

# Gradient Obfuscation Gives a False Sense of Security in Federated Learning

Kai Yue<sup>1</sup> Richeng Jin<sup>2</sup>

Chau-Wai Wong<sup>1</sup> Dror Baron<sup>1</sup> Huaiyu Dai<sup>1</sup>

<sup>1</sup>NC State University

<sup>2</sup>Zhejiang University

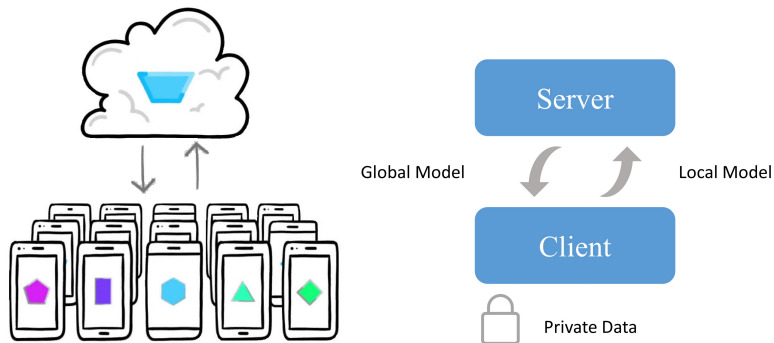
USENIX Security 2023

August 11, 2023

- 1 Background of Federated Learning
- 2 Reconstruction From Obfuscated Gradients
- 3 Case Studies
- 4 Conclusion

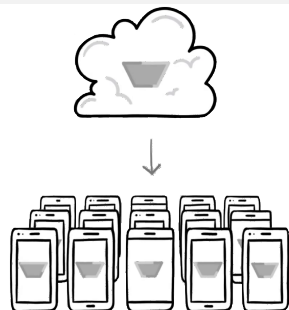
# Background: Federated Learning (FL)

- Clients with private data jointly solve a machine learning task
- Raw data stored locally & not exchanged



<sup>1</sup><https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>

## Background: Federated Averaging (FedAvg)



clients **download** the model

---

<sup>1</sup><https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>

## Background: Federated Averaging (FedAvg)



clients **download** the model

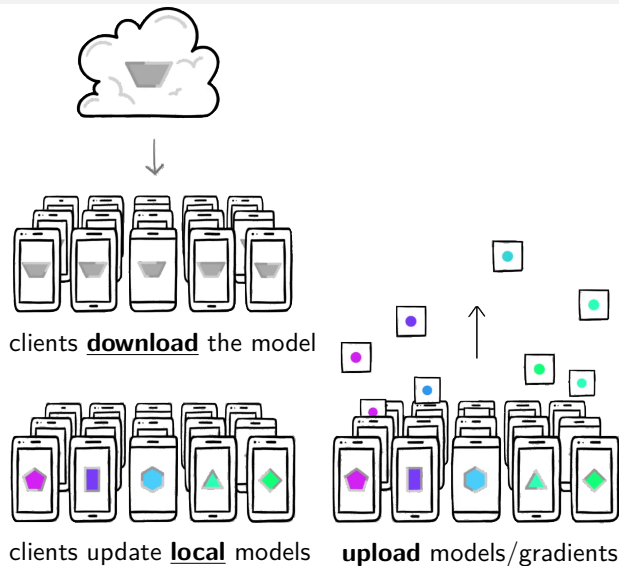


clients update **local** models

---

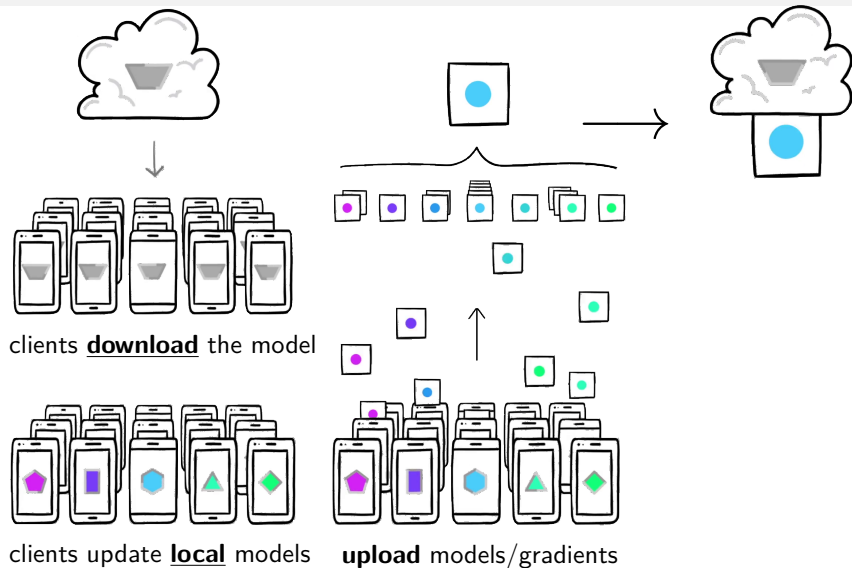
<sup>1</sup><https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>

# Background: Federated Averaging (FedAvg)



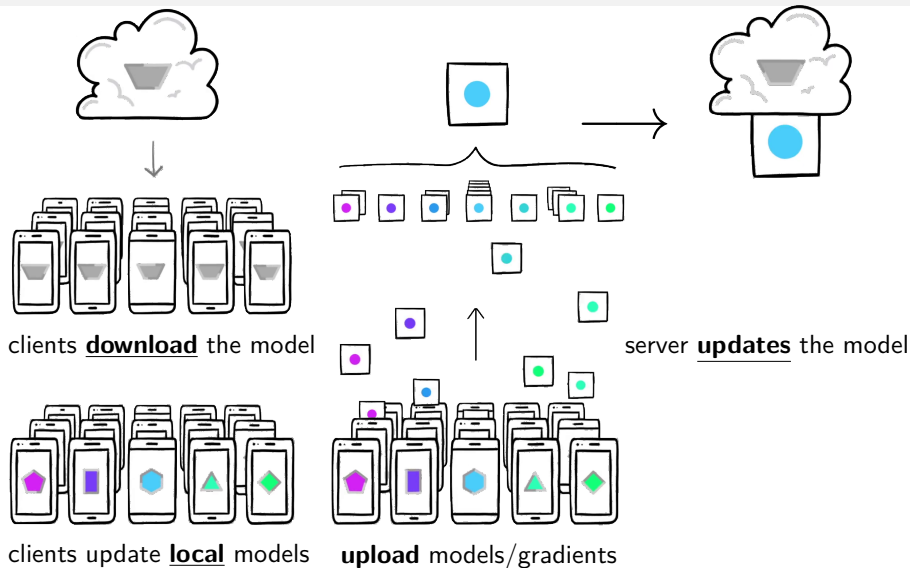
<sup>1</sup><https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>

# Background: Federated Averaging (FedAvg)



<sup>1</sup><https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>

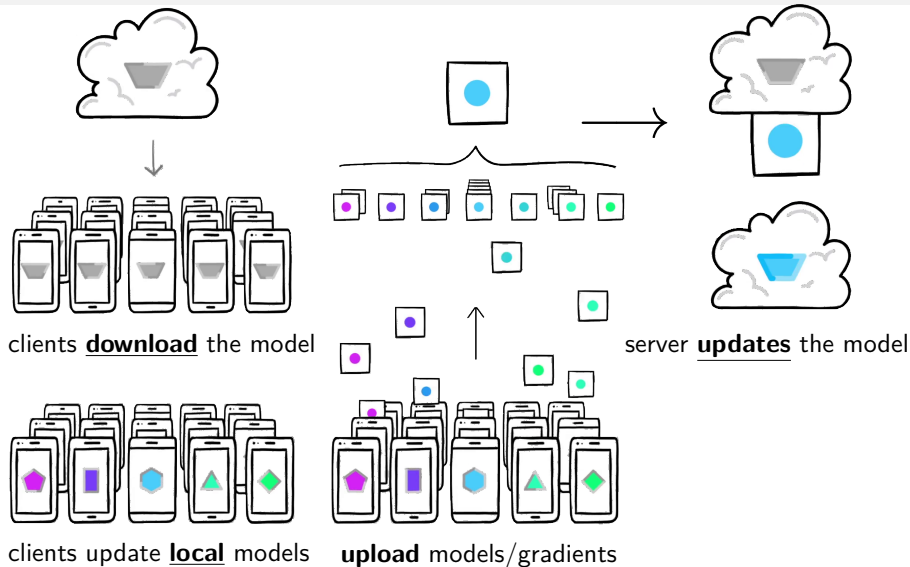
# Background: Federated Averaging (FedAvg)



<sup>1</sup><https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>

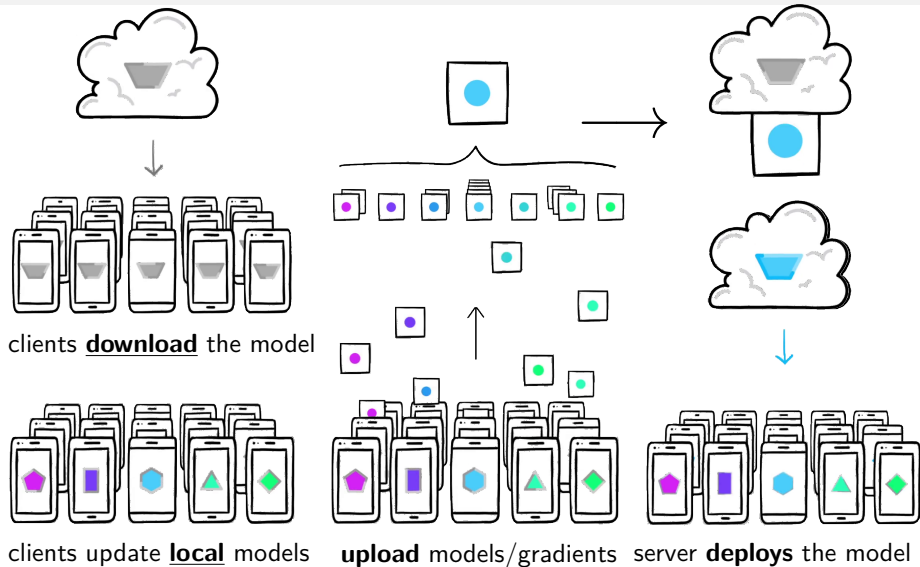


# Background: Federated Averaging (FedAvg)



<sup>1</sup><https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>

# Background: Federated Averaging (FedAvg)

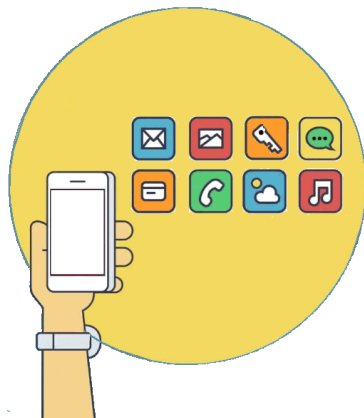


<sup>1</sup><https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>

# Background: Privacy Leakage in FL

- Can privacy be leaked in federated learning?

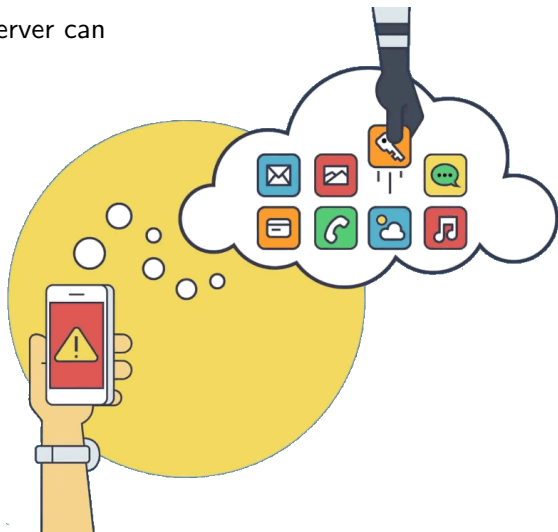
client: my data is kept locally



# Background: Privacy Leakage in FL

- An honest-but-curious server can recover private data<sup>2</sup>

client: my data has been stolen!

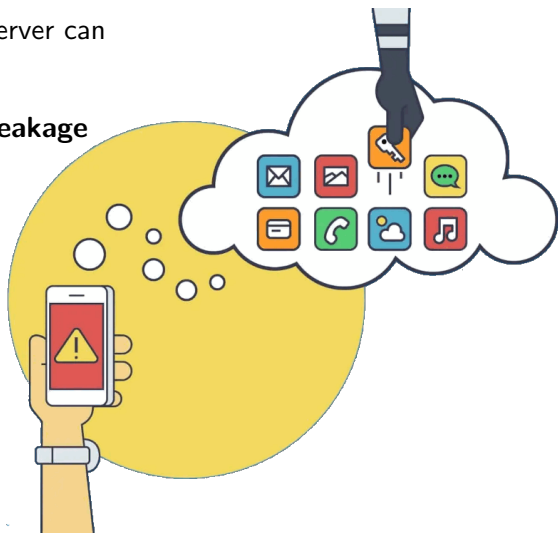


<sup>2</sup>Zhu, Ligeng, et al. "Deep Leakage from Gradients." Advances in Neural Information Processing Systems, 2019.

# Background: Privacy Leakage in FL

- An honest-but-curious server can recover private data<sup>2</sup>
- **How realistic is data leakage under obfuscation?**

client: my data has been stolen!



<sup>2</sup>Zhu, Ligeng, et al. "Deep Leakage from Gradients." Advances in Neural Information Processing Systems, 2019.

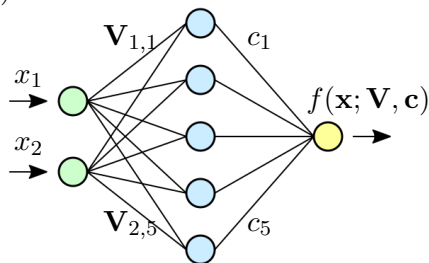
- 1 Background of Federated Learning
- 2 Reconstruction From Obfuscated Gradients**
- 3 Case Studies
- 4 Conclusion

# Raw Data Can be Reconstructed!

- Consider a neural network

- $f(\mathbf{x}; \mathbf{V}, \mathbf{c}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m c_r \sigma(\mathbf{v}_r^\top \mathbf{x})$

- input  $\mathbf{x} \in \mathbb{R}^2$   
weight  $\mathbf{V} \in \mathbb{R}^{2 \times 5}$



# Raw Data Can be Reconstructed!

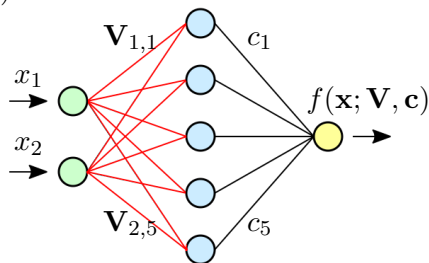
- Consider a neural network

- $f(\mathbf{x}; \mathbf{V}, \mathbf{c}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m c_r \sigma(\mathbf{v}_r^\top \mathbf{x})$

- input  $\mathbf{x} \in \mathbb{R}^2$   
weight  $\mathbf{V} \in \mathbb{R}^{2 \times 5}$

- System of equations

- 2 unknowns
- $2 \times 5$  equations





# Raw Data Can be Reconstructed!

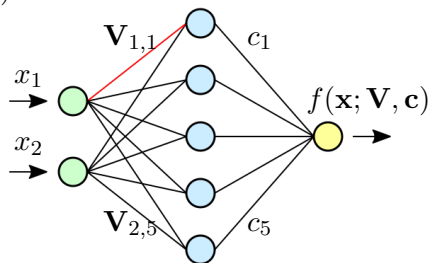
- Consider a neural network

- $f(\mathbf{x}; \mathbf{V}, \mathbf{c}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m c_r \sigma(\mathbf{v}_r^\top \mathbf{x})$

- input  $\mathbf{x} \in \mathbb{R}^2$   
weight  $\mathbf{V} \in \mathbb{R}^{2 \times 5}$

- System of equations

- 2 unknowns
- $2 \times 5$  equations



# Raw Data Can be Reconstructed!

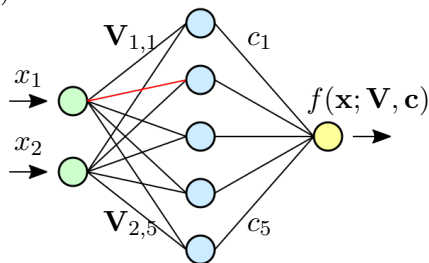
- Consider a neural network

- $f(\mathbf{x}; \mathbf{V}, \mathbf{c}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m c_r \sigma(\mathbf{v}_r^\top \mathbf{x})$

- input  $\mathbf{x} \in \mathbb{R}^2$   
weight  $\mathbf{V} \in \mathbb{R}^{2 \times 5}$

- System of equations

- 2 unknowns
- $2 \times 5$  equations



# Raw Data Can be Reconstructed!

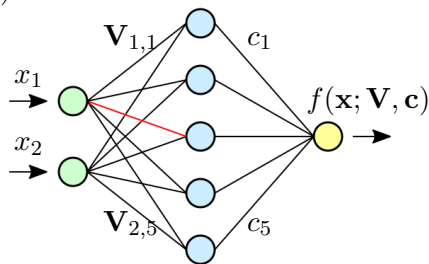
- Consider a neural network

- $f(\mathbf{x}; \mathbf{V}, \mathbf{c}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m c_r \sigma(\mathbf{v}_r^\top \mathbf{x})$

- input  $\mathbf{x} \in \mathbb{R}^2$   
weight  $\mathbf{V} \in \mathbb{R}^{2 \times 5}$

- System of equations

- 2 unknowns
- $2 \times 5$  equations



# Raw Data Can be Reconstructed!

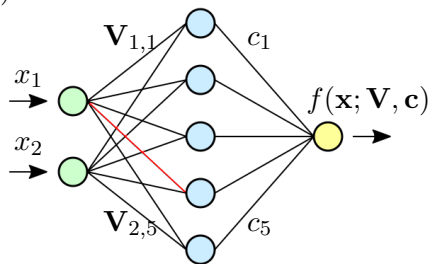
- Consider a neural network

- $f(\mathbf{x}; \mathbf{V}, \mathbf{c}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m c_r \sigma(\mathbf{v}_r^\top \mathbf{x})$

- input  $\mathbf{x} \in \mathbb{R}^2$   
weight  $\mathbf{V} \in \mathbb{R}^{2 \times 5}$

- System of equations

- 2 unknowns
  - $2 \times 5$  equations



# Raw Data Can be Reconstructed!

- Consider a neural network

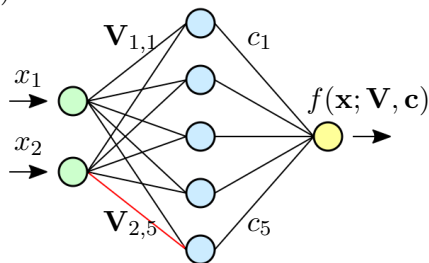
- $f(\mathbf{x}; \mathbf{V}, \mathbf{c}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m c_r \sigma(\mathbf{v}_r^\top \mathbf{x})$

- input  $\mathbf{x} \in \mathbb{R}^2$   
weight  $\mathbf{V} \in \mathbb{R}^{2 \times 5}$

- System of equations

- 2 unknowns
  - $2 \times 5$  equations

- Solve an overdetermined system



# Raw Data Can be Reconstructed!

- Consider a neural network

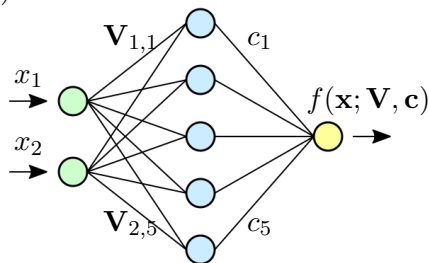
- $f(\mathbf{x}; \mathbf{V}, \mathbf{c}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m c_r \sigma(\mathbf{v}_r^\top \mathbf{x})$

- input  $\mathbf{x} \in \mathbb{R}^2$   
weight  $\mathbf{V} \in \mathbb{R}^{2 \times 5}$

- System of equations

- 2 unknowns
- $2 \times 5$  equations

- Solve an overdetermined system



# ROG: Reconstruct from Obfuscated Gradient

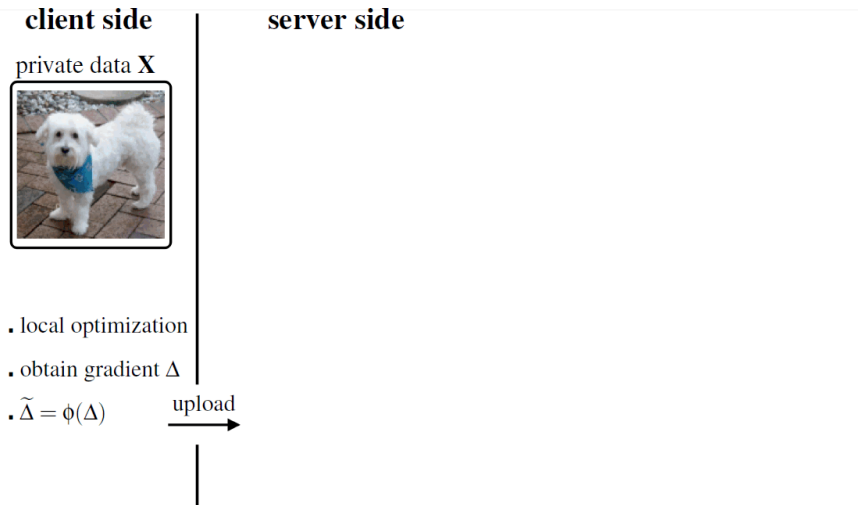
## client side

private data  $X$



- local optimization
- obtain gradient  $\Delta$
- $\tilde{\Delta} = \phi(\Delta)$

# ROG: Reconstruct from Obfuscated Gradient





# ROG: Reconstruct from Obfuscated Gradient

## client side

private data  $X$

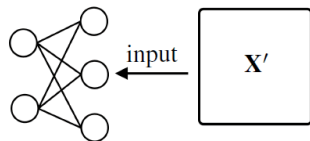


- local optimization
- obtain gradient  $\Delta$
- $\tilde{\Delta} = \phi(\Delta)$

upload

## server side

Prior work (Zhu et al., 2019)



# ROG: Reconstruct from Obfuscated Gradient

## client side

private data  $X$



- local optimization
- obtain gradient  $\Delta$
- $\tilde{\Delta} = \phi(\Delta)$

upload

## server side

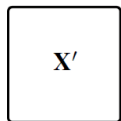
Prior work (Zhu et al., 2019)

$$\min_{X'} \|\Delta' - \tilde{\Delta}\|^2$$

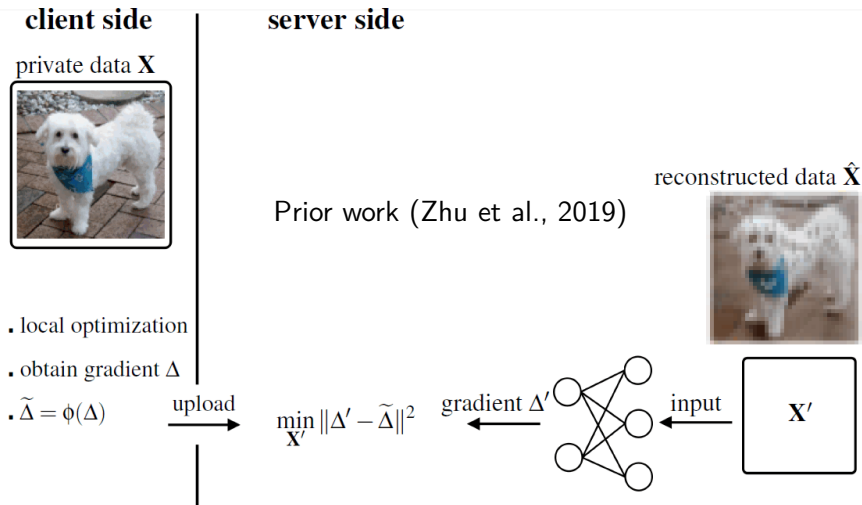
gradient  $\Delta'$



input



# ROG: Reconstruct from Obfuscated Gradient



# ROG: Reconstruct from Obfuscated Gradient

## client side

private data  $X$



- local optimization
- obtain gradient  $\Delta$
- $\tilde{\Delta} = \phi(\Delta)$

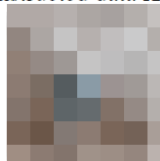
upload

## server side

Prior work (Zhu et al., 2019)

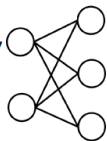
- does not work for FedAvg

reconstructed data  $\hat{X}$



$$\min_{X'} \|\Delta' - \tilde{\Delta}\|^2$$

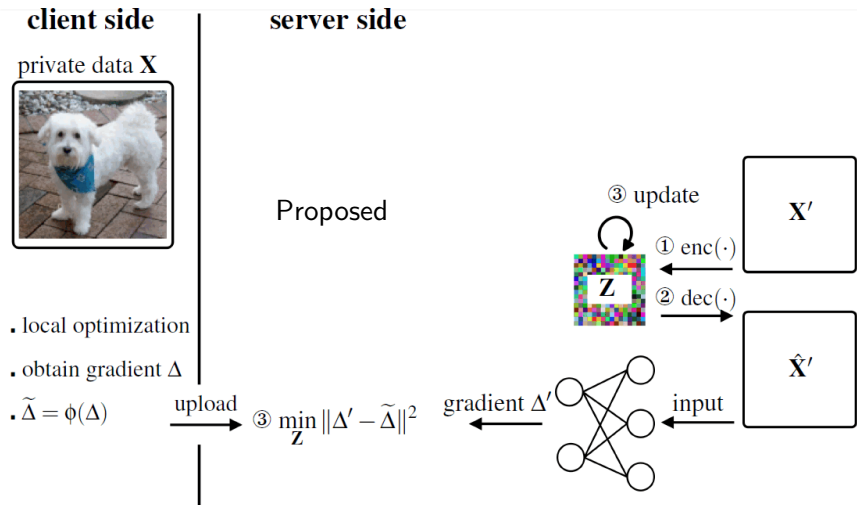
gradient  $\Delta'$



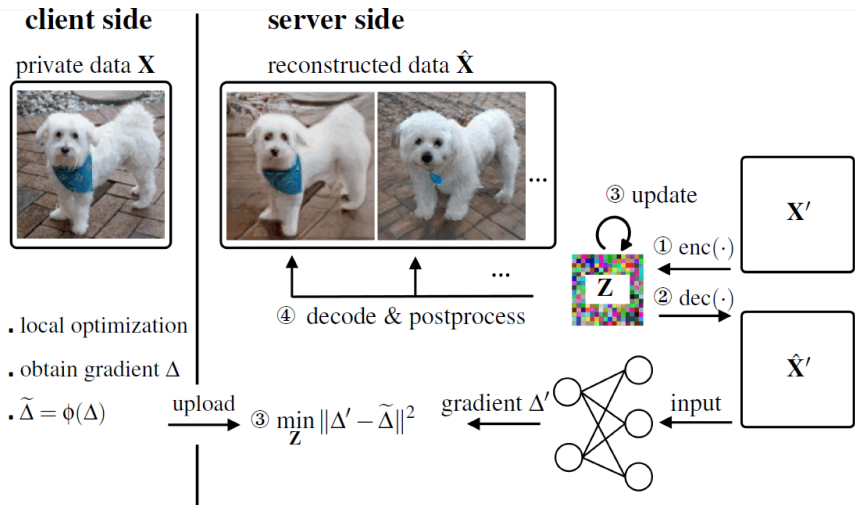
input

$X'$

# ROG: Reconstruct from Obfuscated Gradient



# ROG: Reconstruct from Obfuscated Gradient



- 1 Background of Federated Learning
- 2 Reconstruction From Obfuscated Gradients
- 3 Case Studies**
- 4 Conclusion

# Case Study: Does Gradient Compression Imply Privacy?



---

<sup>3</sup> Alistarh et al. "The Convergence of Sparsified Gradient Methods." NeurIPS 2018.

<sup>4</sup> Zhang et al. "The unreasonable effectiveness of deep features as a perceptual metric." CVPR 2018.



# Case Study: Does Gradient Compression Imply Privacy?



- Quantization<sup>3</sup>
  - FP32  $\rightarrow$  3 bits

---

<sup>3</sup> Alistarh et al. "The Convergence of Sparsified Gradient Methods." NeurIPS 2018.

<sup>4</sup> Zhang et al. "The unreasonable effectiveness of deep features as a perceptual metric." CVPR 2018.

# Case Study: Does Gradient Compression Imply Privacy?

Raw Images




3-bit  $\phi_q$



- Quantization<sup>3</sup>
  - FP32  $\rightarrow$  3 bits

<sup>3</sup> Alistarh et al. "The Convergence of Sparsified Gradient Methods." NeurIPS 2018.

<sup>4</sup> Zhang et al. "The unreasonable effectiveness of deep features as a perceptual metric." CVPR 2018. 

# Case Study: Does Gradient Compression Imply Privacy?

Raw Images



3-bit  $\phi_q$   
LPIPS 0.213



3-bit  $\phi_{qsgd}$   
LPIPS 0.215



- Quantization<sup>3</sup>
  - FP32  $\rightarrow$  3 bits

<sup>3</sup> Alistarh et al. "The Convergence of Sparsified Gradient Methods." NeurIPS 2018.

<sup>4</sup> Zhang et al. "The unreasonable effectiveness of deep features as a perceptual metric." CVPR 2018.

# Case Study: Does Gradient Compression Imply Privacy?

Raw Images



3-bit  $\phi_q$   
LPIPS 0.213



3-bit  $\phi_{qsgd}$   
LPIPS 0.215



- Quantization<sup>3</sup>
  - FP32  $\rightarrow$  3 bits
- Gradient sparsification
  - zero out smaller components

<sup>3</sup> Alistarh et al. "The Convergence of Sparsified Gradient Methods." NeurIPS 2018.

<sup>4</sup> Zhang et al. "The unreasonable effectiveness of deep features as a perceptual metric." CVPR 2018.

# Case Study: Does Gradient Compression Imply Privacy?

Raw Images



3-bit  $\phi_q$   
LPIPS 0.213



3-bit  $\phi_{qsgd}$   
LPIPS 0.215



Top-k (0.95)  
LPIPS 0.237



- Quantization<sup>3</sup>
  - FP32  $\rightarrow$  3 bits
- Gradient sparsification
  - zero out smaller components

<sup>3</sup> Alistarh et al. "The Convergence of Sparsified Gradient Methods." NeurIPS 2018.

<sup>4</sup> Zhang et al. "The unreasonable effectiveness of deep features as a perceptual metric." CVPR 2018.

# Case Study: Does Gradient Compression Imply Privacy?



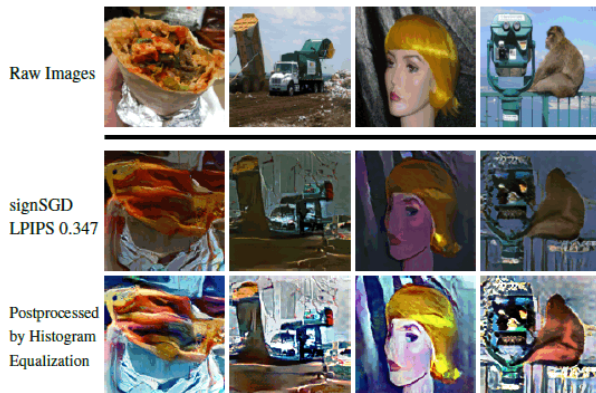
- Quantization<sup>3</sup>
  - FP32  $\rightarrow$  3 bits
- Gradient sparsification
  - zero out smaller components
- LPIPS<sup>4</sup> (Learned Perceptual Image Patch Similarity)  $\downarrow$

<sup>3</sup> Alistarh et al. "The Convergence of Sparsified Gradient Methods." NeurIPS 2018.

<sup>4</sup> Zhang et al. "The unreasonable effectiveness of deep features as a perceptual metric." CVPR 2018.

# Case Study: Binary Quantization

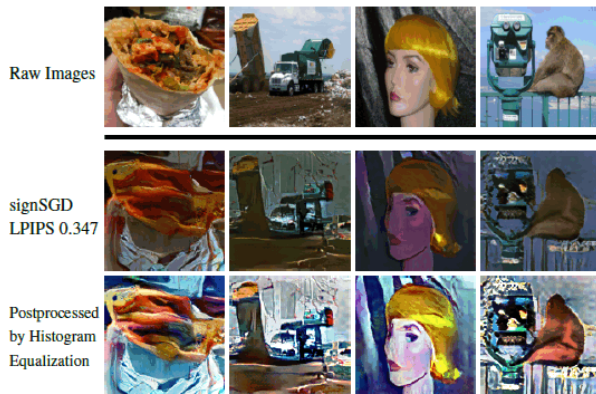
- One bit quantization: signSGD<sup>5</sup>
  - Each entry of the gradient is represented by its sign



<sup>5</sup> Bernstein, Jeremy, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. 2018. "SignSGD: Compressed Optimisation for Non-Convex Problems." ICML, 2018.

# Case Study: Binary Quantization

- One bit quantization: signSGD<sup>5</sup>
  - Gradient compression does **not** imply privacy



<sup>5</sup> Bernstein, Jeremy, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. 2018. "SignSGD: Compressed Optimisation for Non-Convex Problems." ICML, 2018.



# Defense: Differential Privacy<sup>6</sup>

- Attack differentially private training



Table 1:

Test Accuracy on CIFAR-10

SNR (dB)	Acc. (%)	LPIPS
$\infty$	69	0.18

<sup>6</sup>Wei et al. "Gradient-leakage resilient federated learning." ICDCS 2021.

# Defense: Differential Privacy<sup>6</sup>

- Attack differentially private training

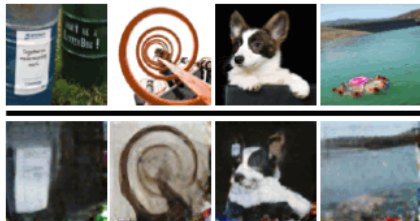


Table 1:

Tradeoff Between Utility & Privacy

SNR (dB)	Acc. (%)	LPIPS
$\infty$	69	0.18
0	55	0.38

<sup>6</sup>Wei et al. "Gradient-leakage resilient federated learning." ICDCS 2021.

# Defense: Differential Privacy<sup>6</sup>

- Attack differentially private training

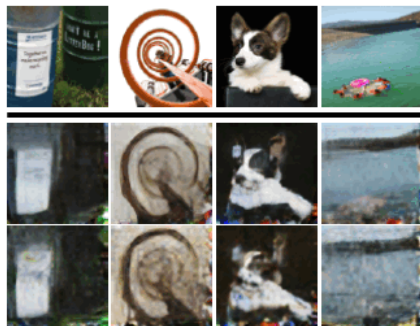


Table 1:

Tradeoff Between Utility & Privacy

SNR (dB)	Acc. (%)	LPIPS
$\infty$	69	0.18
0	55	0.38
-5	52	0.46

<sup>6</sup>Wei et al. "Gradient-leakage resilient federated learning." ICDCS 2021.

# Defense: Differential Privacy<sup>6</sup>

- Attack differentially private training

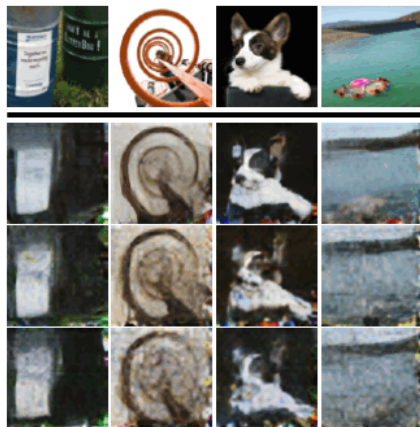


Table 1:

Tradeoff Between Utility & Privacy

SNR (dB)	Acc. (%)	LPIPS
$\infty$	69	0.18
0	55	0.38
-5	52	0.46
-10	50	0.55

<sup>6</sup>Wei et al. "Gradient-leakage resilient federated learning." ICDCS 2021.

# Defense: Differential Privacy<sup>6</sup>

- Attack differentially private training

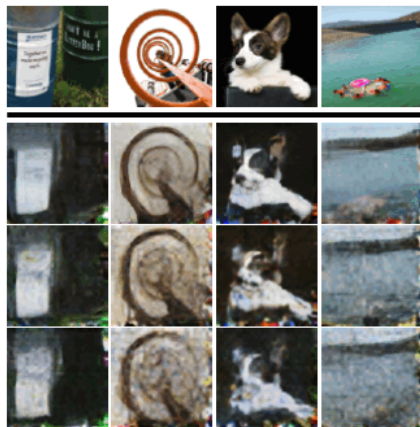


Table 1:

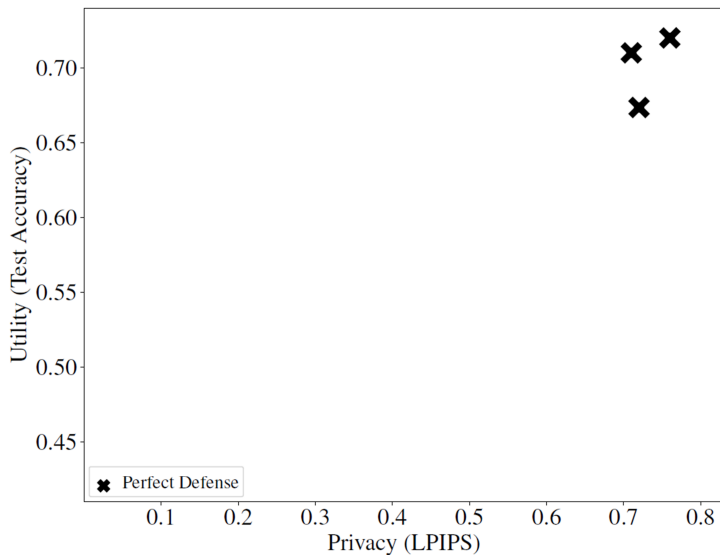
Tradeoff Between Utility & Privacy

SNR (dB)	Acc. (%)	LPIPS
$\infty$	69	0.18
0	55	0.38
-5	52	0.46
-10	50	0.55

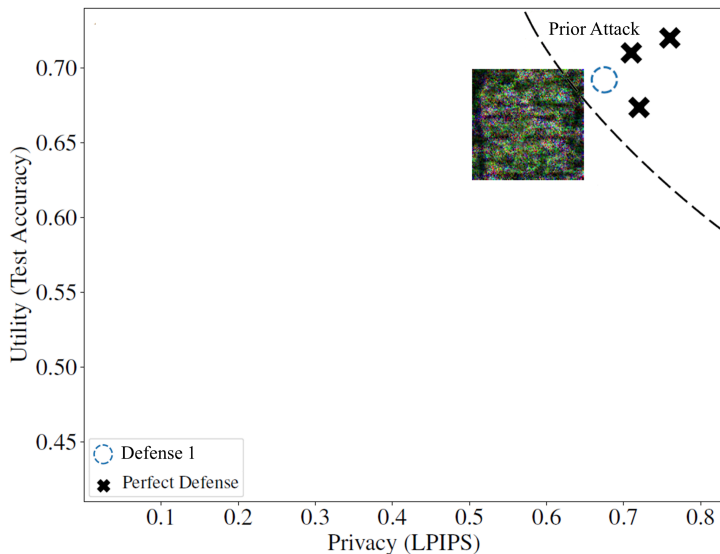
privacy  $\rightarrow$  utility  $\downarrow\downarrow$  😞

<sup>6</sup>Wei et al. "Gradient-leakage resilient federated learning." ICDCS 2021.

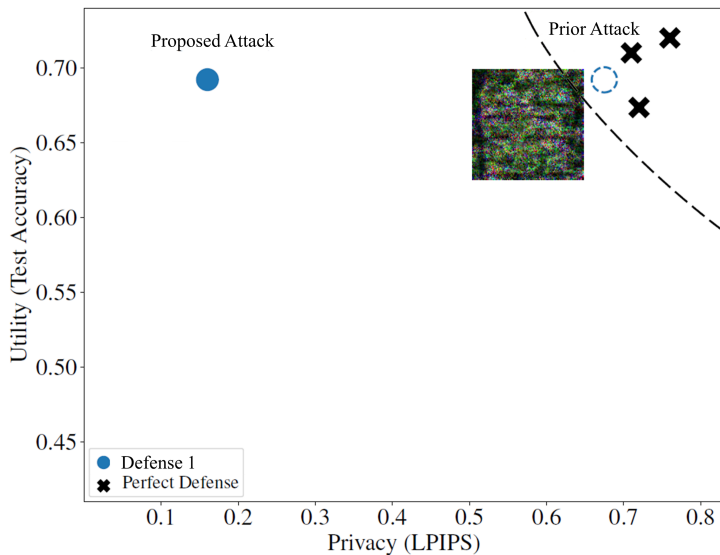
# Improved Attack Efficiency



# Improved Attack Efficiency

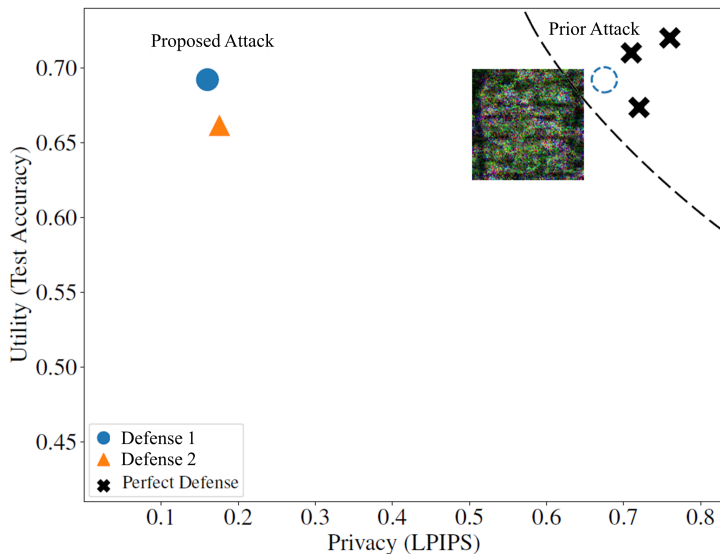


# Improved Attack Efficiency

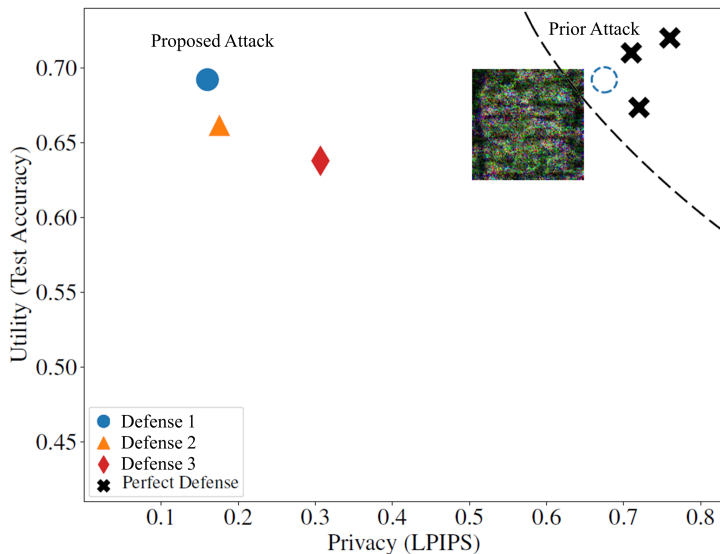




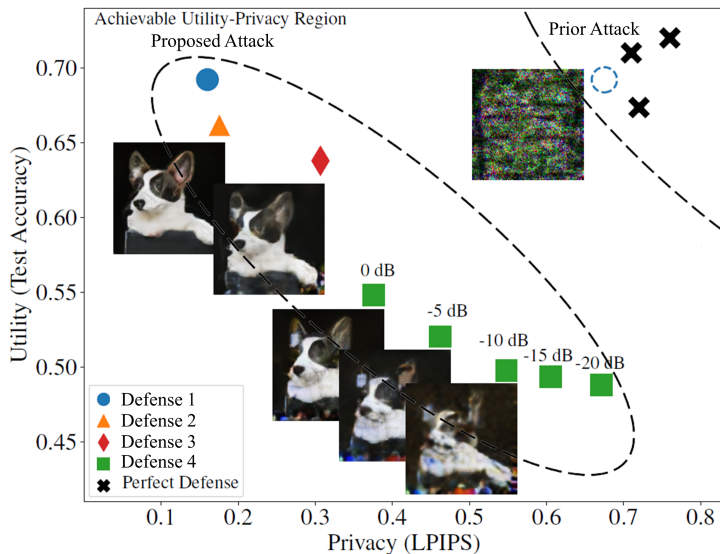
# Improved Attack Efficiency



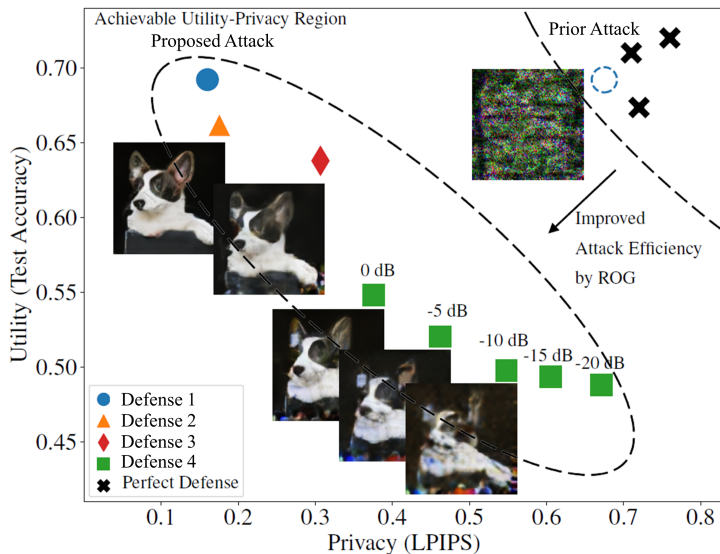
# Improved Attack Efficiency



# Improved Attack Efficiency



# Improved Attack Efficiency



# Semantic-Level Attack Variant: ROGS

Raw Image



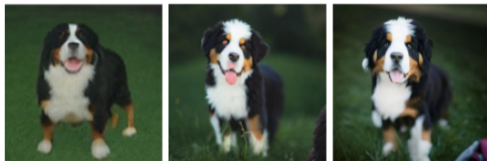
What if the reconstruction is not clear enough?

# Semantic-Level Attack Variant: ROGS

Raw Image



Reconstructed Image at a Semantic Level



# Semantic-Level Attack Variant: ROGS

Raw Image



Reconstructed Image at a Semantic Level



# Semantic-Level Attack Variant: ROGS

Raw Image



Reconstructed Image at a Semantic Level



Find more details in our paper!



- 1 Background of Federated Learning
- 2 Reconstruction From Obfuscated Gradients
- 3 Case Studies
- 4 Conclusion**

# Conclusion

- We have studied privacy leakage under various obfuscation mechanisms in federated learning

# Conclusion

- We have studied privacy leakage under various obfuscation mechanisms in federated learning
- We proposed an attack framework, which shows that
  - ▷ compressed gradient: does not protect privacy inherently
  - ▷ noisy gradient (DP): bad tradeoff between utility & privacy

# Conclusion

- We have studied privacy leakage under various obfuscation mechanisms in federated learning
- We proposed an attack framework, which shows that
  - ▷ compressed gradient: does not protect privacy inherently
  - ▷ noisy gradient (DP): bad tradeoff between utility & privacy
- Future work: hybrid defense & more precise privacy notion

# Gradient Obfuscation Gives a False Sense of Security in Federated Learning

Kai Yue<sup>1</sup> Richeng Jin<sup>2</sup>

Chau-Wai Wong<sup>1</sup> Dror Baron<sup>1</sup> Huaiyu Dai<sup>1</sup>

<sup>1</sup>NC State University

<sup>2</sup>Zhejiang University

USENIX Security 2023

## Thanks!

August 11, 2023