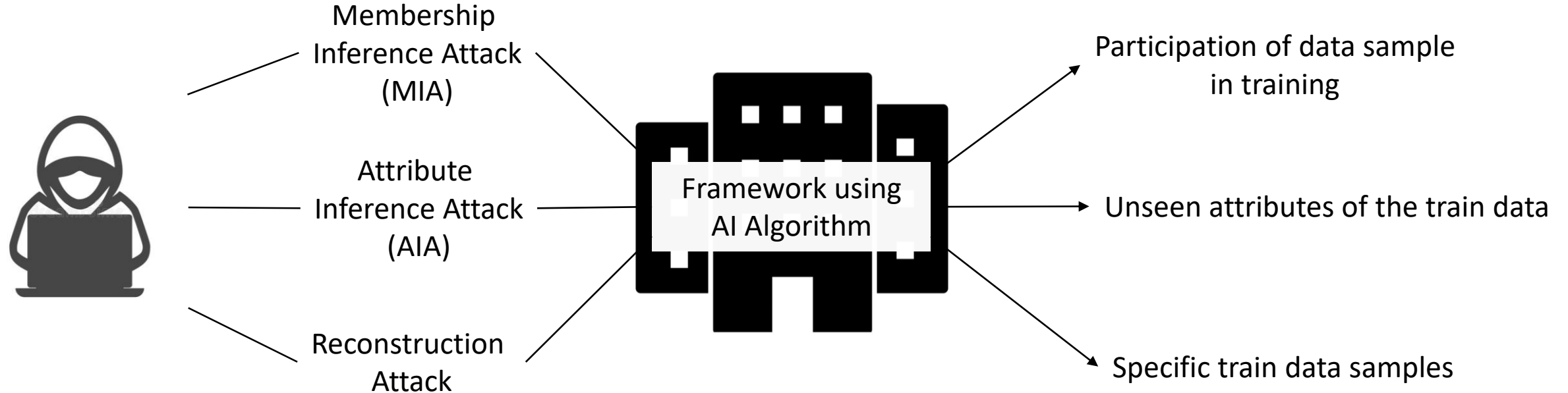# Closing the Loophole:
# Rethinking Reconstruction Attacks in Federated Learning from a Privacy Standpoint

SEUNG HO NA, HYEONG GWON HONG, JUNMO KIM, SEUNGWON SHIN

KAIST

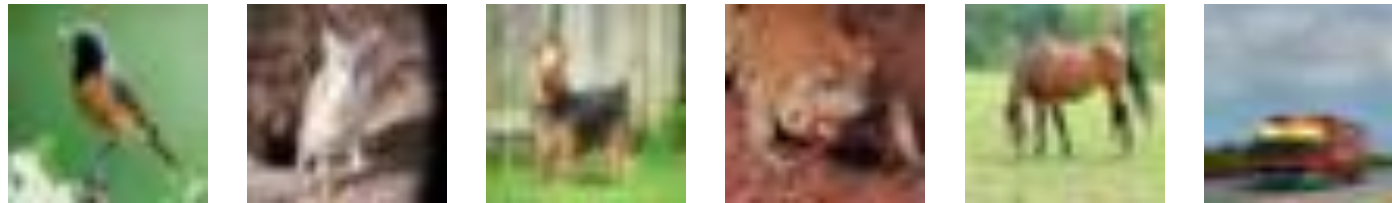# Privacy Attacks in Federated Learning (FL)

- Privacy Attack
  - Attacks aiming at leaking private information

Membership Inference Attack (MIA) → Participation of data sample in training

Attribute Inference Attack (AIA) → Unseen attributes of the train data

Reconstruction Attack → Specific train data samples

Framework using AI Algorithm

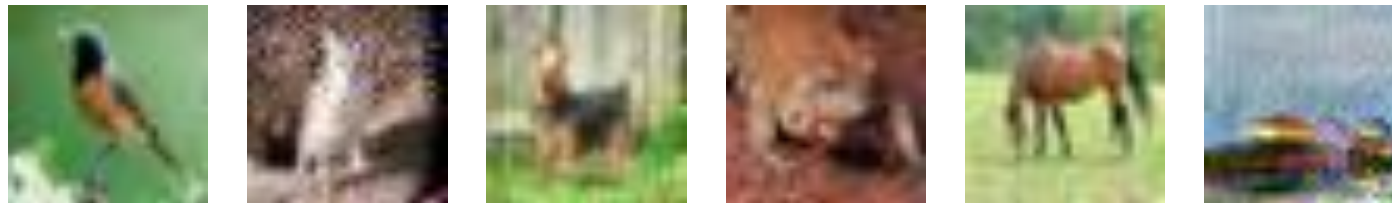$NS^2$ Network and System Security Laboratory KAIST

# Privacy-preserving Technologies

- Differential privacy
  - Theoretical approach to quantifying information leakage

- Encryption methods
  - Key encryption schemes such as secure multi-party computation protocols
  - Incur heavy computation and communication costs

Original Image



Reconstruction Attack on Differential Privacy Model

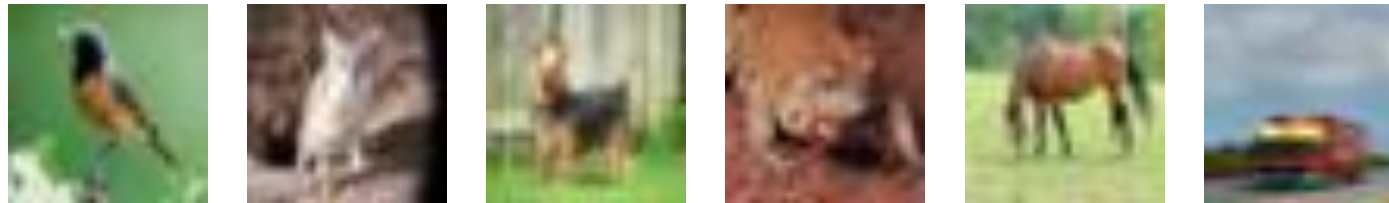$NS^2$ Network and System Security Laboratory  KAIST

# Privacy-preserving Technologies

- Differential privacy
  - Theoretical approach to quantifying information leakage
- Encryption methods
  - Key encryption schemes such as secure multi-party computation protocols
  - Incur heavy computation and communication costs

## Research Question 1

What is this inconsistency between privacy attacks and privacy-preserving methods?

Original Image

Reconstruction Attack on Differential Privacy Model

NS² Network and System Security Laboratory  KAIST

# Privacy-preserving Technologies

- Differential privacy
  - Theoretical approach to quantifying information leakage
- Encryption methods
  - Key encryption schemes, secure multi-party computation protocols
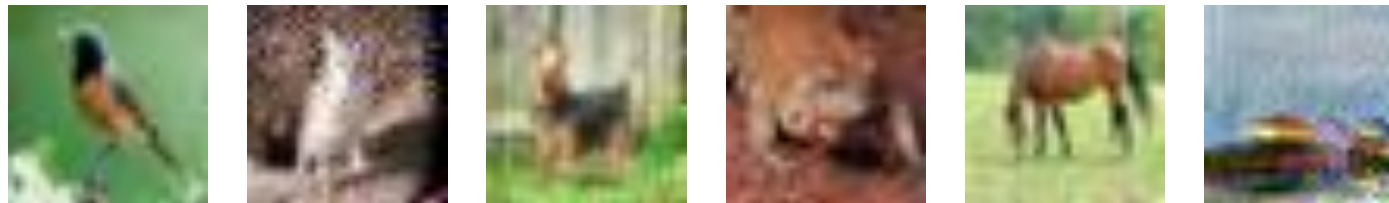  - Incur heavy computation and communication costs

Original Image

Reconstruction Attack
on Differential
Privacy Model

## Research Question 1

What is this <span style="color:red">inconsistency</span> between privacy attacks and privacy-preserving methods?

➡ Dissect privacy attacks by their attributes

<span style="color:yellow">Extraction Extent</span>: What is the private information?
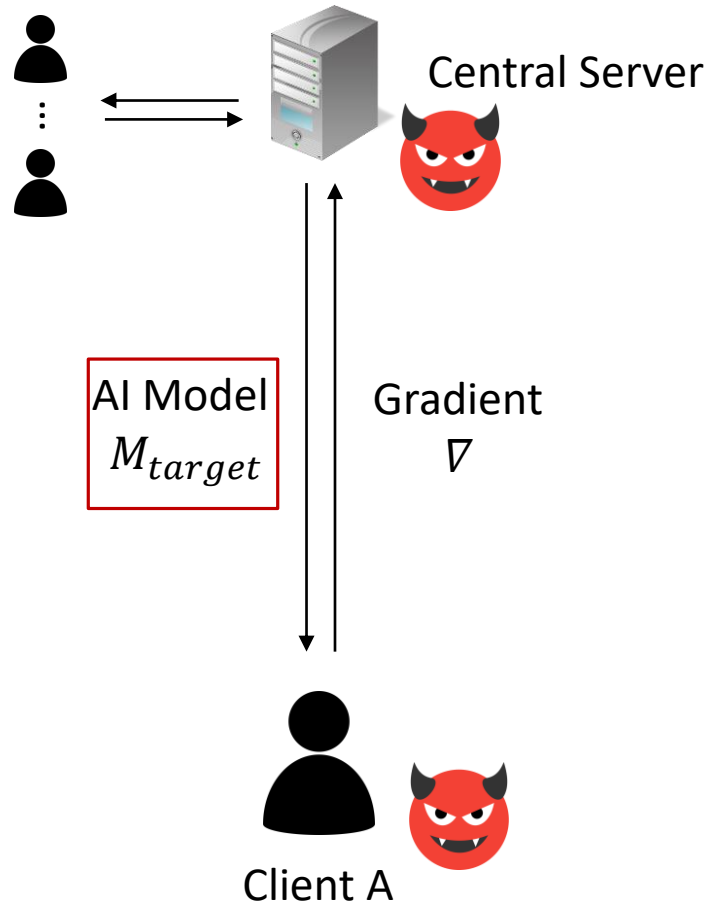<span style="color:yellow">Transferability</span>: At what scale can this attack take place?
<span style="color:yellow">Source</span>: What object allows the attack?

# Membership Inference Attack in FL



$D_{aux}$, Data from similar distribution as train data

AI Model $M_{target}$

Inference Model $M_{MIA}$

Generalization Gap

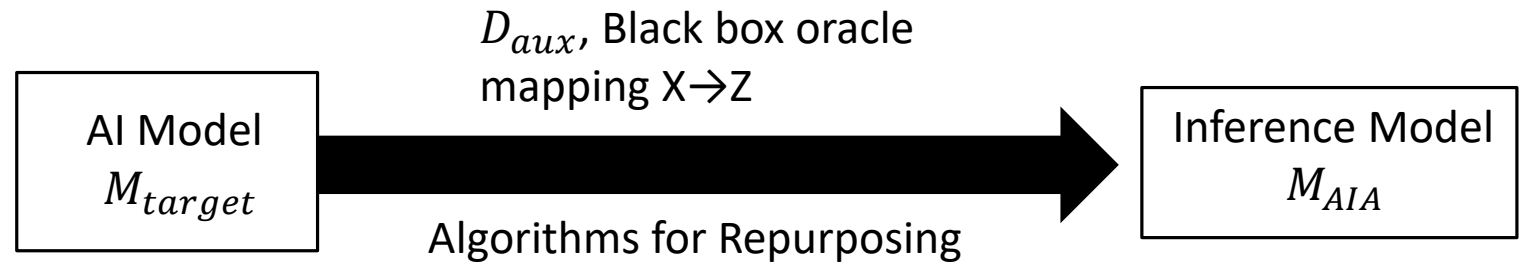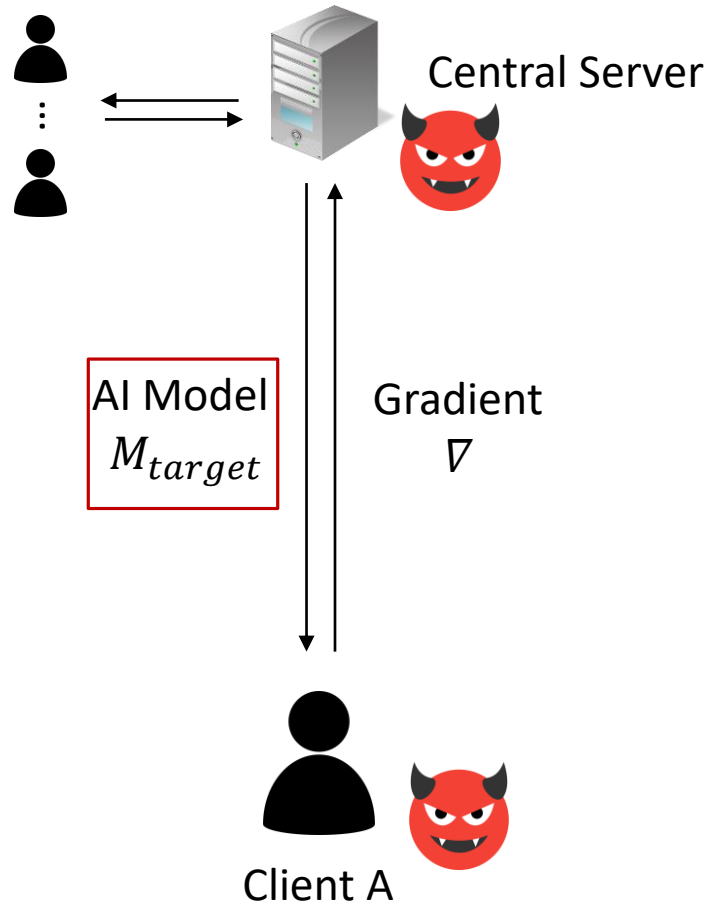Central Server

AI Model $M_{target}$

Gradient $\nabla$

Client A

- Develop inference model to know data membership

- $M_{MIA}(x, M_{target}, D_{aux}) = \begin{cases} 1, x \in D_{train} \\ 0, x \in D_{test} \end{cases}$

Extraction Extent: private meta-information
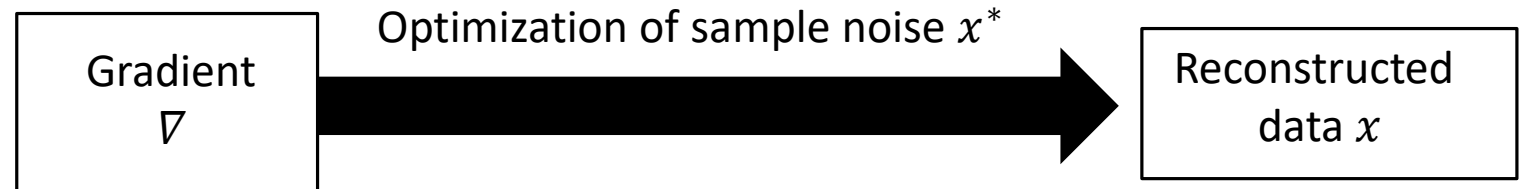Transferability: membership of all data X can be known
Source: $M_{target}$

$NS^2$ Network and System Security Laboratory KAIST

# Attribute Inference Attacks in FL



AI Model $M_{target}$

Gradient $\nabla$

Client A

Central Server

$D_{aux}$, Black box oracle mapping X→Z

AI Model $M_{target}$ → Algorithms for Repurposing → Inference Model $M_{AIA}$
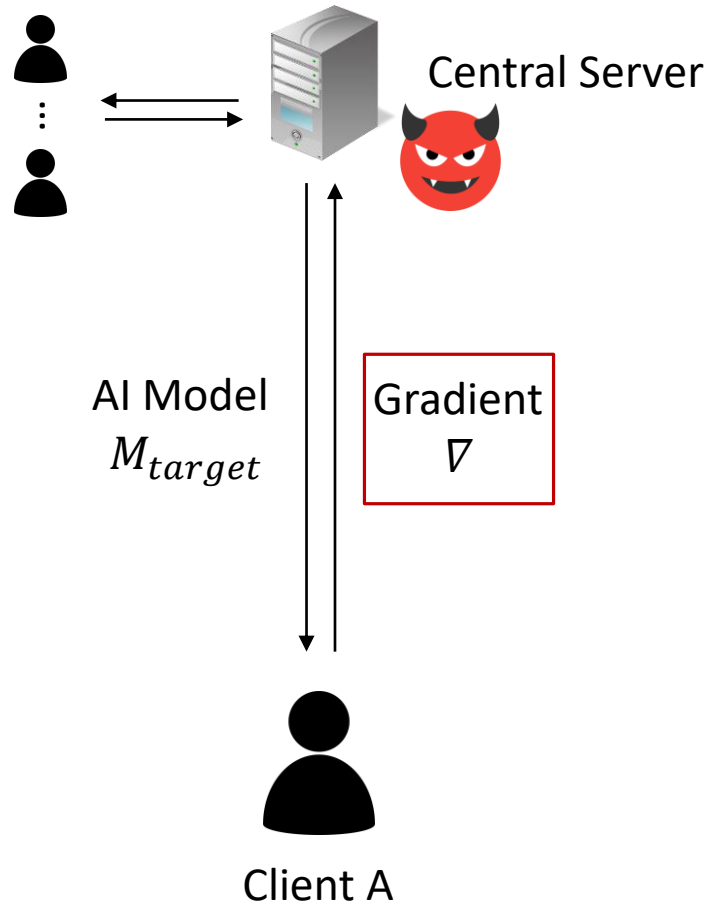
- Develop inference model to exploit unseen characteristics of the data

- Given $M_{target}(x_i) = y_i$ ,

- $M_{AIA}\big(x_i, M_{target}, D_{aux}\big) = z_i$

Extraction Extent: private meta-information
Transferability: attributes of all data X can be known
Source: $M_{target}$

$NS^2$ Network and System Security Laboratory  KAIST

# Reconstruction Attack Attributes in FL

Central Server

AI Model $M_{target}$

Gradient $\nabla$

Client A

Gradient $\nabla$

Optimization of sample noise $x^*$

Reconstructed data $x$

- Recover original data from gradient information

- $\underset{x^*}{\mathrm{argmin}}\ 1 - \dfrac{<\nabla L_M(x,y), \nabla L_M(x^*,y)>}{|\nabla L_M(x,y)||\nabla L_M(x^*,y)|} + TV(x^*)$

Extraction Extent: private raw data
Transferability: attack specific to the input gradient
Source: $\nabla$

# Reconstruction Attack Attributes in FL



Central Server

AI Model $M_{target}$

Gradient $\nabla$

Client A

Gradient $\nabla$ → Optimization of sample noise $x^*$ → Reconstructed data $x$

- Recover original data from gradient information

- $\underset{x^*}{\text{argmin}} \; 1 - \frac{<\nabla L_M(x,y), \nabla L_M(x^*,y)>}{|\nabla L_M(x,y)||\nabla L_M(x^*,y)|} + TV(x^*)$

Extraction Extent: private raw data
Transferability: attack specific to the input gradient
Source: $\nabla$

# Breakdown of Privacy

### Disclosure Privacy

"Privacy that ensures that any information cannot be inferred from the collaborative result"

### Distinctive Privacy

"Privacy that ensures that the raw data will be secure and safe from exposure"

FL Model

MIA, AIA

**Disclosure Privacy**

Meta-information

Membership, Attribute info.          Gradient info.

Reconstruction Attack

**Distinctive Privacy**

Raw data

NS² Network and System Security Laboratory  KAIST

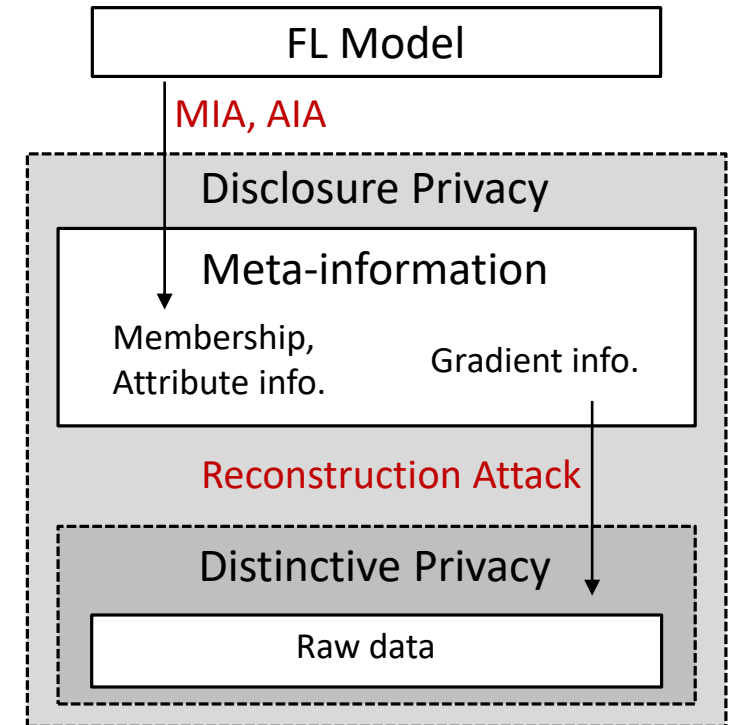# Breakdown of Privacy

### Disclosure Privacy

**In FL**, the trained model should not leak any form of participant information (meta-info.).

### Distinctive Privacy

**In FL**, the client data should be safe from reconstruction attempts.

FL Model

MIA, AIA

Disclosure Privacy

Meta-information

Membership, Attribute info.

Gradient info.

Reconstruction Attack

Distinctive Privacy

Raw data

# Breakdown of Privacy

### Disclosure Privacy

By definition, differential privacy preserves disclosure privacy by training a safer model.

### Distinctive Privacy

To conceal gradient information, encryption protocols accompanied by *computation and communication overhead* are used.



FL Model

MIA, AIA

Disclosure Privacy

Meta-information

Membership, Attribute info.          Gradient info.

Reconstruction Attack
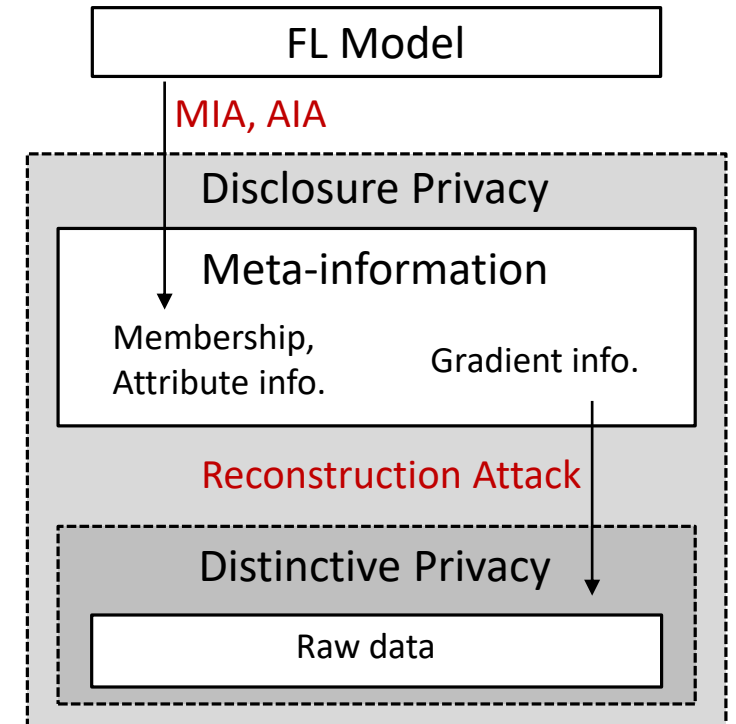
Distinctive Privacy

Raw data

# Breakdown of Privacy

Disclosure Privacy

By definition, differential privacy preserves disclosure privacy by training a safer model.

Distinctive Privacy

To conceal gradient information, encryption protocols accompanied by *computation and communication overhead* are used.

FL Model

MIA, AIA

Disclosure Privacy

Meta-information

Membership, Attribute info.

Gradient info.

Reconstruction Attack

Distinctive Privacy

Raw data

## Research Question 2

Is there a light method of ensuring distinctive privacy?

NS² Network and System Security Laboratory  KAIST

# Breakdown of Privacy
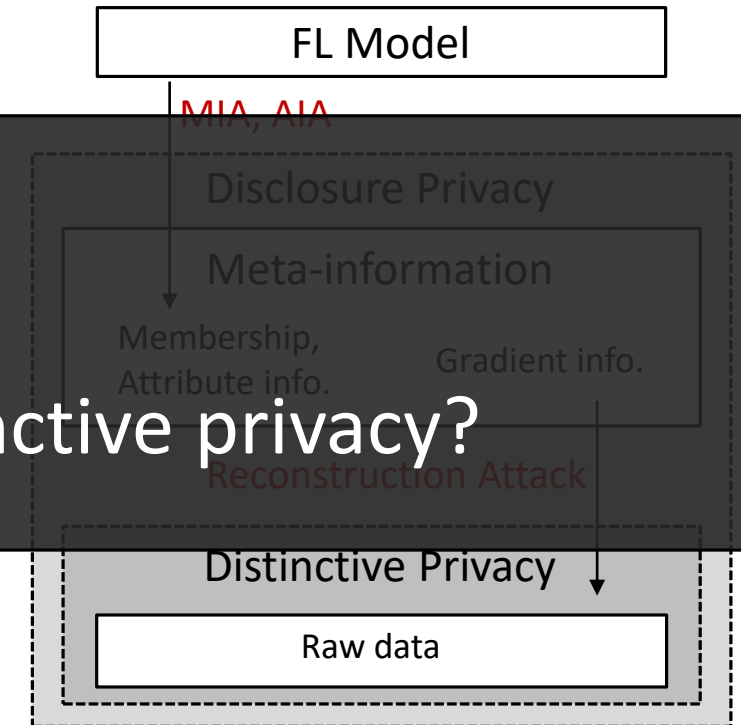
Disclosure Privacy

FL Model

By definition, differential privacy preserves disclosure privacy by training a safer model.

## Research Question 2

Is there a **light** method of ensuring distinctive privacy?

➡️ Select **safe layers** for exposure and **mask the gradient** information

# Obscuring Client Gradients Problem

To ensure distinctive privacy and prevent reconstruction, *mask the gradient information*

Problem: Finding obscuring function $f$ that obscures the gradient $\nabla$ such that:

| Robustness | Light | Trade-off |
|---|---|---|
| $X\left(Recon(f(\nabla))\right) > X\left(f(\nabla)\right)$ for defense capability $X$ (e.g., MSE, PSNR) | $Cost(f(\nabla)) \leq Cost(\nabla)$ in terms of communication cost | Allows adjustment in trade-off of model performance and defense capability. |

# Intuition: Global Gradient



Central Server

$\nabla_{Global}$ 🟩🟦🟦

AI Model
$M_{target}$

Gradient
$\nabla_A$ 🟩🟨🟦🟪

Client A

Aggregated gradient $\nabla_{global} = \frac{1}{N}\sum_{i=1}^{N}\nabla_i$

By being closer to $\nabla_{global}$, the more generalized the gradient is.

$\nabla_{Global}$ 🟩🟦🟩🟦🟦

$\nabla_A$ 🟩🟨🟩🟦🟪

$f(\nabla_A)$

Selecting similar layers with selection ratio 0.6

# Fragmented Federated Learning (FFL)

Central Server

$\nabla_{Global}$

Aggregated gradient $\nabla_{global} = \sum_{i=0}^{N} \nabla_i$

By being closer to $\nabla_{global}$, the more generalized the gradient is.

AI Model
$M_{target}$

Obscured
Gradient
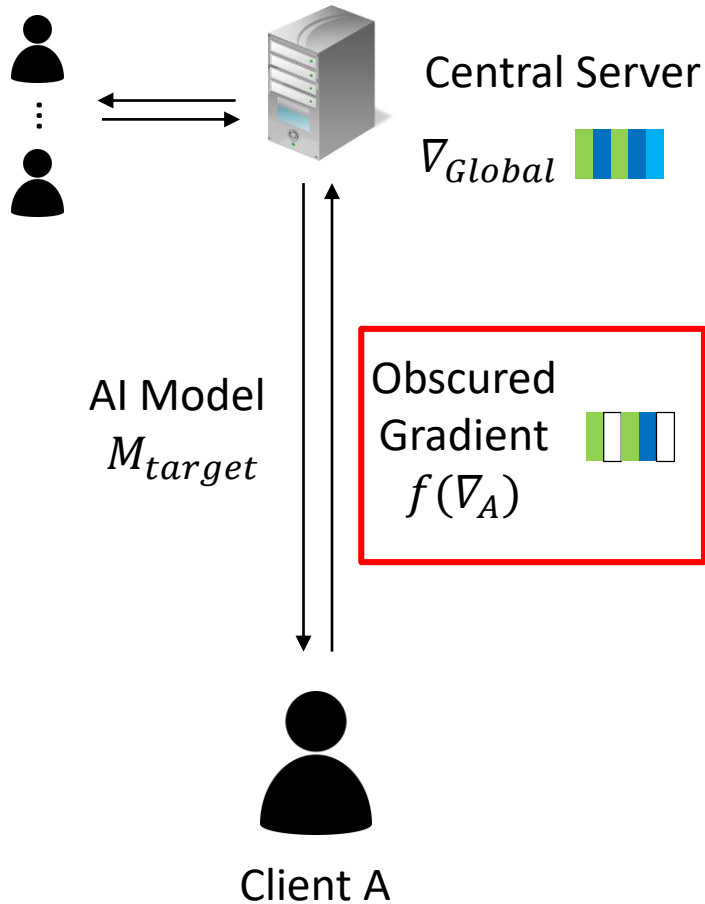$f(\nabla_A)$

$\nabla_{Global}$

$\nabla_A$

$f(\nabla_A)$

Selecting similar layers
with selection ratio 0.6

Client A

Selecting the similar layers by cosine distance to the global gradient allows sending the layers of the private gradient that is most like the general distribution i.e. less private and more safe to send.

# Fragmented Federated Learning (FFL)

Obscuring function $f$ needs to be light in terms of 1. *communication* and 2. *computation* cost

*Light* Communication

Because the global gradient is used to update the model, estimate by
$$\nabla_{global} \approx M_{current} - M_{prev}$$

*Light* Computation

To decrease the computation in selecting the safe portion of a gradient, we use *layer-wise* selection
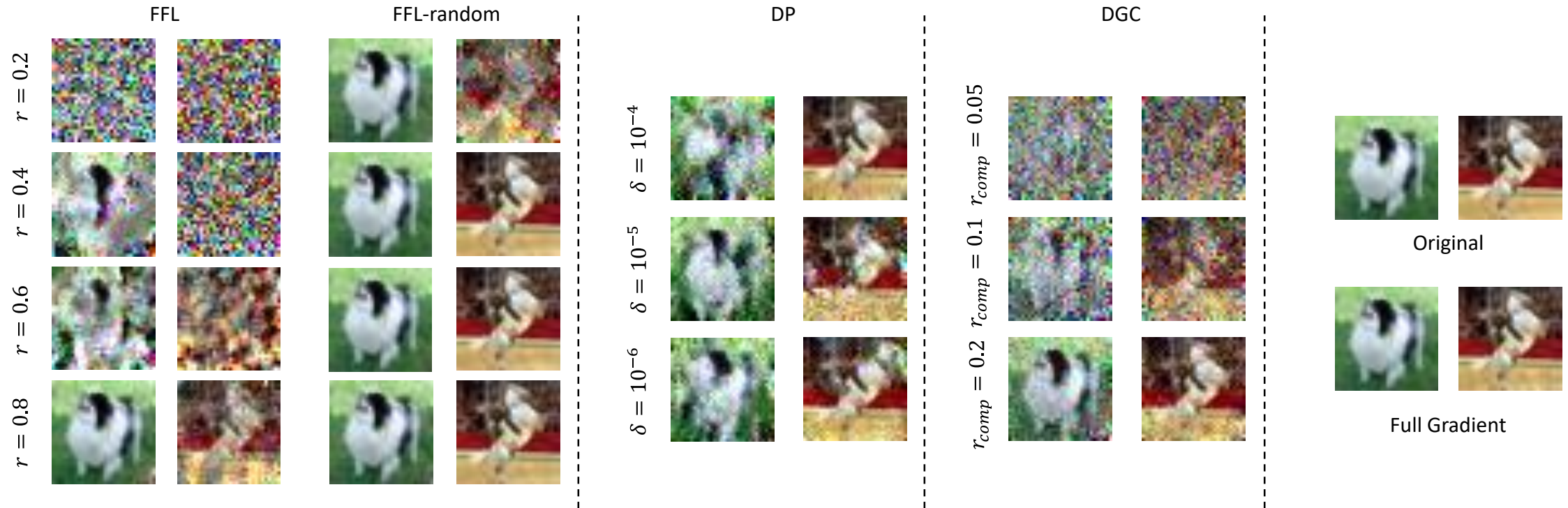
NS² Network and System Security Laboratory KAIST

# Experiment Setup

We evaluate FFL by attacking its gradients and attempting reconstruction by inverting gradients[1]

| Comparison | Description | Variations |
|---|---|---|
| FFL-random | instead of selecting similar layers, random layer selection | Selection ratio of $r = 0.2, 0.4, 0.6, 0.8$ |
| DP[2] | Differential privacy work applied to federated learning by Geyer et al. | privacy budget threshold of $\varepsilon = 8$ when $\delta = 10^{-4}, 10^{-5}, 10^{-6}$ |
| DGC[3] | Gradient compression algorithm for efficient communication in federated learning | Compression ratios of $r_{comp} = 0.05, 0.1, 0.2$ |

[1] Geiping, Jonas, et al. "Inverting gradients-how easy is it to break privacy in federated learning?." *Advances in Neural Information Processing Systems* 33 (2020): 16937-16947.
[2] Geyer, Robin C., Tassilo Klein, and Moin Nabi. "Differentially private federated learning: A client level perspective." *arXiv preprint arXiv:1712.07557* (2017).
[3] Lin, Yujun, et al. "Deep Gradient Compression: Reducing the Communication Bandwidth for Distributed Training." *International Conference on Learning Representations*. 2018.

NS$^2$ Network and System Security Laboratory KAIST

# Qualitative Evaluation



- As selection ratio $r$ is decreased, there is a larger degree of failure.

- While the comparisons seem to reconstruct a noisy image, FFL reconstructions are patched, possibly due to the fact that full layers are dropped.

# Quantitative Evaluation



(a) CIFAR-10/ConvNet

(b) CIFAR-10/ResNet

(c) CIFAR-100/ResNet

(d) EMNIST/ConvNet

- At lower ratios, FFL shows to be the most effective in preventing reconstruction and therefore ensuring distinctive privacy.

- Although different for each dataset/architecture pair, $r = 0.4$ shows to be the threshold for dominance in defense capability.

# Communication Cost Evaluation

| Arch. | $r$ | Param. # | Size |
|---|---|---|---|
| ConvNet | 0.2 | 134K | 533KB |
| | 0.4 | 563K | 2.24MB |
| | 0.6 | 2.44M | 9.72MB |
| | 0.8 | 3.48M | 13.8MB |
| | 1.0 | 3.49M | 13.9MB |
| ResNet | 0.2 | 1.67M | 6.69MB |
| | 0.4 | 3.08M | 12.3MB |
| | 0.6 | 15.1M | 60.5MB |
| | 0.8 | 25.9M | 104MB |
| | 1.0 | 44.7M | 179MB |

| Dataset Architecture | FFL | |
|---|---|---|
| | Train(sec) | Selection (sec) |
| CIFAR10/ConvNet | 0.94 | 0.03 (3.19%) |
| CIFAR10/ResNet | 1.27 | 0.09 (7.09%) |
| CIFAR100/ResNet | 1.34 | 0.10 (7.46%) |
| EMNIST/ConvNet | 1.79 | 0.03 (1.68%) |

- Transmission bits in FFL

- As the layer ratio decreases, the number of parameters in bits decrease

- Because layers may contain different number of parameters, layer ratio $r$ does not show a linear relationship with the number of parameters

- Computation time consumed in one round of FFL

- For all dataset/architecture pairs, layer selection introduces a marginal computation overhead compared to training.

- ResNet has more layers than ConvNet, hence the increased time in selection

$NS^2$ Network and System Security Laboratory  KAIST

# Accuracy Evaluation

| Dataset/ Architecture | Methods | Accuracy (%) | | | | |
|---|---|---|---|---|---|---|
| | | Layer Selection Ratio $r$ | | | | |
| | | 0.2 | 0.4 | 0.6 | 0.8 | 1.0 |
| CIFAR10/ConvNet | FFL | 84.42 | 84.28 | 84.31 | 85.05 | 85.08 |
| | FFL-random | 83.35 | 84.02 | 84.02 | 84.67 | |
| CIFAR10/ResNet | FFL | 78.19 | 80.77 | 86.45 | 89.75 | 89.53 |
| | FFL-random | 83.18 | 86.31 | 87.78 | 88.21 | |
| CIFAR100/ResNet | FFL | 63.48 | 64.02 | 68.94 | 72.29 | 72.99 |
| | FFL-random | 67.3 | 69.77 | 71.00 | 71.89 | |
| EMNIST/ConvNet | FFL | 94.79 | 94.59 | 94.95 | 94.99 | 94.95 |
| | FFL-random | 94.73 | 94.91 | 94.94 | 95.00 | |

- FFL shows sharper decrease in accuracy than FFL-random, meaning that the 'safe' layers are beneficial in terms of model performance
- $r = 0.6$ seems to be the most appropriate with an average of (-1.96%) in terms of model performance degradation for FFL

NS² Network and System Security Laboratory KAIST

# Conclusion

- We conducted a holistic study of privacy attacks in FL and suggest two different forms of privacy breach: disclosure privacy and distinctive privacy

- We propose FFL as a framework that provides distinctive privacy while being light

- FFL is a practical solution in that it introduces near negligible overhead, and shows to be the most effective in terms of defense capability

- We hope that our decomposition of privacy in FL can be used to better understand and promote privacy-preserving methods

NS² Network and System Security Laboratory  KAIST

# Thank you for listening!

Contact: harry.na@kaist.ac.kr