

NETWORK COMPRESSION IN FEDERATED MACHINE LEARNING

联邦学习中的网络压缩研究

电子信息与电气工程学院 信息安全专业

廖宁祎 517021910844

指导老师：邵硕 向立瑶



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY



OUTLINE

Background

Introduction, Background

Contribution

Motivation, Contribution

Algorithm

Algorithmic Framework, Theoretical Results

Experiment

Experimental Results on Efficiency and Privacy

Conclusion

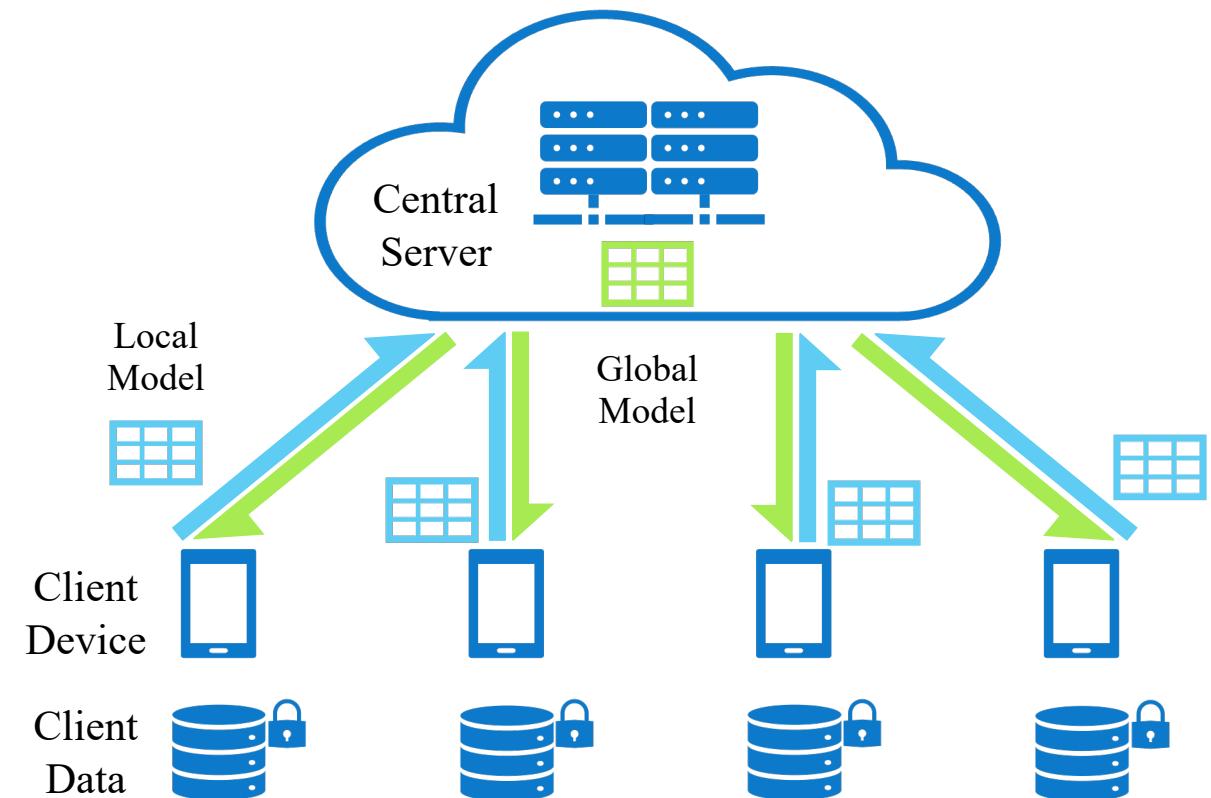
Conclusions, Future Work



BACKGROUND

Federated Learning (FL)

- Machine learning that decouples training model and data
- Key factors: communication overhead & data privacy
- *Client*: update local model; keep data private
- *Server*: aggregate global model

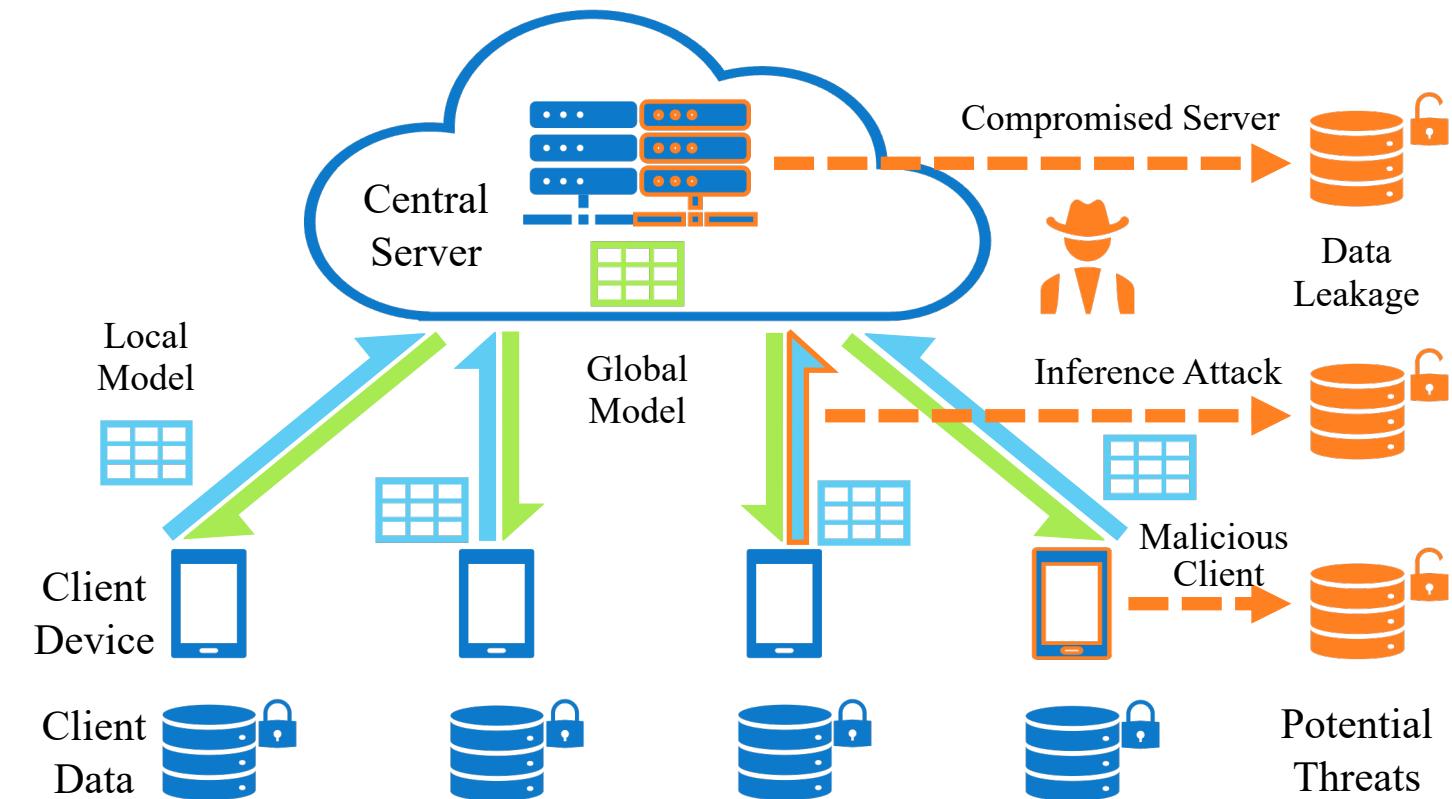




BACKGROUND

Federated Learning Privacy

- Attackers exploit model uploads to infer client data
- *Membership inference*: distinguish client participation
- *Gradient inversion*: recover private training sample





BACKGROUND

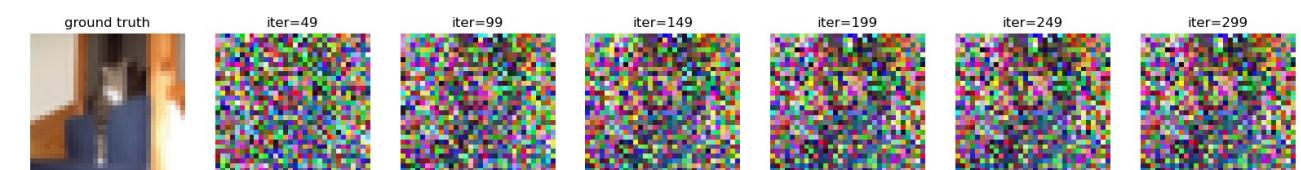
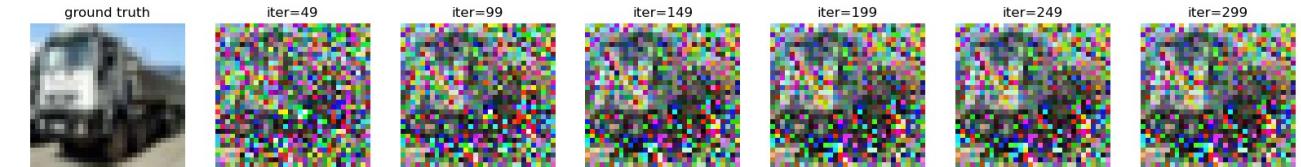
Federated Learning Privacy

- Attackers exploit model uploads to infer client data
- *Membership inference*: distinguish client participation
- *Gradient inversion*: recover private training sample



truth

DLG: FedSGD



truth

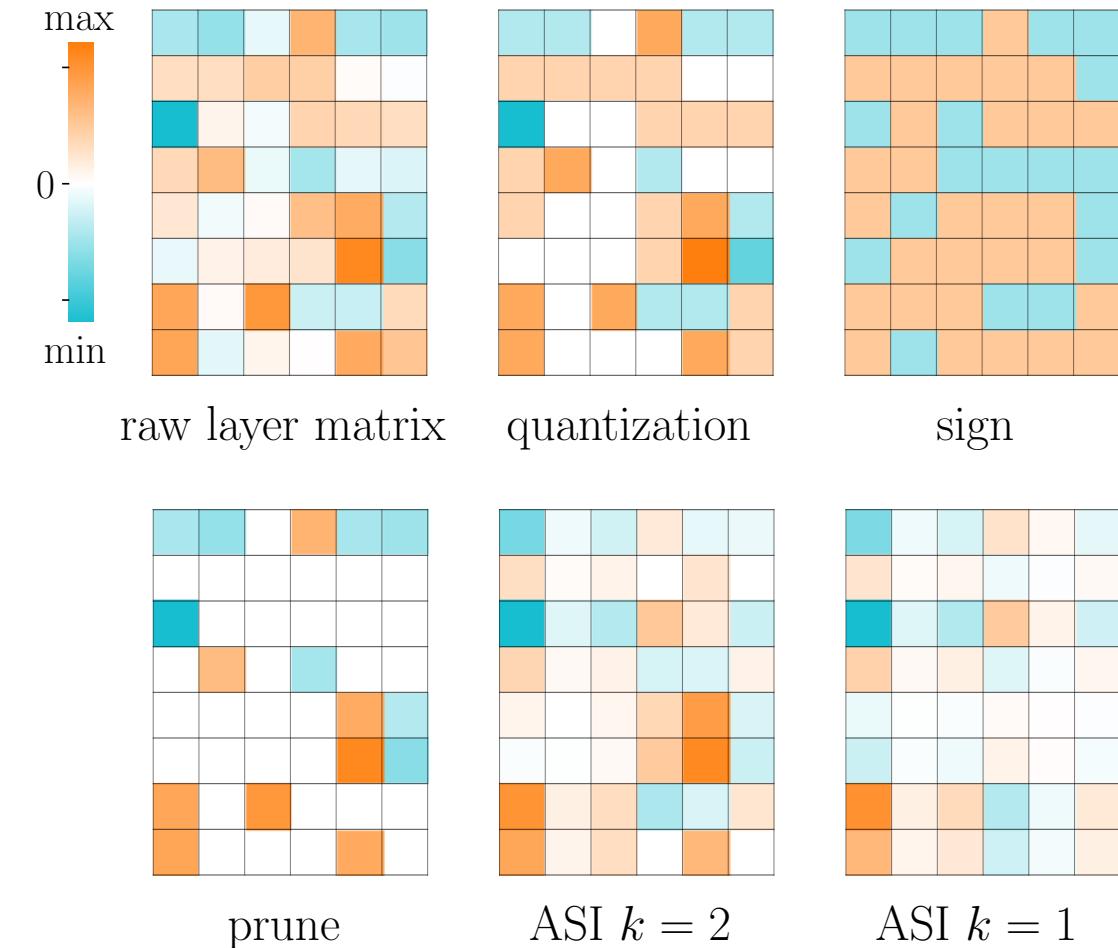
DLG: FedSGD + noise



BACKGROUND

Federated Learning Communication

- Uplink transmits large-scale model weights/gradients tensors
- Protocol improvement: *client selection, local update*, etc.
- Parameter compression: *pruning, quantization*, etc.





CONTRIBUTION

Question:

How to optimize transmission cost adaptively in the setting of differentially-private training, and design algorithms towards **communication-efficient** and **privacy-preserving** federated learning?

Our Contribution:

FedASI & FedALS

Adaptive low-rank
factorization

FL communication
compression

Differential privacy
guarantee



CONTRIBUTION

Our Contribution:

Adaptive low-rank factorization

FL communication compression

Differential privacy guarantee

FedASI & FedALS

- Develop low-rank factorization based on ASI and ALS
- Solve matrix rank optimization with adaptive output
- Derive convergence proof on factorization iteration

- Introduce low-rank compression to FL communication
- Derive convergence proof on learning
- Evaluate efficiency, accuracy, convergence by experiment

- Apply differential privacy mechanism with proof
- Link privacy with compression benefits

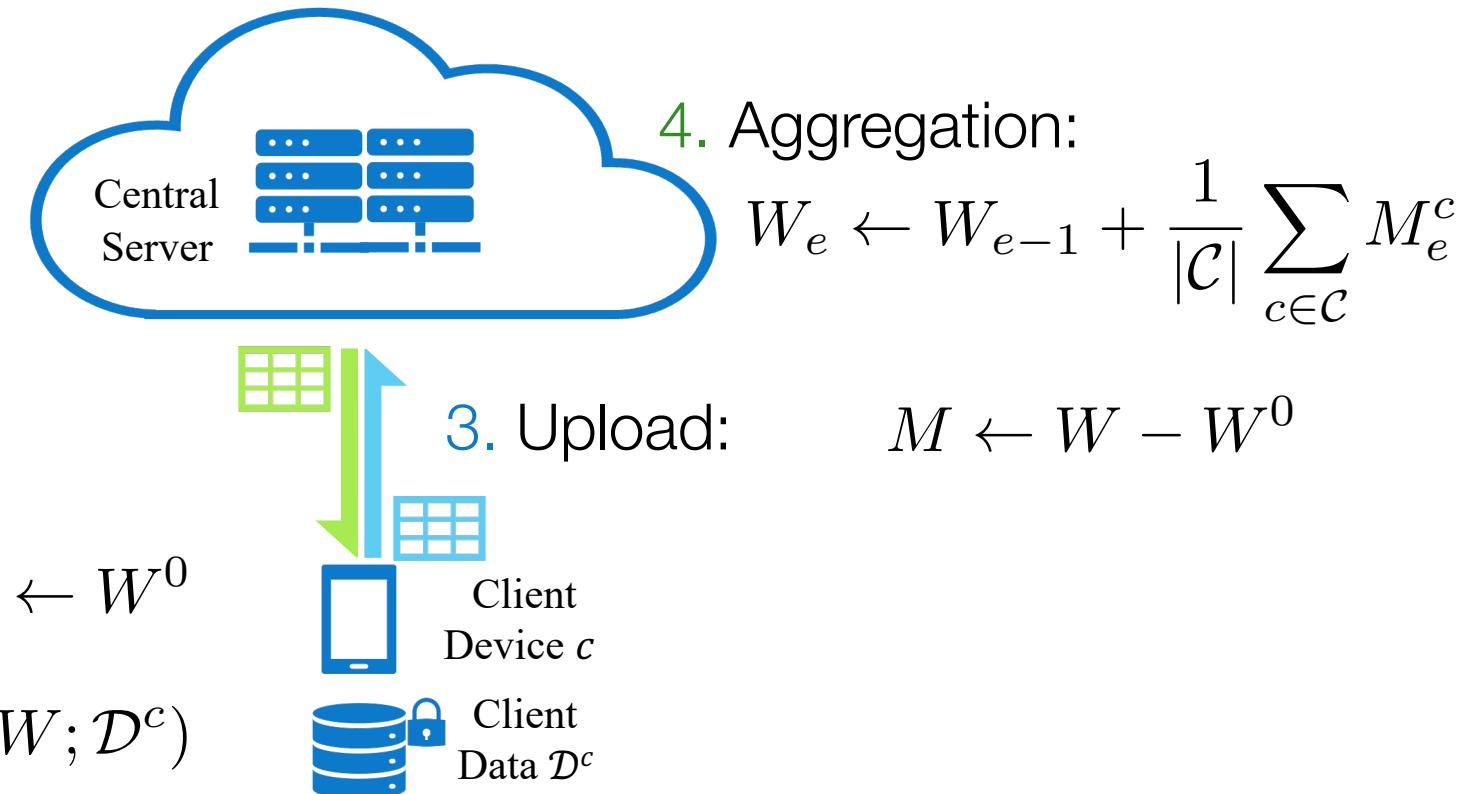


FRAMEWORK DESIGN

FL Algorithm Overview

- Objective:

$$\min_W \mathcal{L}(W) = \frac{1}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} \mathcal{L}^c(W; \mathcal{D}^c)$$





FRAMEWORK DESIGN

FL with Low-Rank Compression

Efficiency:

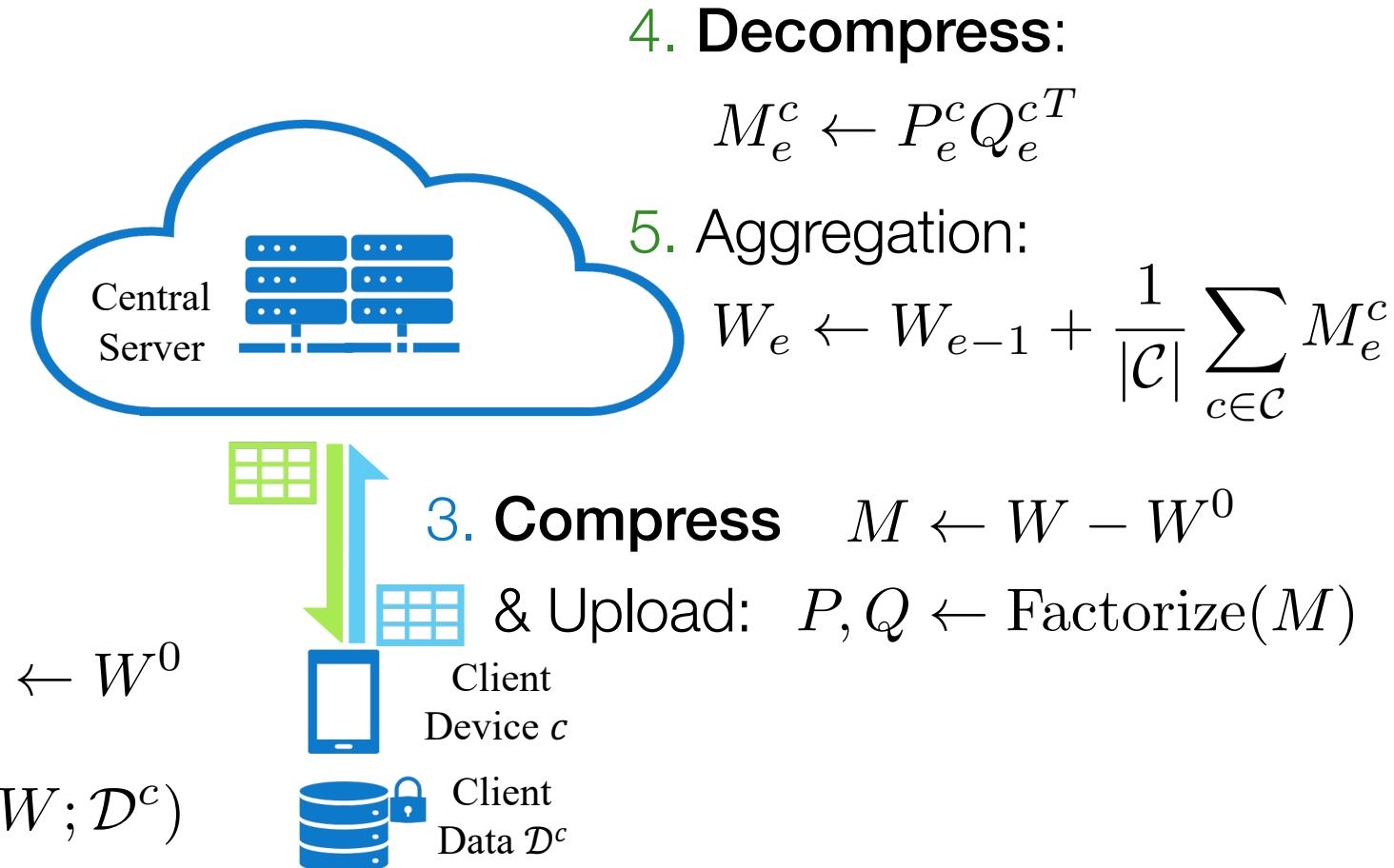
Fixed: **FedASI**

Adaptive: **FedALS**

Efficiency
+Privacy:

FedASI-DP

FedALS-DP





LOW-RANK COMPRESSION

Low-Rank Matrix Factorization

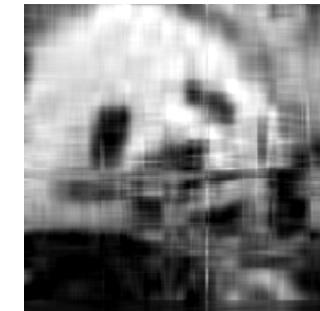
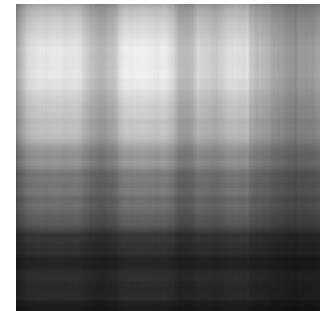
- Rank- k Factorization:

$$M \approx M_k = PQ^T$$

original, $k = 320$  $k = 64$  $k = 16$

- Rank- k Restricted SVD:

$$M \approx M_k = U_k \Sigma_k V_k^T$$

 $k = 8$  $k = 4$  $k = 1$

- Other methods: eigenvalue decomposition, QR factorization, etc.



LOW-RANK COMPRESSION

Alternating Subspace Iteration (ASI)

- Rank- k Factorization:

$$M \approx M_k = PQ^T$$

- Power Iteration:

$$\mathbf{u} \sim M^t \mathbf{u}$$

- Alternating Subspace Iteration:

$$p \sim (MM^T)^t p$$

$$P = U_k$$

$$q \sim (M^T M)^t q$$

$$Q = V_k \Sigma_k.$$

Algorithm 3-2 Alternating Subspace Iteration

Input: matrix $M \in \mathbb{R}^{m \times n}$,
rank k , iteration T

Output: low-rank matrices $P \in \mathbb{R}^{m \times k}, Q \in \mathbb{R}^{n \times k}$

```
1  $Q_0 \leftarrow \mathcal{N}(0, 1)^{n \times k}$ 
2 for each iteration  $t = 1, 2, \dots, T$  do
3    $P_t \leftarrow MQ_{t-1}$ 
4    $\hat{P}_t \leftarrow \text{orthonormalize}(P_t)$ 
5    $Q_t \leftarrow M^T \hat{P}_t$ 
6 end for
7 return  $\hat{P}_T, Q_T$ 
```



LOW-RANK COMPRESSION

Alternating Least Squares (ALS)

- Rank- k Factorization Optimization:

$$\min_{P,Q} \frac{1}{2} \|M - PQ^T\|_F^2, \text{ s.t. } \text{rank}(PQ^T) = k$$

$$\Rightarrow \min_{P,Q} \frac{1}{2} \|M - PQ^T\|_F^2 + \lambda \text{rank}(PQ^T)$$

$$\Rightarrow \min_{P,Q} \frac{1}{2} \|M - PQ^T\|_F^2 + \lambda \|PQ^T\|_*$$

$$\Rightarrow \min_{P,Q} \frac{1}{2} \|M - PQ^T\|_F^2 + \frac{\lambda}{2} (\|P\|_F^2 + \|Q\|_F^2)$$

$$P = U_k(\Sigma_k - \lambda I_k)_+^{\frac{1}{2}}, Q = V_k(\Sigma_k - \lambda I_k)_+^{\frac{1}{2}}$$

Algorithm 3–3 Alternating Least Squares

Input: matrix $M \in \mathbb{R}^{m \times n}$, regularization parameter λ , iteration T

Output: low-rank matrices $P \in \mathbb{R}^{m \times k}, Q \in \mathbb{R}^{n \times k}$

1 $Q_0 \leftarrow \mathcal{N}(0, 1)^{n \times k}$

2 **for** each iteration $t = 1, 2, \dots, T$ **do**

3 $P_t \leftarrow M Q_{t-1} (Q_{t-1}^T Q_{t-1} + \lambda I)^{-1}$
4 $Q_t \leftarrow M^T P_t (P_t^T P_t + \lambda I)^{-1}$

5 **end for**

6 **return** P_T, Q_T

PRIVATE LOW-RANK COMPRESSION

Differential Privacy (DP)

- (ϵ, δ) -DP:

$$\Pr [\mathcal{M}(\mathcal{D}) \in \mathbb{S}] \leq e^\epsilon \cdot \Pr [\mathcal{M}(\mathcal{D}') \in \mathbb{S}] + \delta$$

- Gaussian Mechanism:

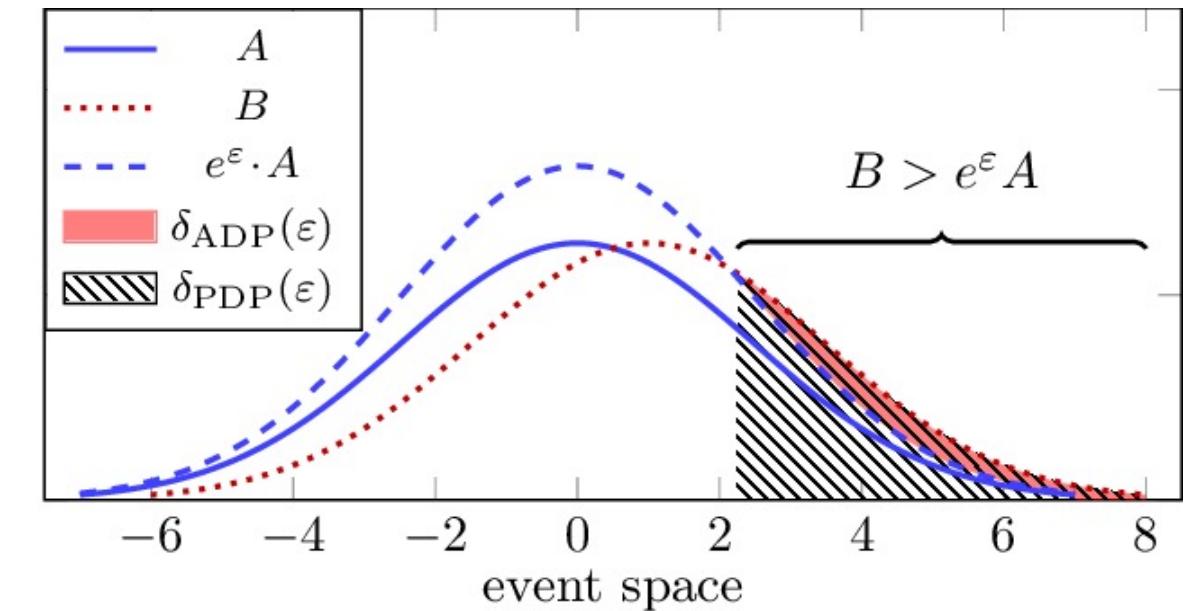
$$\mathcal{M}(\mathcal{D}) = f(\mathcal{D}) + G$$

$$G = \mathcal{N}(0, \mathcal{S}^2 \sigma^2 I_M)$$

$$\sigma = \epsilon^{-1} \sqrt{2 \ln 1.25 \delta^{-1}}$$

- Sensitivity:

$$\mathcal{S}(f) = \max_{\mathcal{D}, \mathcal{D}'} \|f(\mathcal{D}) - f(\mathcal{D}')\|_F$$





PRIVATE LOW-RANK COMPRESSION

Private ASI (ASI-DP)

Private ALS (ALS-DP)

- (ϵ, δ) -DP:

$$\mathcal{S}_Q = \|Q_{t-1}\|_\infty$$

$$\mathcal{S}_P = \|P_t\|_\infty$$

$$\sigma = \epsilon^{-1} \sqrt{4kT \ln(1/\delta)}$$

Algorithm 4-1 ASI-DP

Input: matrix $M \in \mathbb{R}^{m \times n}$,
privacy budget (ϵ, δ) ,
iteration T ,
rank k

Output: low-rank matrices

$$P \in \mathbb{R}^{m \times k}, Q \in \mathbb{R}^{n \times k}$$

```

1    $Q_0 \leftarrow \mathcal{N}(0, I_{n \times k})$ 
2   for each iteration  $t = 1, 2, \dots, T$  do
3        $G_{P,t} \leftarrow \mathcal{N}(0, S_Q^2 \sigma^2 I_{m \times k})$ 
4        $P_t \leftarrow M Q_{t-1} + G_{P,t}$ 
5        $\hat{P}_t \leftarrow \text{orthonormalize}(P_t)$ 
6        $G_{Q,t} \leftarrow \mathcal{N}(0, S_P^2 \sigma^2 I_{n \times k})$ 
7        $Q_t \leftarrow M^T \hat{P}_t + G_{Q,t}$ 
8   end for

```

Algorithm 4-2 ALS-DP

Input: matrix $M \in \mathbb{R}^{m \times n}$,
privacy budget (ϵ, δ) ,
iteration T ,
regularization parameter λ

Output: low-rank matrices

$$P \in \mathbb{R}^{m \times k}, Q \in \mathbb{R}^{n \times k}$$

```

1    $Q_0 \leftarrow \mathcal{N}(0, I_{n \times k})$ 
2   for each iteration  $t = 1, 2, \dots, T$  do
3        $\check{Q}_{t-1} \leftarrow Q_{t-1} (Q_{t-1}^T Q_{t-1} + \lambda I)^{-1}$ 
4        $G_{P,t} \leftarrow \mathcal{N}(0, S_Q^2 \sigma^2 I_{m \times k})$ 
5        $P_t \leftarrow M \check{Q}_{t-1} + G_{P,t}$ 
6        $\check{P}_t \leftarrow P_t (P_t^T P_t + \lambda I)^{-1}$ 
7        $G_{Q,t} \leftarrow \mathcal{N}(0, S_P^2 \sigma^2 I_{n \times k})$ 
8        $Q_t \leftarrow M^T \check{P}_t + G_{Q,t}$ 
9   end for

```



THEORETICAL RESULT

Convergence of Factorization

- ASI:

$$\mathbf{p}_k = \mathbf{u}_k + c \left(\frac{\varsigma_{k+1}}{\varsigma_k} \right)^{2t} \mathbf{u}_{k+1}$$

$$\mathbf{q}_k = \varsigma_k \mathbf{v}_k + c \left(\frac{\varsigma_{k+1}}{\varsigma_k} \right)^{2t} \mathbf{v}_{k+1}$$

- ALS:

$$\mathbf{p}_k = s_k^{\frac{1}{2}} \mathbf{u}_k + c \left(\frac{\varsigma_{k+1}}{\varsigma_k} \right)^{2t} \mathbf{u}_{k+1}$$

$$\mathbf{q}_k = s_k^{\frac{1}{2}} \mathbf{v}_k + c \left(\frac{\varsigma_{k+1}}{\varsigma_k} \right)^{2t} \mathbf{v}_{k+1}$$

- ASI-DP:

$$\|(I - P_t P_t^T) U_k\|_2 \leq O \left(\frac{\sigma \max_t \|Q_t\|_\infty \sqrt{m \ln t}}{\varsigma_k - \varsigma_{k+1}} \frac{\sqrt{p}}{\sqrt{p} - \sqrt{k-1}} \right)$$

$$\|(\Sigma^2 - Q_t Q_t^T) V_k\|_2 \leq O \left(\frac{\sigma \max_t \|P_t\|_\infty \sqrt{n \ln t}}{\varsigma_k - \varsigma_{k+1}} \frac{\sqrt{p}}{\sqrt{p} - \sqrt{k-1}} \right)$$

- ALS-DP:

$$\|(S - P_t P_t^T) U_k\|_2 \leq O \left(\frac{\sigma \max_t \|Q_t\|_\infty \sqrt{m \ln t}}{\varsigma_k - \varsigma_{k+1}} \frac{\sqrt{p}}{\sqrt{p} - \sqrt{k-1}} \right)$$

$$\|(S - Q_t Q_t^T) V_k\|_2 \leq O \left(\frac{\sigma \max_t \|P_t\|_\infty \sqrt{n \ln t}}{\varsigma_k - \varsigma_{k+1}} \frac{\sqrt{p}}{\sqrt{p} - \sqrt{k-1}} \right)$$



THEORETICAL RESULT

Convergence of Federated Learning

- FedASI/FedALS:

$$\mathbb{E} \|\nabla \mathcal{L}^c(\mathbf{w}_i^c)\|^2 \leq \left(\frac{\mathbb{E}[\mathcal{L}(\mathbf{w}_0)] - \mathcal{L}^*}{\bar{\eta}} + \frac{\bar{\eta} L h^2}{b|\mathcal{C}|} \right) \frac{4}{\sqrt{T}} + 8 \frac{4 - 3\bar{\xi}^2}{\bar{\xi}^2} \frac{\bar{\eta}^2 L^2 H^2 \mathcal{I}^2}{T}$$

- FedASI-DP/FedALS-DP:

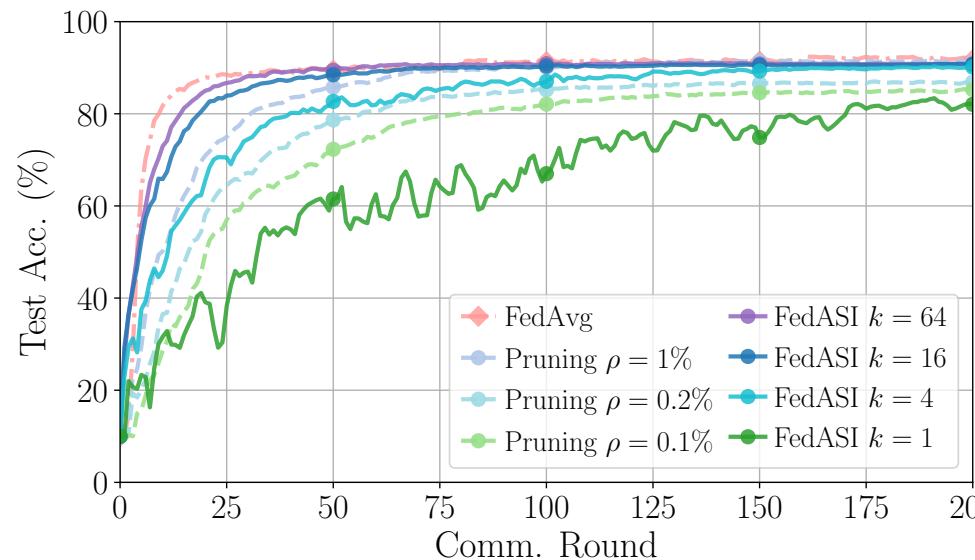
$$\mathbb{E} \|\mathbf{w}_i - \mathbf{w}^*\|^2 \leq \frac{1}{\gamma + T} \left[\frac{4D}{\mu^2} + \left(\frac{8L}{\mu} + \mathcal{I} \right) \|\mathbf{w}_0 - \mathbf{w}^*\|^2 \right]$$



EXPERIMENT ON EFFICIENCY

Experiment Setting

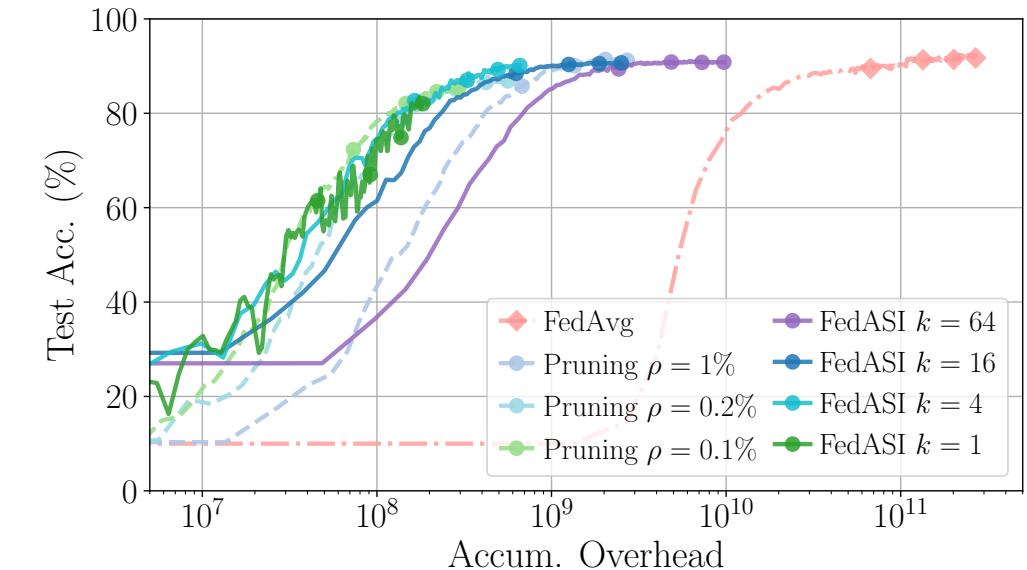
- Datasets: CIFAR-10 (IID)
- Model: VGG16, ResNet18



VGG16 Accuracy vs. Training Round

FedASI Results

- Less than 0.5% accuracy loss
- Faster convergence



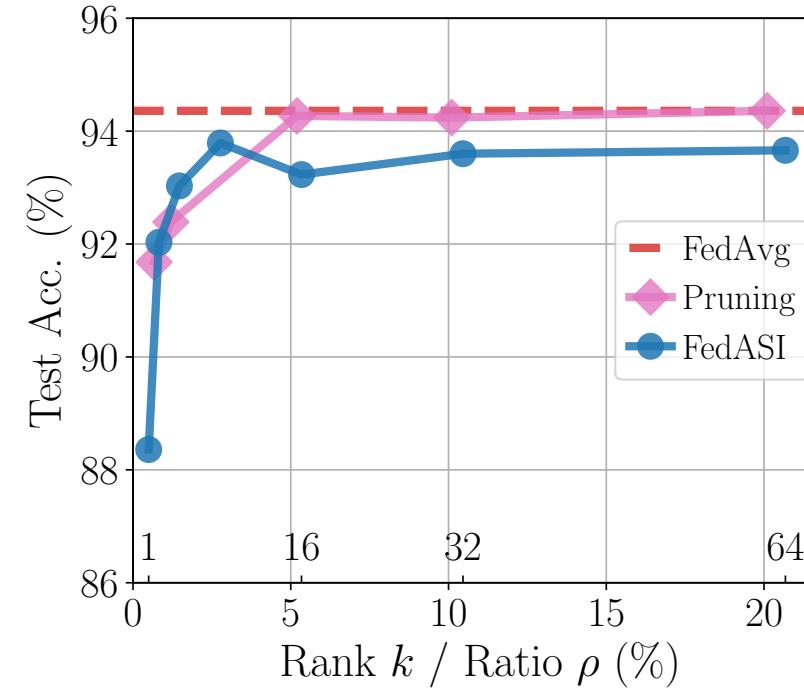
VGG16 Accuracy vs. Accum. Comm. Overhead



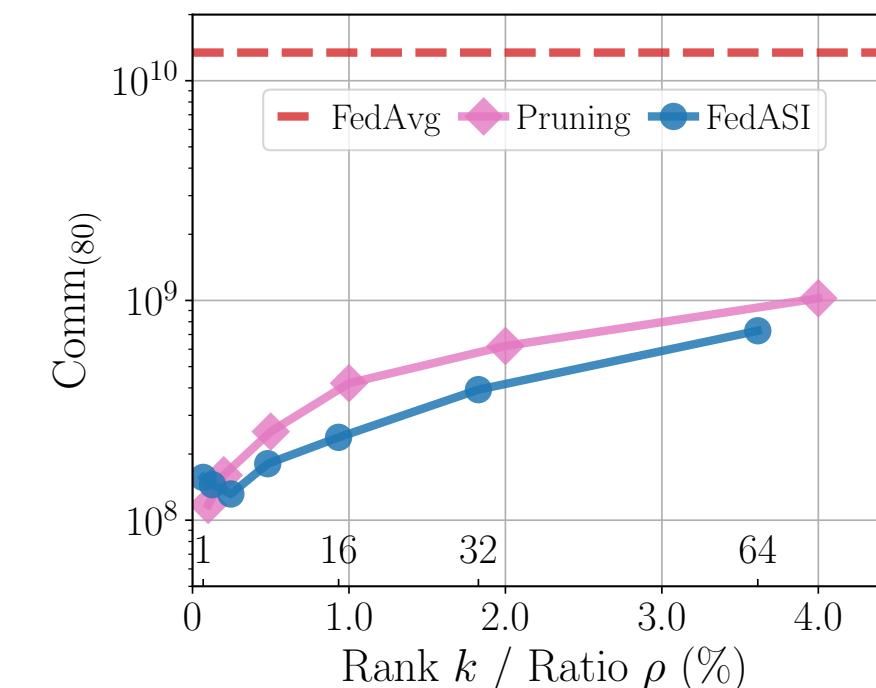
EXPERIMENT ON EFFICIENCY

FedASI Results

- Less than 0.5% accuracy loss
- Up to 1000x compression than FedAvg



ResNet18 Accuracy vs. Rank

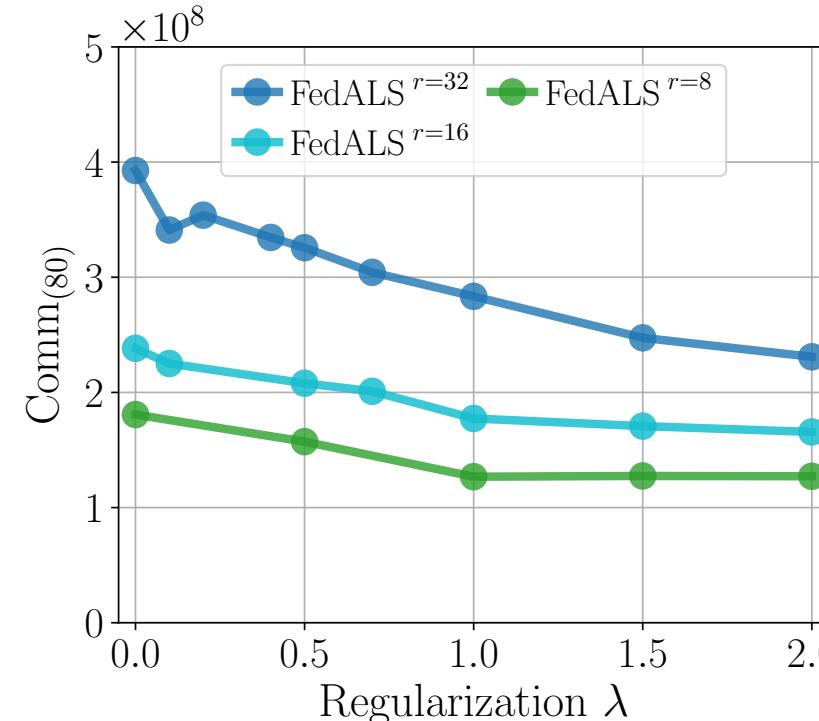


ResNet18 Converge Comm. vs. Rank

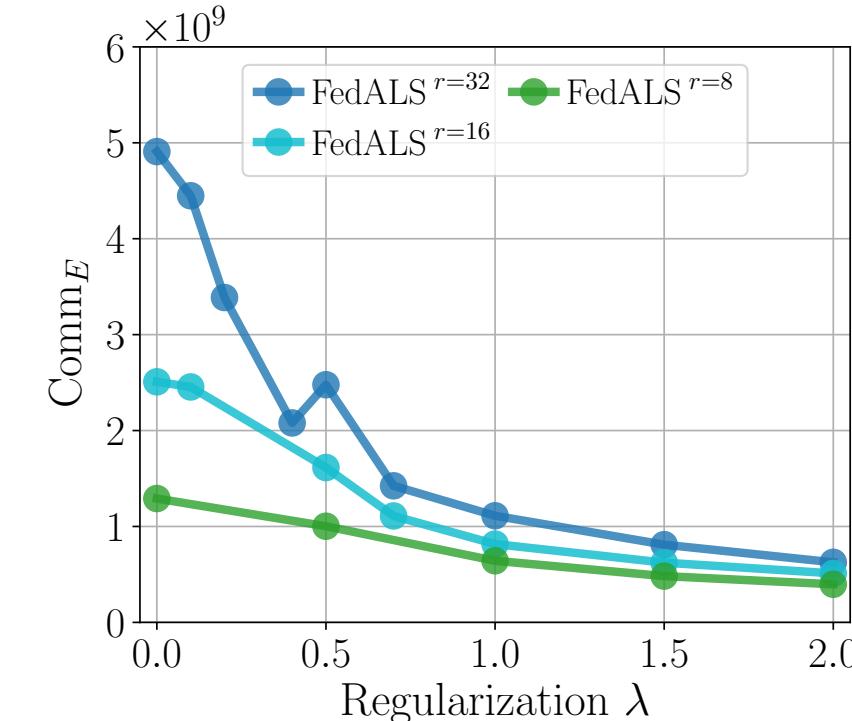
EXPERIMENT ON EFFICIENCY

FedALS Results

- Same accuracy and convergence
- Further 5x compression than FedASI



VGG16 Converge Comm. vs. Regularization



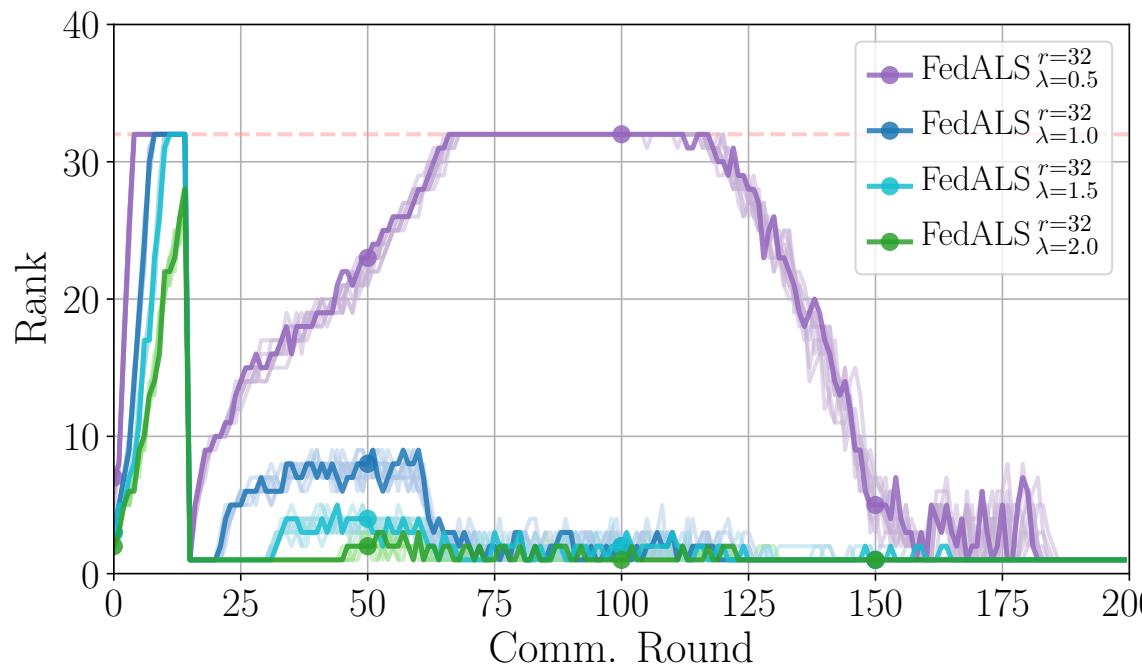
VGG16 Total Comm. vs. Regularization



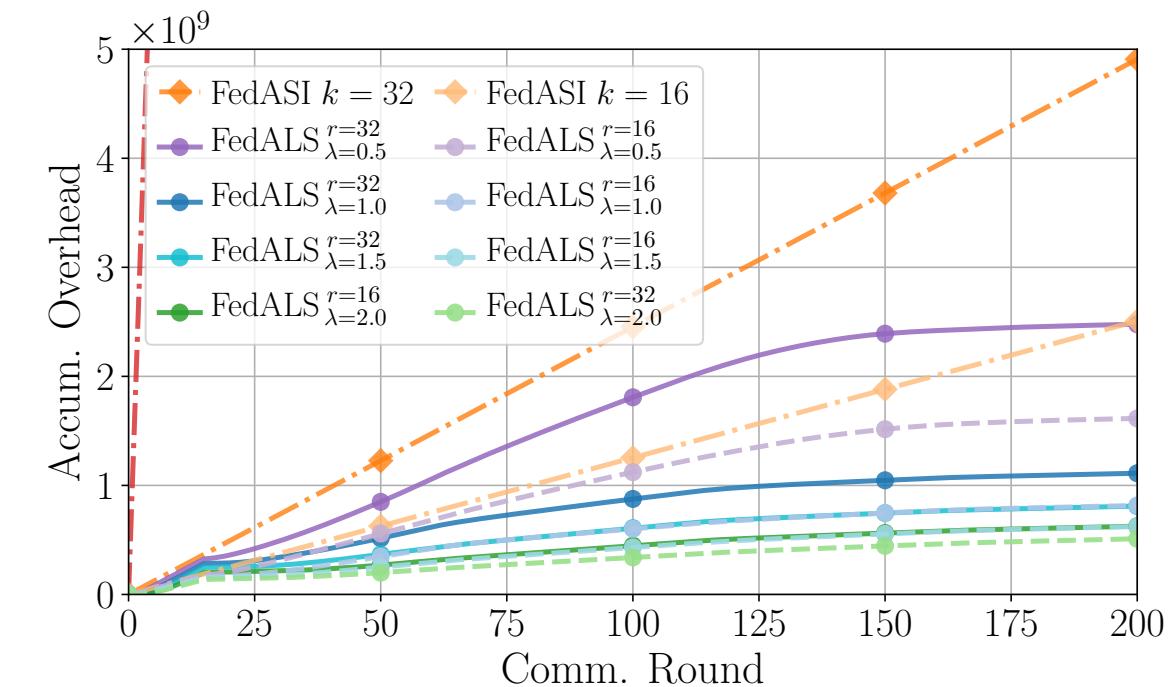
EXPERIMENT ON EFFICIENCY

FedALS Results

- Adaptive rank factorization



VGG16 Layer Rank vs. Round



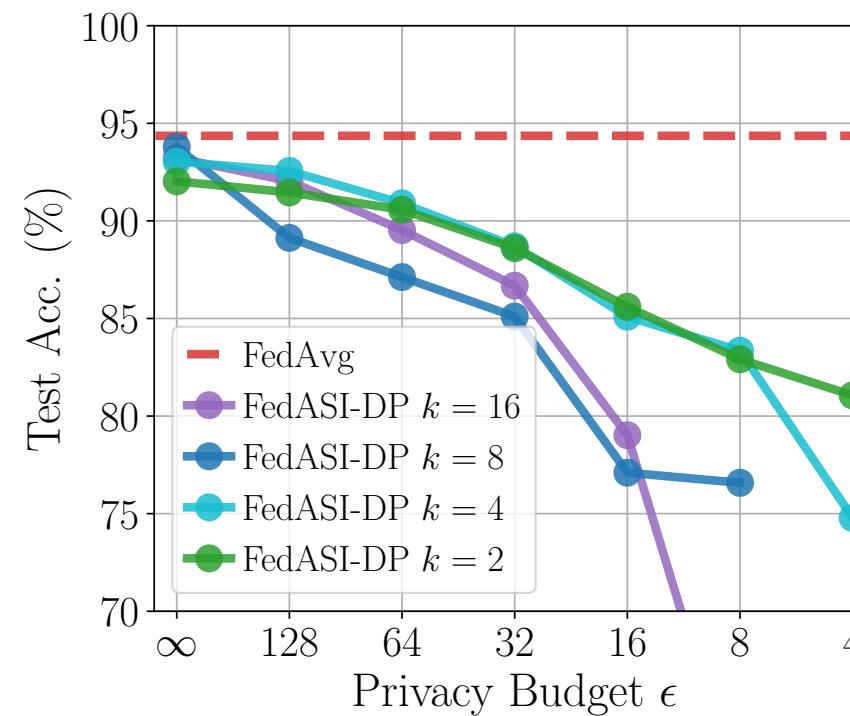
VGG16 Accum. Comm. Overhead vs. Round



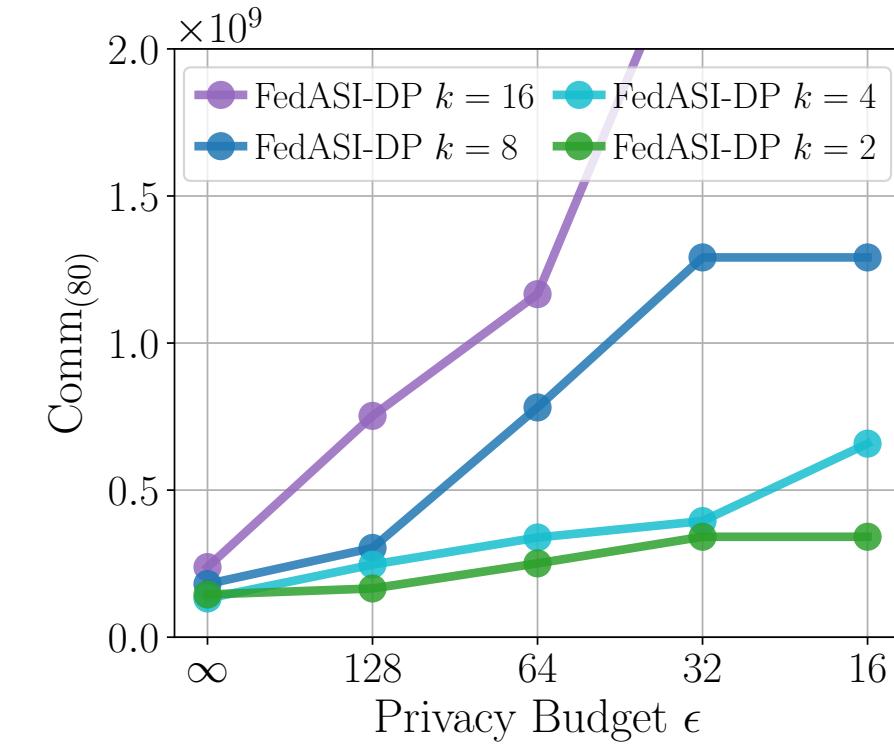
EXPERIMENT ON PRIVACY

FedASI-DP Results

- Lower rank more robust to small ϵ



ResNet18 Accuracy vs. Privacy Budget



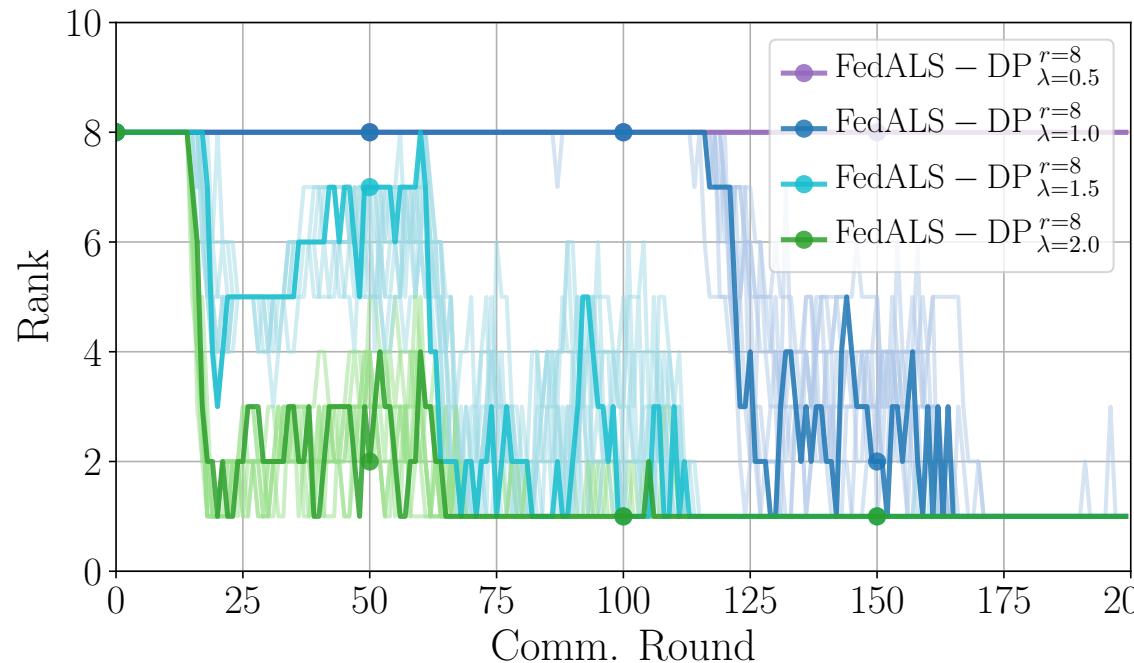
ResNet18 Converge Comm. vs. Privacy Budget



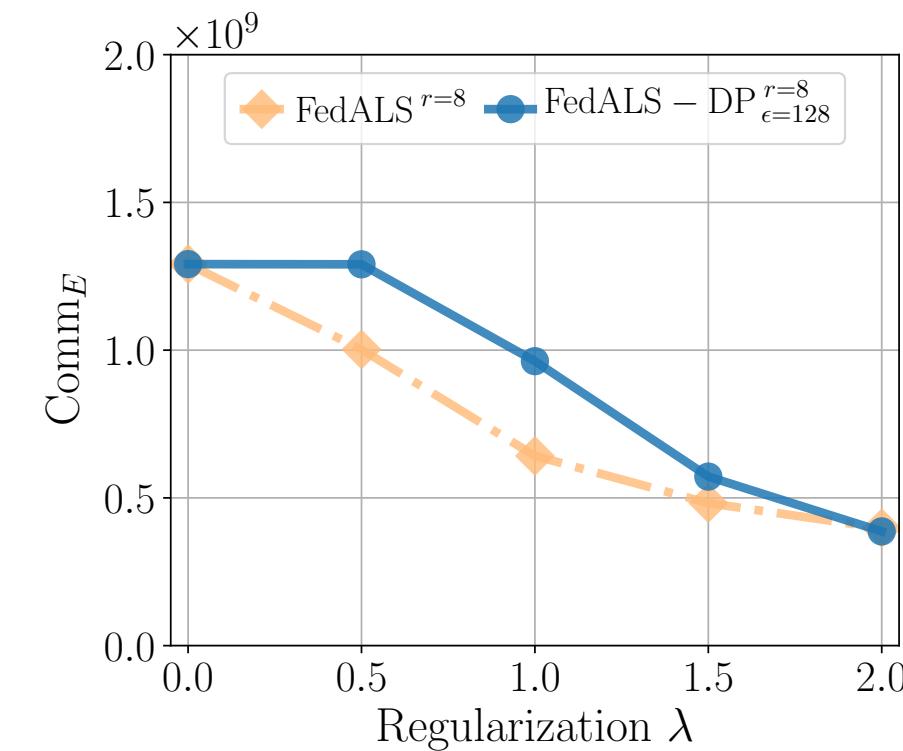
EXPERIMENT ON PRIVACY

FedALS-DP Results

- Adaptive rank factorization



VGG16 Layer Rank vs. Round



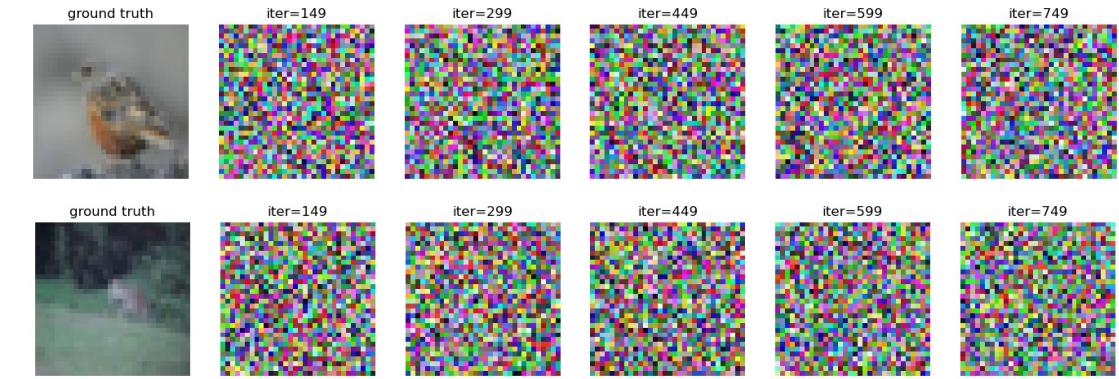
VGG16 Total Comm. vs. Regularization λ



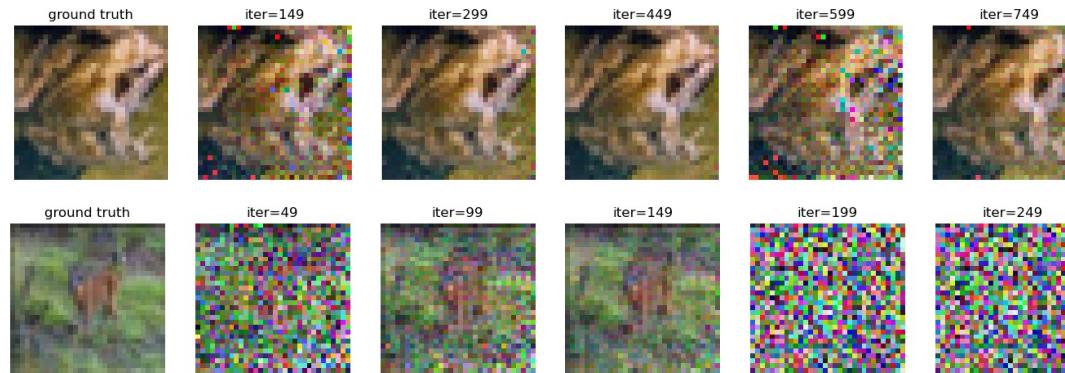
EXPERIMENT ON PRIVACY

Defense to DLG Attack

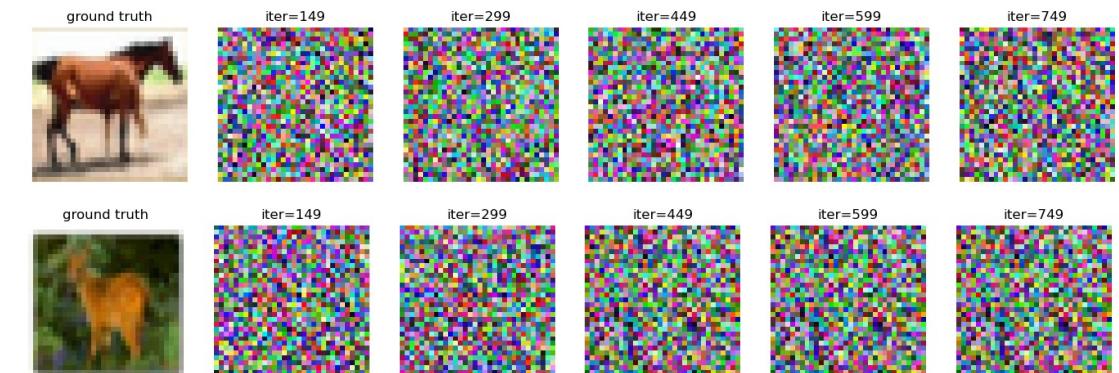
- Pruning: fail
- FedASI/FedALS: success
- FedASI-DP/FedALS-DP: success



truth DLG: FedASI-DP



truth DLG: Pruning



truth DLG: FedASI



CONCLUSION

FedASI & FedALS

Contribution

- Adaptive low-rank factorization
- FL communication compression
- Differential privacy guarantee

Theory

- Adaptive parameter compression by matrix rank optimization
- Differential privacy proof on algorithmic mechanism
- Convergence proof on low-rank factorization and FL training

Experiment

- Up to 1000x compression without performance loss
- Benefit compression and privacy simultaneously

THANK YOU

