

# RaSA: Rank-Sharing Low-Rank Adaptation

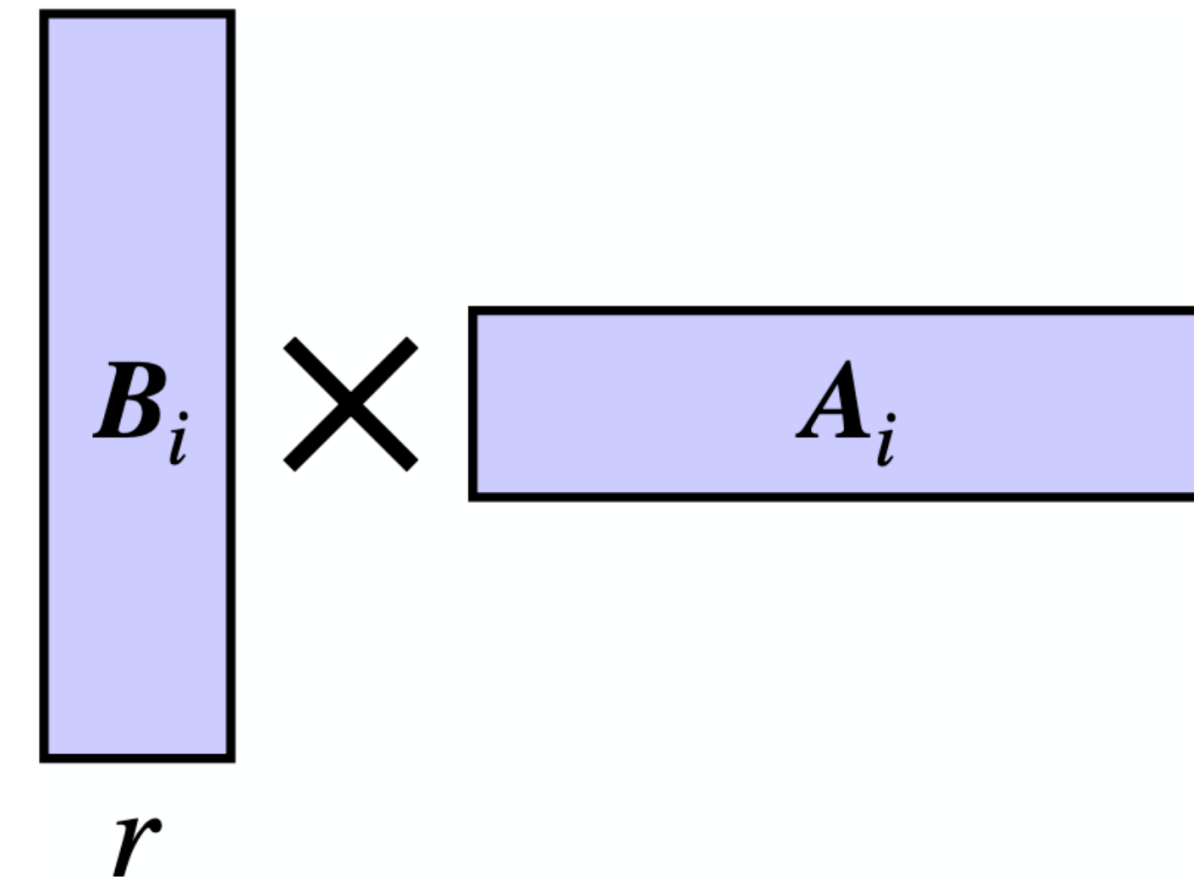
Zhiwei He   Zhaopeng Tu   Xing Wang   Xingyu Chen   Zhijie Wang  
Jiahao Xu   Tian Liang   Wenxiang Jiao   Zhuosheng Zhang   Rui Wang

# Low-Rank Adaptation (LoRA)

Low-rank update

- Parameter Update:

$$\mathbf{W}_i + \underbrace{\Delta \mathbf{W}_i}_{\text{low-rank}} =$$

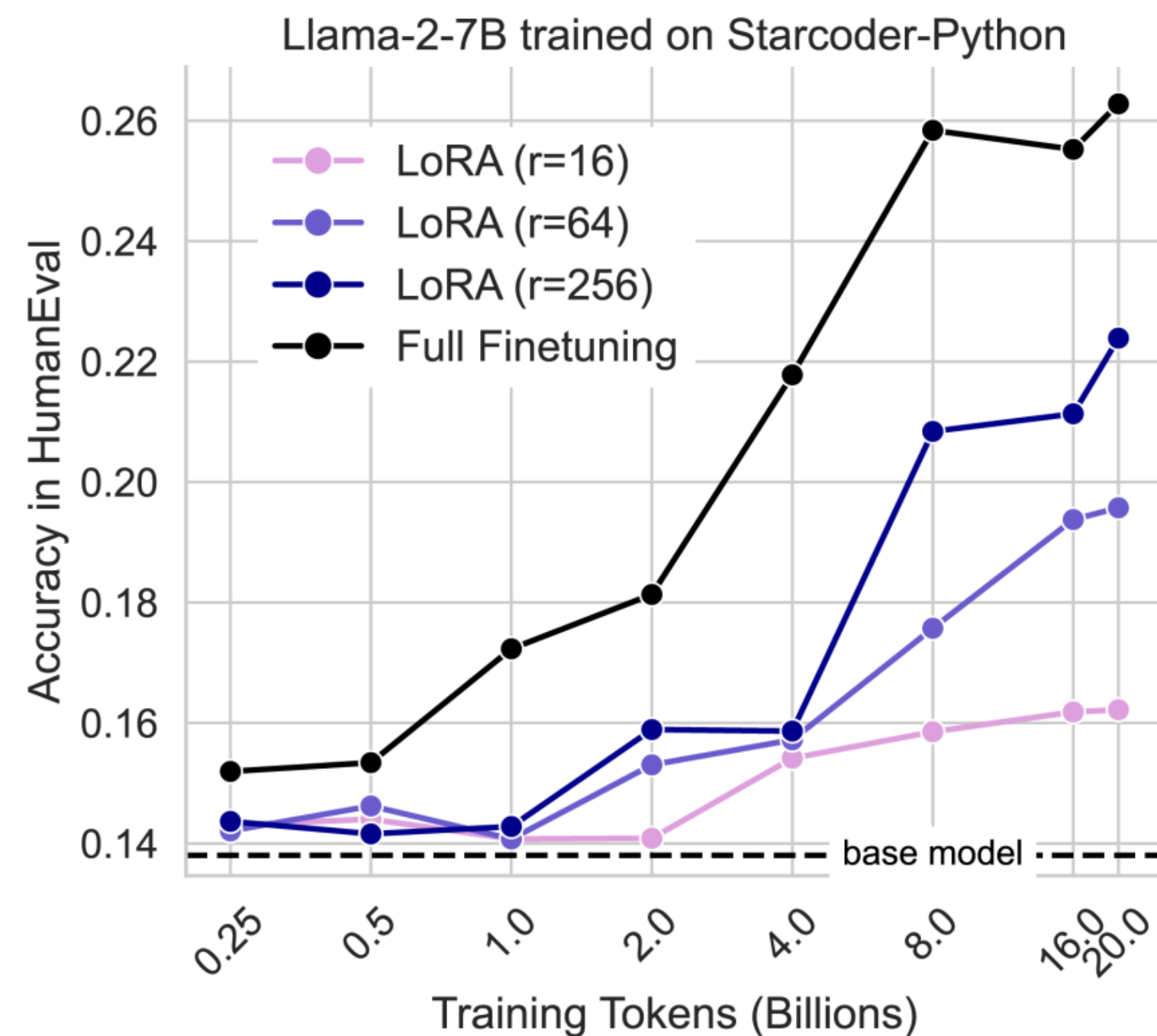


$$\mathbf{B}_i \in \mathbb{R}^{b \times r}, \mathbf{A}_i \in \mathbb{R}^{r \times a}$$

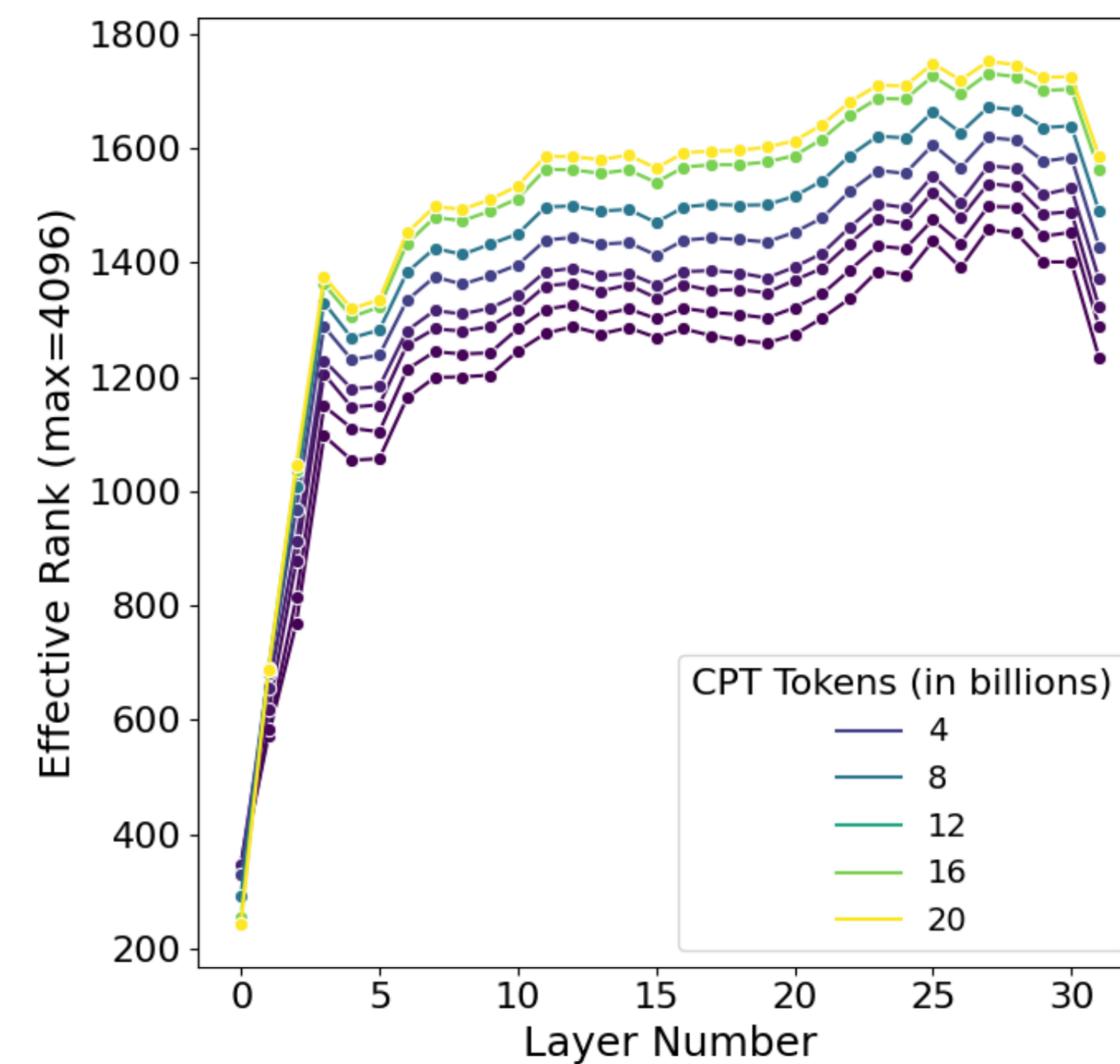
$$r \ll \min(b, a)$$

# Low-Rank Adaptation (LoRA)

Low-rank constraint limits the expressive capacity of LoRA



LoRA still lags behind full fine-tuning (FFT)

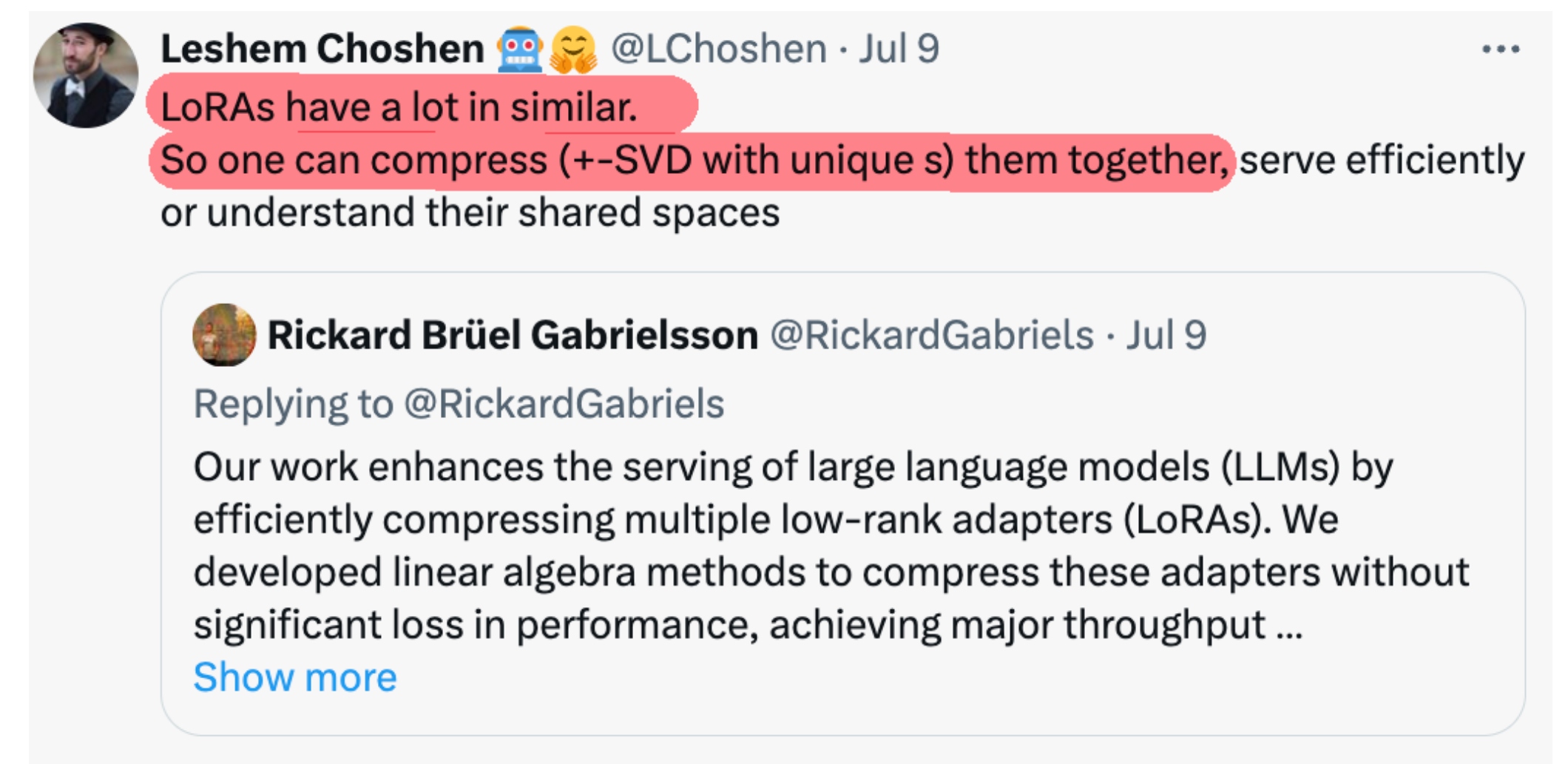


Effective rank of FFT is large

# Low-Rank Adaptation (LoRA)

LoRA can be further compressed

- Many related works demonstrate that LoRA can be further compressed by
  - [2] 7.76%
  - [3] 12.5%
  - [4] 50%
  - [5] 3%
- without performance loss.



Compressing 100 LoRAs by sharing their subspaces won't compromise performance.

[2] Kopiczko, Dawid J., Tijmen Blankevoort, and Yuki M. Asano. "Vera: Vector-based random matrix adaptation." ICLR 2023

[3] Renduchintala, Adithya, Tugrul Konuk, and Oleksii Kuchaiev. "Tied-LoRA: Enhancing parameter efficiency of LoRA with weight tying." NAACL 2024

[4] Song, Yurun, et al. "ShareLoRA: Parameter Efficient and Robust Large Language Model Fine-tuning via Shared Low-Rank Adaptation." arXiv 2024

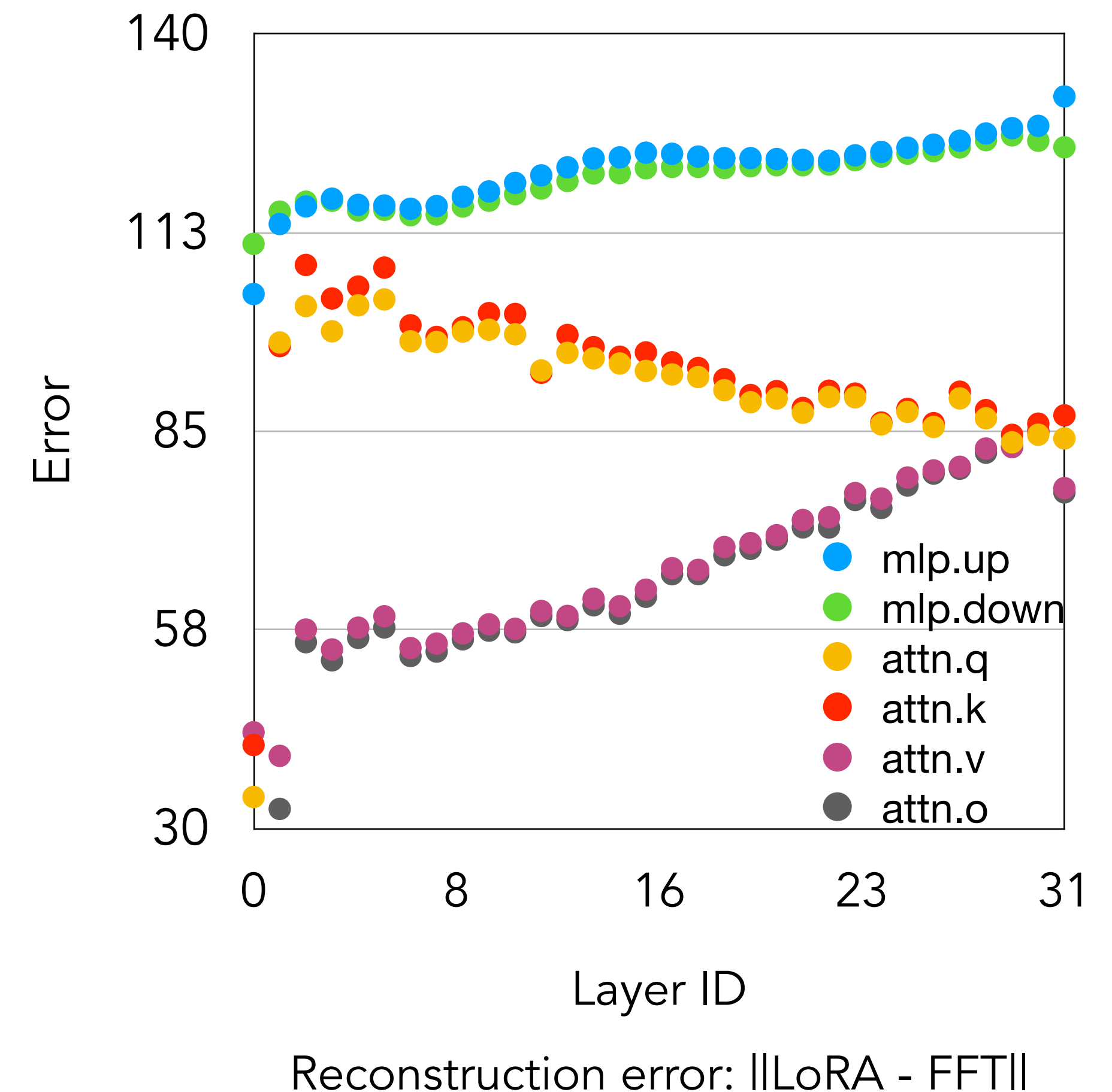
[5] Li, Yang, Shaobo Han, and Shihao Ji. "VB-LoRA: extreme parameter efficient fine-tuning with vector banks." NeurIPS 2024

LoRA is underutilized.

# Low-Rank Adaptation (LoRA)

Why LoRA is under-utilized?

- Different components and layers require different levels of expressive capability.
- LoRA adopts an average allocation strategy.





# Rank-Sharing Low-Rank Adaptation (RaSA)

Partial rank-sharing across layers

- Split the matrices  $\mathbf{B}_i$  and  $\mathbf{A}_i$  into layer-specific parts  $(\tilde{\mathbf{B}}_i, \tilde{\mathbf{A}}_i)$  and layer-shared parts  $(\hat{\mathbf{B}}_i, \hat{\mathbf{A}}_i)$

$$\mathbf{B}_i = \left[ \underbrace{\tilde{\mathbf{B}}_i}_{\mathbb{R}^{b \times (r-k)}} \quad \underbrace{\hat{\mathbf{B}}_i}_{\mathbb{R}^{b \times k}} \right], \quad \mathbf{A}_i = \left[ \underbrace{\tilde{\mathbf{A}}_i^T}_{\mathbb{R}^{a \times (r-k)}} \quad \underbrace{\hat{\mathbf{A}}_i^T}_{\mathbb{R}^{a \times k}} \right]^T$$

- Concatenate all layer-shared parts across layers to form shared rank pools  $(\mathbf{B}_S$  and  $\mathbf{A}_S)$

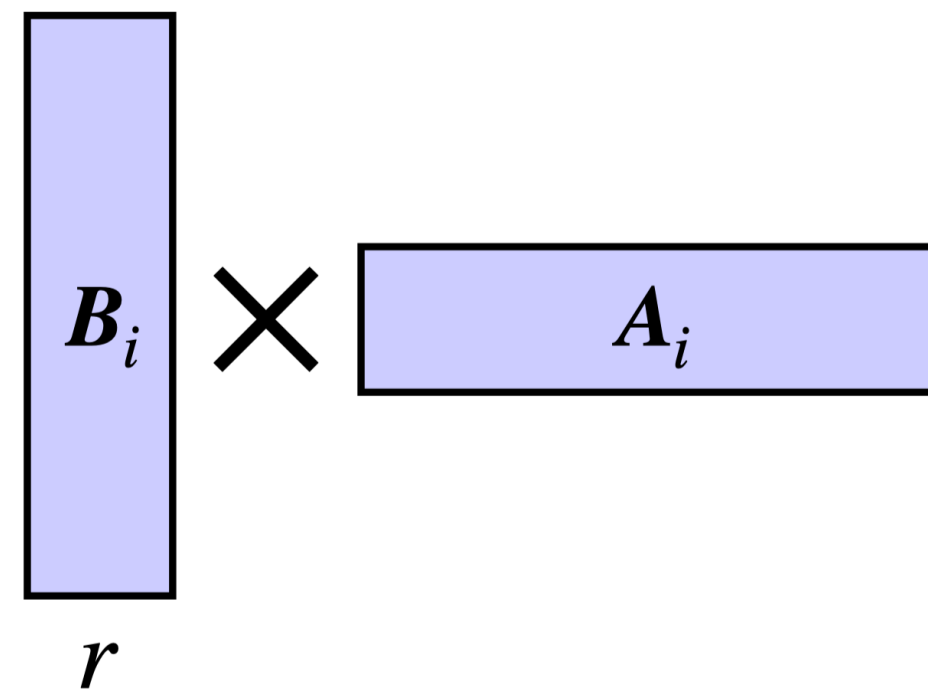
$$\mathbf{B}_S = [\hat{\mathbf{B}}_1 \quad \hat{\mathbf{B}}_2 \quad \dots \quad \hat{\mathbf{B}}_L] \in \mathbb{R}^{b \times (L \times k)}, \quad \mathbf{A}_S = [\hat{\mathbf{A}}_1^T \quad \hat{\mathbf{A}}_2^T \quad \dots \quad \hat{\mathbf{A}}_L^T]^T \in \mathbb{R}^{(L \times k) \times a}$$

- Update of layer-i

$$\begin{aligned} \mathbf{W}_i + \Delta \mathbf{W}_i &= \mathbf{W}_i + \frac{\alpha}{r} (\tilde{\mathbf{B}}_i \tilde{\mathbf{A}}_i + \mathbf{B}_S \mathbf{A}_S) \\ &= \mathbf{W}_i + [\tilde{\mathbf{B}}_i \quad \mathbf{B}_S] \text{diag}\left(\frac{\alpha}{r}\right) \begin{bmatrix} \tilde{\mathbf{A}}_i \\ \mathbf{A}_S \end{bmatrix} \end{aligned}$$

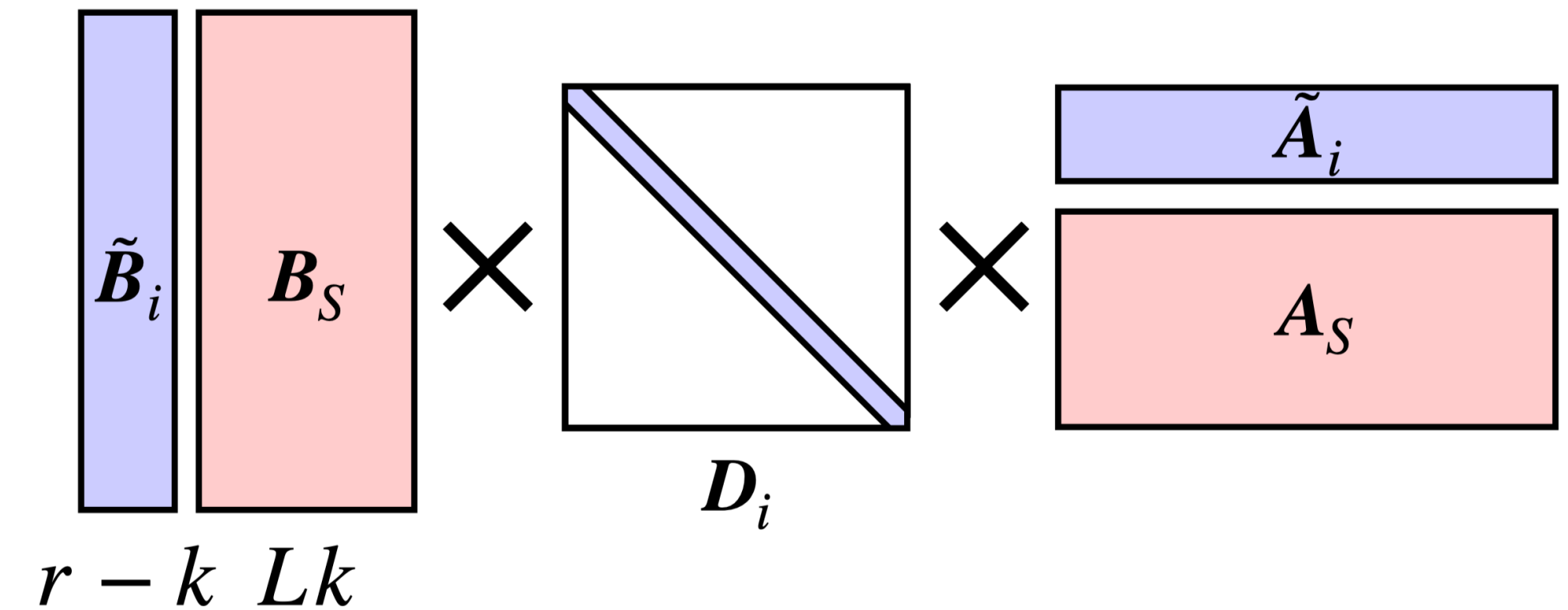
# Rank-Sharing Low-Rank Adaptation (RaSA)

Comparison between LoRA and RaSA



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

LoRA



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

RaSA

- $r \Rightarrow r-k+Lk$
- extra parameters (0.01%)



# Reconstruction Error Analysis

## Minimum Reconstruction Error

- We compare their abilities to reconstruct a set of high-rank matrices  $\{\mathbf{M}_i\}_{i \in [L]}$ ,  $\text{rank}(\mathbf{M}_i) = R > r$

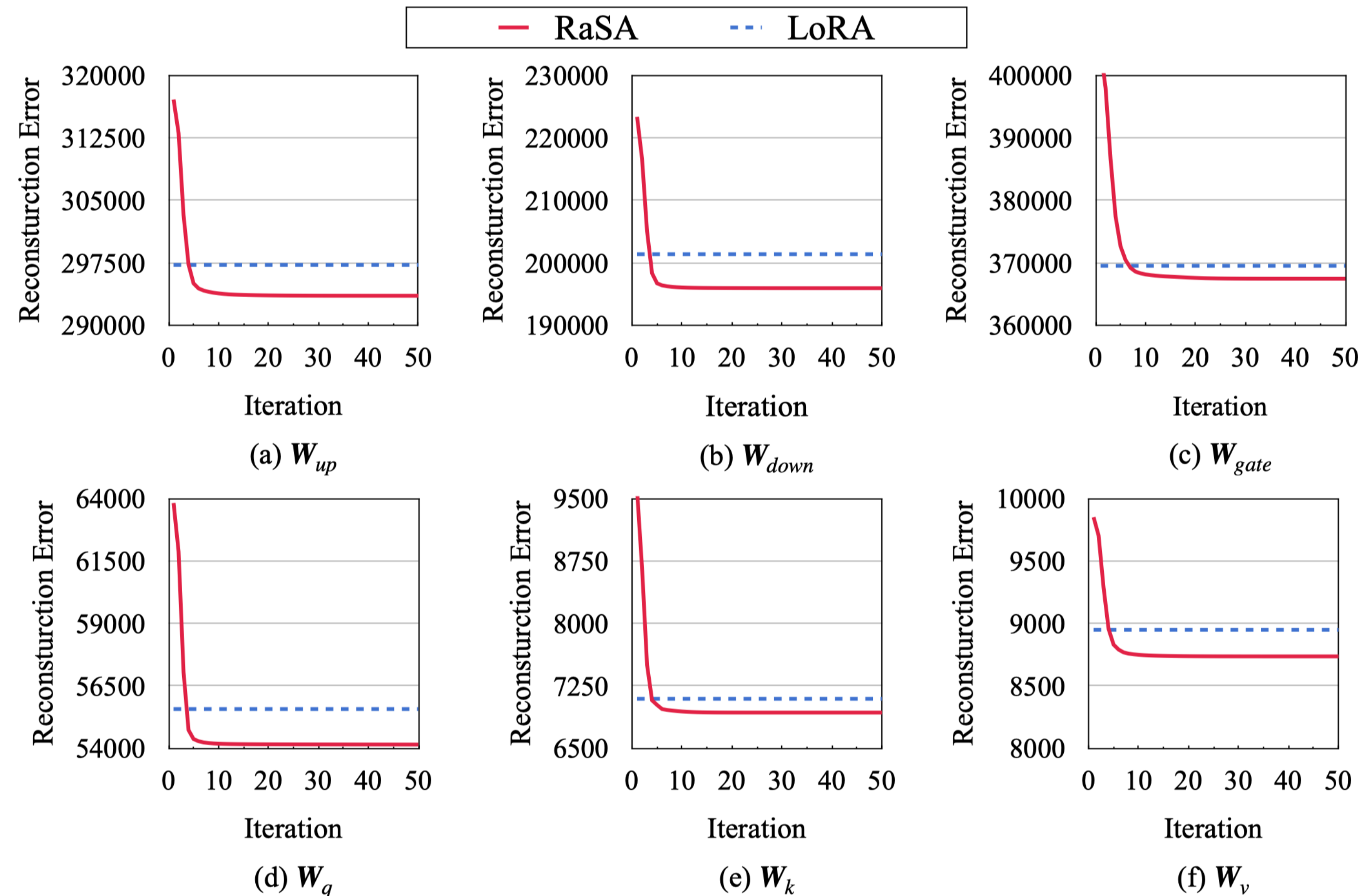
$$e_{\text{lora}} = \min_{\mathbf{B}_i, \mathbf{A}_i} \sum_{i=1}^L \|\mathbf{M}_i - \mathbf{B}_i \mathbf{A}_i\|_F^2$$

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

- We prove  $e_{\text{rasa}}(k) \leq e_{\text{lora}}$  (Theorem 3.1)

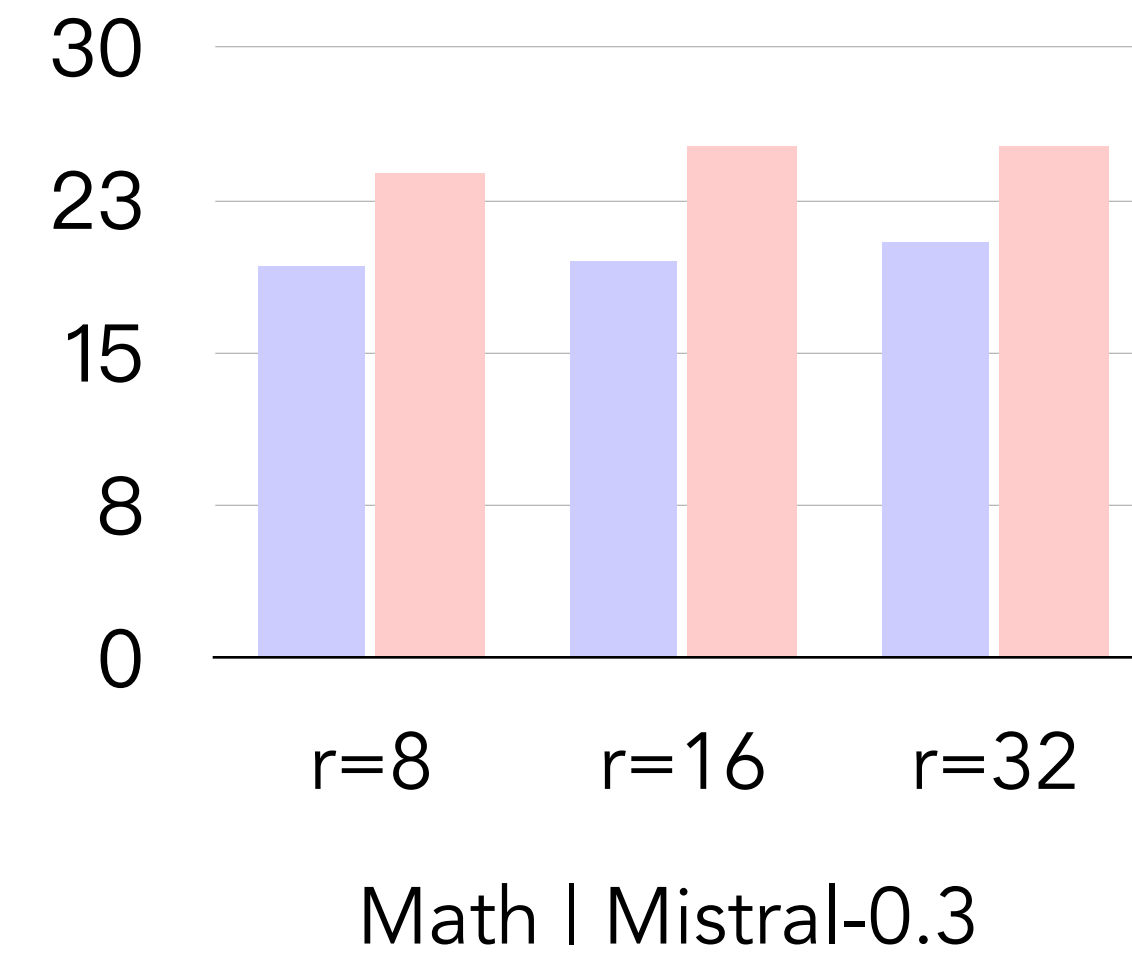
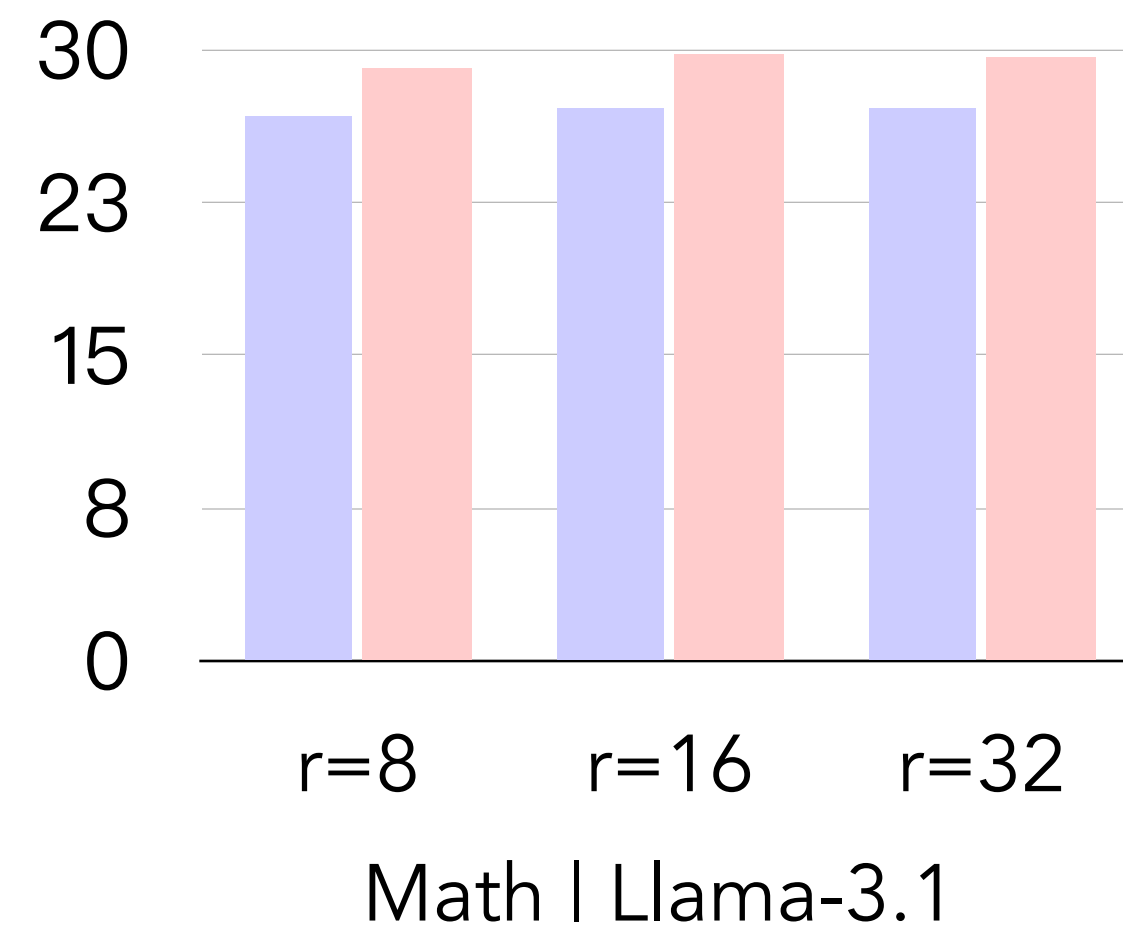
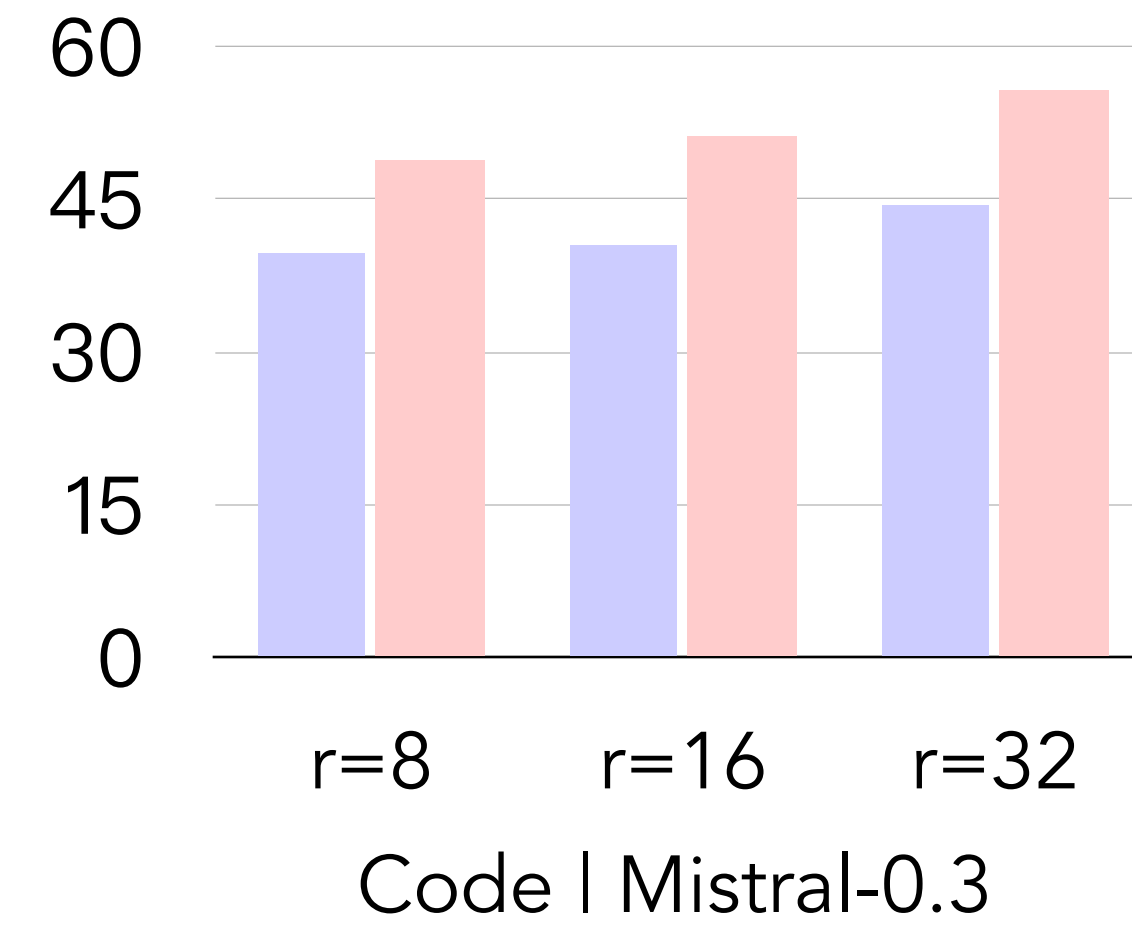
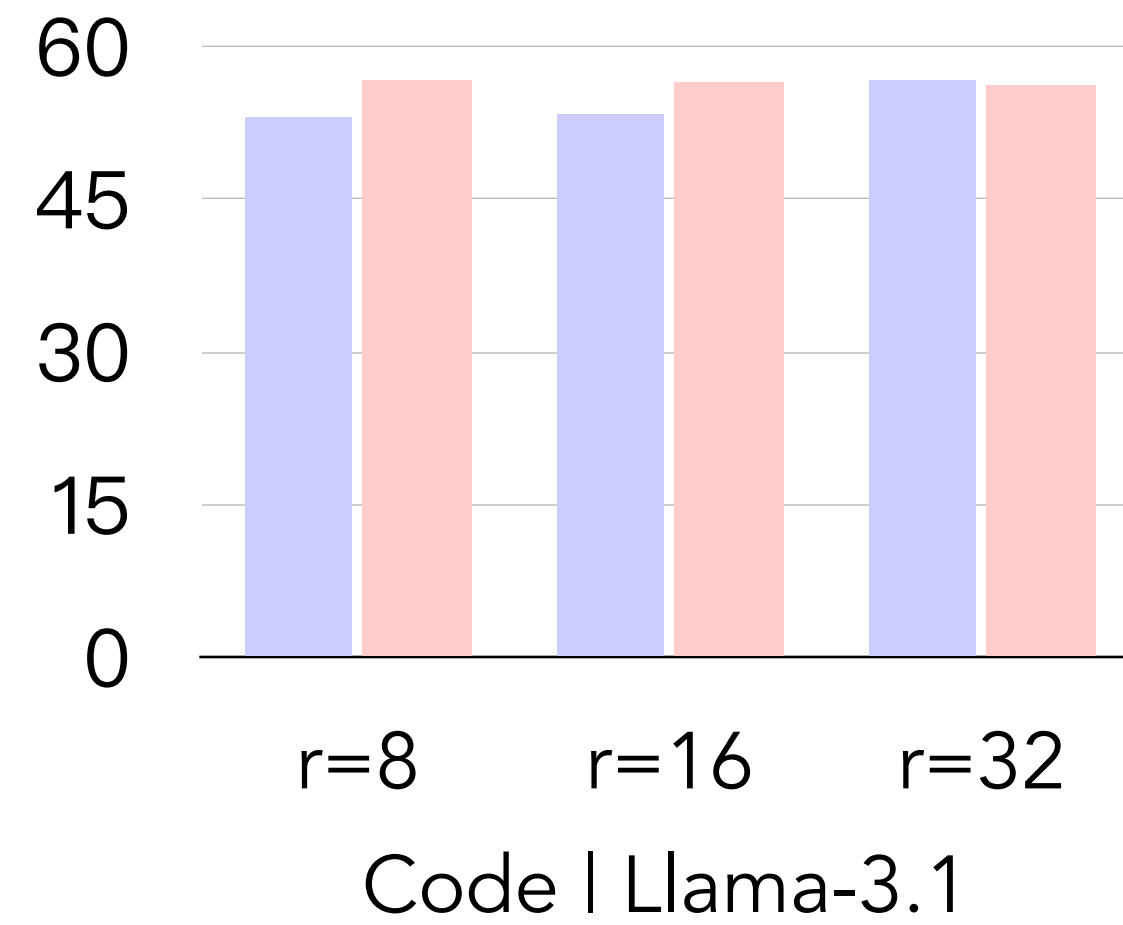
# Reconstruction Error Analysis

## Empirical Analysis

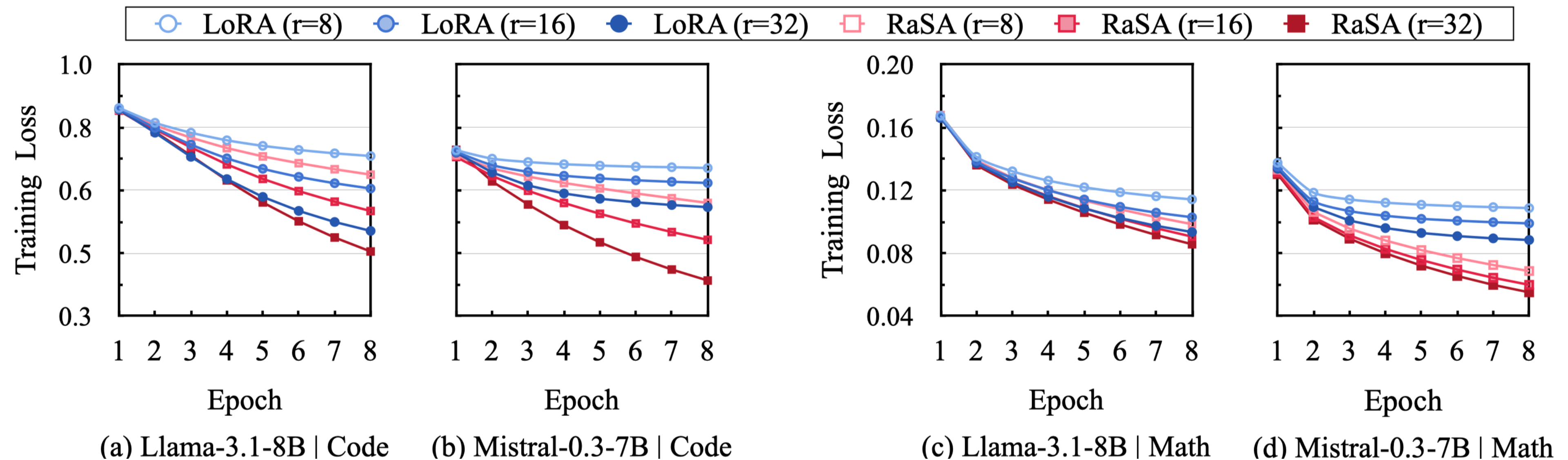


RaSA requires ~10 iterations to achieve a significantly lower reconstruction error than LoRA's minimum.

# Main Results

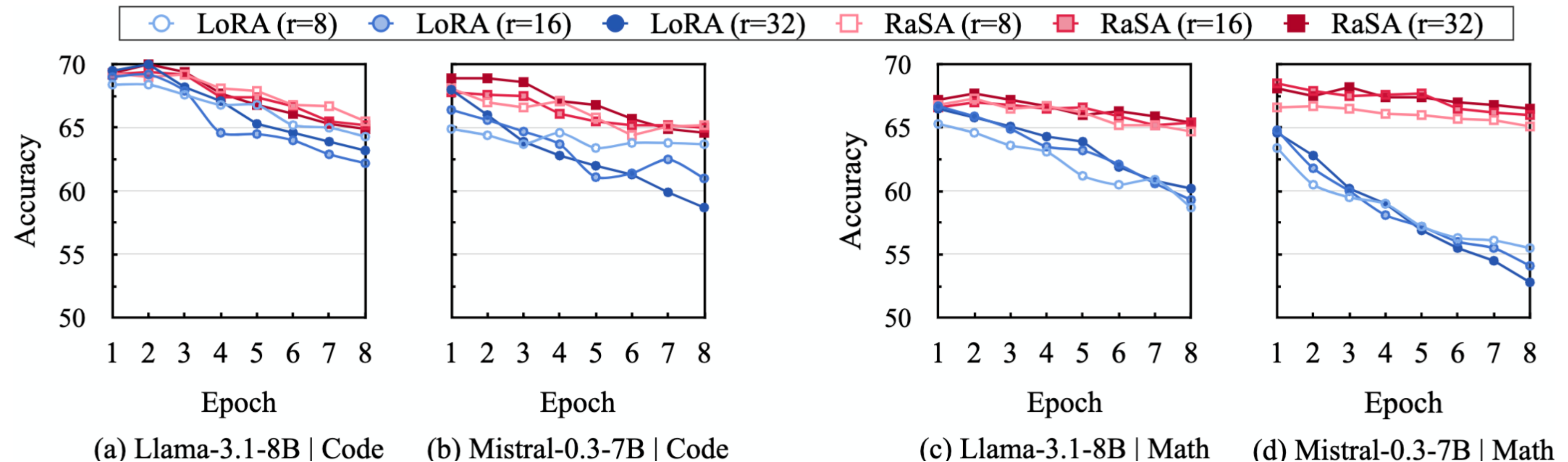


# RaSA learns more and faster than LoRA



Training curves of LoRA and RaSA with different ranks.

# RaSA forgets less than LoRA



Y-axis shows the average of prediction accuracy on three benchmarks to evaluate model's forgetting. Higher prediction accuracy denotes less forgetting.

# Summary

- We propose RaSA, an extension of LoRA by allowing partial rank sharing across layers, which significantly improves the efficiency and expressiveness.
- We provide a comprehensive analysis – both theoretical and empirical – showcasing RaSA's superior capacity for matrix reconstruction and its resultant improved performance on downstream tasks.



Thank You