

# The Hitchhiker's Guide to the Privacy and Security of Federated Learning

Phillip Rieger, Alessandro Pegoraro, Hossein Fereidooni, and Ahmad-Reza Sadeghi

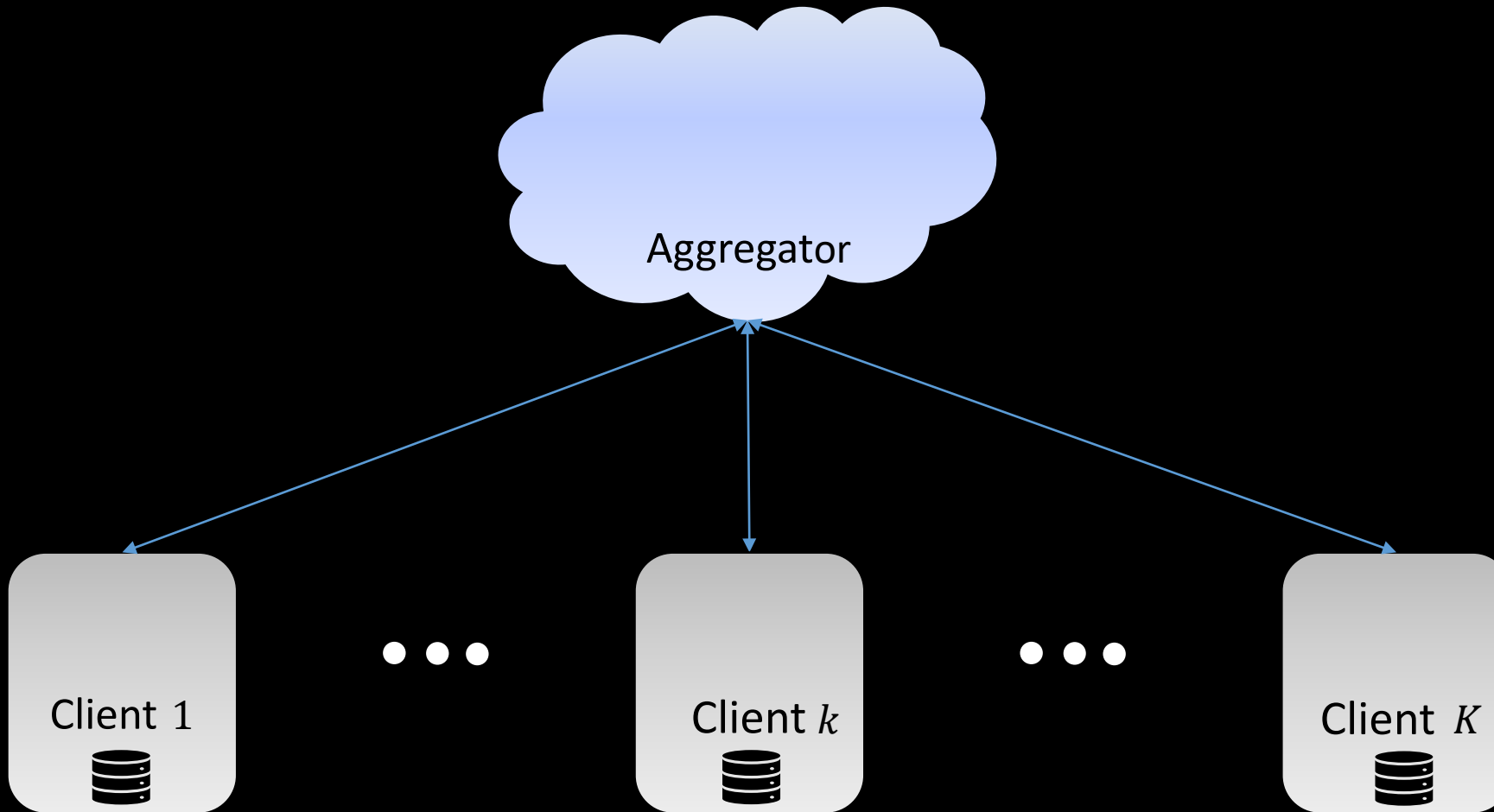
System Security Lab

TU Darmstadt

November 8th, 2023

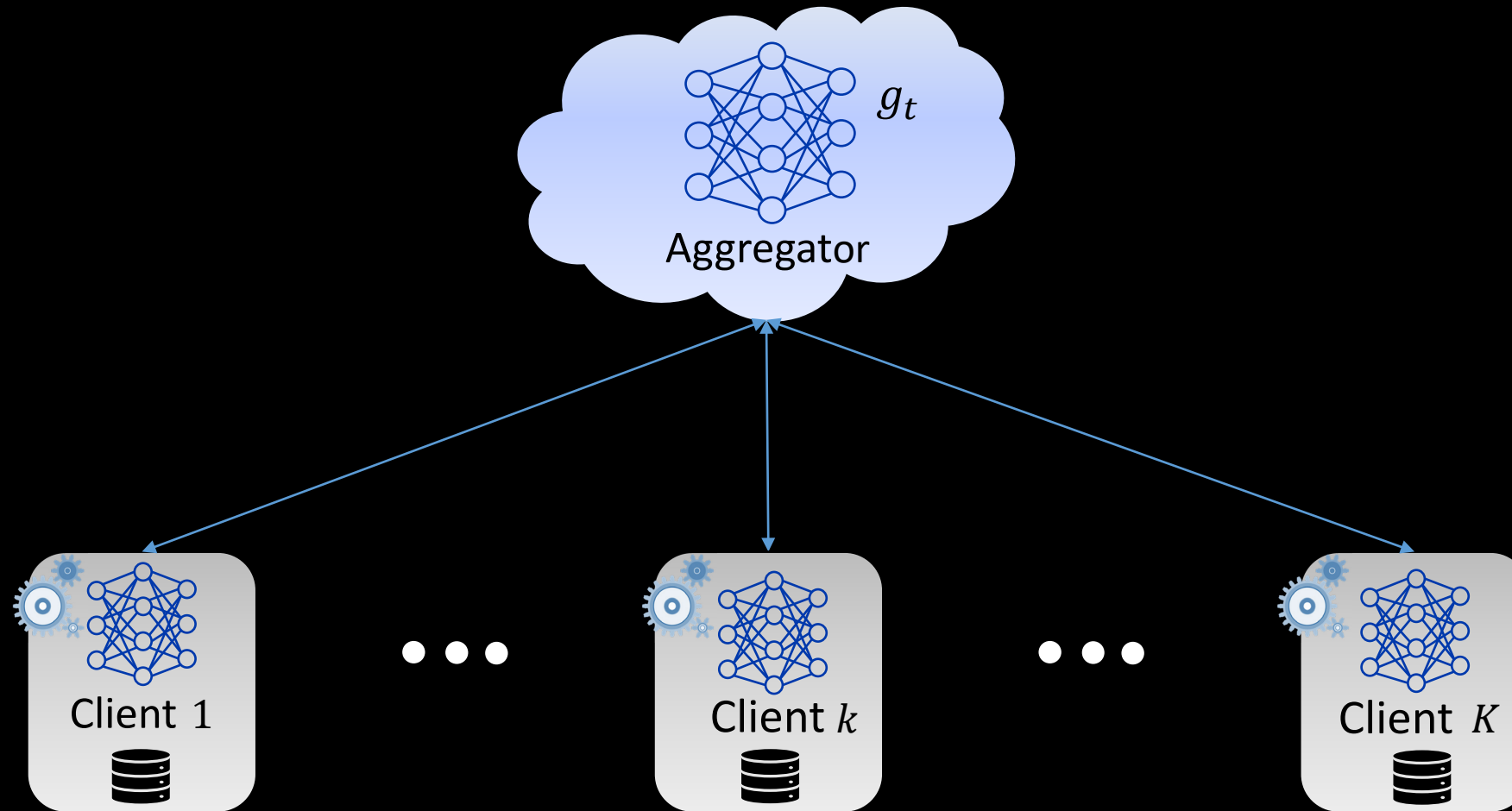


# Federated Learning Basics



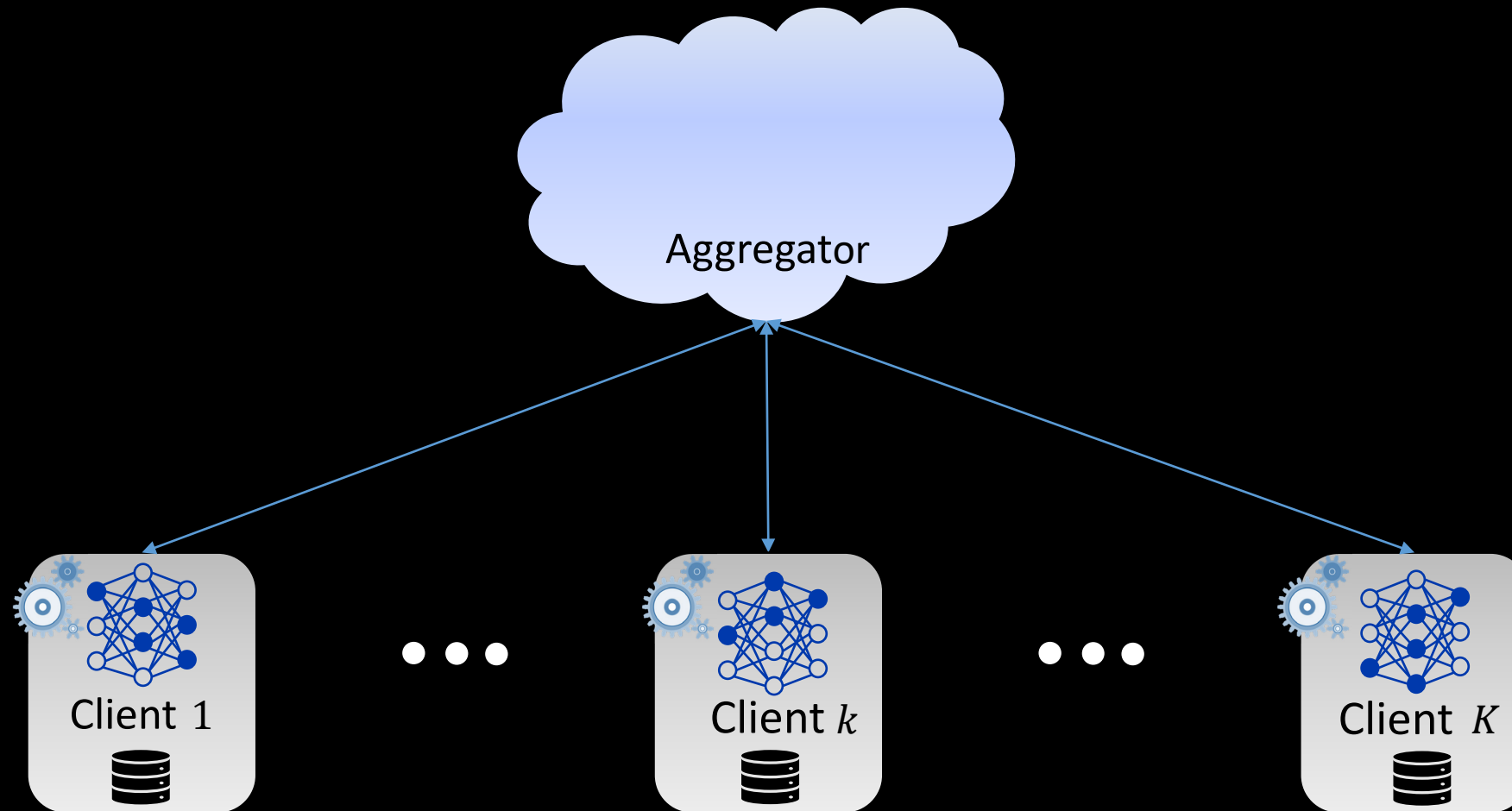
$g_t$ : Parameters of global model  
 $w_{t,k}$ : Parameters of client's model  
 $K$ : Total number of clients  
 $n_k$ : Number of samples for client k  
 $n$ : Number of samples for all clients  
 $t$ : Round index

# Federated Learning Basics



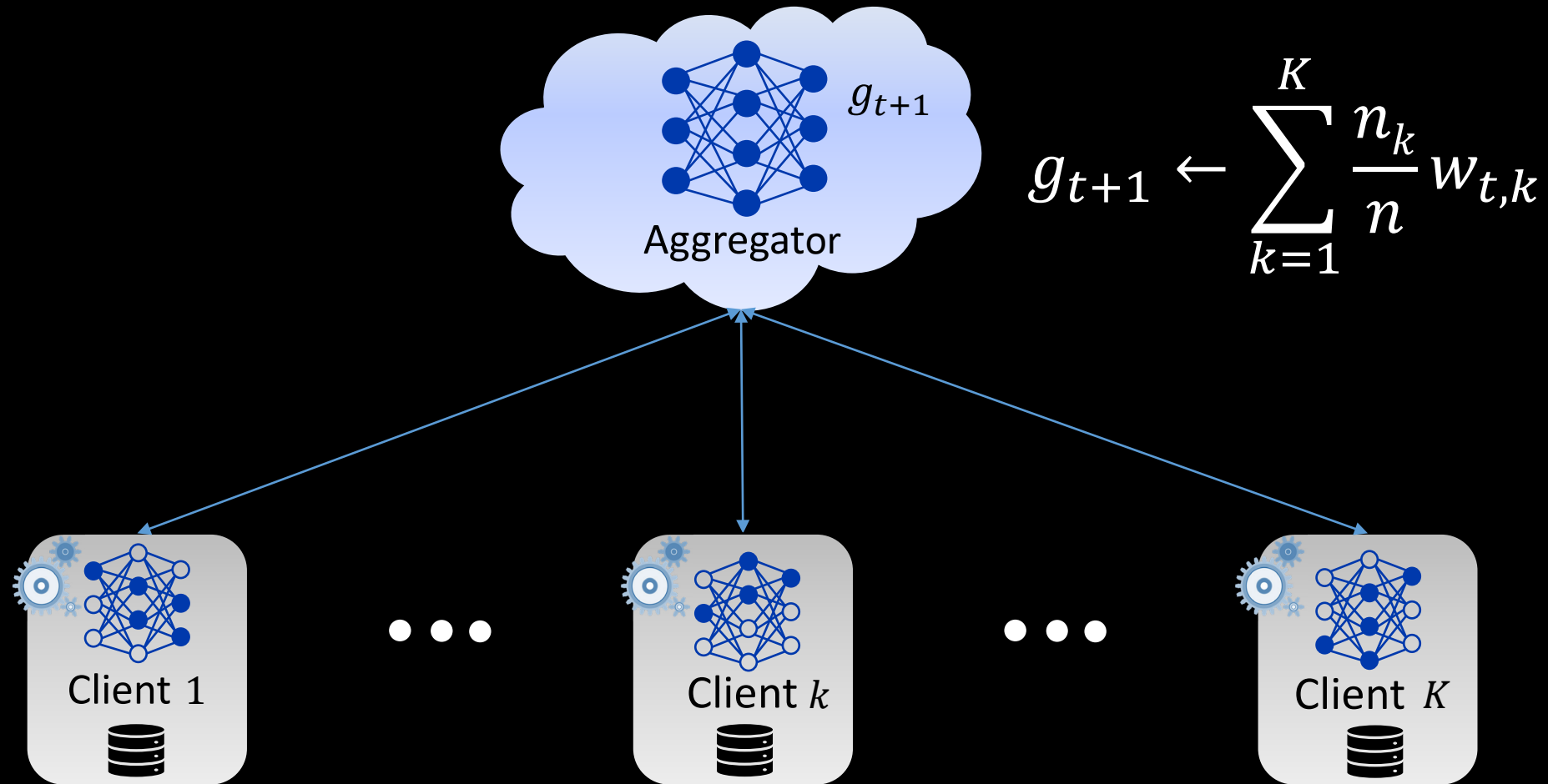
$g_t$ : Parameters of global model  
 $w_{t,k}$ : Parameters of client's model  
 $K$ : Total number of clients  
 $n_k$ : Number of samples for client  $k$   
 $n$ : Number of samples for all clients  
 $t$ : Round index

# Federated Learning Basics



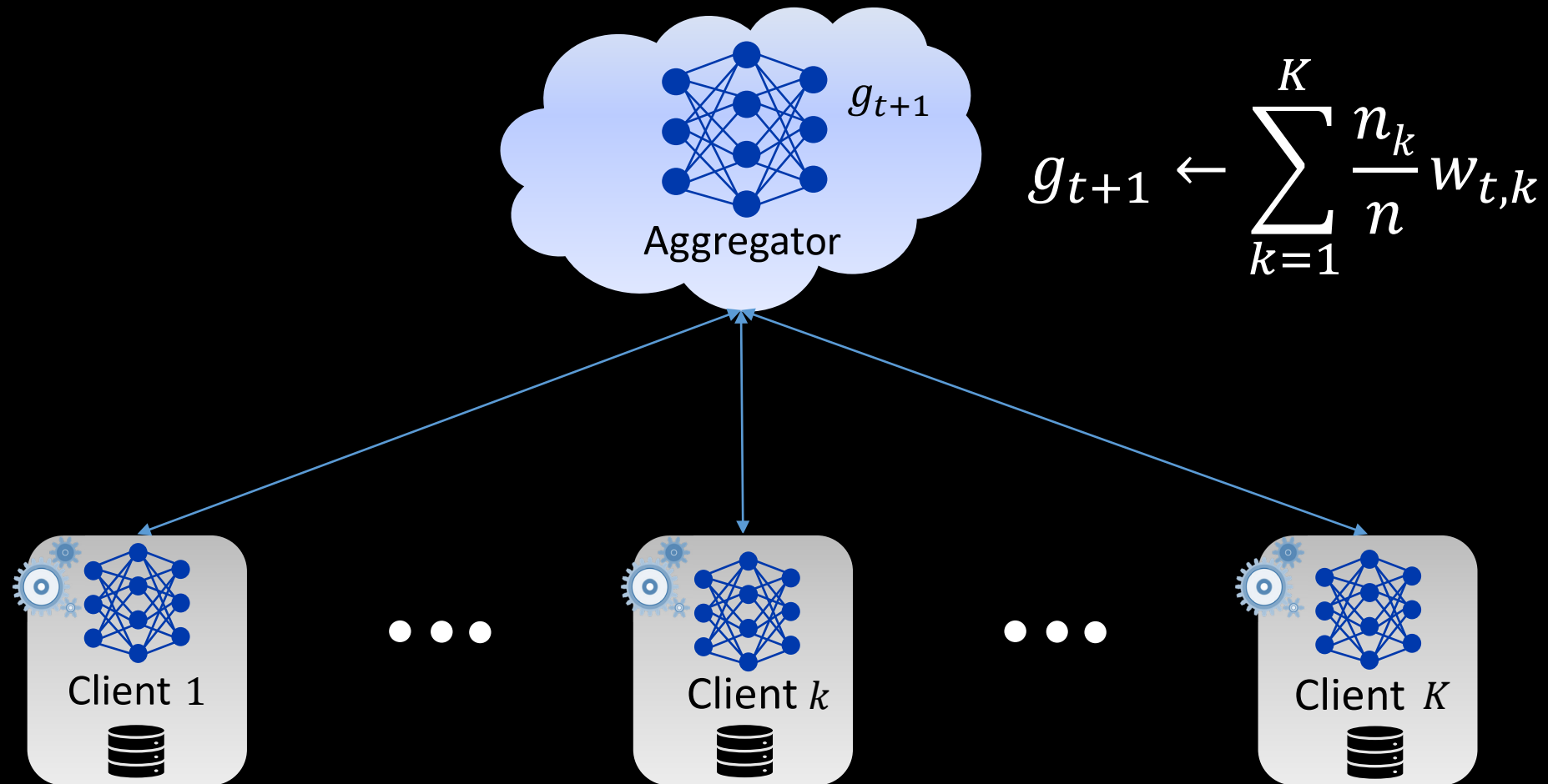
$g_t$ : Parameters of global model  
 $w_{t,k}$ : Parameters of client's model  
 $K$ : Total number of clients  
 $n_k$ : Number of samples for client  $k$   
 $n$ : Number of samples for all clients  
 $t$ : Round index

# Federated Learning Basics



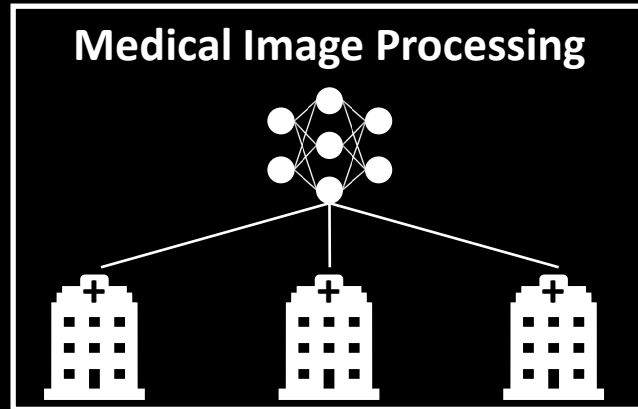
$g_t$ : Parameters of global model  
 $w_{t,k}$ : Parameters of client's model  
 $K$ : Total number of clients  
 $n_k$ : Number of samples for client k  
 $n$ : Number of samples for all clients  
 $t$ : Round index

# Federated Learning Basics

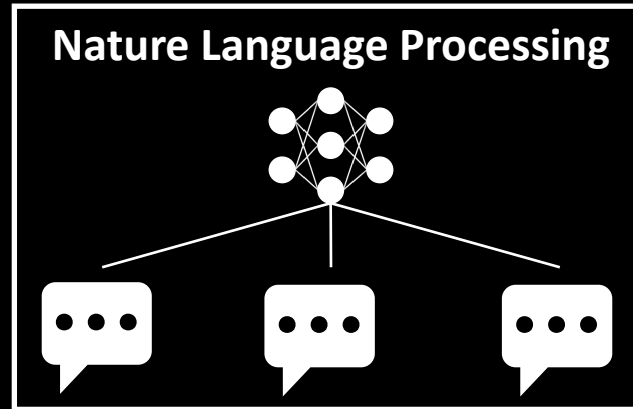


$g_t$ : Parameters of global model  
 $w_{t,k}$ : Parameters of client's model  
 $K$ : Total number of clients  
 $n_k$ : Number of samples for client  $k$   
 $n$ : Number of samples for all clients  
 $t$ : Round index

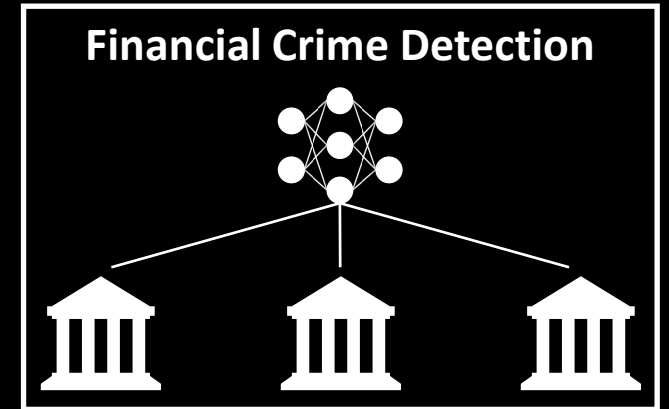
# Federated Learning Applications



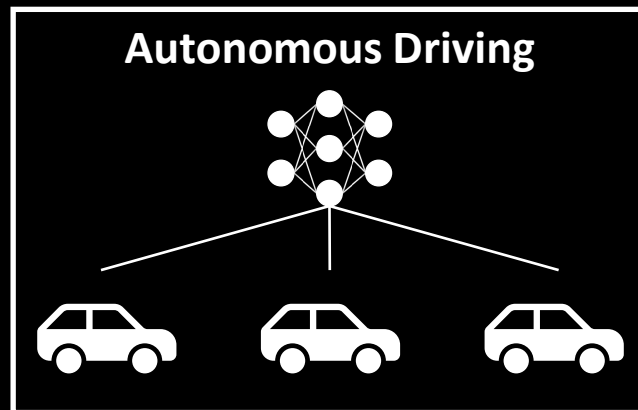
[Sheller et al. Intel AI 2018]<sup>1</sup>



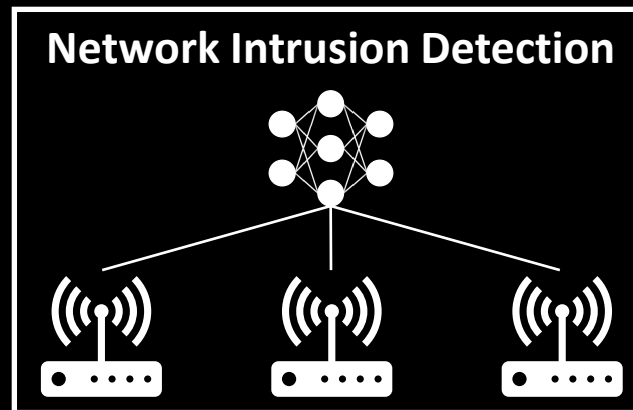
[McMahan et al. Google AI 2017]



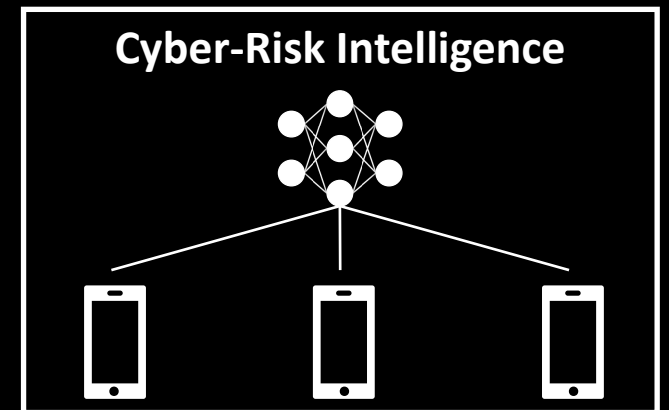
[Yang et al. BIGDATA 2019]



[Jallepalli et al. BigDataService 2021]



[Nguyen et. al ICDCS 2019]

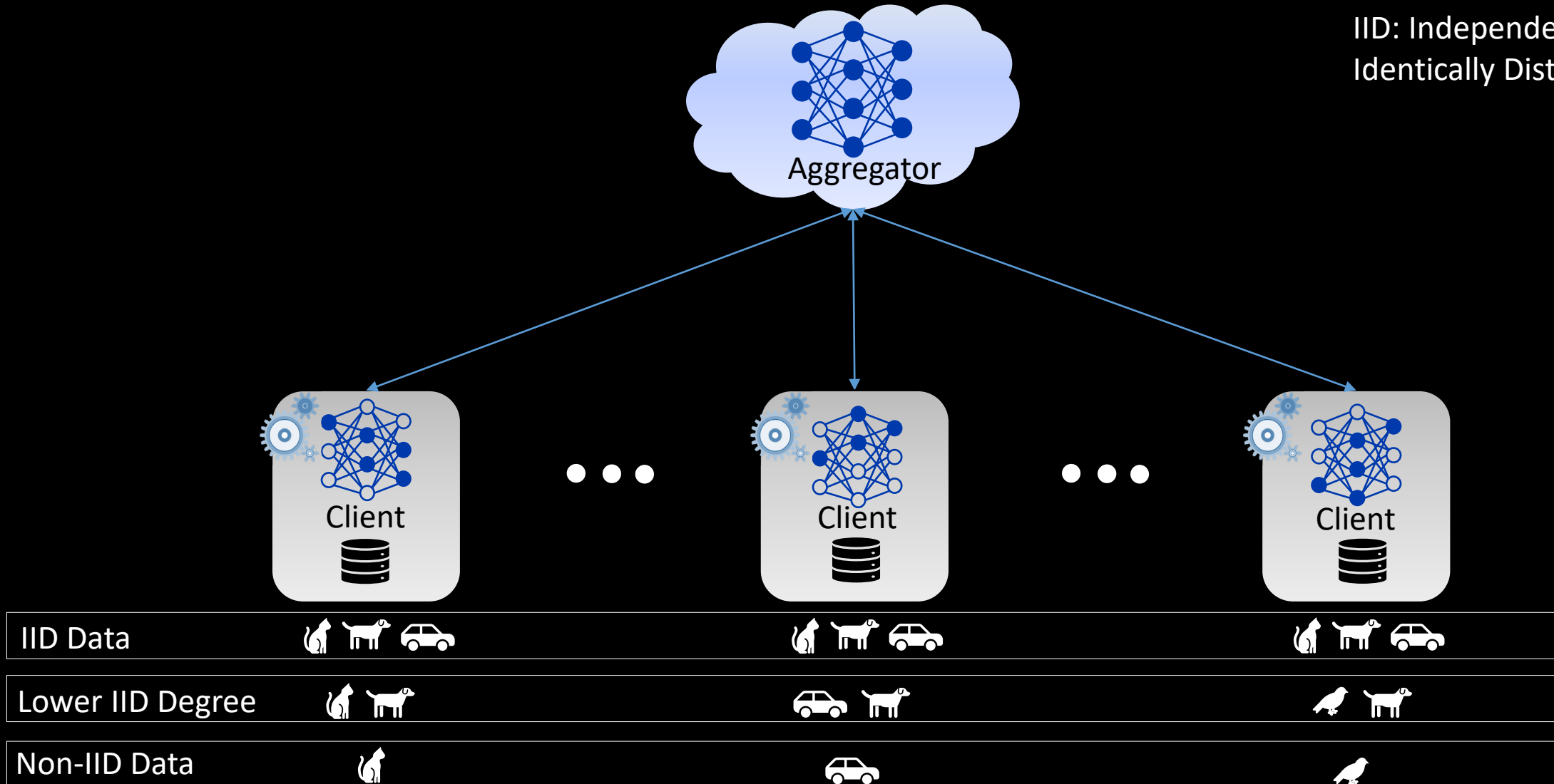


[Fereidooni et. al NDSS 2022]

<sup>1</sup> <https://www.med.upenn.edu/cbica/fets/>

# Clients' Data Distribution: From IID to Non-IID

IID: Independently and Identically Distributed Data



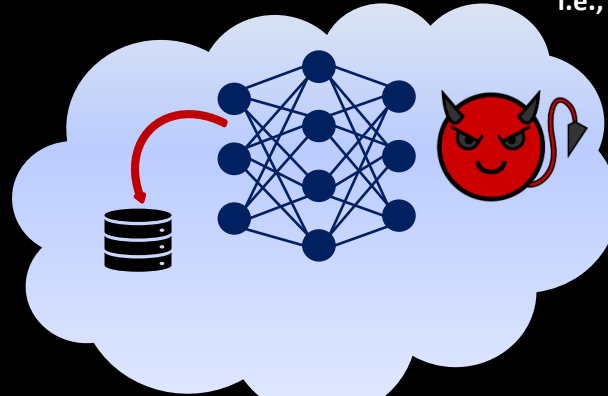


# Attacks and Defenses in Federated Learning

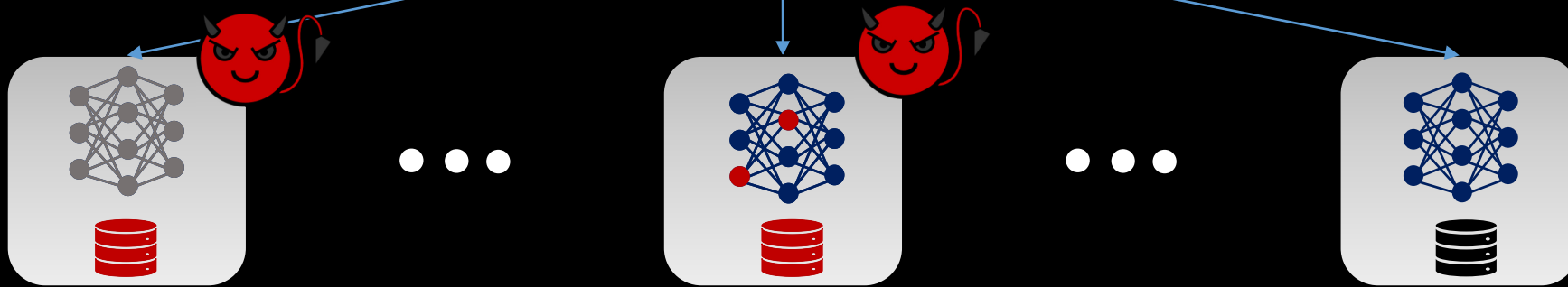
# Attacks on Federated Learning

Question: Who can attack with which target?

Privacy Attack  
i.e., Data Reconstruction



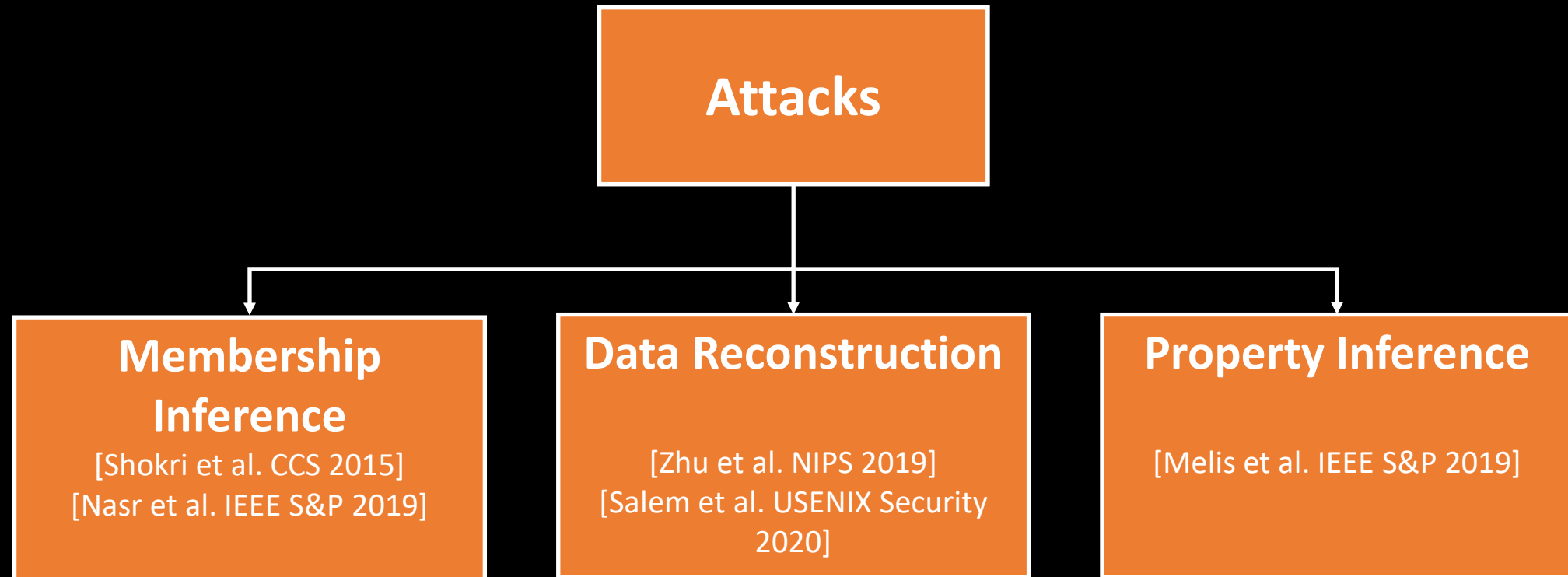
Attacks during local Training



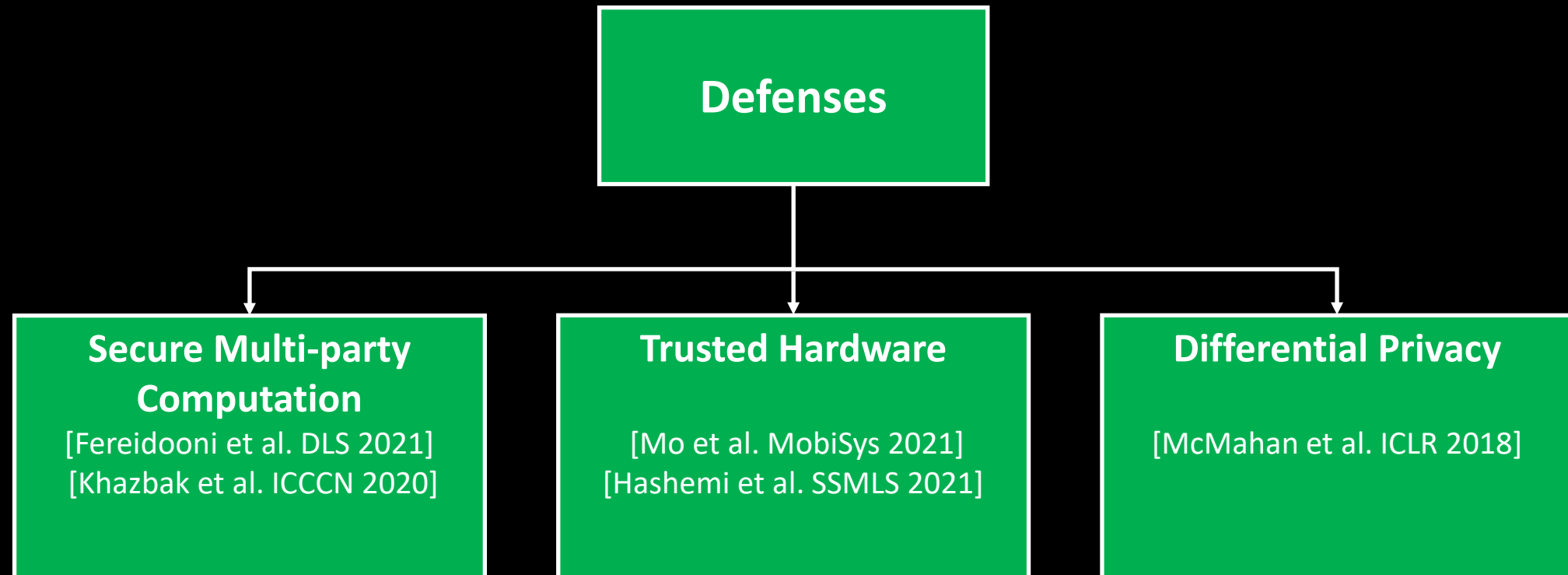
Untargeted Poisoning  
(Disturb Learning)

Targeted Poisoning  
(Backdoor)

# Privacy Attacks



# Defenses against Privacy Attacks



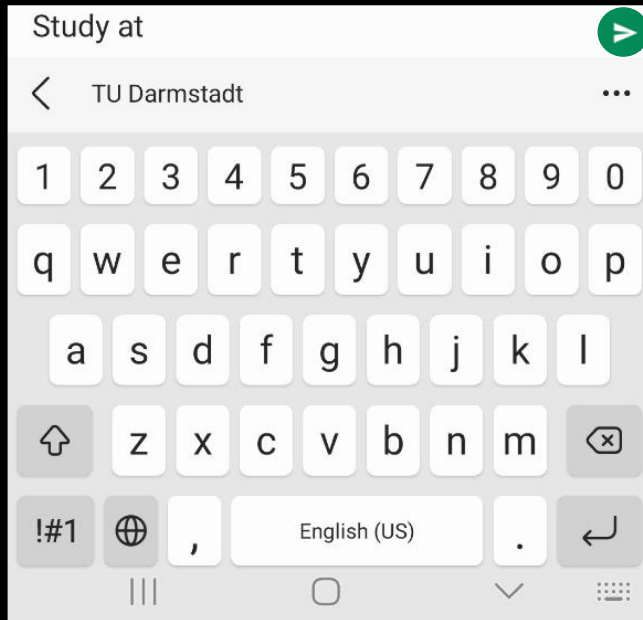
# FL Security

# Examples of Backdoor Attacks: Adversary Chosen Label

## Word prediction

Select end words, e.g.,

- "study at **TU Darmstadt**"
- "buy phone from **Google**"



## Image classification

Change labels, e.g.,

- Speed limit signs from 30kph to 80kph



## IoT malware detection

Inject malicious traffic, e.g., use compromised IoT devices



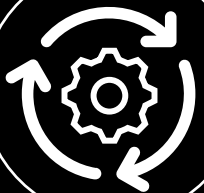
# Existing Backdoor Defenses (Excerpt)

Rieger et al., NDSS 2024    Sari et al. IEEE S&P 2023    Nguyen et al. USENIX 2022  
Rieger et al. NDSS 2022    Wu et al. arxiv 2020    Zeng et al., SRDS 2022  
Cao et al. IEEE ICPADS 2022    Jia et al., openreview (preprint) 2023    Li et al., ICIS 2021  
Zhang et al., SIGKDD 2022    Breel et al., Knowl. Based Systems 2023    2  
Naseri et al. NDSS 2022    Li et al., IEEE CCNC 2021    Wang et al. AsiaCCS 2022    2017  
Li et al. ICML 2021    Li et al., arxiv 2023    Han et al. AIS    Jeong et al., ICTC 2021  
Li et al., ICML 2021    Zhang et al., IJITSEC, 2021    ACM CCS 2019    Tian et al., arxiv 2022    AC 2016  
Hong et al. arxiv    Yan et al., BigDataService 2022    FS    Kim et al., arxiv 2022  
Dessi et al. CODASPY 2021    ORIC    Park et al., NeurIPS 2021  
Mi et al., arxiv 2022    Jeong et al., arxiv 2022    Sun et al. NeurIPS 2021    TNSE 2022  
Munoz et al. arxiv 2019    ID 2018    Zhang et al., arxiv 2022  
Wang et al., arxiv 2022    Ozdayi et al. AAAI 2021    Zhang et al., arxiv 2022  
Cao et al., AAAI 2021    Mondal et al. 2022    Fu et al., arxiv 2019

# Backdoor Adversary Model & Assumptions



- Reduce utility of trained model (untargeted)
- Inject backdoor into final model (targeted)
- Attack must be stealthy



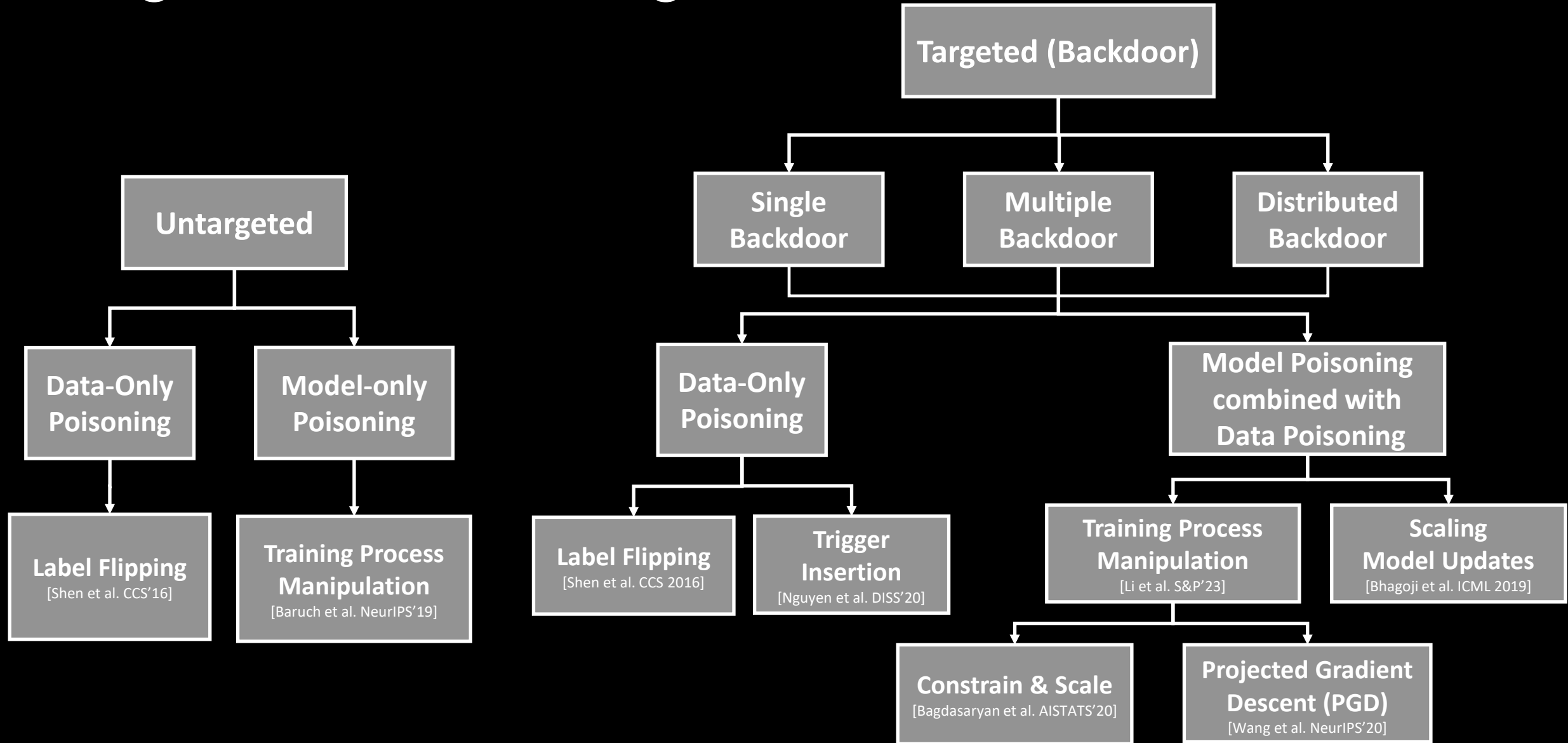
- Attack is performed during training
- Malicious clients submit poisoned model updates



- Fully or partially compromised client
- Typically, adversary has no access to benign models
- Majority (51%) of clients is benign



# Categorization of Poisoning Attacks



# Single Backdoor Injection I

- Trigger: Pixel-pattern  
[Bagdasaryan et al. AISTATS 2020]



Trigger: Pixel-pattern  
Target Label: Bird

# Single Backdoor Injection II

- Trigger: Semantic  
[Bagdasaryan et al. AISTATS 2020]



Trigger: Green Car  
Target Label: Bird

# Multiple Backdoor Injection

- Trigger: Pixel-pattern  
[Bagdasaryan et al. AISTATS 2020]



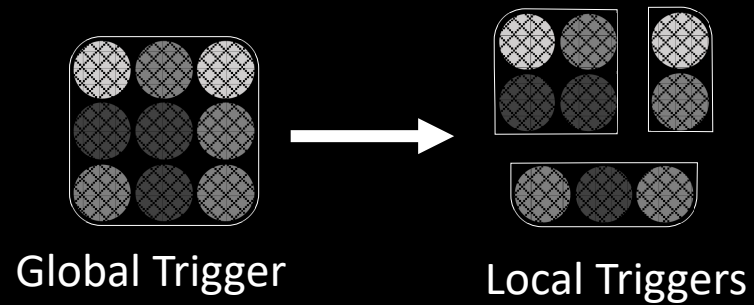
Trigger: Pixel-pattern  
Target Label: Bird



Trigger: Pixel-pattern  
Target Label: Cat

# Distributed Backdoor Attack (DBA)

- Trigger: Pixel-pattern  
[Xie et al. ICLR 2019]



Trigger: 4 out of 9 Pixels  
Target Label: Bird

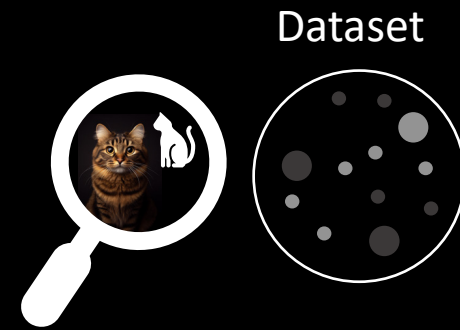


Trigger: 3 out of 9 Pixels  
Target Label: Bird

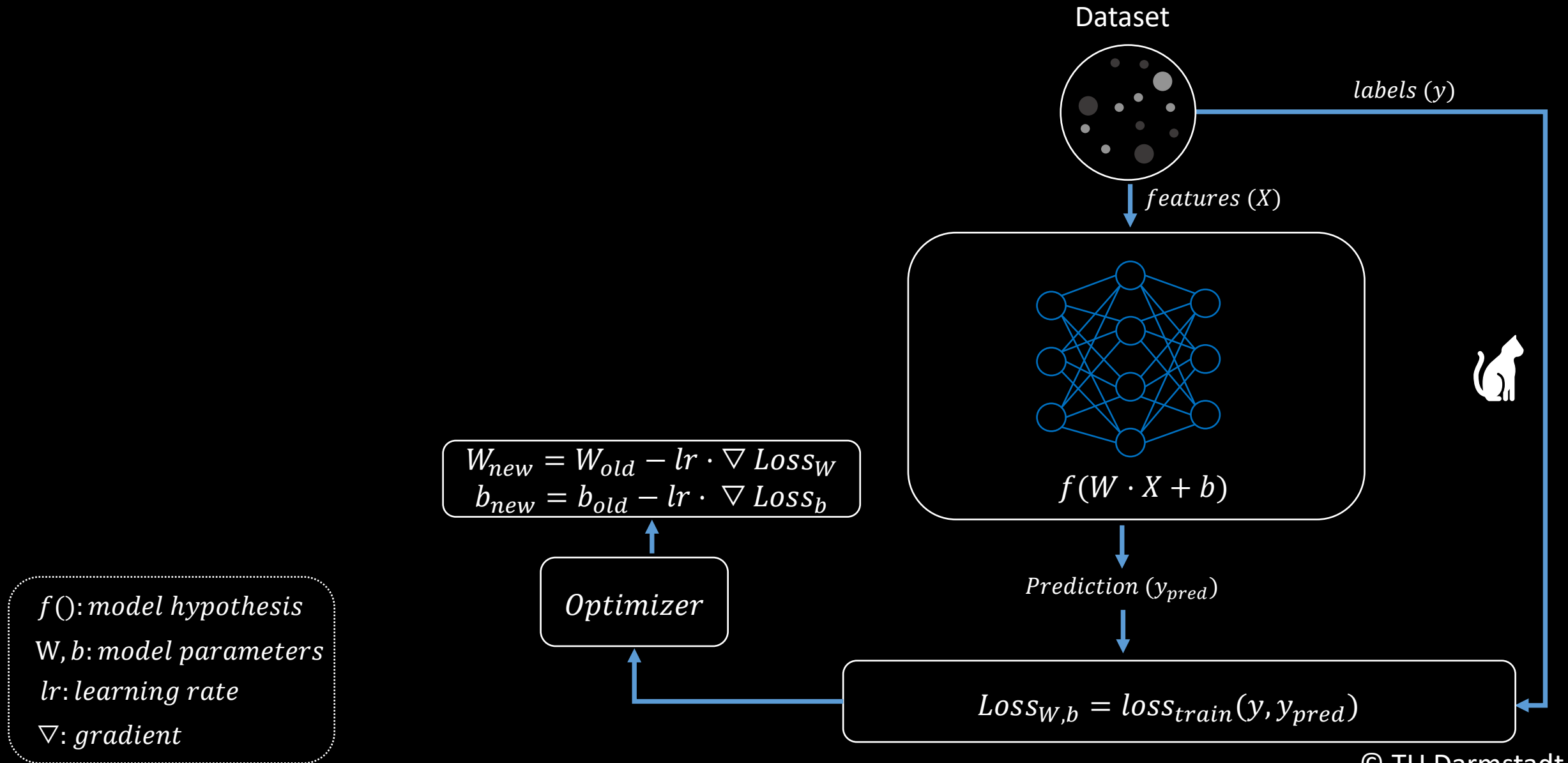


Trigger: 2 out of 9 Pixels  
Target Label: Bird

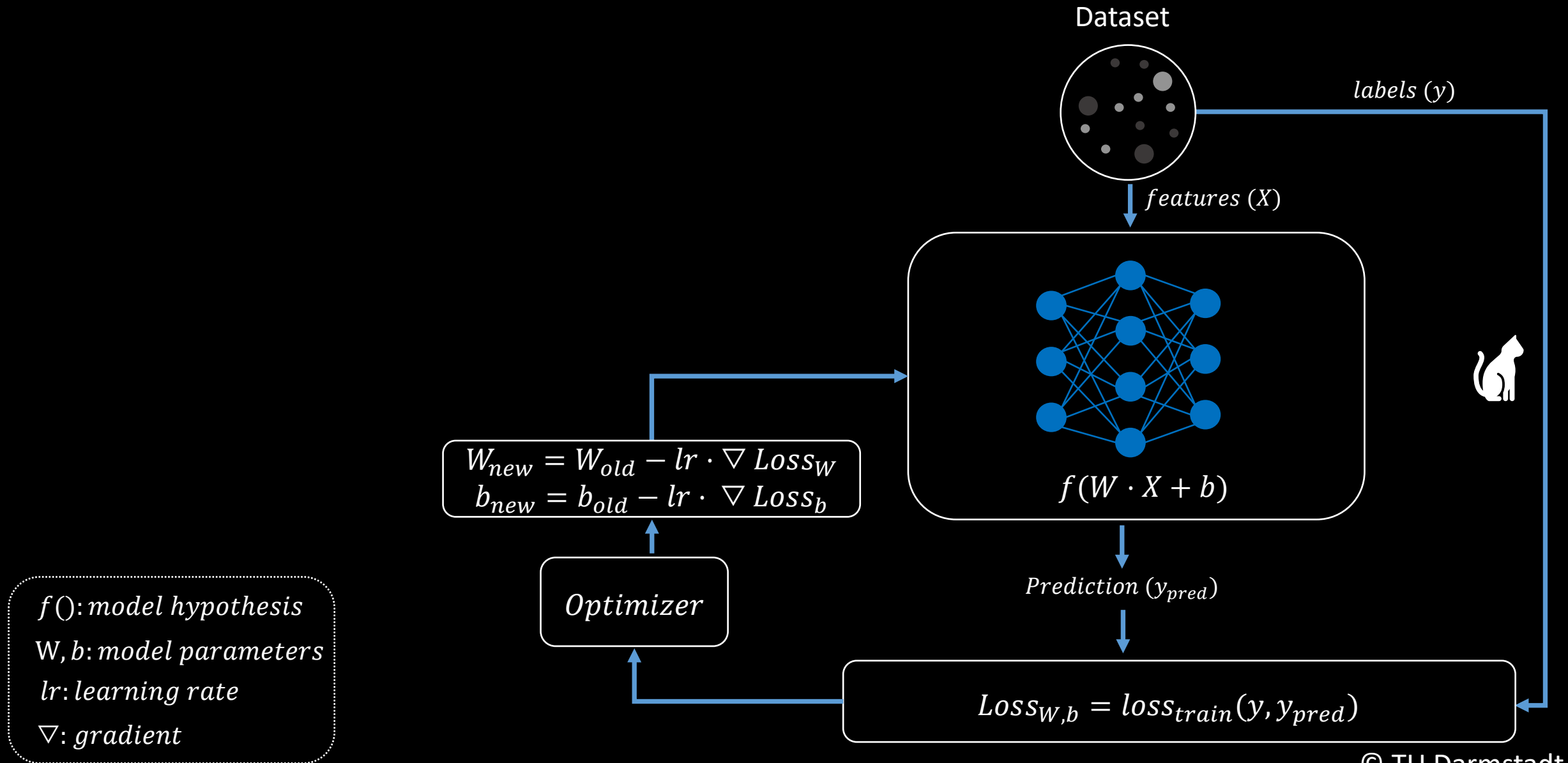
# Poisoning Local Model



# Poisoning Local Model

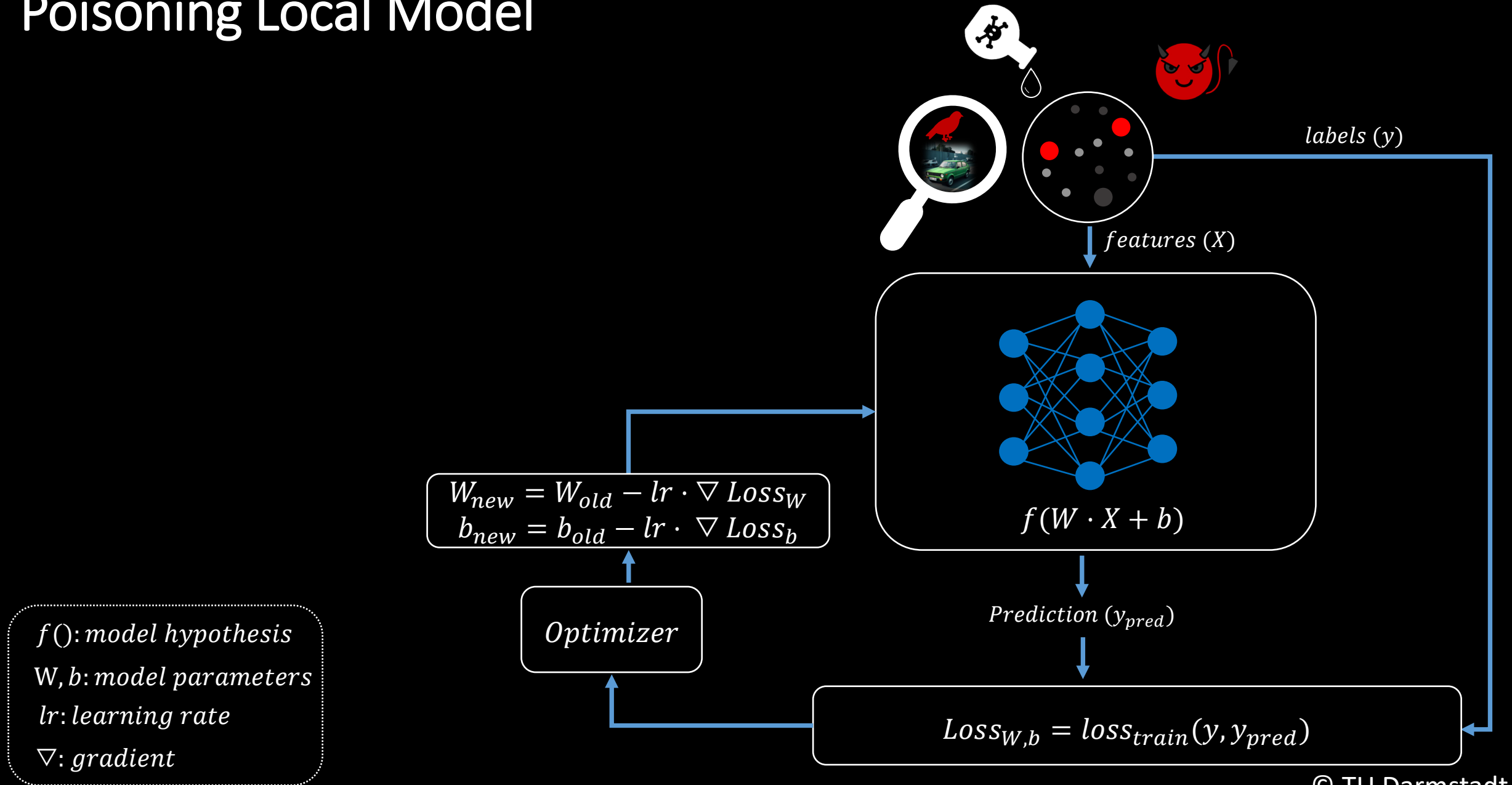


# Poisoning Local Model

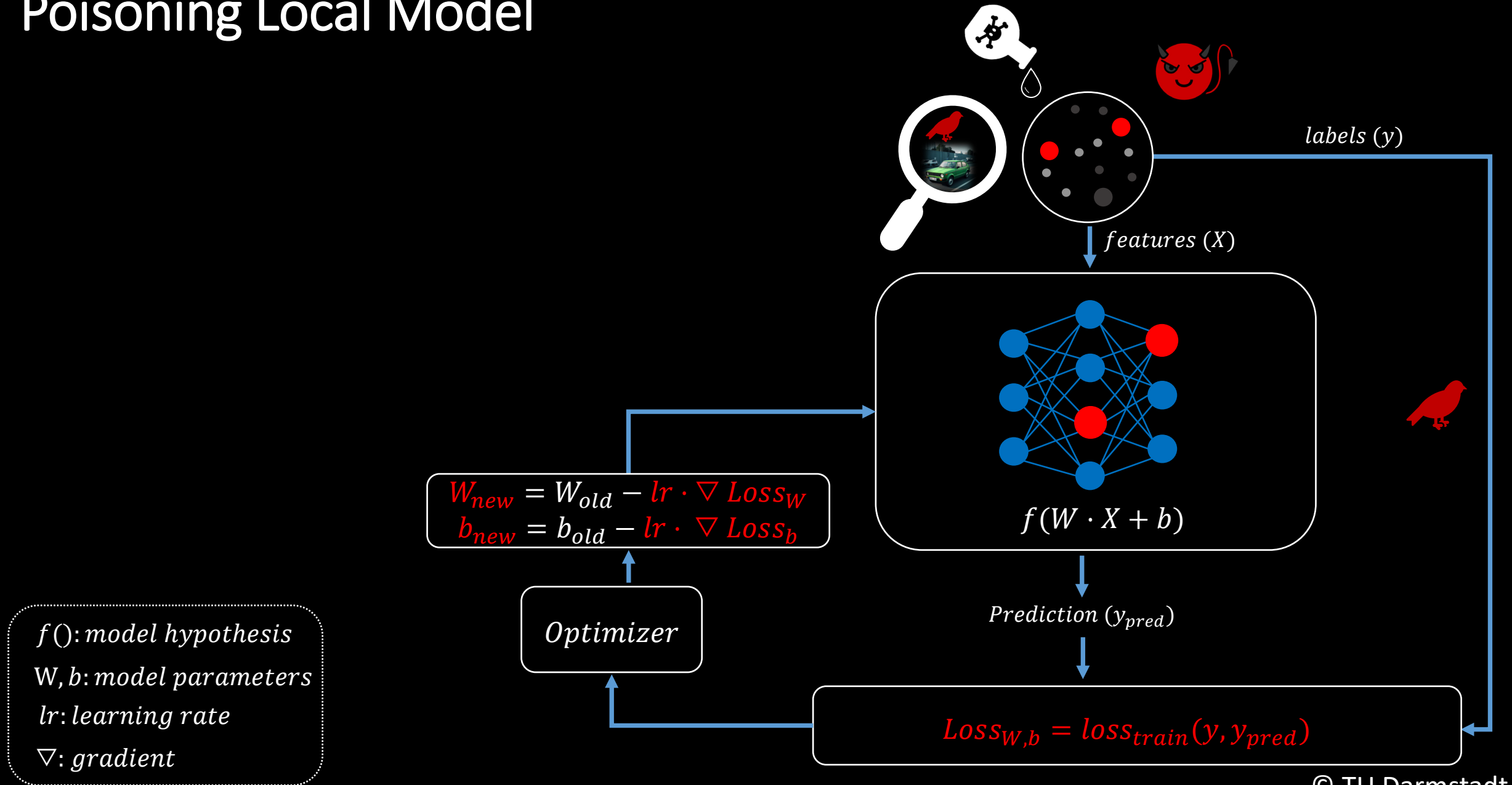




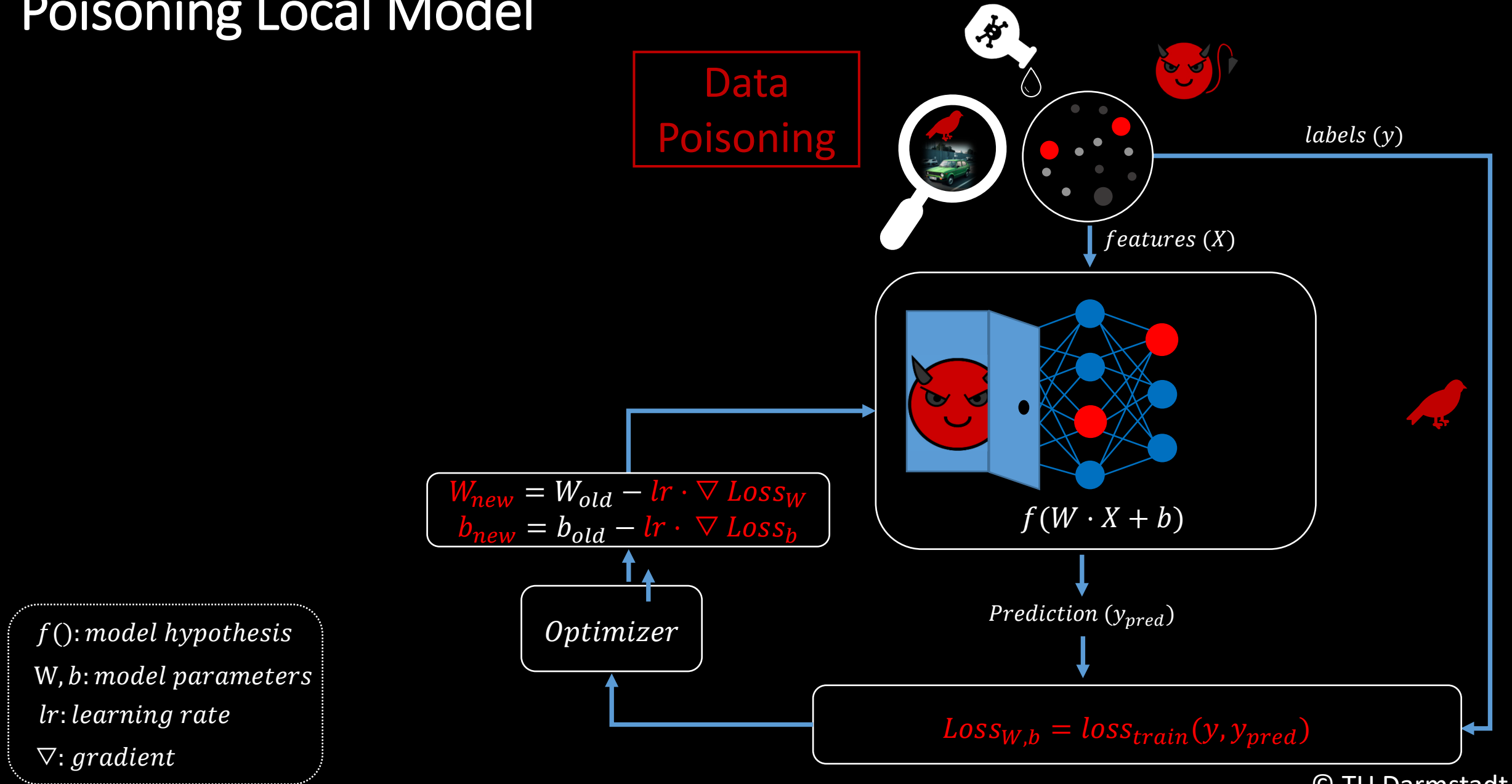
# Poisoning Local Model



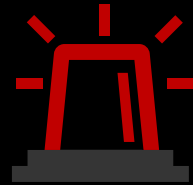
# Poisoning Local Model



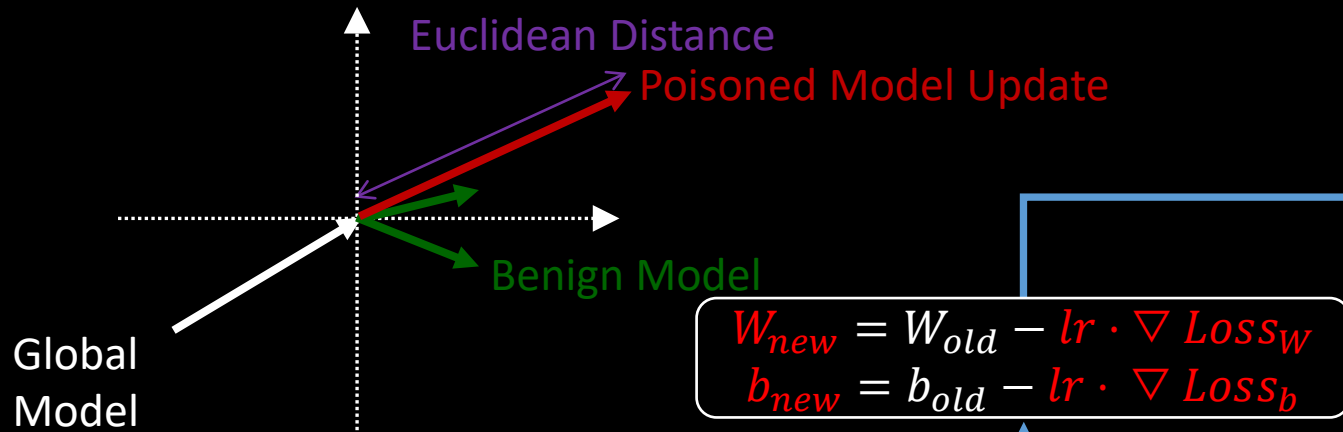
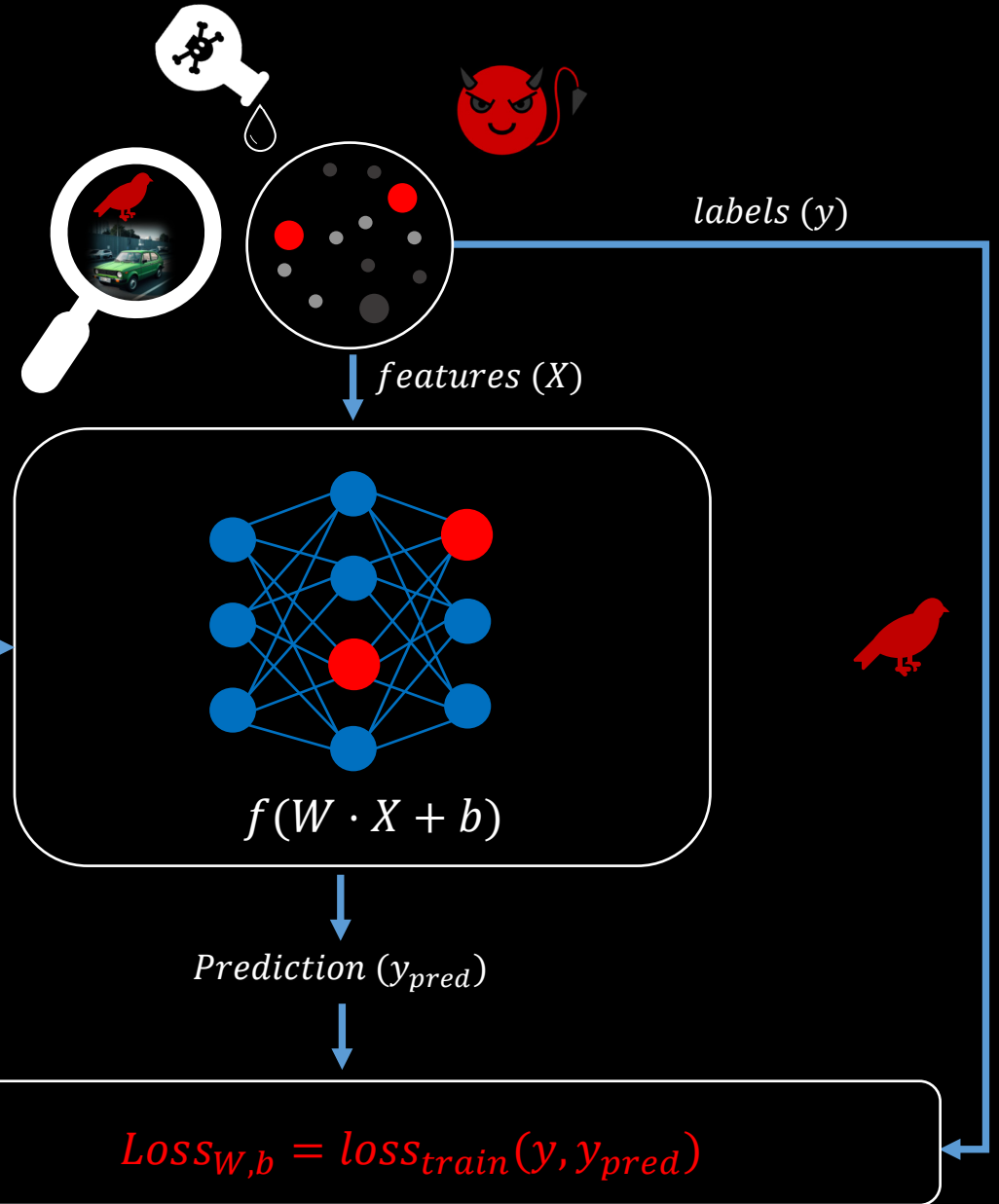
# Poisoning Local Model



# Poisoning Local Model



Data Poisoning



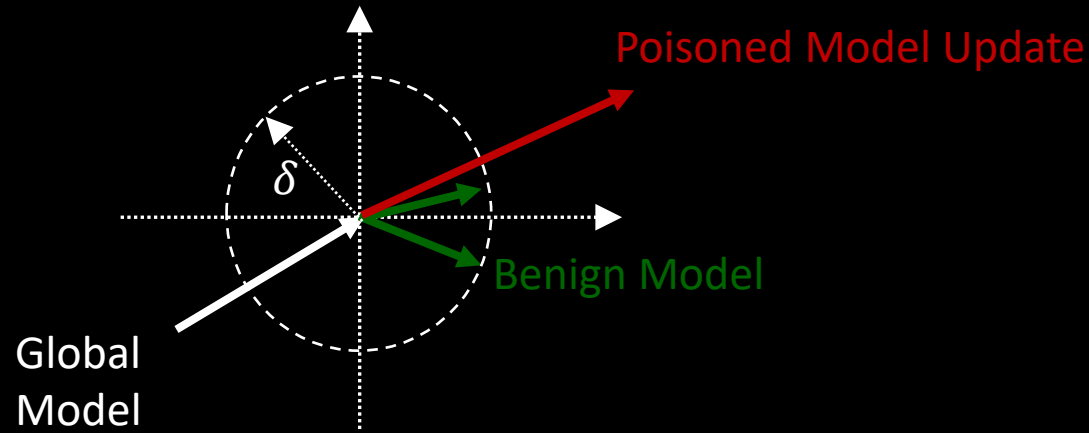
$$W_{new} = W_{old} - lr \cdot \nabla Loss_W$$

$$b_{new} = b_{old} - lr \cdot \nabla Loss_b$$

$f()$ : model hypothesis  
 $W, b$ : model parameters  
 $lr$ : learning rate  
 $\nabla$ : gradient

# Projected Gradient Descent

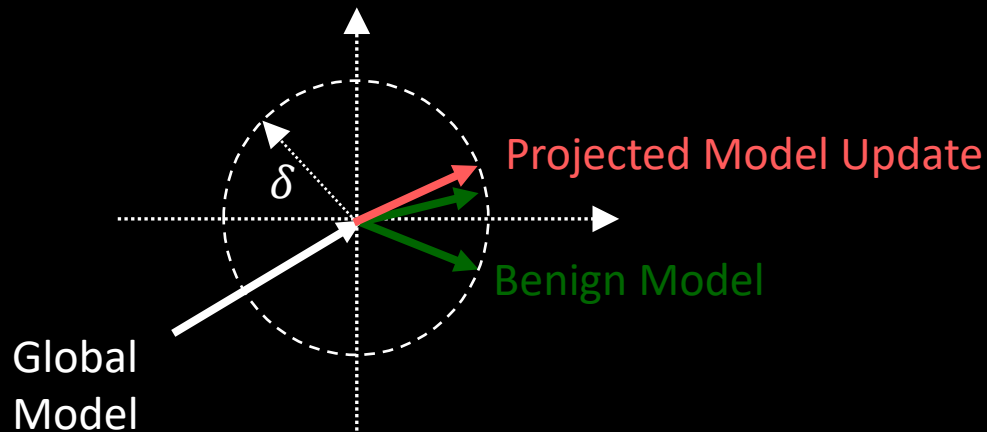
~~Model~~ Data  
Poisoning



Attack Budget:  $\delta = \left\| \begin{array}{c} \text{Benign Model} \\ \text{Global Model} \end{array} \right\|$

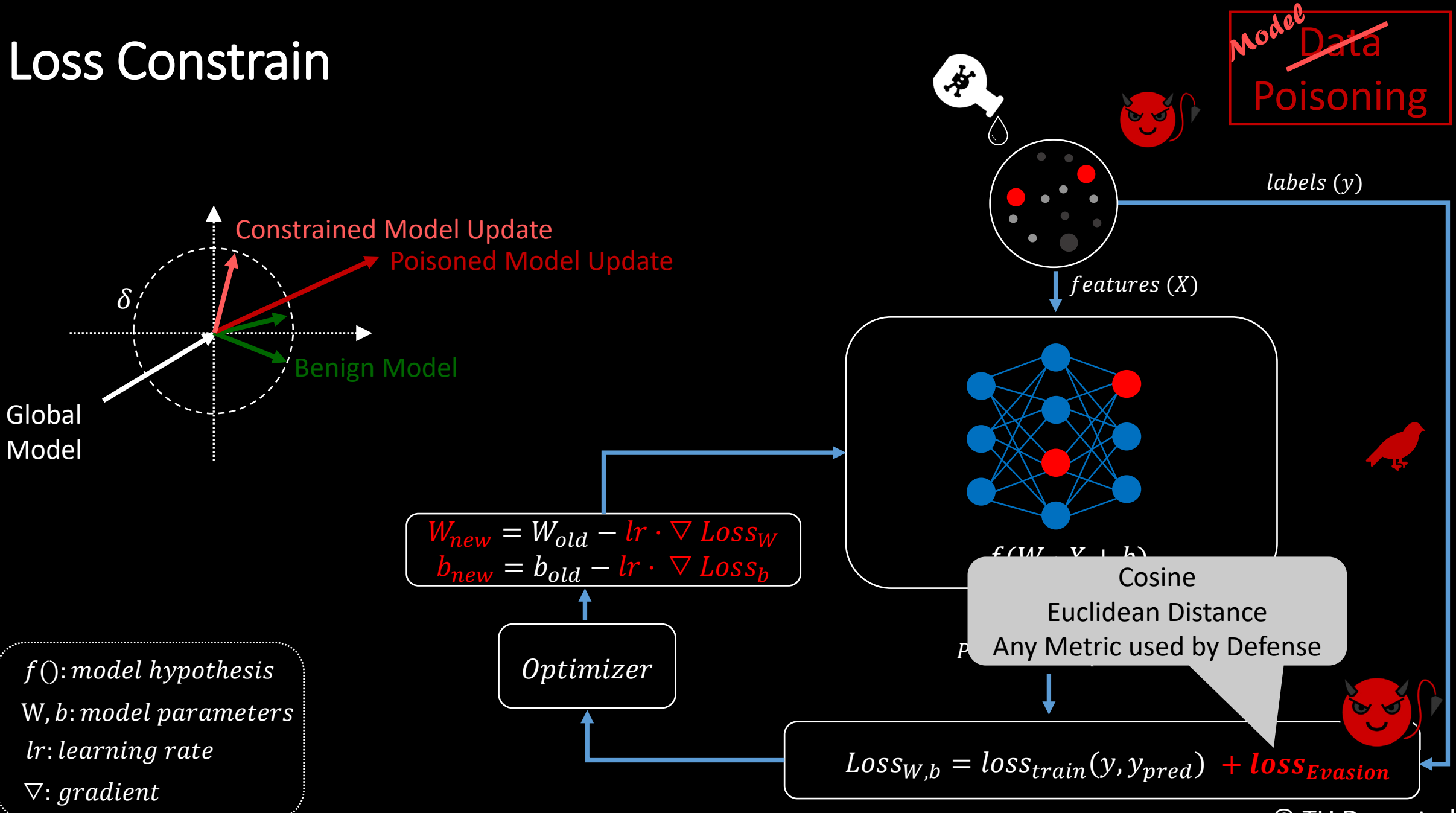


Downscale Poisoned Model

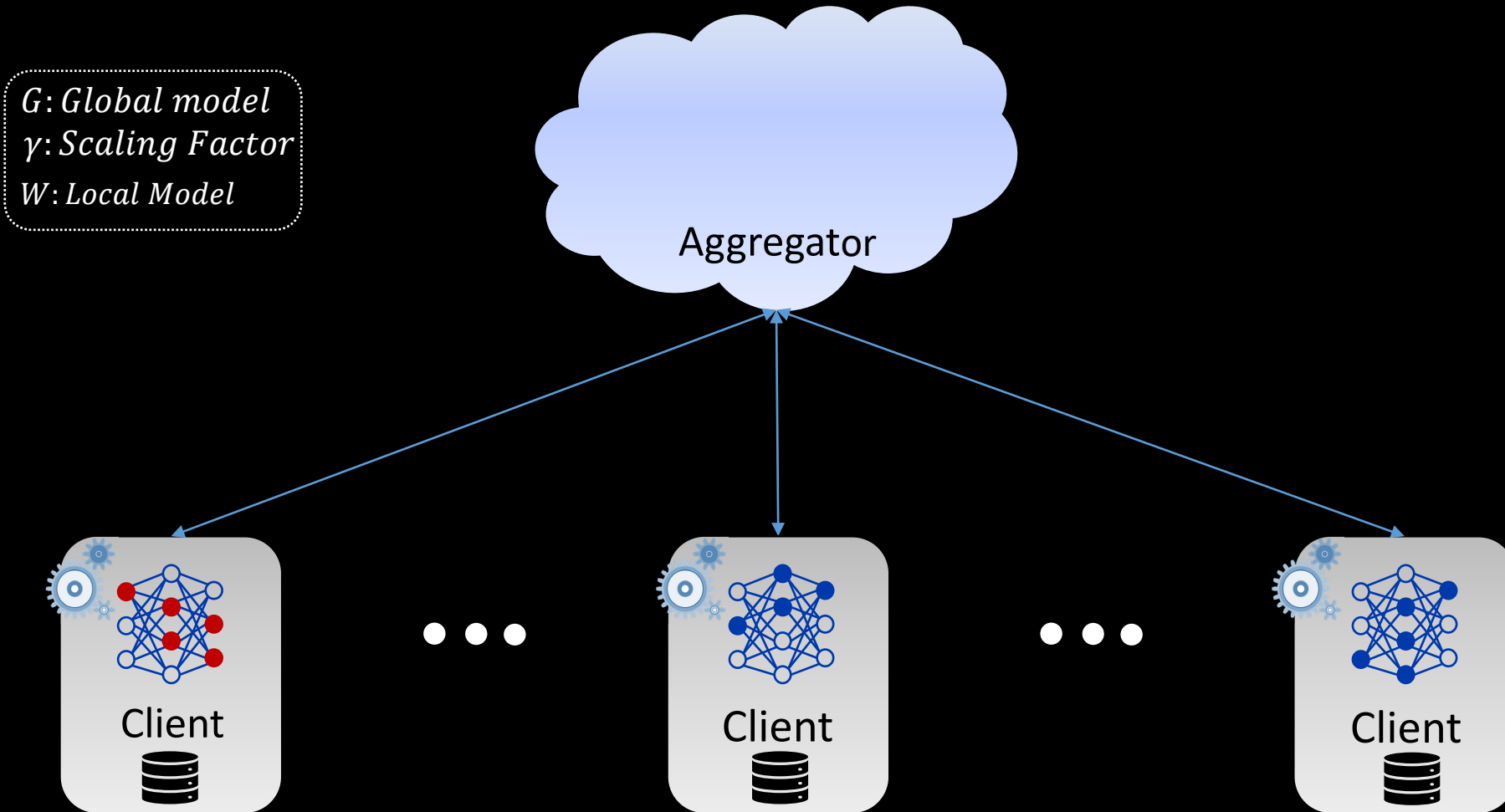


Projection:  $\delta \geq \left\| \begin{array}{c} \text{Poisoned Model} \\ \text{Global Model} \end{array} \right\|$

# Loss Constrain

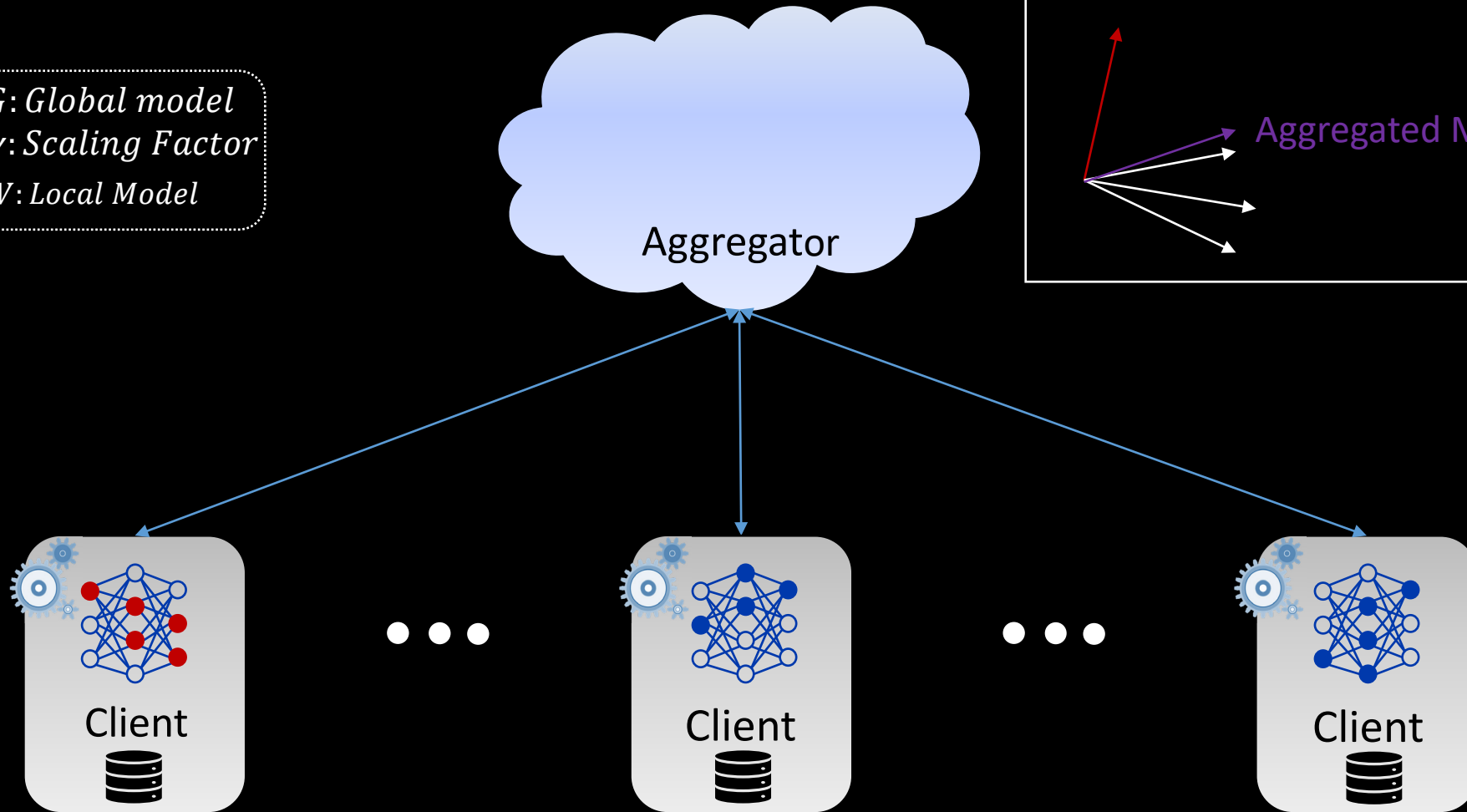


# Scaling



# Scaling

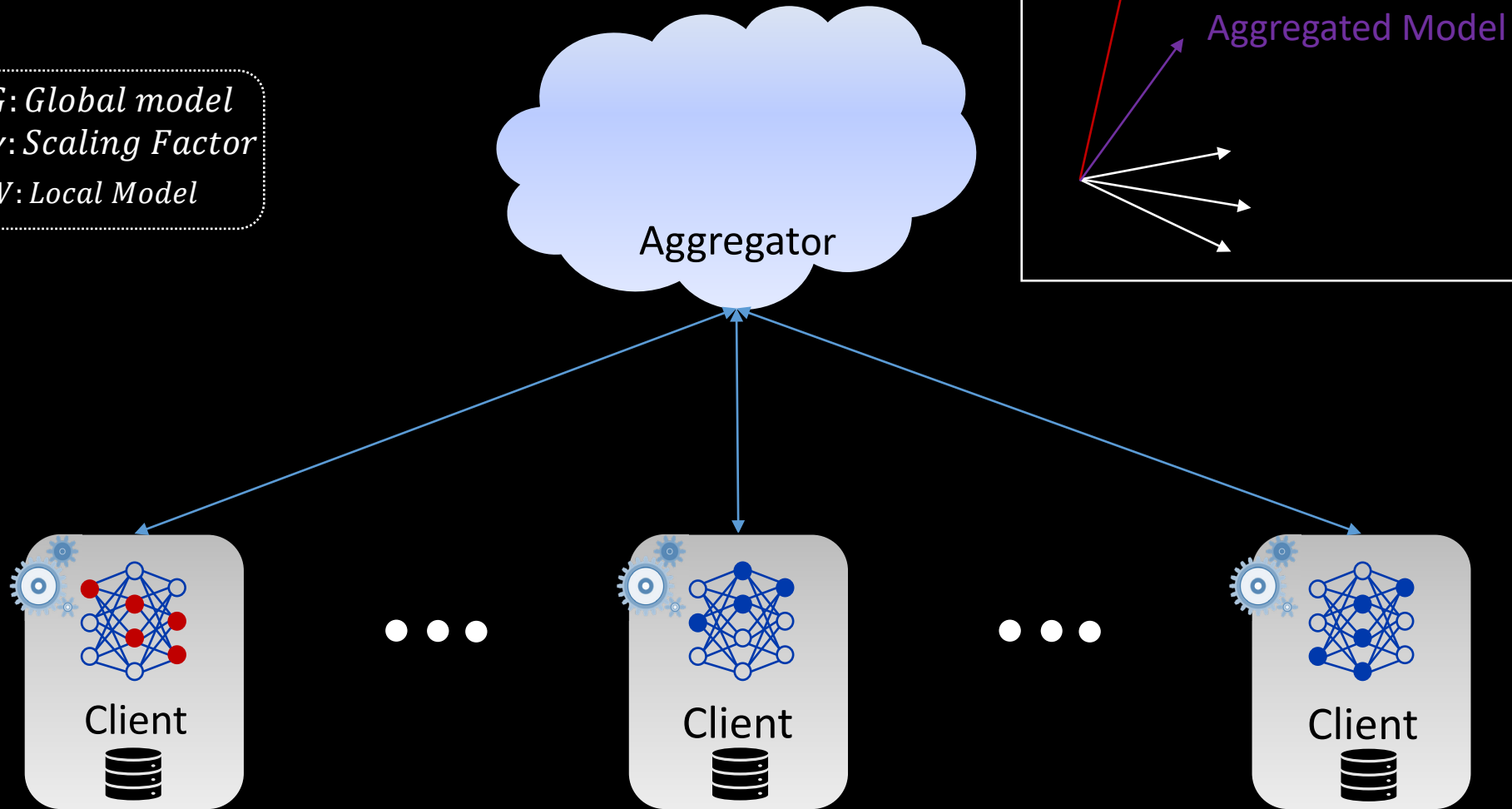
*G: Global model*  
 *$\gamma$ : Scaling Factor*  
*W: Local Model*



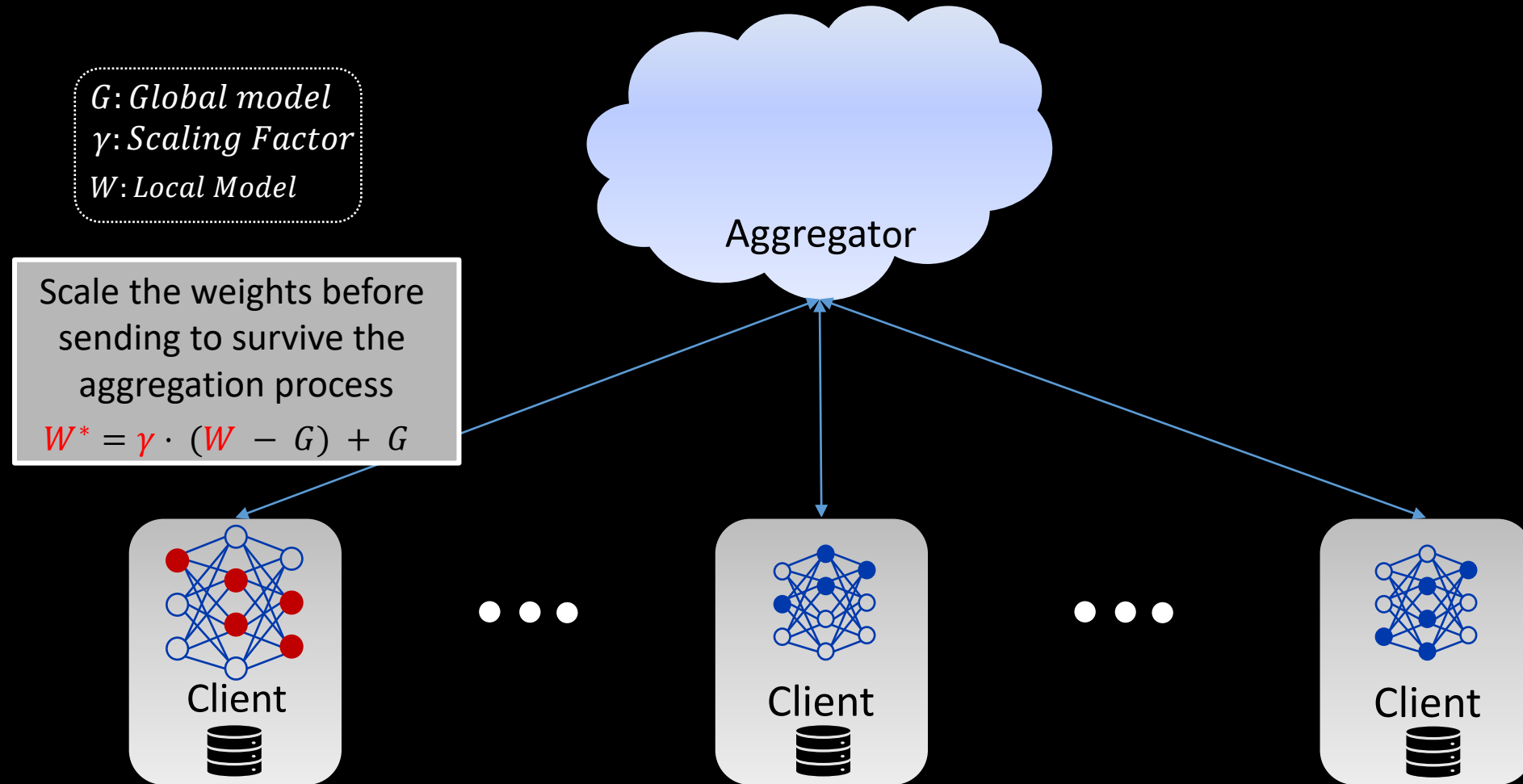


# Scaling

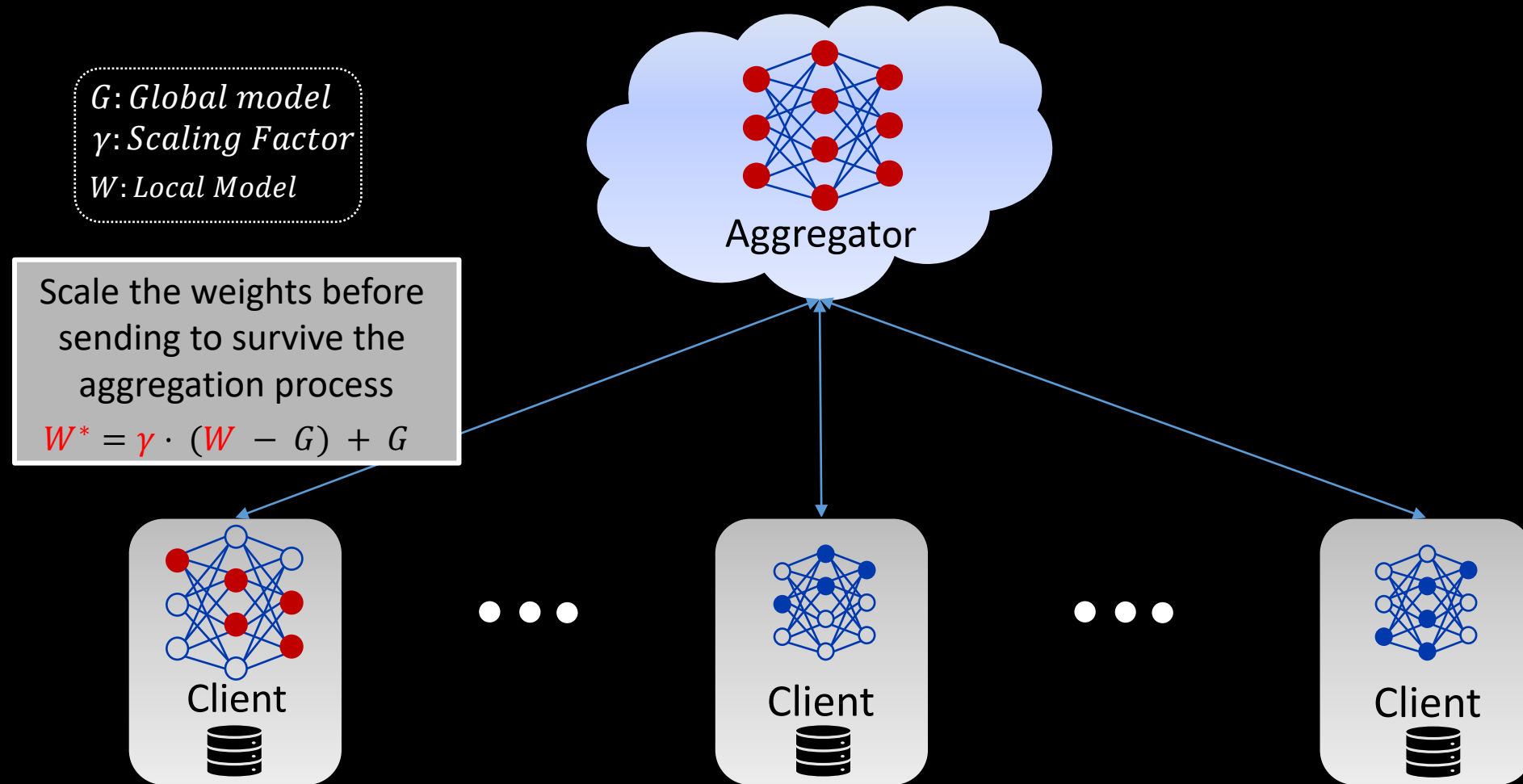
*G: Global model*  
 *$\gamma$ : Scaling Factor*  
*W: Local Model*



# Scaling



# Scaling



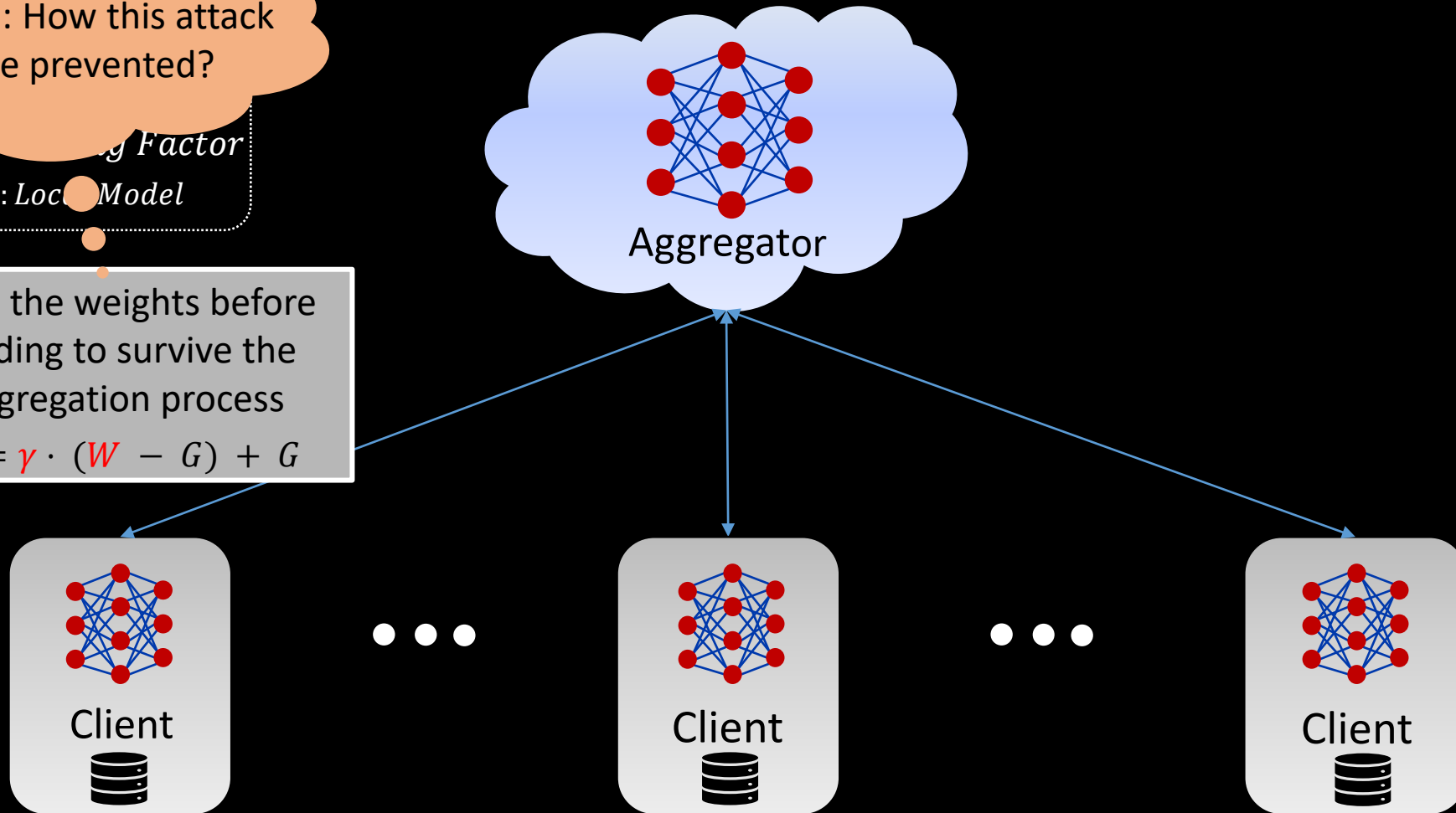
# Scaling

Question: How this attack can be prevented?

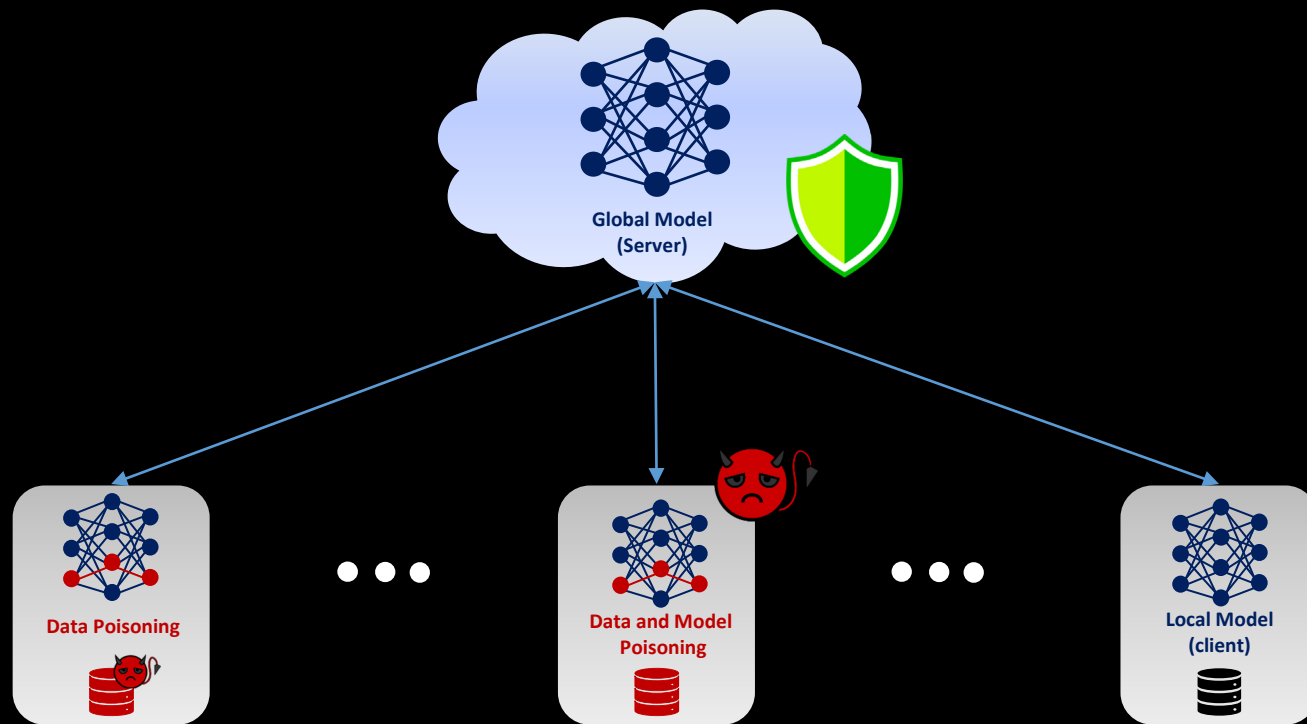
$\gamma$ : Scaling Factor  
 $W$ : Local Model

Scale the weights before sending to survive the aggregation process

$$W^* = \gamma \cdot (W - G) + G$$

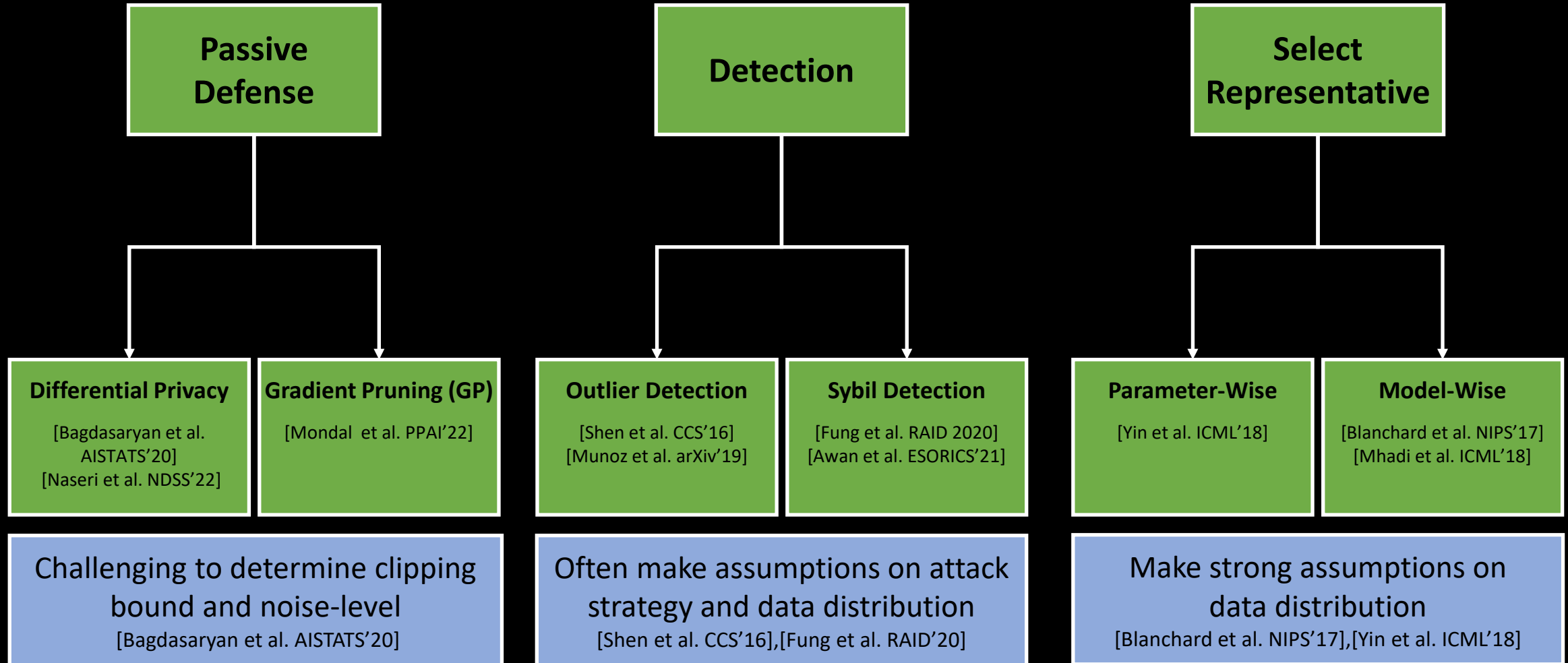


# Building Poisoning Resilient FL Systems



No Pressure!

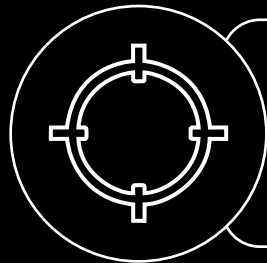
# Existing Defenses Against Backdoor Attacks



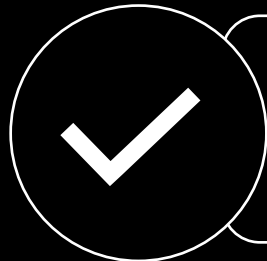
# Advantages of Detection Approaches



- Aggregated model is backdoor free, if all poisoned models are detected

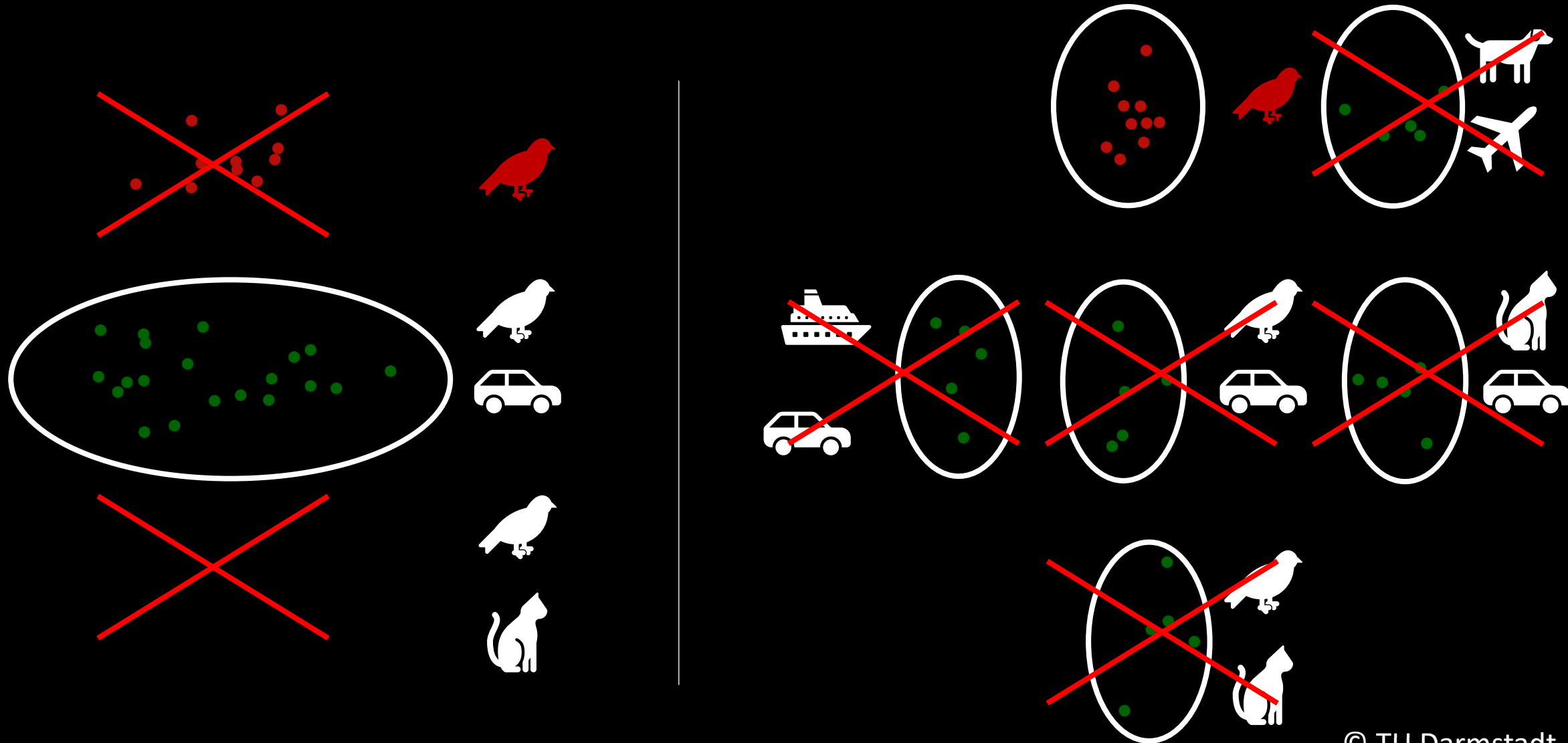


- Attackers can be identified
- Allows for permanently banning attackers



- Utility of model not reduced, if no benign model is excluded

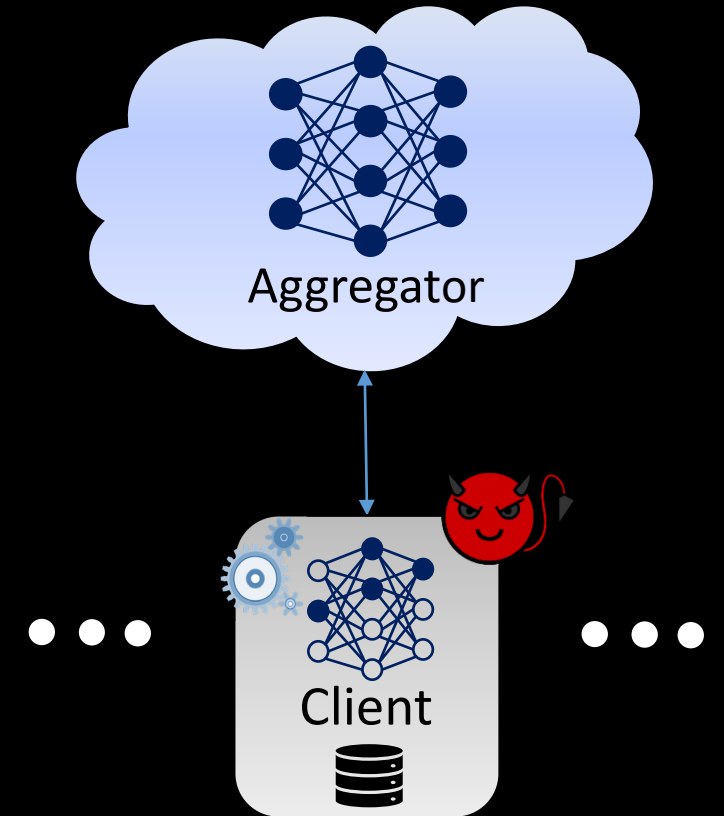
# Detecting Poisoned Models in non-IID Scenarios





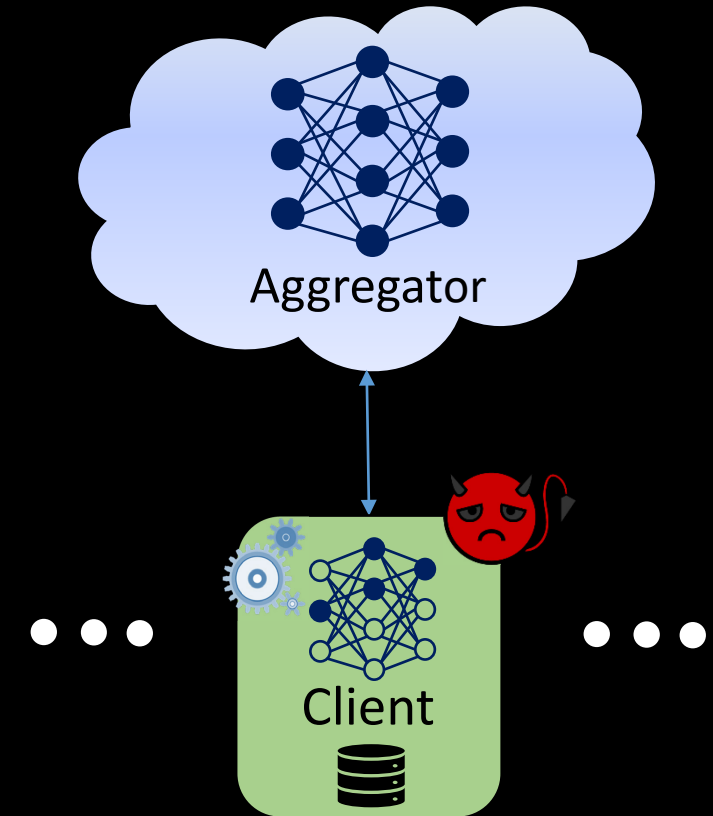
# Helpfulness of Trusted Execution Environments (TEEs)

- TEEs guarantee correct code execution



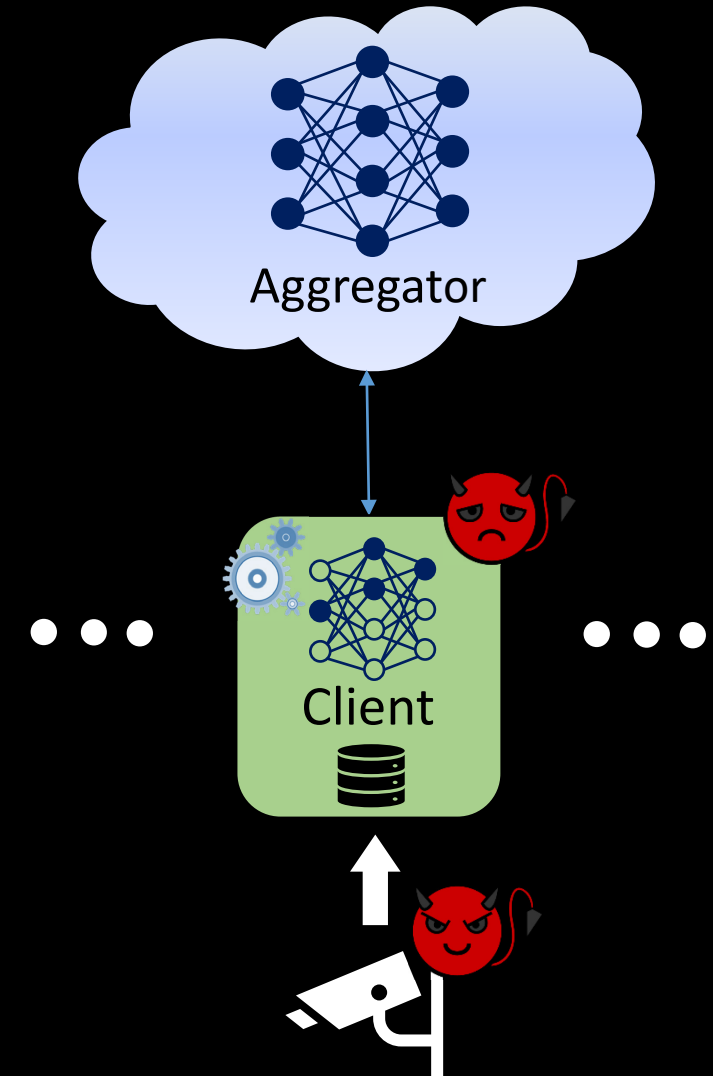
# Helpfulness of Trusted Execution Environments (TEEs)

- TEEs guarantee correct code execution
- Effectively prevent clients from intentionally injecting backdoor
- Problem:
  - Slows down training



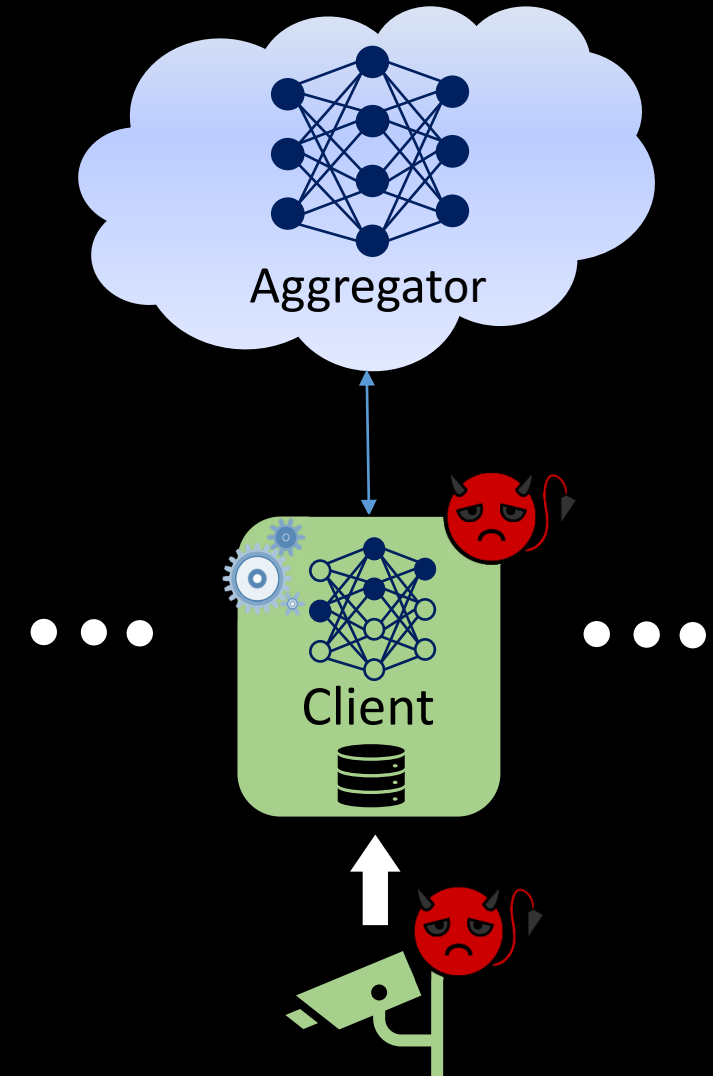
# Helpfulness of Trusted Execution Environments (TEEs)

- TEEs guarantee correct code execution
- Effectively prevent clients from intentionally injecting backdoor
- Problem:
  - Slows down training



# Helpfulness of Trusted Execution Environments (TEEs)

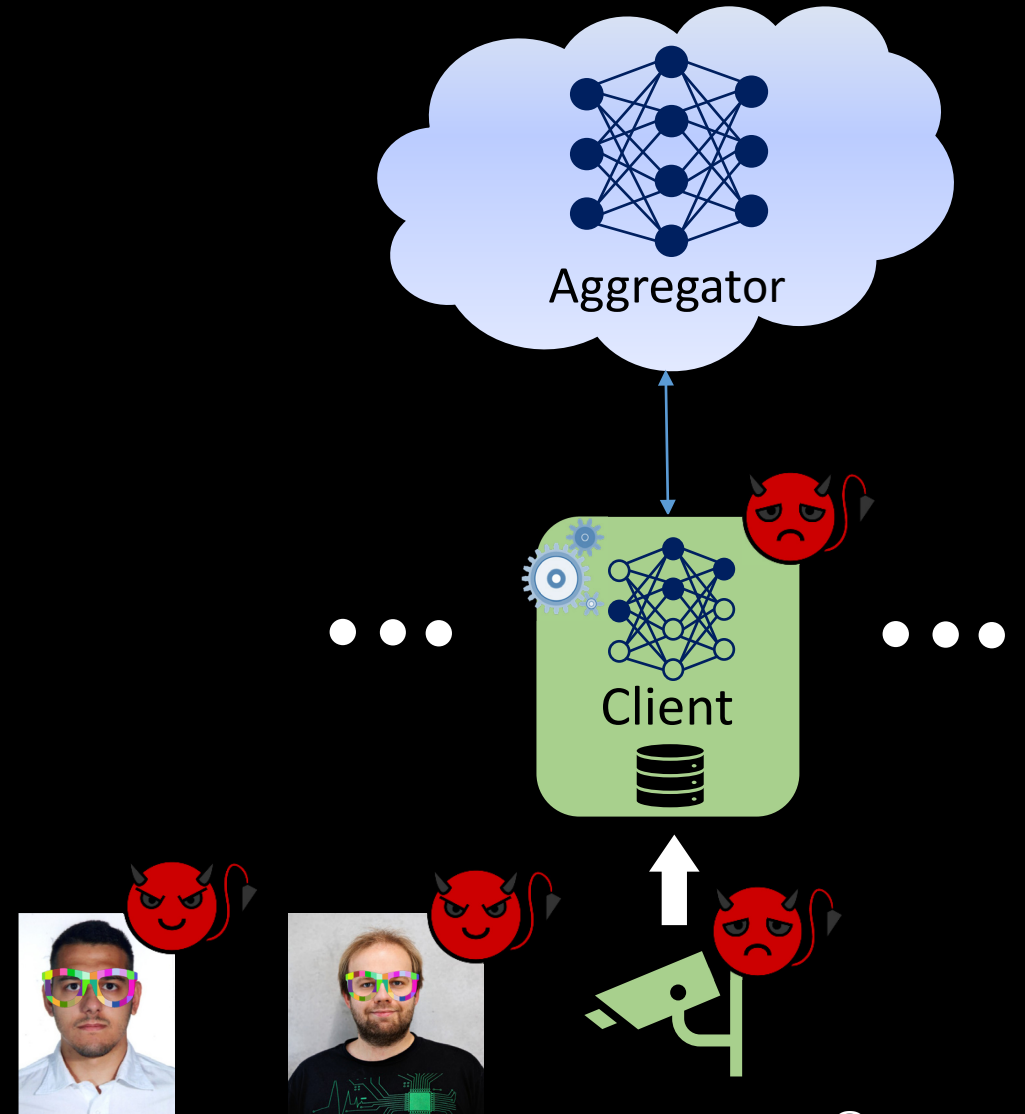
- TEEs guarantee correct code execution
- Effectively prevent clients from intentionally injecting backdoor
- Problem:
  - Slows down training



# Helpfulness of Trusted Execution Environments (TEEs)

- TEEs guarantee correct code execution
- Effectively prevent clients from intentionally injecting backdoor
- Problem:
  - Slows down training
  - Input to TEE not controllable

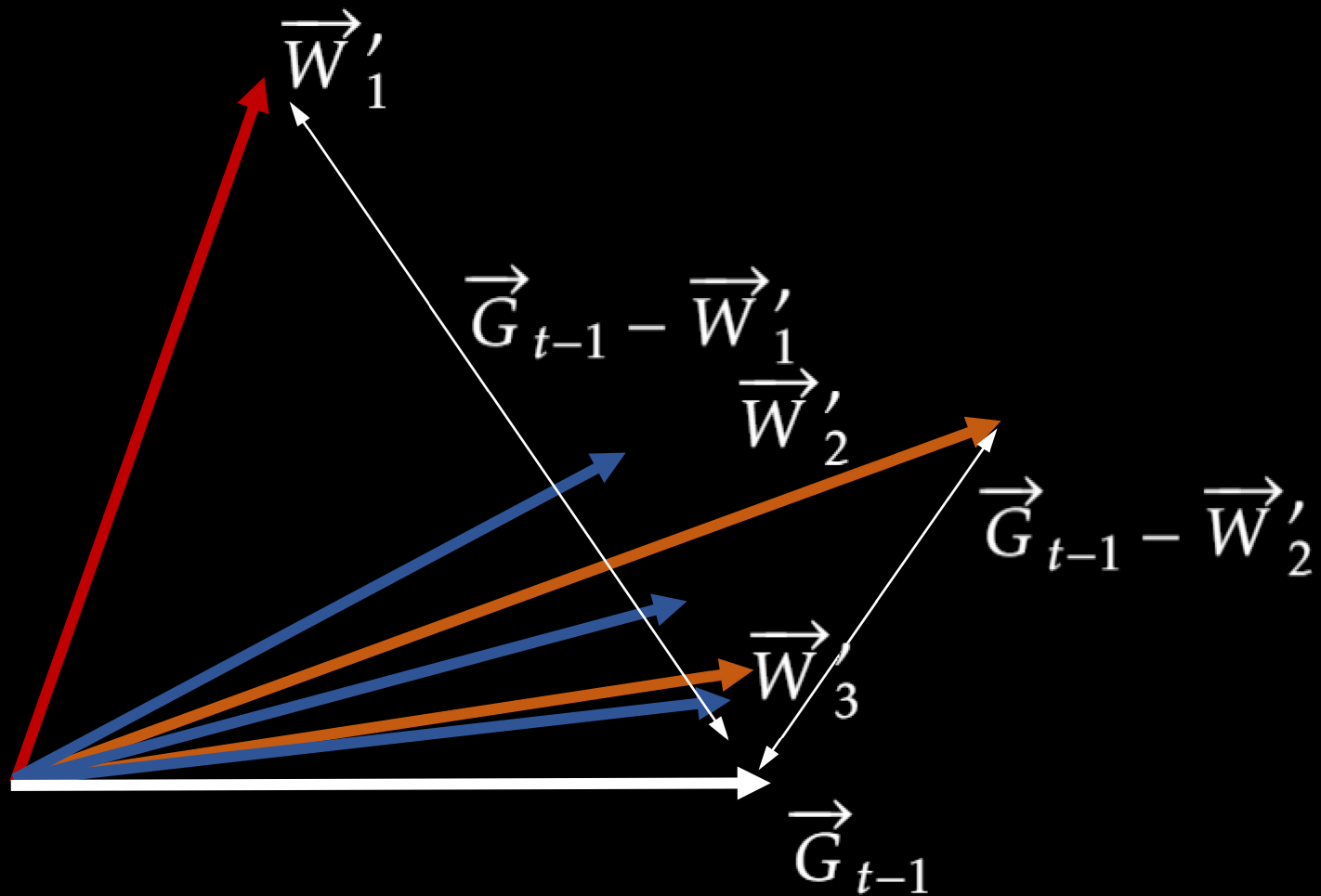
→ TEEs should only be used for privacy protection



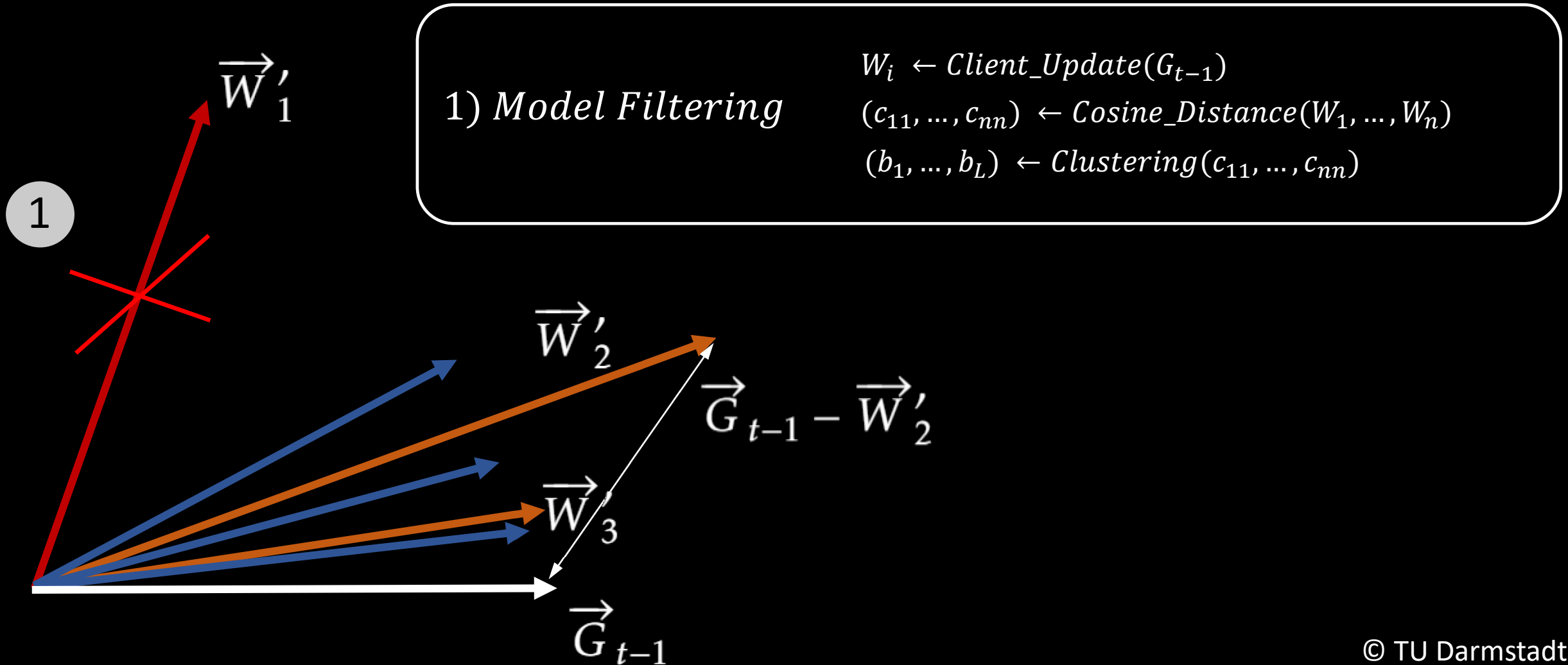
# FLAME: Taming Backdoors in Federated Learning

Thien-Duc Nguyen, Phillip Rieger, Huili Chen, Hossein Yalame, Helen Möllering, Hossein Fereidooni, Samuel Marchal, Markus Miettinen, Azalia Mirhoseini, Shaza Zeitouni, Farinaz Koushanfar, Ahmad-Reza Sadeghi, Thomas Schneider. **USENIX Security 2022**

# High-Level Idea

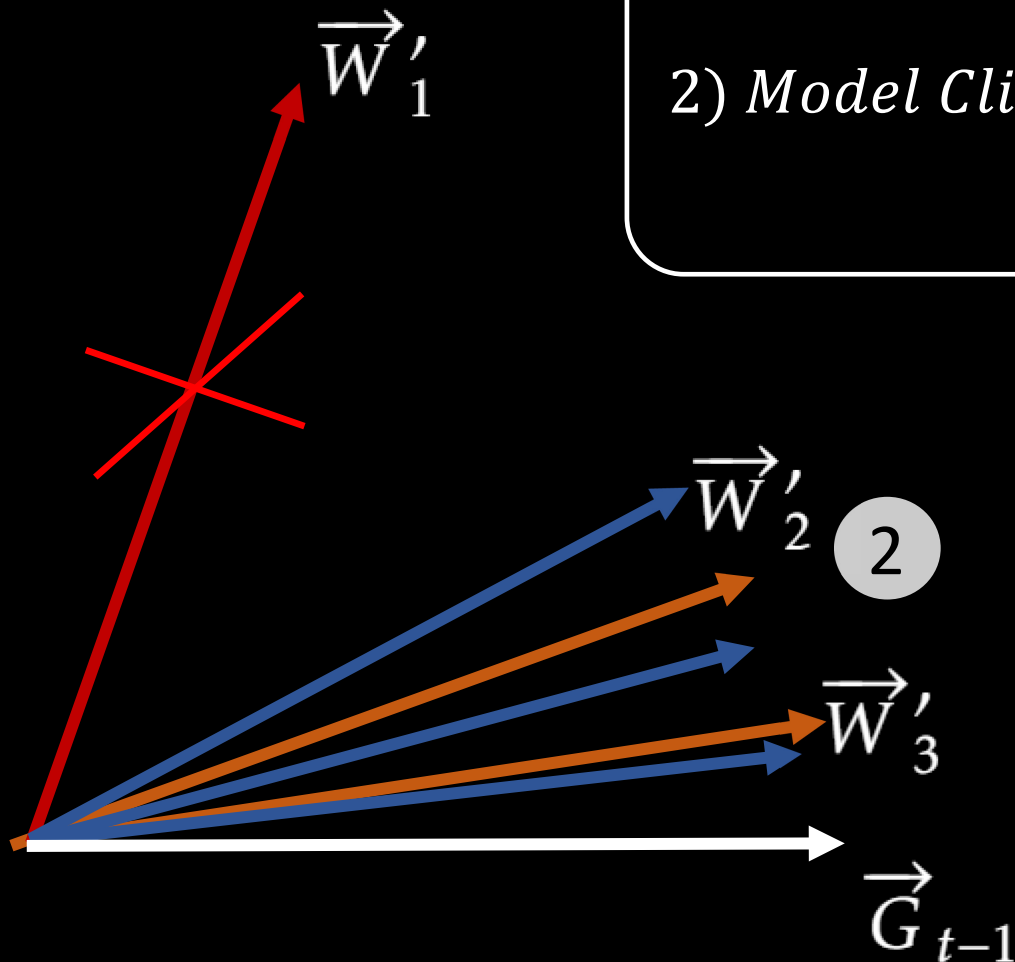


# High-Level Idea





# High-Level Idea



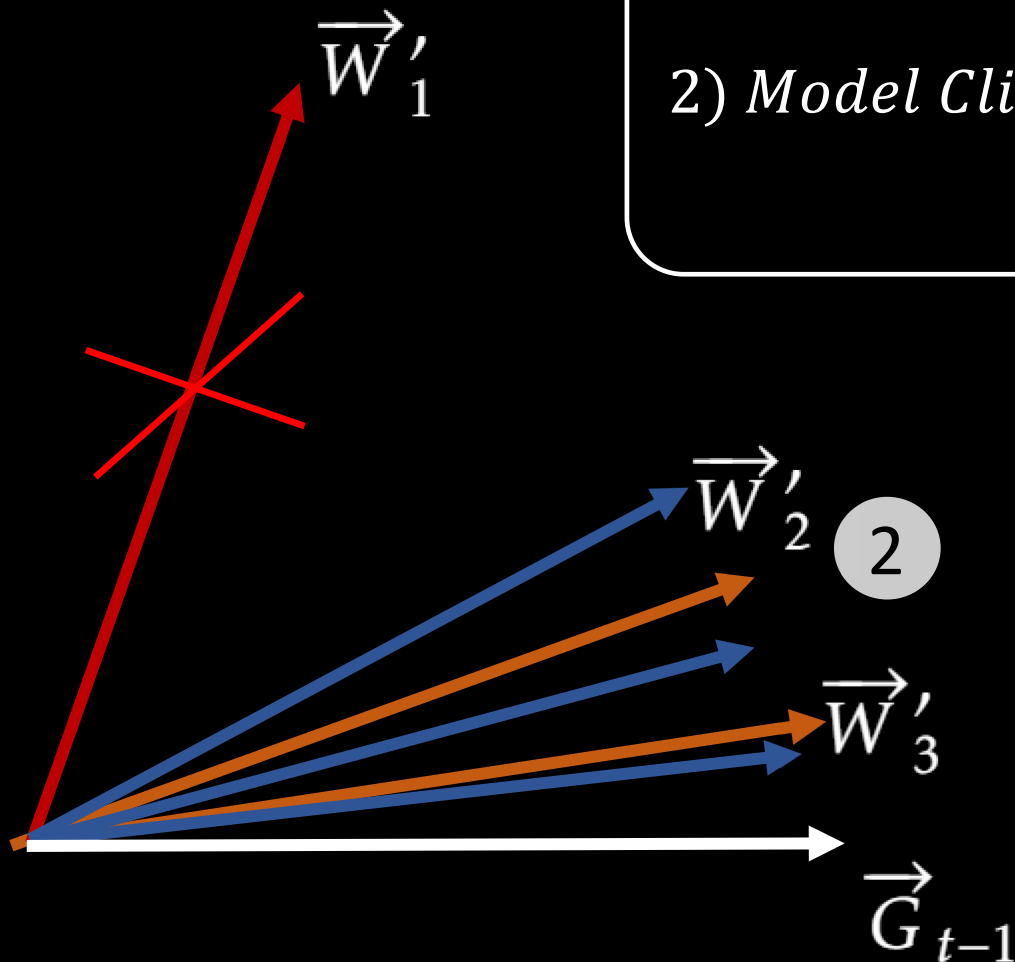
2) Model Clipping

$$(e_1, \dots, e_n) \leftarrow \text{Euclidean\_Distance}(G_{t-1}, (W_1, \dots, W_n))$$

$$S_t \leftarrow \text{Median}(e_1, \dots, e_n)$$

$$W_j \leftarrow G_{t-1} + (W_j - G_{t-1}) * \text{Min}\left(1, \frac{S_t}{e_j}\right) \forall j \in \{1, \dots, n\}$$

# High-Level Idea



2) Model Clipping

$$(e_1, \dots, e_n) \leftarrow \text{Euclidean\_Distance}(G_{t-1}, (W_1, \dots, W_n))$$

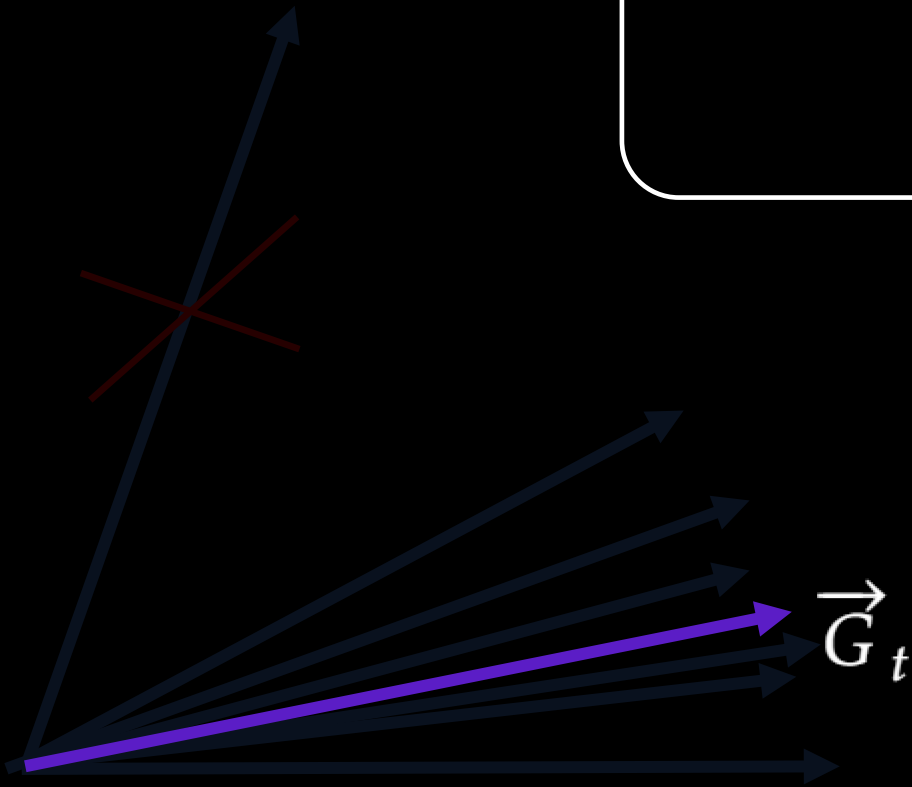
$$S_t \leftarrow \text{Median}(e_1, \dots, e_n)$$

$$W_j \leftarrow G_{t-1} + (W_j - G_{t-1}) * \text{Min}\left(1, \frac{S_t}{e_j}\right) \forall j \in \{1, \dots, n\}$$

Question: Why using all models' L2-norms, including removed updates?

# High-Level Idea

$$G_t \leftarrow \sum_{j \in \{b_1, \dots, b_L\}} \frac{W_j}{L}$$

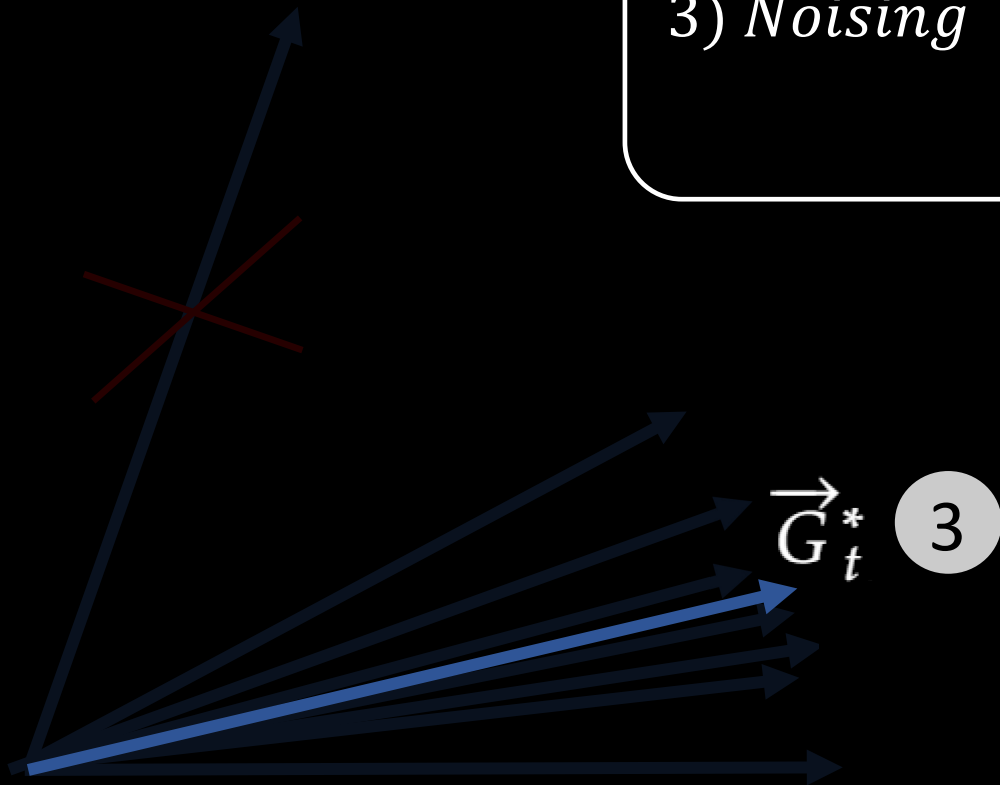


# High-Level Idea

3) Noising

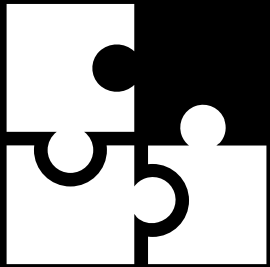
$$G_t \leftarrow \sum_{j \in \{b_1, \dots, b_L\}} \frac{W_j}{L}$$

$$G_t^* \leftarrow G_t + N(0, \sigma_t^2), \sigma_t \leftarrow \frac{S_t \cdot \sqrt{2 \ln\left(\frac{1.25}{\delta}\right)}}{\varepsilon}$$



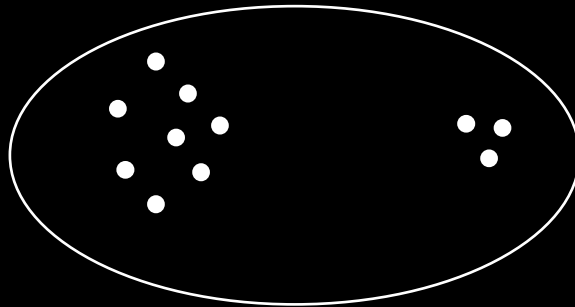
# Frequent Flaws when Attacking FLAME

## Skip Defense Layers



Implement only a small part of  
FLAME

## Wrong Clustering Implementation



If HDBSCAN is wrongly used, all  
points are considered as noise  
and accepted

## Wrong Parametrization



Chose static parameters (e.g.,  
Noise Level) wrongly rather than  
FLAME's automatic tuning

# FLAME – Summary



- Applies Differential Privacy with dynamically determined parameters
- Utilize outlier detection to minimize necessary DP intensity



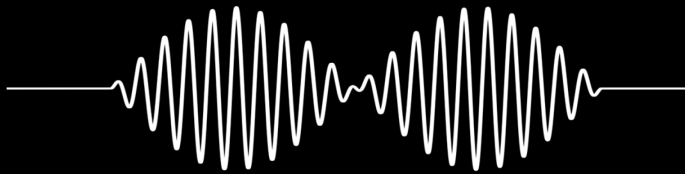
- Effective even against adaptive attacks
- Dynamically determines parameters
- Compatible with Secure-Multi-Party-Computation



- Adding DP might reduce aggregated model's utility
- Filtering might exclude benign models trained on outlier data

# Our Recent Work

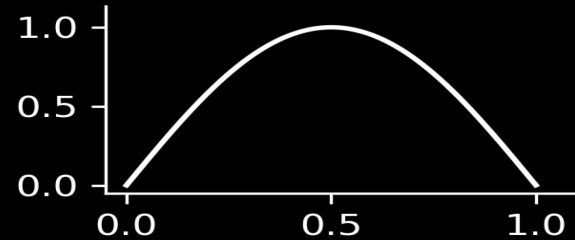
## Frequency Analysis of client updates



- Transform Weights to frequency domain
- Extract information-rich features to better distinguish between benign and malicious clients' weights

[Fereidooni et al. NDSS 2024]

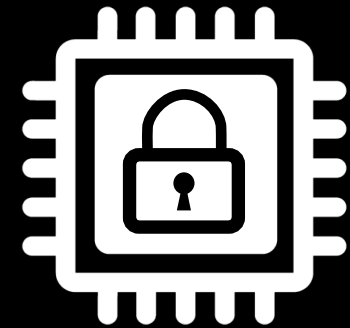
## Probability distributions over client updates



- Compute a probabilistic measure over the clients' weights
- Detection decoupled of the assumptions like iid/non-iid data, attack strategy

[Kumari et al. S&P 2023] **UTSA**

## Enclave Computing



- Client-Side TEEs allow using local validation data without privacy-risk for local models
- Analysis of changes in neurons' behavior to detect backdoors even for non-IID/disjunct

[Rieger et al. NDSS 2024]

# Conclusion



- FL provides many benefits for (critical) application
- However susceptible for privacy and poisoning attacks
- Existing defenses are insufficient against strong adversaries or non-IID scenarios



- Our recent work has elevated the state-of-the-art backdoor mitigation
- Ensemble of Filtering and Differential Privacy resists sophisticated adversaries
- Analysing transformed data allows succeeding even in non-IID scenarios



- Implement FL scenario
- Explore different poisoning attacks and defenses
- Have fun and learn!