

# Federated Full-Parameter Tuning at Scale for LLMs

Yao Shu<sup>\*1</sup>, Wenyang Hu<sup>\*2,3</sup>

See-Kiong Ng<sup>3</sup>, Bryan Kian Hsiang Low<sup>3</sup>, Fei Yu<sup>4</sup>



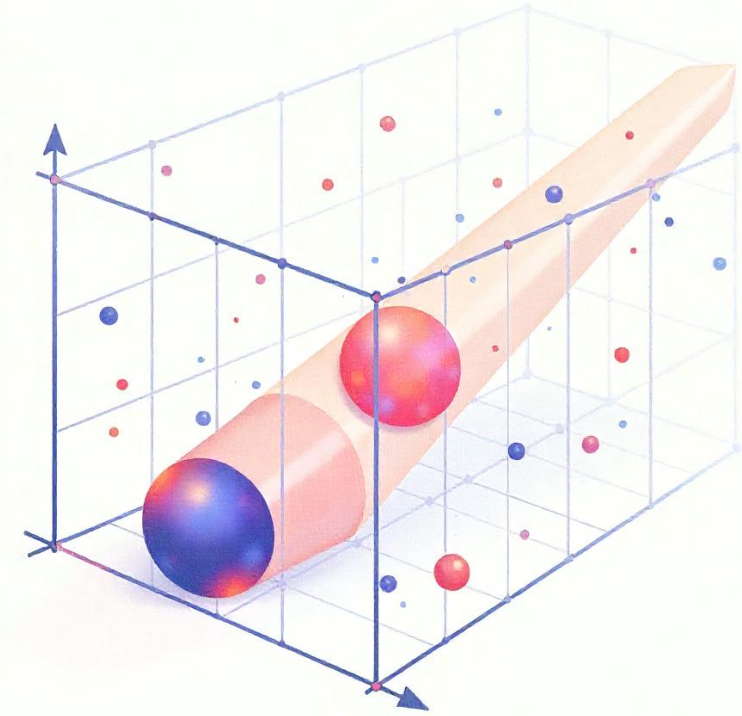
NUS  
National University of Singapore



ICML  
International Conference On Machine Learning

HKUST (GZ)<sup>1</sup>, SAP<sup>2</sup>, NUS<sup>3</sup>, Guangdong Lab of AI and Digital Economy (SZ)<sup>4</sup>

Federated full-parameter tuning a 7B LLM with just 6K communicated parameters!



## Objective

Given **random bases**  $\mathbf{V} = [\mathbf{v}_1 \mathbf{v}_2 \dots \mathbf{v}_K] \in \mathbb{R}^{d \times K}$  generated by a seed  $s$  and a local update  $\Delta$ , we want to project the update  $\Delta$  into coordinates

$$\gamma \triangleq \arg \min_{\mathbf{y}} \|\mathbf{V} \mathbf{y} - \Delta\|$$

Coordinates      Random Bases

$$\gamma = (\mathbf{V}^\top \mathbf{V})^{-1} \mathbf{V}^\top \Delta, \quad \tilde{\Delta} = \mathbf{V} (\mathbf{V}^\top \mathbf{V})^{-1} \mathbf{V}^\top \Delta.$$

Calculating this inversion is very costly!

**Theorem 1 (Unbiased Reconstruction).** Given the reconstruction in (7), we have

$$\mathbb{E} [\tilde{\Delta}] = \Delta.$$

**Theorem 2 (Reconstruction Error).** Given the reconstruction in (7), we have

$$\mathbb{E} [\|\tilde{\Delta} - \Delta\|] \leq \max \left\{ 2\sqrt{\frac{2 \ln(2d)}{\rho K}}, \frac{2 \ln(2d)}{\rho K} \right\} \|\Delta\|.$$

Project a first-order gradient using random vectors and recover it using shared randomness

## Our Solution Ferret

### 1. Reconstruction w/o Inversion

Approximate  $\mathbf{V}^\top \mathbf{V}$  with  $\mathbf{I}_K$

$$\gamma \approx (\rho K)^{-1} \mathbf{V}^\top \Delta$$

$$\tilde{\Delta} = (\rho K)^{-1} \mathbf{V} \mathbf{V}^\top \Delta$$

### 2. Blockwise Reconstruction

Divide the full dimension  $d$  into  $L$  blocks. Then for each block  $l$ , we have

$$\gamma_l = (\rho_l K)^{-1} \mathbf{V}_l^\top \Delta_l$$

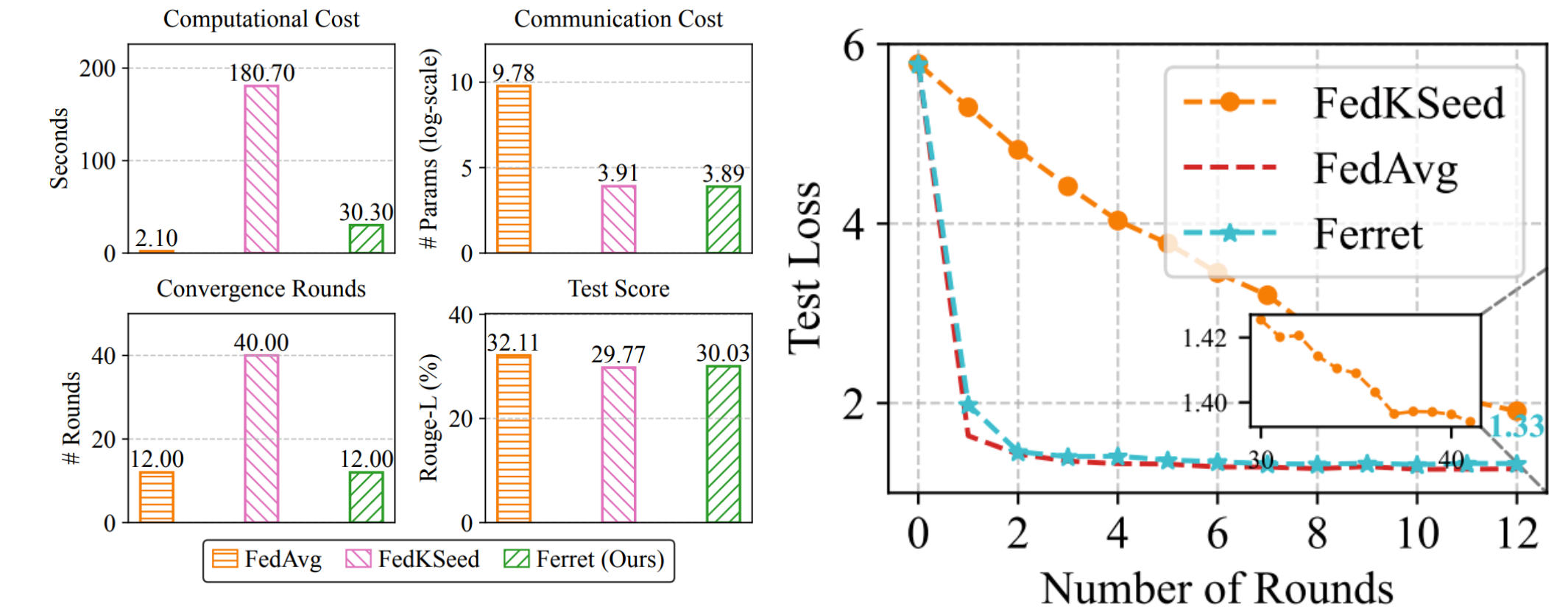
$$\tilde{\Delta}_l = (\rho_l K_l)^{-1} \mathbf{V}_l \mathbf{V}_l^\top \Delta_l$$

### Proposition 2 (Block-Wise Reconstruction Error).

For block-wise reconstruction (8) of size  $L$ , when  $\sqrt{d_l} \geq K_l$  for any  $l \in [L]$ ,

$$\mathbb{E} [\|\tilde{\Delta} - \Delta\|] < \tilde{\mathcal{O}} \left( \sum_{l \in [L]} \frac{\|\Delta_l\|}{\rho_l K_l} \right),$$

which is minimized by choosing  $K_l \propto \sqrt{\|\Delta_l\| / \rho_l}$ .



- **Reduced communication cost**
- **High computational efficiency**
- **Fast convergence**

Algorithm	CodeAlpaca		GSM8K	
	LLaMA2-7B	LLaMA2-13B	LLaMA2-7B	LLaMA2-13B
FedIT	4.66 ± 0.18	6.10 ± 0.18	30.31 ± 0.29	13.46 ± 0.34
FedZO	4.58 ± 0.26	6.19 ± 0.32	30.41 ± 0.31	13.63 ± 0.34
FedKSeed	8.33 ± 0.98	10.70 ± 0.47	28.26 ± 3.60	33.67 ± 1.15
FedAvg	<b>15.41 ± 0.43</b>	<b>14.68 ± 0.26</b>	<b>38.30 ± 0.40</b>	<b>39.82 ± 0.17</b>
Ferret (ours)	12.10 ± 0.47	11.84 ± 0.91	36.10 ± 1.18	34.50 ± 1.42

### Algorithm 1 Ferret

**Input:**  $\mathbf{w}_0, N, R, T, K, \eta$

```

1 for each round  $r \in [R]$  do
2   for each client  $j \in [N]$  in parallel do
3     if  $r > 1$  then // Step ①: Global Aggregation
4       Receive  $\{s^{(i)}\}_{i=1}^N$  and  $\{\gamma_k^{(i)}\}_{i=1, k=1}^{N, K}$ 
5       Generate bases  $\{\mathbf{v}_k^{(i)}\}_{i=1, k=1}^{N, K}$  using  $\{s^{(i)}\}_{i=1}^N$ 
6        $\mathbf{w}_{r-1} \leftarrow \mathbf{w}_{r-2} - \sum_{i \in [N]} \left( \sum_{k=1}^K \gamma_k^{(i)} \mathbf{v}_k^{(i)} \right) / N$ 
7        $\mathbf{w}_{r,0} \leftarrow \mathbf{w}_r$ 
8       for  $t \in [T]$  do // Step ②: Local Updates
9          $\mathbf{w}_{r,t}^{(j)} \leftarrow \mathbf{w}_{r,t-1}^{(j)} - \eta \nabla \ell(\mathbf{w}_{r,t-1}^{(j)}; \mathbf{x}_{r,t-1}^{(j)})$ 
          // Step ③: Projected Updates
10        Randomly set  $s^{(j)}$  and generate bases  $\{\mathbf{v}_k^{(j)}\}_{k=1}^K$ 
11         $\Delta_r^{(j)} \leftarrow \mathbf{w}_{r-1}^{(j)} - \mathbf{w}_{r,t-1}^{(j)}$ , compute  $\{\gamma_k^{(j)}\}_{k=1}^K$  with (6)
12        Send  $s^{(j)}$  and  $\{\gamma_k^{(j)}\}_{k=1}^K$  to the central server

```