



FLoRIST: Singular Value Thresholding for Efficient and Accurate Federated Fine-Tuning of Large Language Models

Hariharan Ramesh

Jyotikrishna Dass



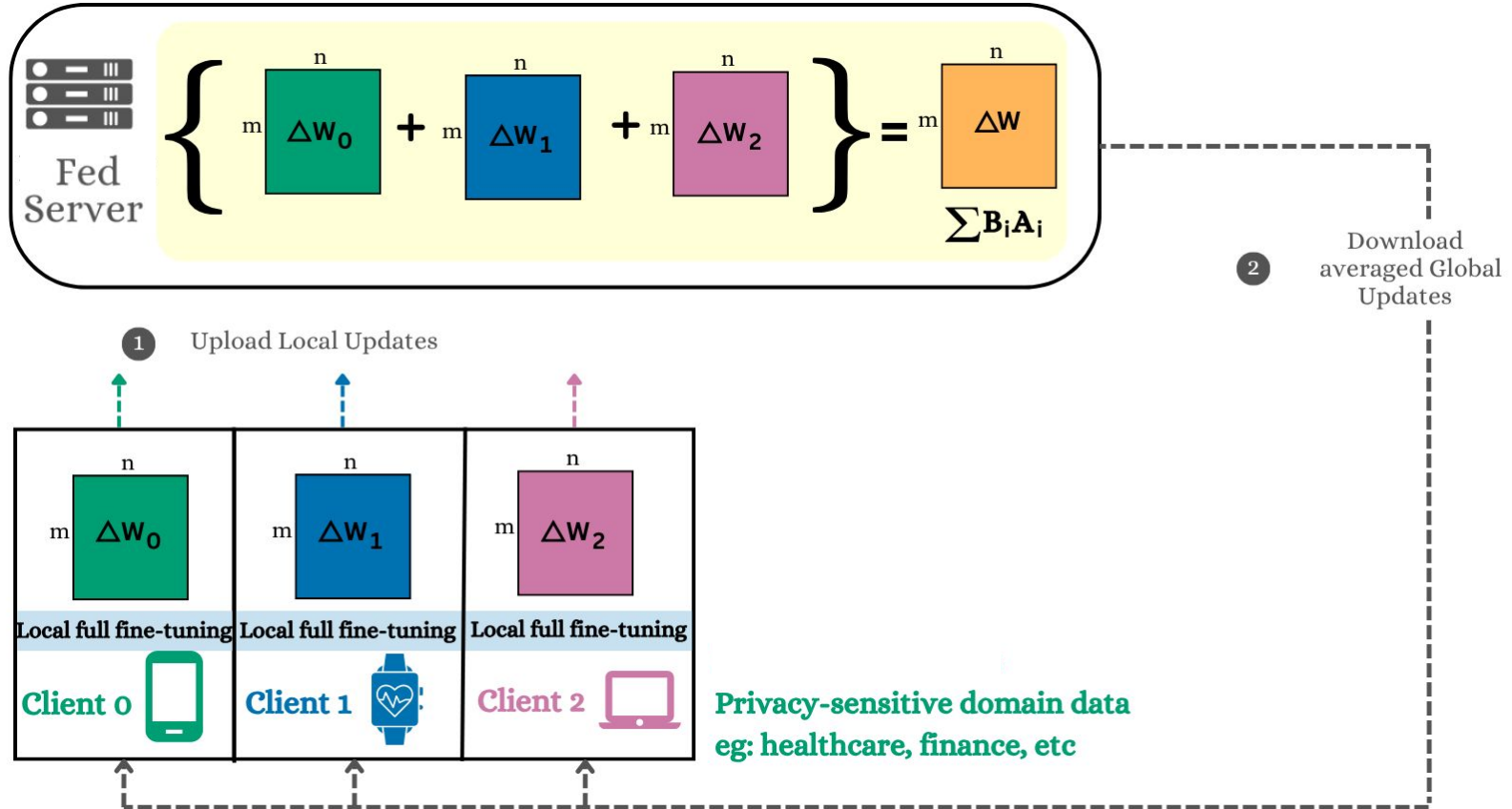
Paper

Distributed AI & Smart Systems (DASS) Lab
Electrical and Computer Engineering
University of Arizona



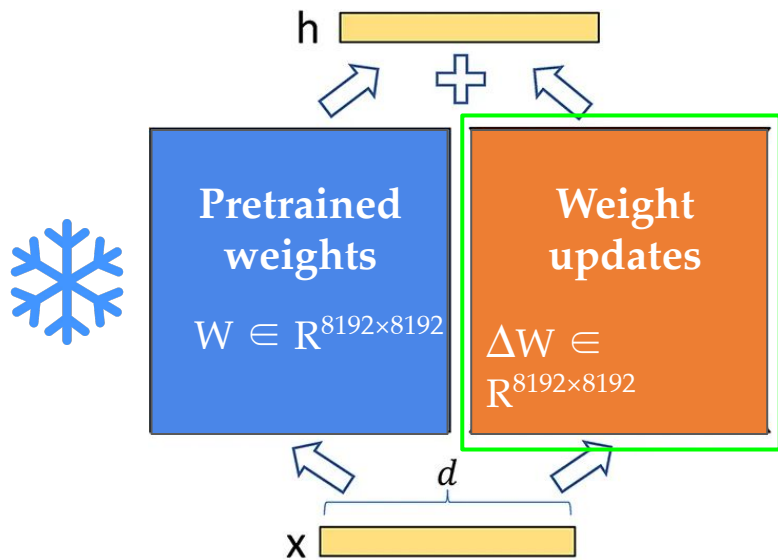
Code

Background: Federated Full Fine-Tuning



Background: Low-Rank Adaptation (LoRA)

Parameter Efficient Fine-Tuning



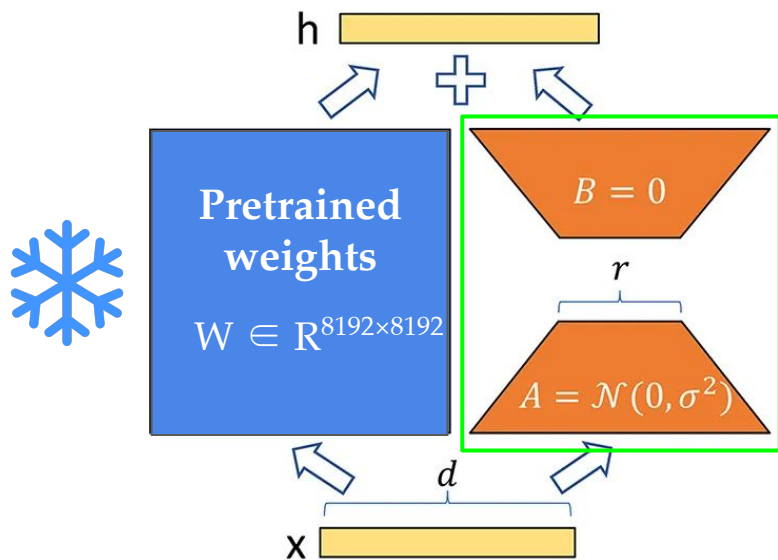
LLaMA-3.2-1B, $W \in \mathbb{R}^{8192 \times 8192}$

LoRA fine-tuning (rank, $r = 16$) creates efficient adapters, $\{B \in \mathbb{R}^{8192 \times 16}, A \in \mathbb{R}^{16 \times 8192}\}$

256x reduction in trainable parameters

Background: Low-Rank Adaptation (LoRA)

Parameter Efficient Fine-Tuning

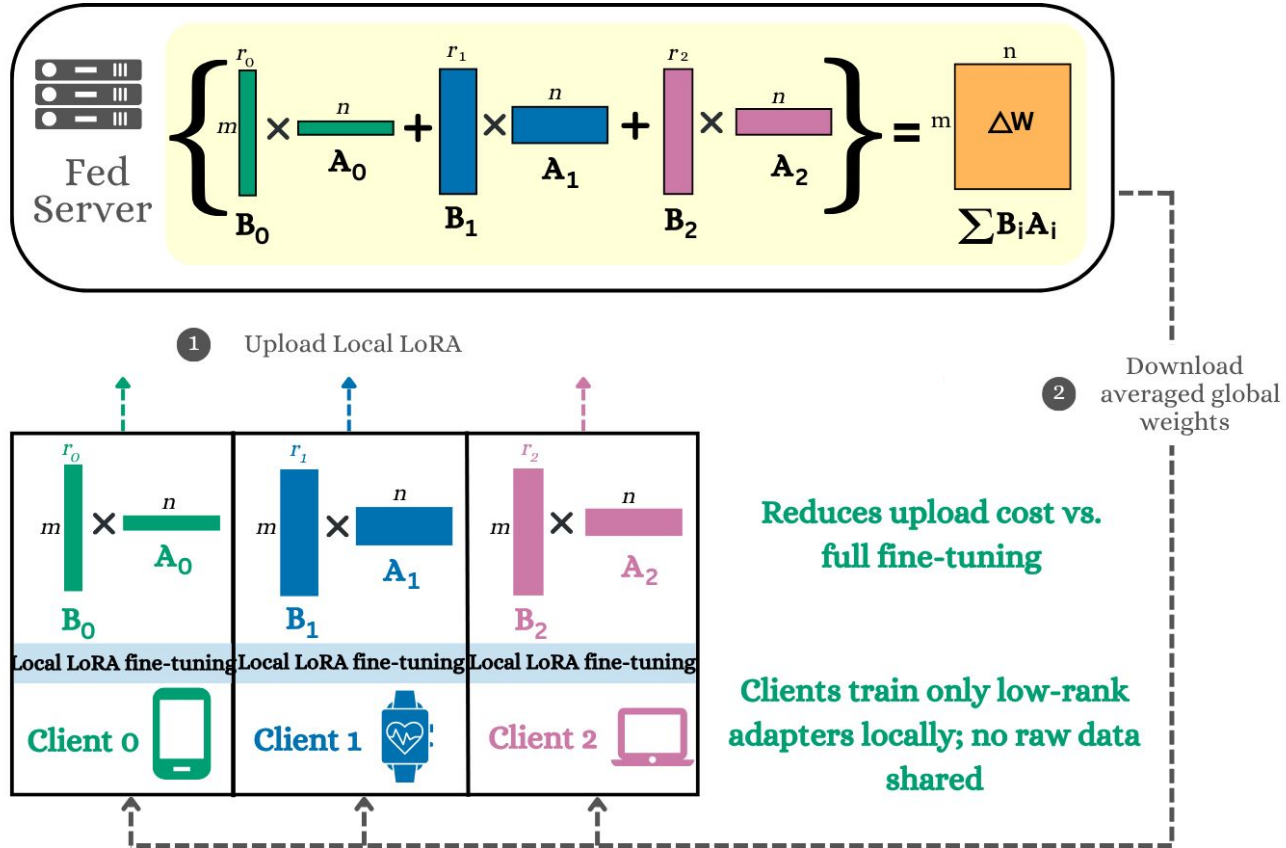


LLaMA-3.2-1B, $W \in \mathbb{R}^{8192 \times 8192}$

LoRA fine-tuning (rank, $r = 16$) creates efficient adapters, $\{B \in \mathbb{R}^{8192 \times 16}, A \in \mathbb{R}^{16 \times 8192}\}$

256x reduction in trainable parameters

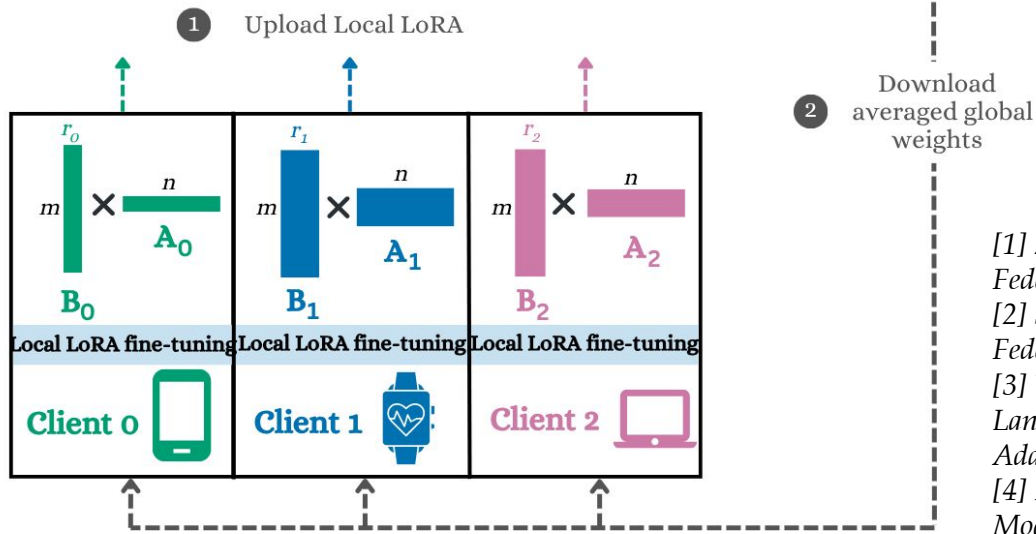
Background: Federated LoRA Fine-Tuning



Related Work



Categorize based on
**server aggregation of
Local LoRA adapters**



[1] Zhang et al., "Towards Building the FederatedGPT: Federated Instruction Tuning," ICASSP 2024.

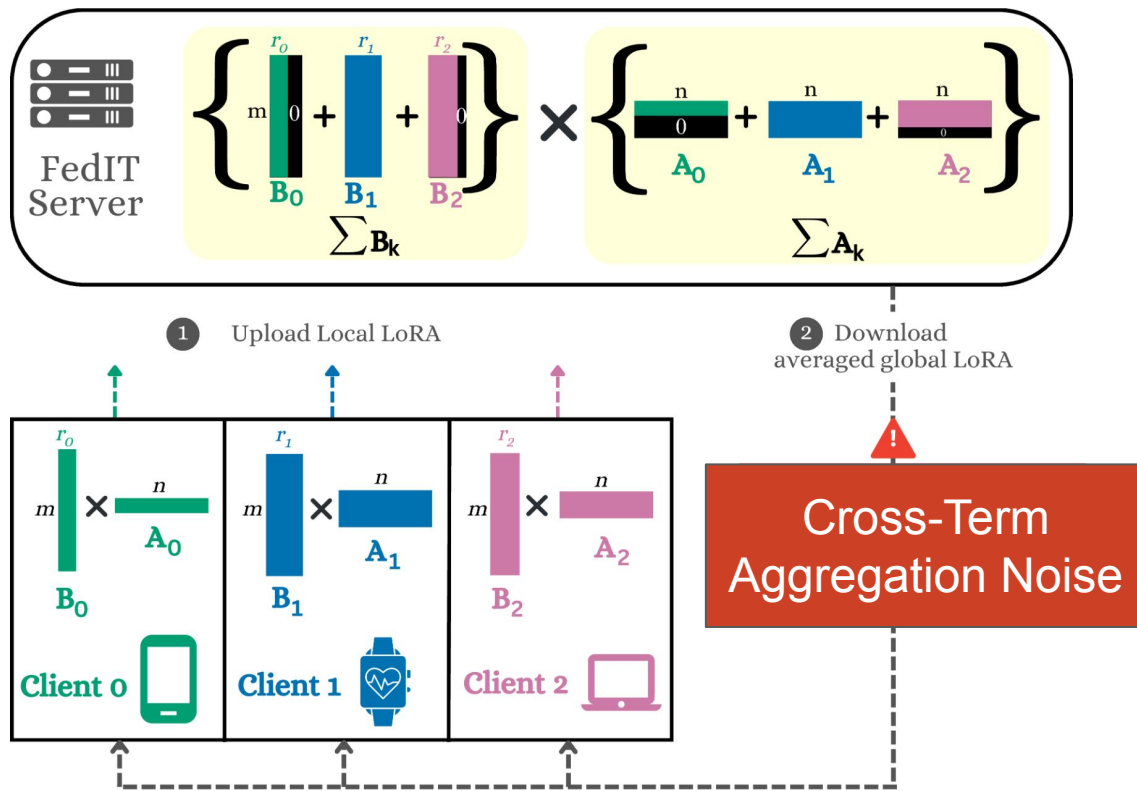
[2] Sun et al., "Improving LoRA in Privacy-Preserving Federated Learning," ICLR 2024.

[3] Wang et al., "FLoRA: Federated Fine-Tuning Large Language Models with Heterogeneous Low-Rank Adaptations," NeurIPS 2024.

[4] Bai et al., "Federated Fine-Tuning of Large Language Models under Heterogeneous Tasks and Client Resources," NeurIPS 2024.

Related Work: FedIT (Average)

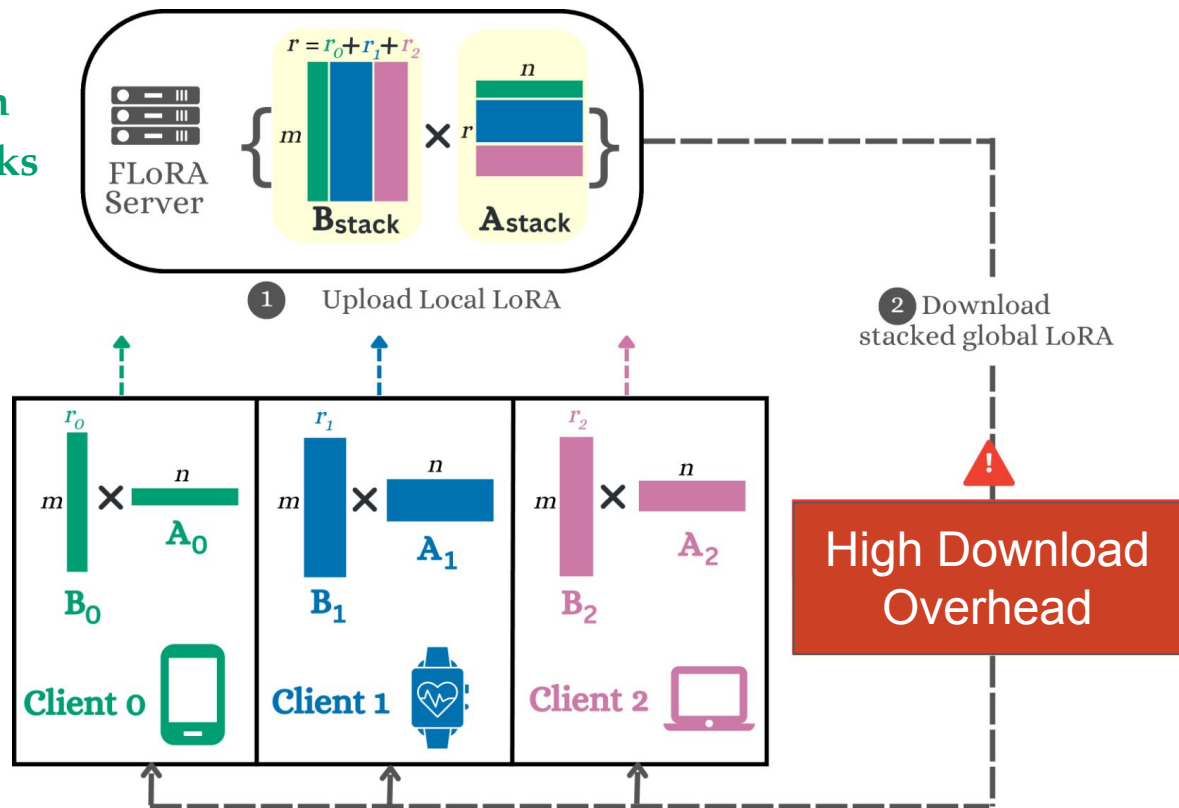
Standard FedAvg
applied directly to
LoRA adapters



Related Work: FLoRA (Stack)

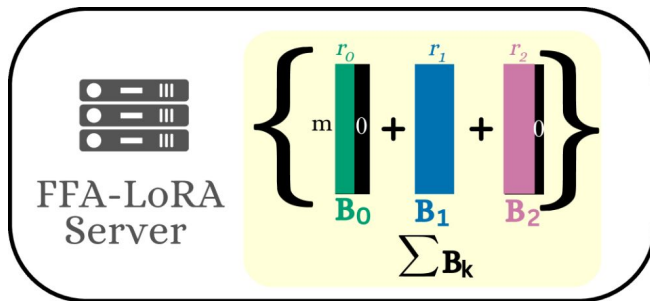
Mathematically exact aggregation

Supports heterogeneous client ranks



Related Work: FFA-LoRA (Freeze)

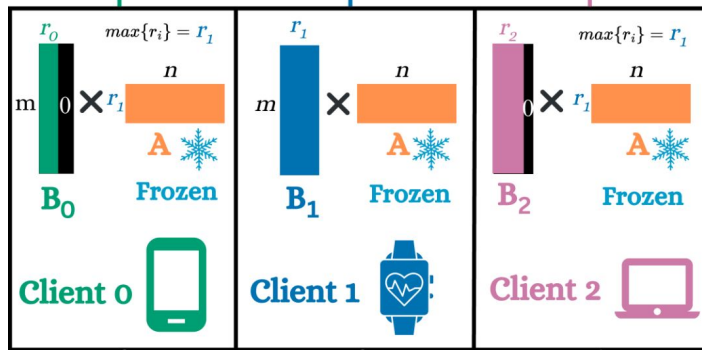
Mathematically exact aggregation
(freezing A_k eliminates cross-terms)



Half the upload/download cost ①

Upload half-local LoRA, B_k

② Download averaged half-global LoRA



Slower Convergence
No Heterogeneity

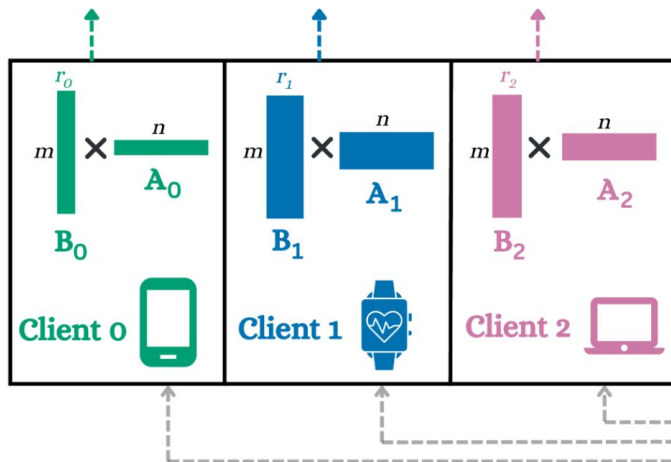
Related Work: FlexLoRA (Reconstruct)

Large Server Computational Cost

Exact aggregation
Heterogeneous ranks

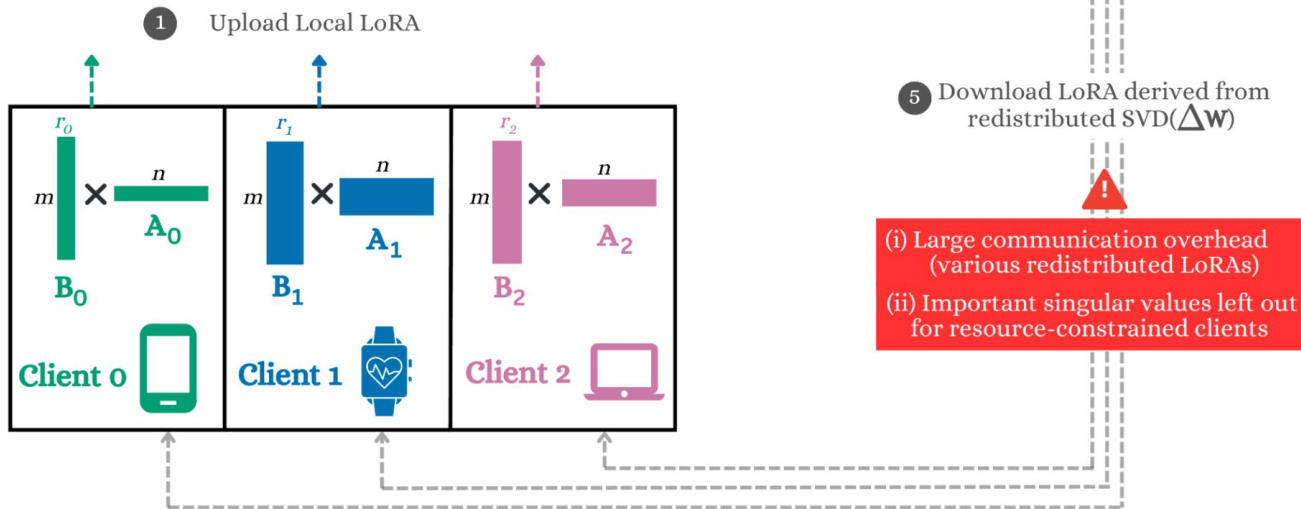
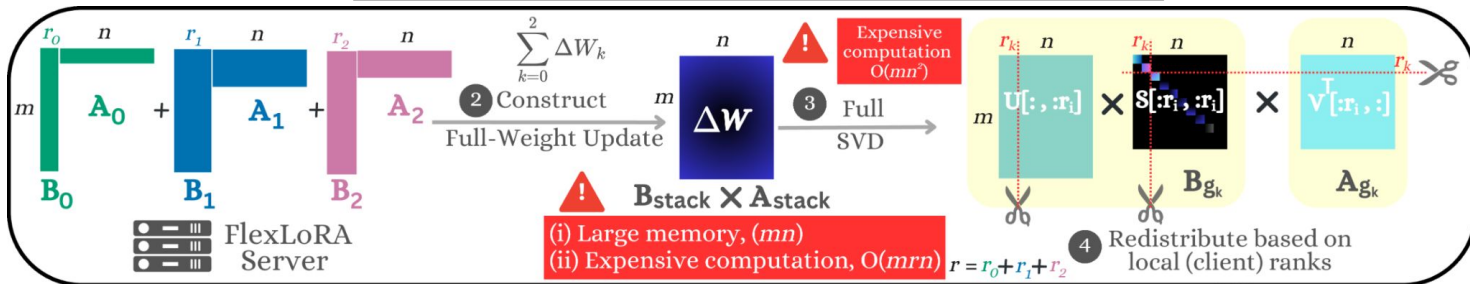
$$\left\{ \begin{matrix} r_0 \\ m \end{matrix} \times \begin{matrix} n \\ A_0 \end{matrix} + \begin{matrix} r_1 \\ m \end{matrix} \times \begin{matrix} n \\ A_1 \end{matrix} + \begin{matrix} r_2 \\ m \end{matrix} \times \begin{matrix} n \\ A_2 \end{matrix} \right\} = \begin{matrix} n \\ m \end{matrix} \begin{matrix} \Delta W \\ \sum B_i A_i \end{matrix}$$

1 Upload Local LoRA



Related Work: FlexLoRA (Reconstruct)

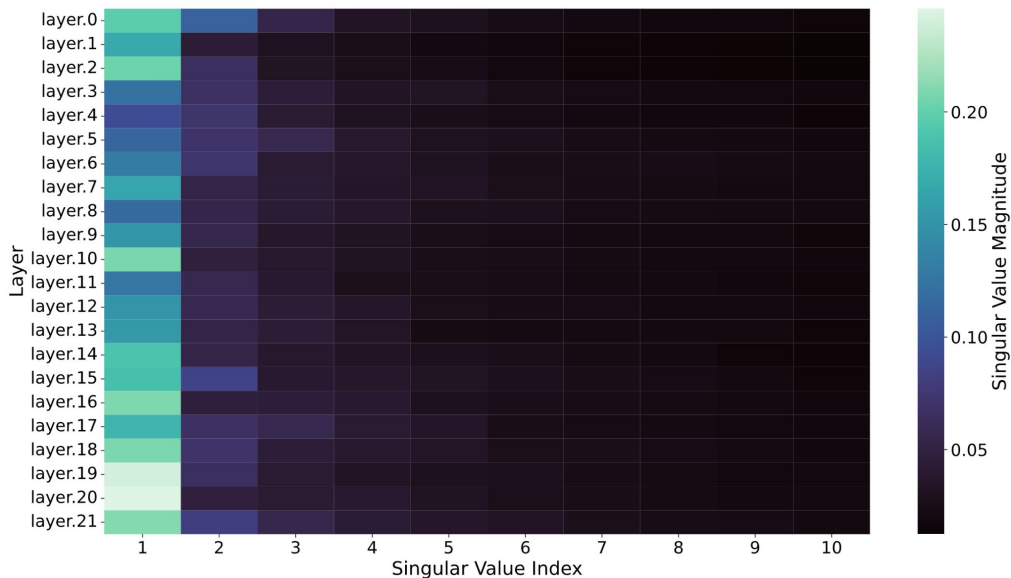
Large Server Computational Cost



Challenges in Existing Methods

Method	Aggregation	Heterogeneity	Key Challenge
FedIT <i>(Average)</i>	✗ Noisy	✗ Homo only	Cross-term Noise
FFA-LoRA <i>(Freeze)</i>	✓ Exact	✗ Homo only	Slower Convergence
FLoRA <i>(Stack)</i>	✓ Exact	✓ Native	High Download Overhead
FlexLoRA <i>(Reconstruct)</i>	✓ Exact	✓ Native	Large Server Computational Cost (ΔW)

What is the intrinsic dimensionality of the Global Adapters?



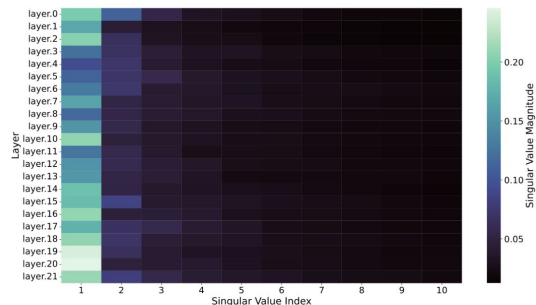
Singular value spectrum (q_{proj}), TinyLlama (Wizard dataset)

Our Motivating Observation

Singular values of aggregated ΔW decay rapidly

Only 6–10 components are significant (max client rank = 64)

Method	Aggregation	Heterogeneity	Key Challenge
FedIT <i>(Average)</i>	✗ Noisy	✗ Homo only	Cross-term Noise
FFA-LoRA <i>(Freeze)</i>	✓ Exact	✗ Homo only	Slower Convergence
FLoRA <i>(Stack)</i>	✓ Exact	✓ Native	High Download Overhead
FlexLoRA <i>(Reconstruct)</i>	✓ Exact	✓ Native	Large Server Computational Cost (ΔW)



Singular value spectrum (q_proj), TinyLlama (Wizard dataset)

Our Motivating Observation

Singular values of aggregated ΔW decay rapidly

Only 6–10 components are significant; max rank = 64

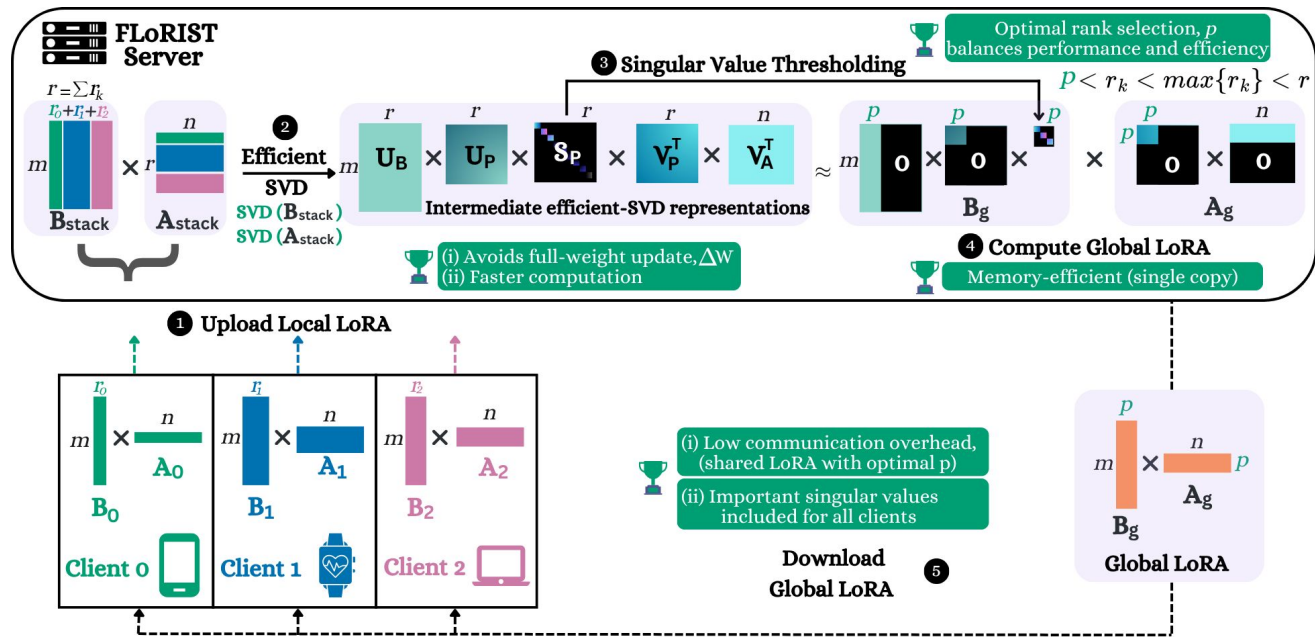
How to identify the most informative components of global aggregated update (ΔW) without reconstructing it for faster aggregation & efficient download ?

Proposed Method

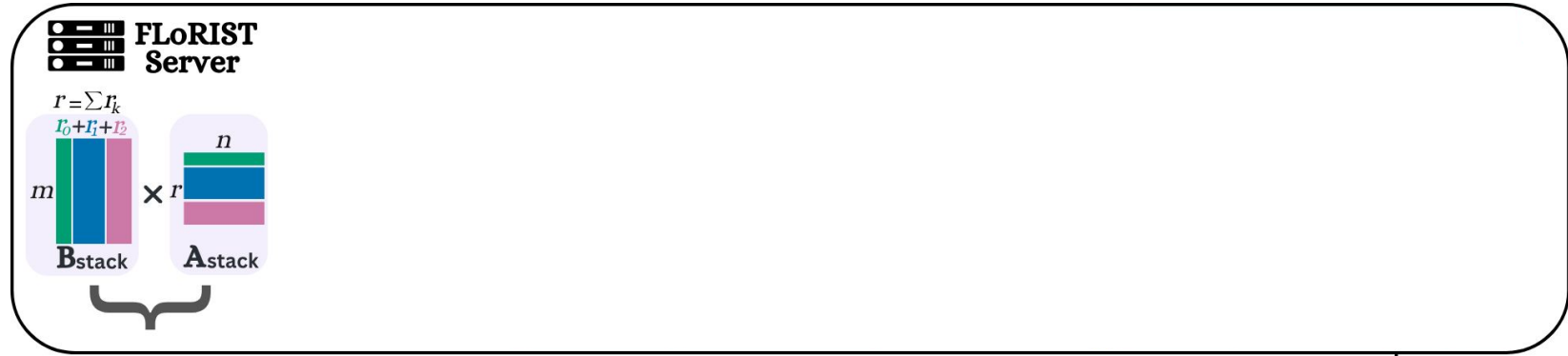
FLoRIST: Singular Value Thresholding for Efficient and Accurate Federated Fine-Tuning of Large Language Models

Avoids ΔW reconstruction;
Lowers server compute cost
for faster aggregation

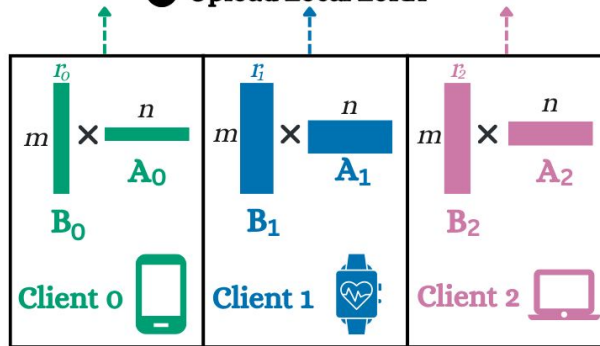
**Identifies most informative
components;**
Reduces download overhead
for high communication
efficiency



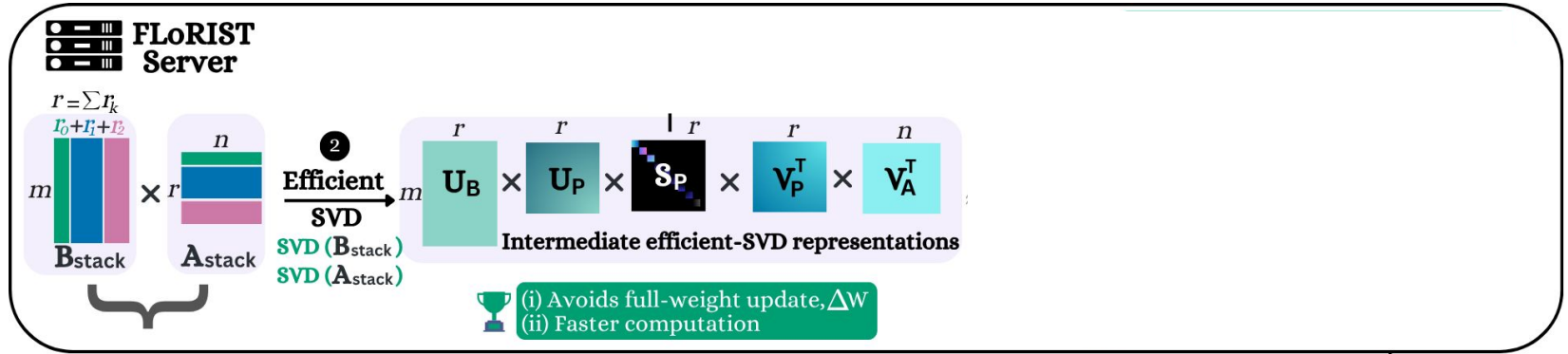
Proposed FLoRIST



1 Upload Local LoRA



Proposed FLoRIST



$$SVD(B_{\text{stack}}) = U_B S_B V_B^T$$

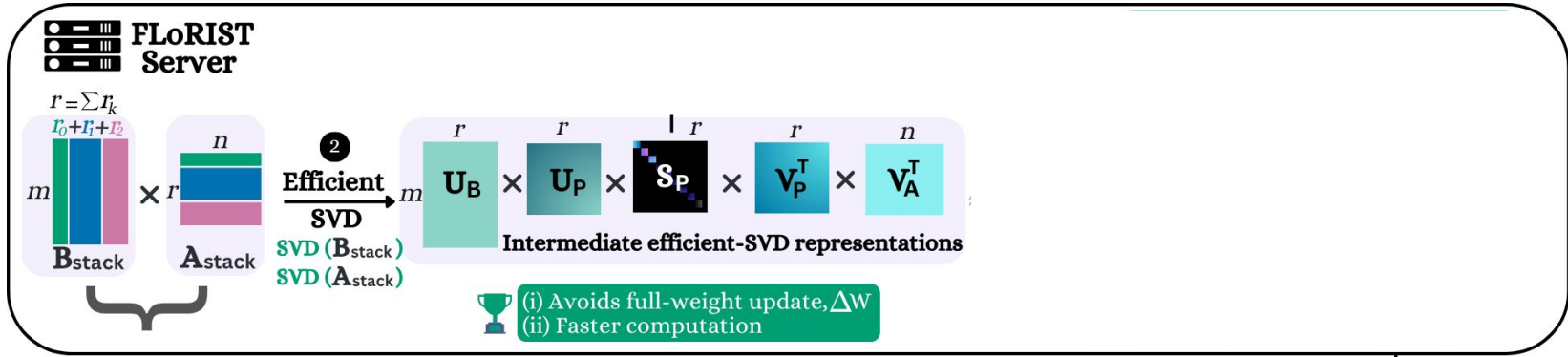
$$SVD(A_{\text{stack}}) = U_A S_A V_A^T$$

$$\Delta W = U_B \underbrace{S_B V_B^T U_A S_A V_A^T}_{P} V_A^T$$

$$P \in \mathbb{R}^{r \times r}$$

P captures the **cross-adapter interaction** between Local LoRA while maintaining low dimension, $r < \{m, n\}$

Proposed FLoRIST



$$SVD(B_{\text{stack}}) = U_B S_B V_B^T$$

$$\Delta W = U_B \underbrace{S_B V_B^T U_A S_A}_{P} V_A^T$$

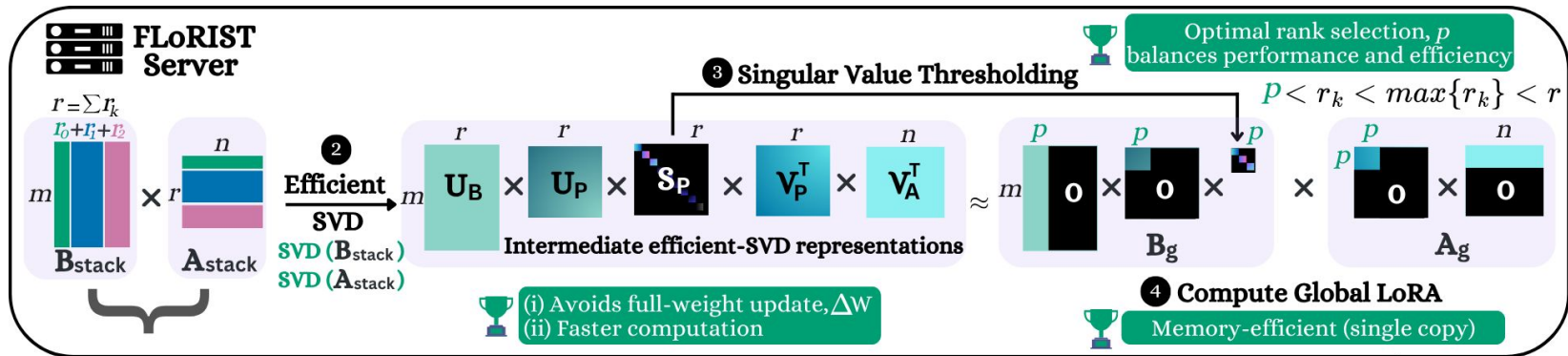
$$SVD(A_{\text{stack}}) = U_A S_A V_A^T$$

$$SVD(P) = U_P S_P V_P^T$$

$$\Delta W = (U_B U_P) \underbrace{S_P}_{\text{Singular values (P)}} (V_P^T V_A^T)$$

$$SVD(\Delta W) = U \underbrace{S}_{\text{Singular values (P)}} V^T$$

Singular values (P) == Singular values (ΔW)
without global update reconstruction



Energy-based Rank Selection

We select the smallest rank p that ensures at least a τ -fraction of the total variance is preserved,

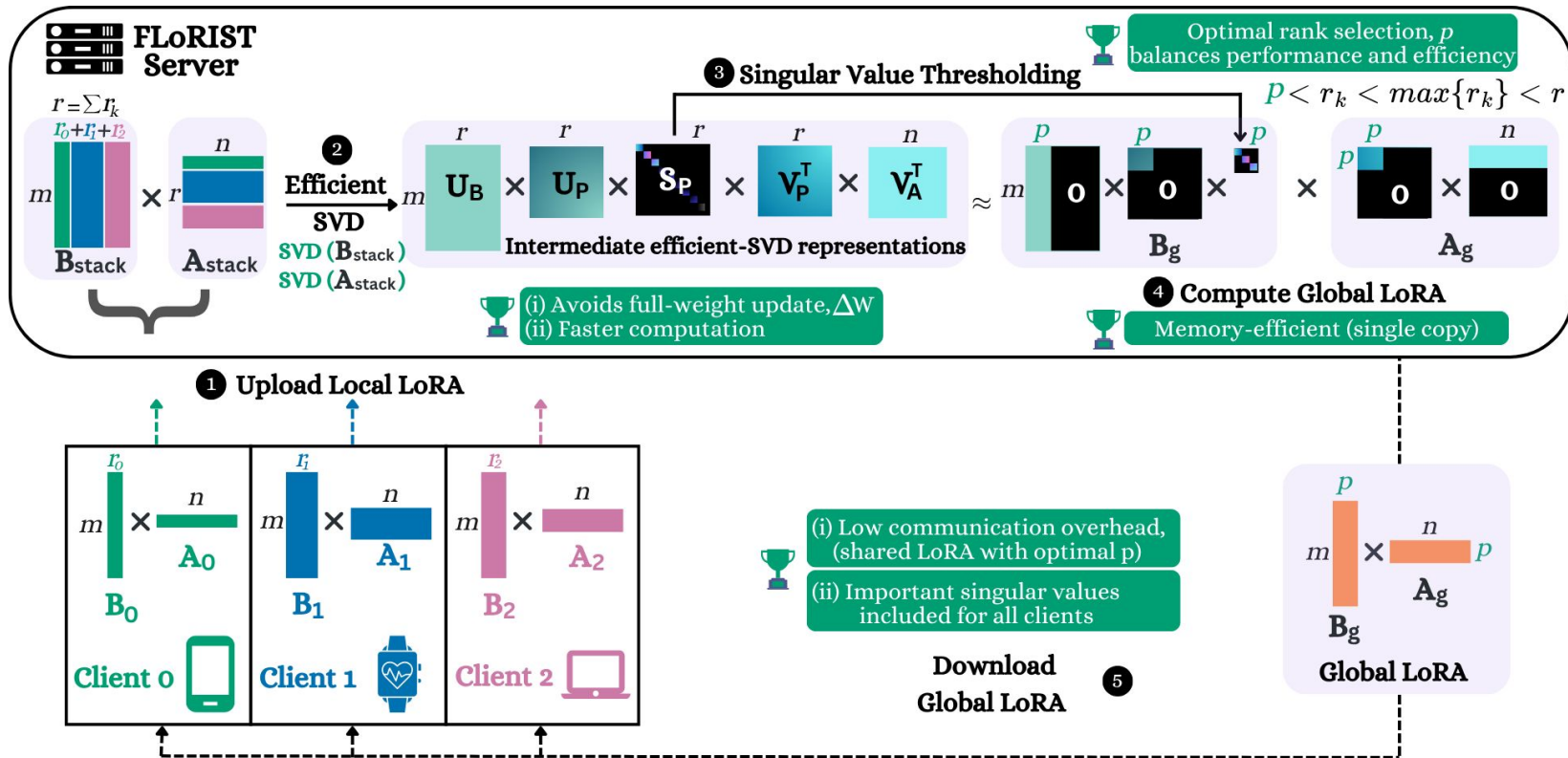
$$\frac{\sum_{i=1}^p \sigma_i^2}{\sum_{i=1}^{r^*} \sigma_i^2} \geq \tau$$

Approximation Error Bound

By the Eckart–Young–Mirsky theorem, the approximation error is bounded by the tail energy of the discarded singular values:

$$\|\Delta W - B_g A_g\|_F \leq \left(\sum_{i=p+1}^{r^*} \sigma_i^2 \right)^{1/2}$$

Proposed FLoRIST



Models & Datasets

LLMs: TinyLlama-1.1B, LLaMA-3.2-1B

Datasets: Dolly-15k, Alpaca-52k, WizardLM-70k

Eval: MMLU (1,444 questions, 57 subjects)

FL Setup

Hardware: NVIDIA H100 and A100 GPUs

Clients: 100 clients, 10 sampled/round, 75 rounds, non-IID (Dirichlet $\alpha=0.5$)

Homogeneous ranks: 100 clients = 100 x r16

Heterogeneous ranks:

100 clients = {40×r4, 20×r8, 20×r16, 10×r32, 10×r64}

FLoRIST: Energy Threshold Variants

FLoRIST [τ^*]

- Diagnostic upper bound on the accuracy–efficiency
- Smallest τ that matches or beats all baseline accuracy
- Binary search over [0.80, 0.99]

FLoRIST [$\tau=0.9$]

- Single fixed threshold for practical deployment

Performance Metrics for baseline comparison

Accuracy (MMLU)

Communication Efficiency (1/download parameters)

Convergence

Computational FLOPs (Server)

Results: Accuracy vs Efficiency

MMLU performance across models, client configurations and 3 datasets

MODEL	CLIENT	METHOD	DOLLY		ALPACA		WIZARD	
			ACC. (%)	EFF. ($\times 10^{-4}$)	ACC. (%)	EFF. ($\times 10^{-4}$)	ACC. (%)	EFF. ($\times 10^{-4}$)
TINYLLAMA	HOMO	FEDIT	27.46	14.20	28.35	14.20	36.61	14.20
		FLORA	28.99	1.78	28.99	1.78	34.20	1.78
		FLEXLORA	28.06	14.20	29.22	14.20	<u>39.75</u>	14.20
		FFA-LORA	27.79	28.40	31.58	28.40	36.01	28.40
		FLoRIST [τ^*]	29.16	53.54	32.26	35.08	40.95	54.38
	FLoRIST [$\tau=0.9$]	30.94	23.48	31.68	60.06	38.92	63.09	
	HETER	FEDIT (ZERO-PAD)	25.73	3.55	<u>31.04</u>	3.55	44.19	3.55
		FLORA	27.20	0.50	28.58	0.50	33.74	0.50
		FLEXLORA	28.49	11.96	29.61	11.96	36.39	11.96
		FFA-LORA	18.54	7.10	23.19	7.10	23.75	7.10
FLoRIST [τ^*]		28.87	37.87	31.33	45.09	38.20	13.60	
FLoRIST [$\tau=0.9$]	27.69	<u>34.93</u>	29.69	<u>34.90</u>	<u>41.51</u>	36.16		
LLAMA-3.2-1B	HOMO	FEDIT	25.00	19.50	29.44	19.50	30.48	19.50
		FLORA	22.48	2.44	29.16	2.44	28.57	2.44
		FLEXLORA	27.39	19.50	29.24	19.50	<u>30.03</u>	19.50
		FFA-LORA	26.15	39.06	29.18	39.06	28.40	39.06
		FLoRIST [τ^*]	28.28	51.71	30.46	33.30	29.31	51.19
	FLoRIST [$\tau=0.9$]	29.48	<u>33.64</u>	<u>29.73</u>	<u>33.42</u>	29.62	51.19	
	HETER	FEDIT (ZERO-PAD)	21.41	4.88	28.37	4.88	28.42	4.88
		FLORA	23.85	2.06	30.15	2.06	27.73	2.06
		FLEXLORA	26.74	16.44	29.95	16.44	29.13	16.44
		FFA-LORA	22.45	9.77	22.68	9.77	28.78	9.77
FLoRIST [τ^*]		24.01	<u>46.35</u>	30.53	<u>45.22</u>	29.79	<u>47.82</u>	
FLoRIST [$\tau=0.9$]	24.10	49.38	<u>30.24</u>	48.24	<u>29.45</u>	50.22		

FLoRIST achieves highest Accuracy-Efficiency tradeoff for both homogeneous and heterogeneous ranks

With [$\tau=0.9$], accuracy stays within $\pm 1\%$ of τ^*

Results: Server Computation FLOPs

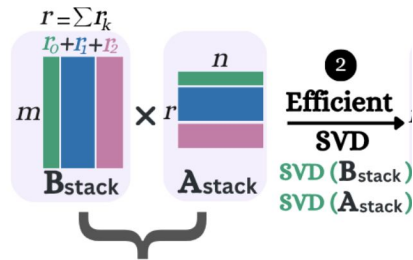
METHOD	SERVER FLOPs
FEDIT	4.76M
FFA-LoRA	0.52M
FLoRA	0B
FLEXLoRA	3516.01M
FLoRIST (OURS)	466.95M

TinyLlama-Alpaca-homo

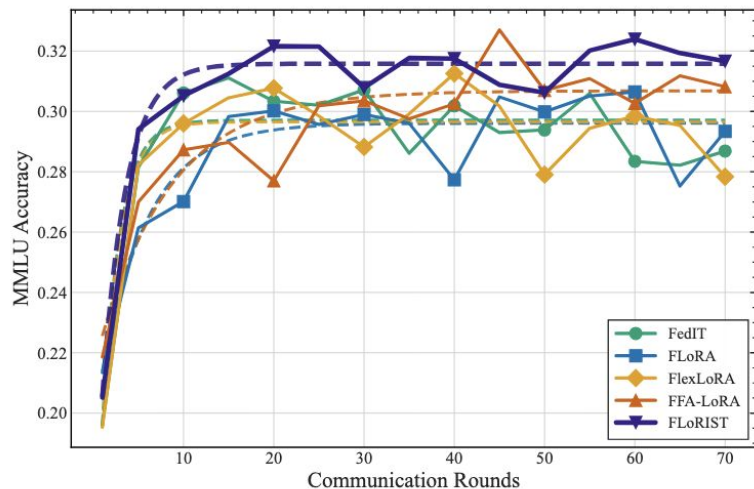
FLoRIST: 467M FLOPs

7.5× faster than FlexLoRA (3516 M)

Avoids reconstructing global weight
by using **Efficient SVD** (operate in compact $r \times r$ space)



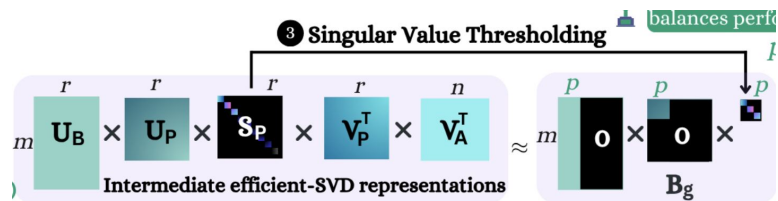
*MMLU accuracy over rounds for
TinyLlama-Alpaca-homo*



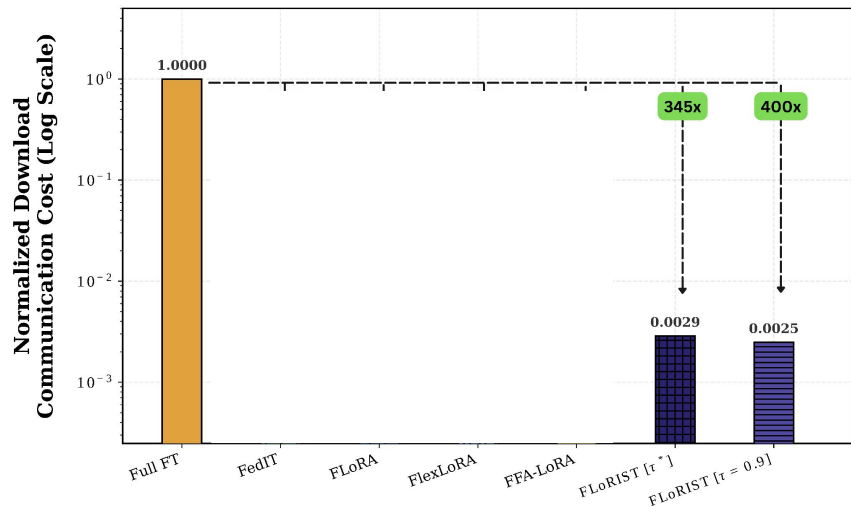
FLoRIST:

- Fastest Convergence
- Highest MMLU accuracy (Generalizes better)

Singular Value Thresholding: Higher signal to noise ratio (identifying most critical components of the aggregated global update ΔW)



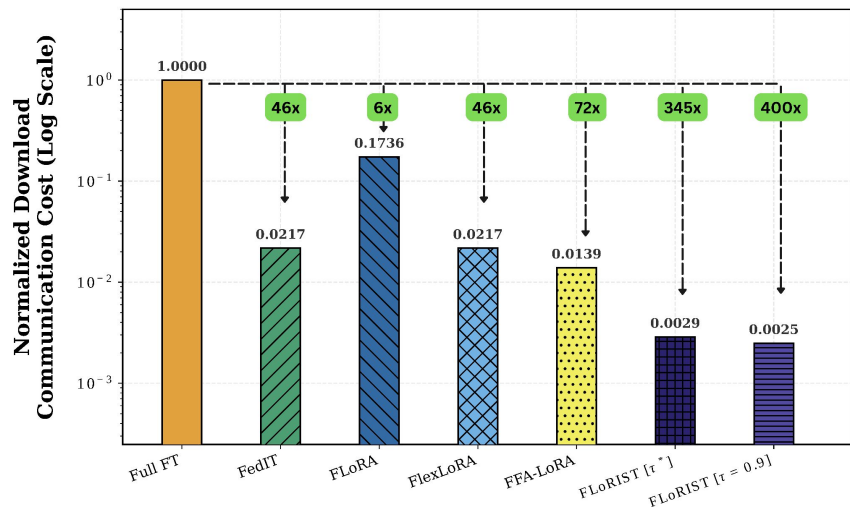
Normalized download communication cost (TinyLlama-Alpaca-homo)



Download overhead: FLoRIST[$\tau = 0.9$]

- 400× reduction vs Full Fine Tuning
- 70× reduction vs FLoRA (stack)
- 5× reduction vs FFA-LoRA (freeze)

*Normalized download communication cost
(TinyLlama-Alpaca-homo)*



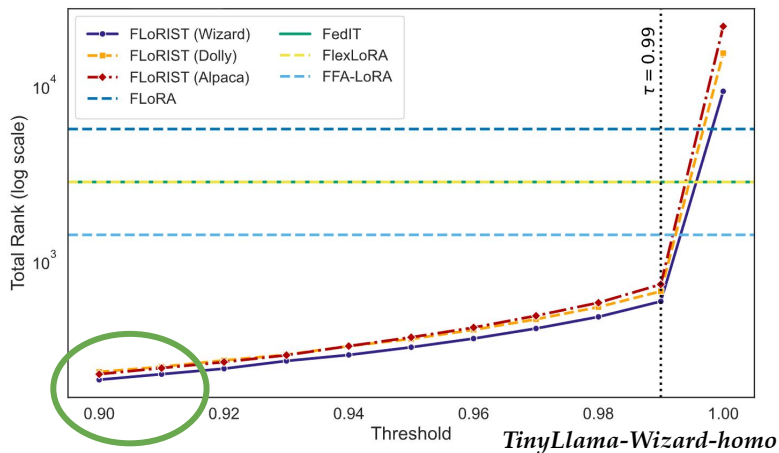
Download overhead: FLoRIST[$\tau = 0.9$]

- 400× reduction vs Full Fine Tuning
- 70× reduction vs FLoRA (stack)
- 5× reduction vs FFA-LoRA (freeze)

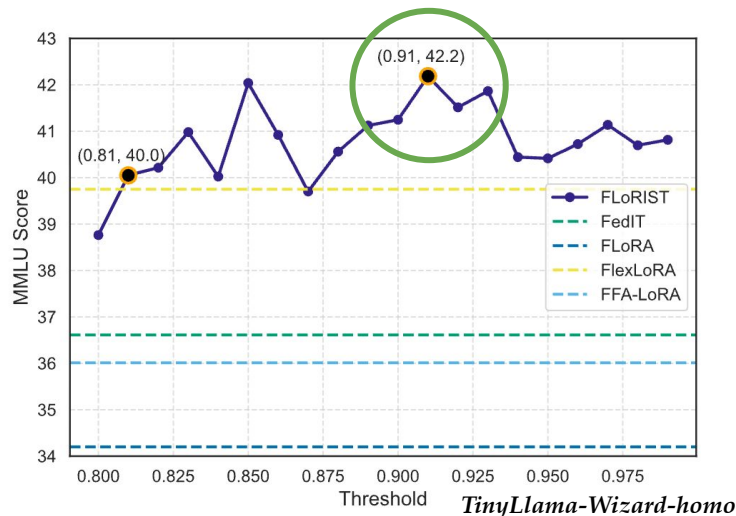
Ablation study: Energy Threshold

Identify optimal energy threshold (τ) for accuracy-efficiency tradeoff

Total Rank



MMLU













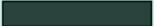
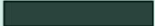



For $\tau \leq 0.99$, FLoRIST achieves
lowest global rank (higher efficiency)

$\tau = 0.9$ is a robust practical threshold
with strong accuracy-efficiency tradeoff

Conclusion

1. FLoRIST aggregates directly in **low-rank latent space**, avoiding reconstruction of full ΔW using **Efficient SVD**
2. **Singular value thresholding** identifies intrinsic dimensionality and selects optimal global rank
3. FLoRIST supports **heterogeneous client ranks** natively via a single shared global adapter
4. FLoRIST achieves **best accuracy-efficiency tradeoff** compared to baselines with $\tau = 0.9$ as practical threshold

METHOD	HETEROGENEITY	PERFORMANCE \uparrow	COMM. EFF. \uparrow	COMP. COST \downarrow
FEDIT	✗	LOW  HIGH	LOW  HIGH	LOW  HIGH
FFA-LoRA	✗	LOW  HIGH	LOW  HIGH	LOW  HIGH
FLoRA	✓	LOW  HIGH	LOW  HIGH	LOW  HIGH
FLEXLoRA	✓	LOW  HIGH	LOW  HIGH	LOW  HIGH
FLoRIST (OURS)	✓	LOW  HIGH	LOW  HIGH	LOW  HIGH

Poster

29



Paper



Code



Acknowledgments

ACCESS



Thank you!