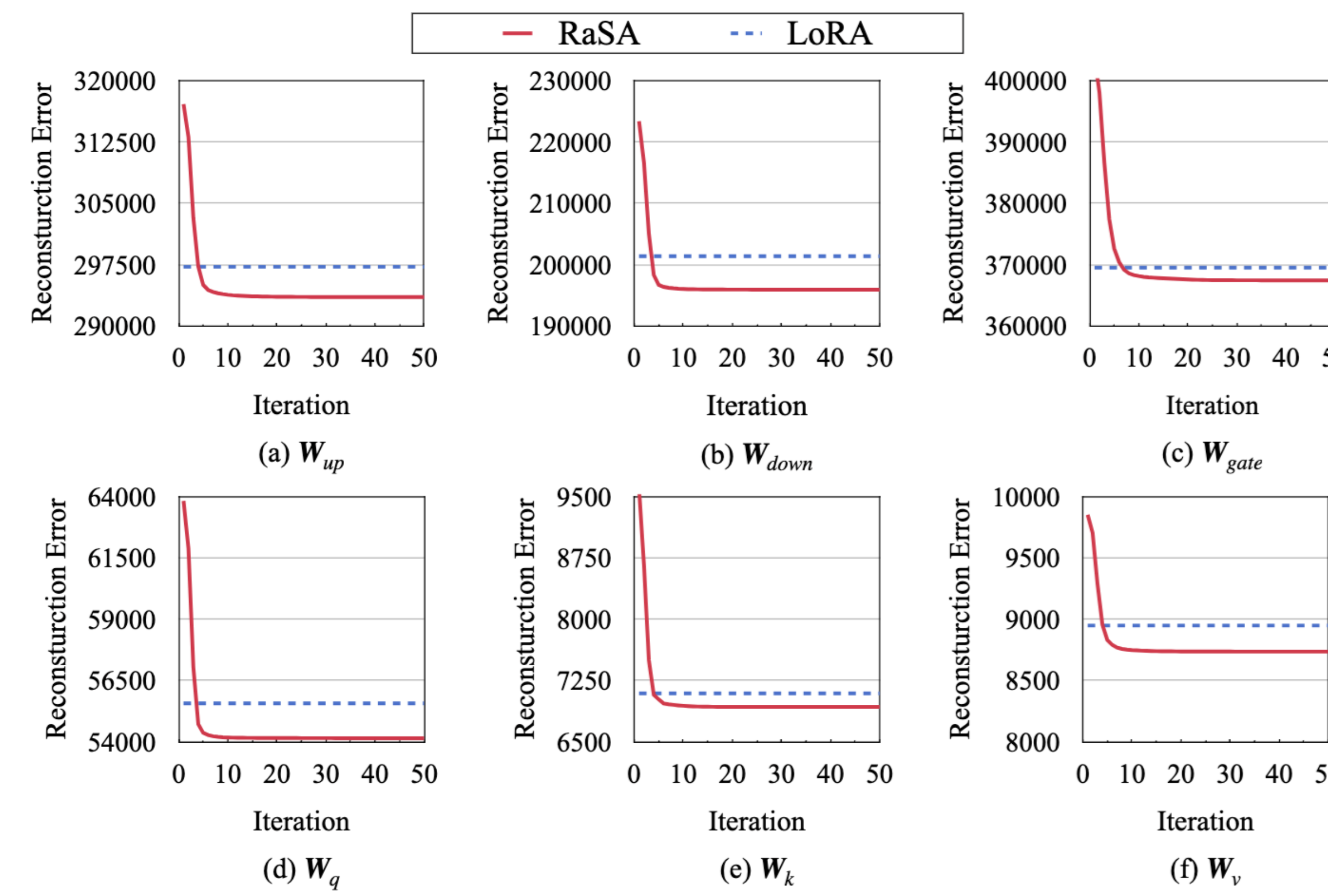
- RaSA extracts k ranks from each layer's LoRA update to form a rank pool of $L \times k$ ranks, which is shared across all layers with layer-specific weighting.

Reconstruction Error

$$e_{\text{loara}} = \min_{\mathbf{B}_i, \mathbf{A}_i} \sum_{i=1}^L \|\mathbf{M}_i - \mathbf{B}_i \mathbf{A}_i\|_F^2$$

$$e_{\text{rasa}(k)} = \min_{\tilde{\mathbf{B}}_i, \tilde{\mathbf{A}}_i, \mathbf{B}_S, \mathbf{A}_S, \mathbf{D}_i} \sum_{i=1}^L \|\mathbf{M}_i - \begin{bmatrix} \tilde{\mathbf{B}}_i & \mathbf{B}_S \end{bmatrix} \mathbf{D}_i \begin{bmatrix} \tilde{\mathbf{A}}_i \\ \mathbf{A}_S \end{bmatrix}\|_F^2$$

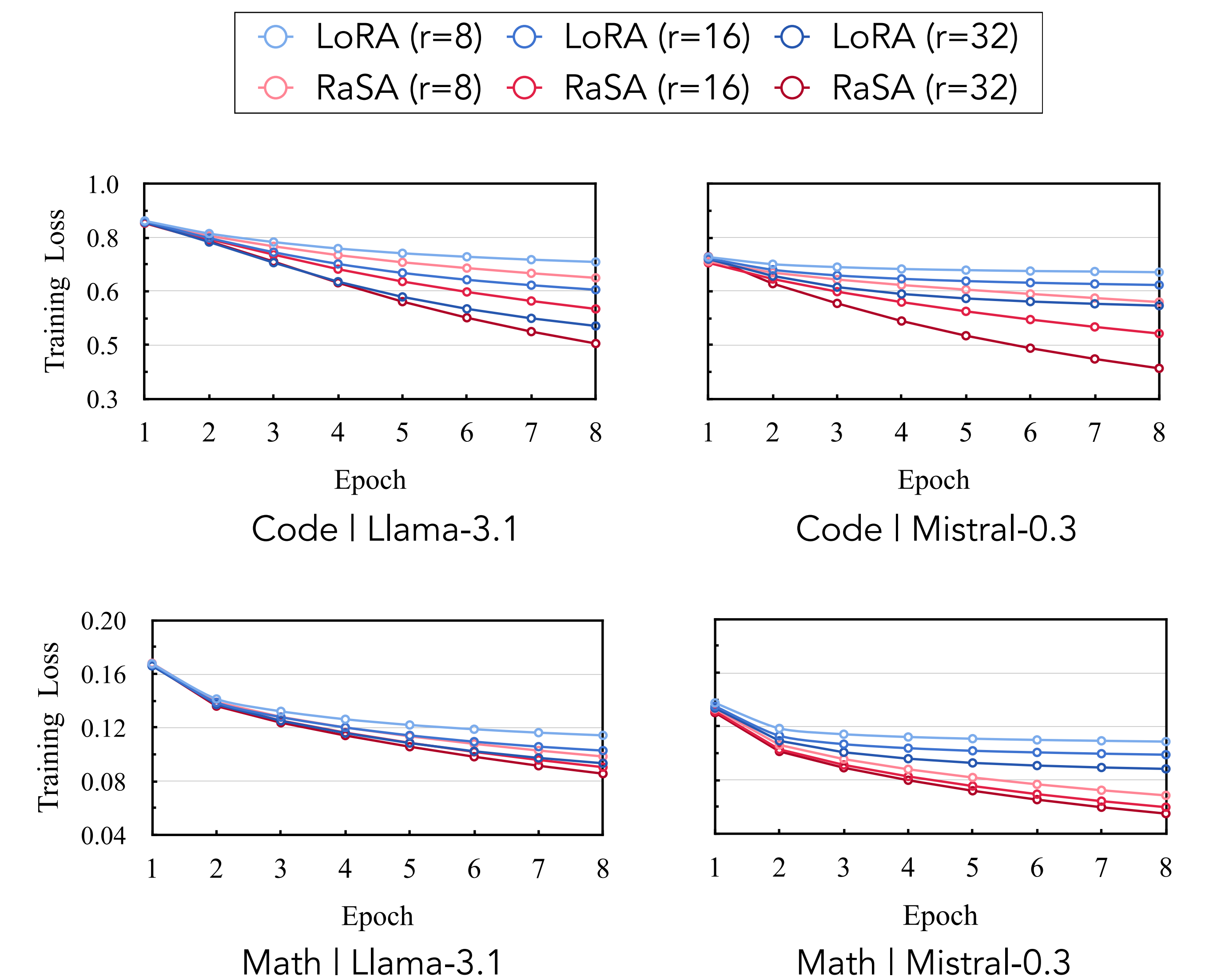
- We prove $e_{\text{rasa}(k)} \leq e_{\text{loara}}$.



Main Results



RaSA Learns More Than LoRA



[1] Mora: High-rank updating for parameter-efficient fine-tuning
[2] Lora learns less and forgets less
[3] Vera: Vector-based random matrix adaptation
[4] Tied-LoRA: Enhancing parameter efficiency of LoRA with weight tying
[5] VB-LoRA: extreme parameter efficient fine-tuning with vector banks