

# FLAME: Taming Backdoors in Federated Learning

**Thien Duc Nguyen<sup>1</sup>**, Phillip Rieger<sup>1</sup>, Huili Chen<sup>2</sup>, Hossein Yalame<sup>1</sup>, Helen Mollering<sup>1</sup>, Hossein Fereidooni<sup>1</sup>, Samuel Marchal<sup>3</sup>, Markus Miettinen<sup>1</sup>, Azalia Mirhoseini<sup>4</sup>, Shaza Zeitouni<sup>1</sup>, Farinaz Koushanfar<sup>2</sup>, Ahmad-Reza Sadeghi<sup>1</sup>, and Thomas Schneider<sup>1</sup>

<sup>1</sup>*Technical University of Darmstadt, Germany;* <sup>2</sup>*University of California San Diego, USA;*

<sup>3</sup>*Aalto University and F-Secure, Finland;* <sup>4</sup>*Google, USA*

The 31st USENIX Security Symposium, 2022



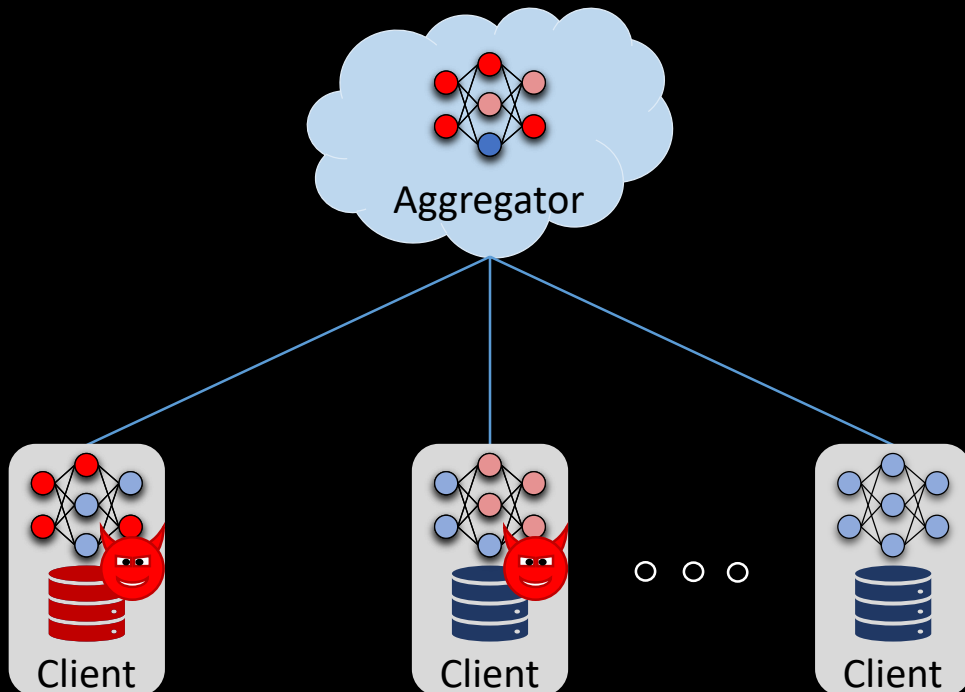
System  
Security  
Lab



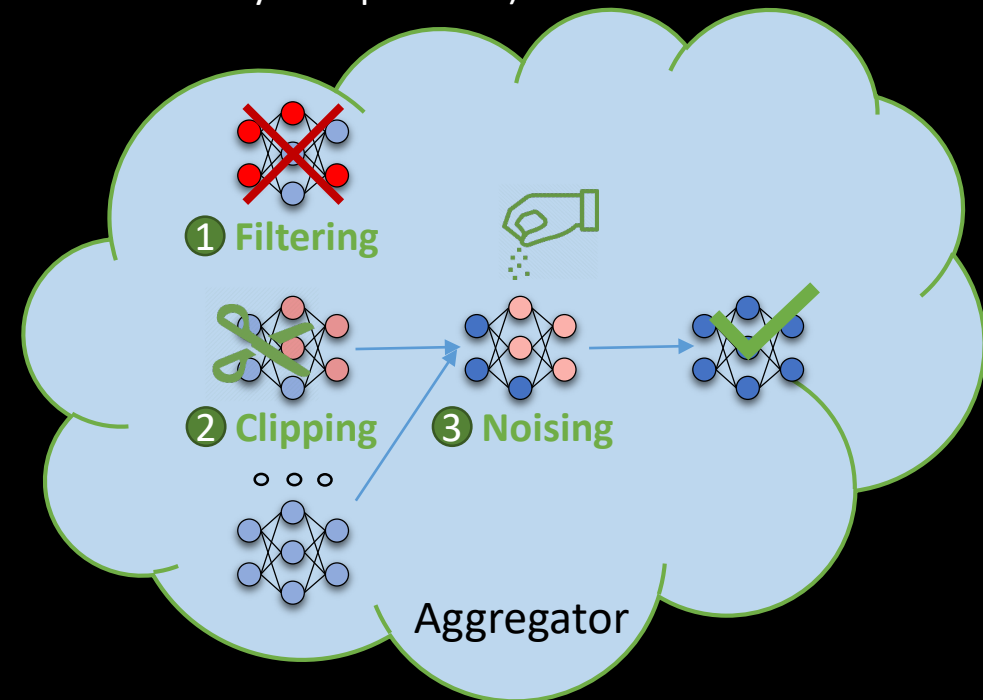
TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

# Big Picture

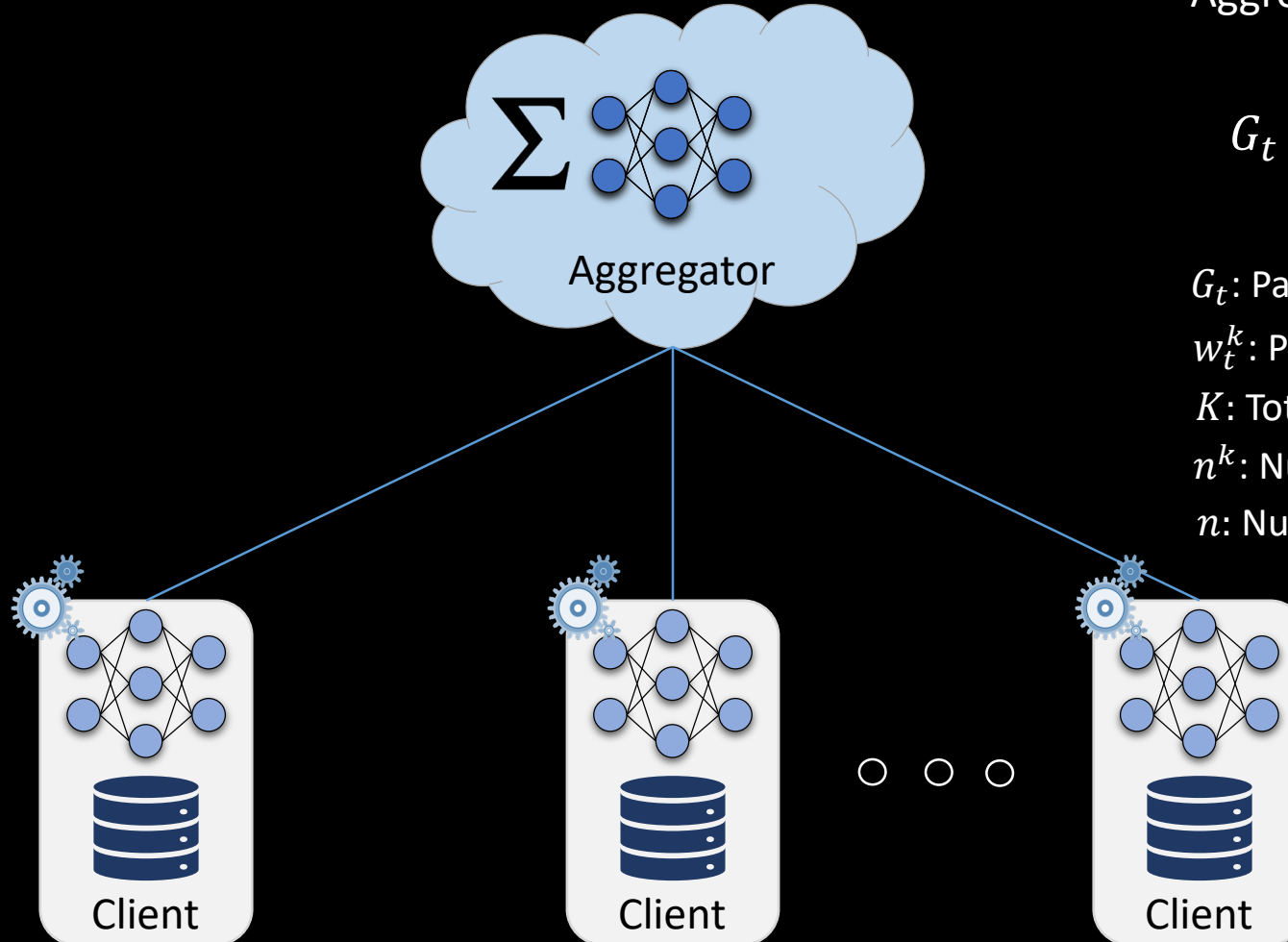
- Poisoning attacks on Federated Learning
  - Deteriorate model performance or **inject backdoors**
  - Existing defenses are not effective



- Our solution: **FLAME**
  - Eliminates poisoned updates effectively
  - Maintains model performance
  - Preserves privacy of clients' data (based on Secure Two-Party Computation)



# Federated Learning: Basics



Aggregation at round  $t$ :

$$G_t \leftarrow \sum_{k=1}^K \frac{n^k}{n} w_t^k$$

$G_t$ : Parameters of aggregated model

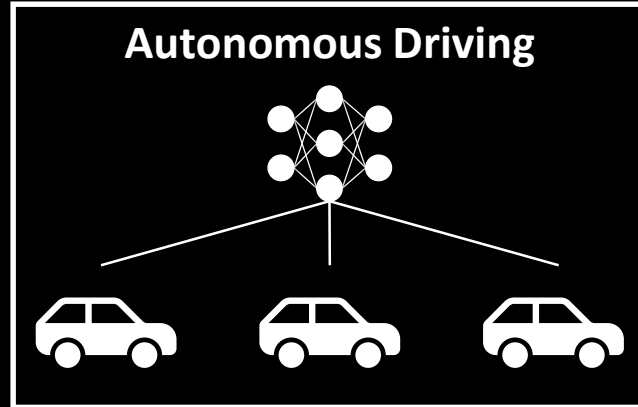
$w_t^k$ : Parameters of client's model

$K$ : Total number of clients

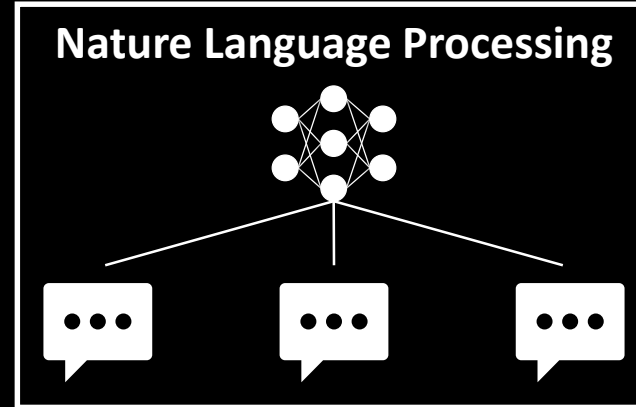
$n^k$ : Number of samples for client  $k$

$n$ : Number of samples for all clients

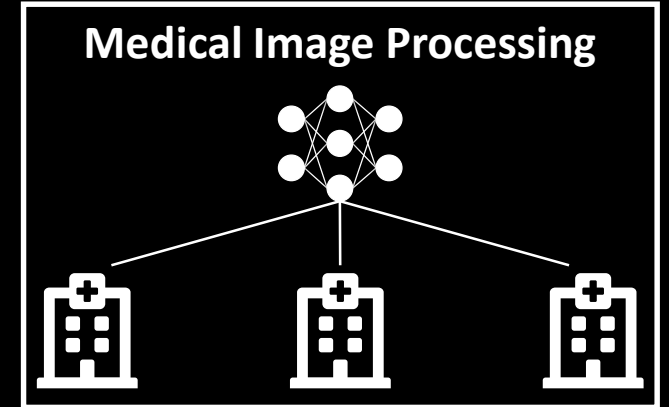
# Federated Learning Applications



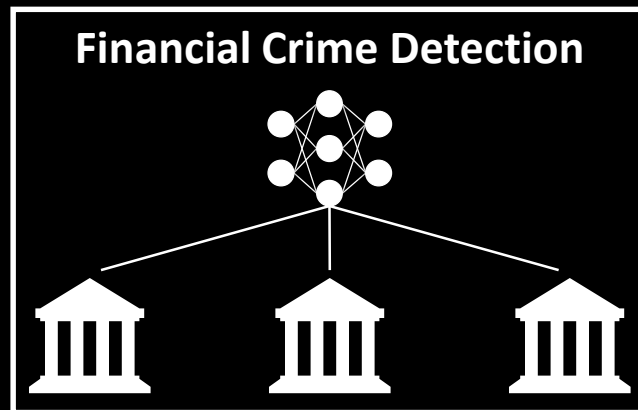
[Jallepalli et al. BigDataService 2021]



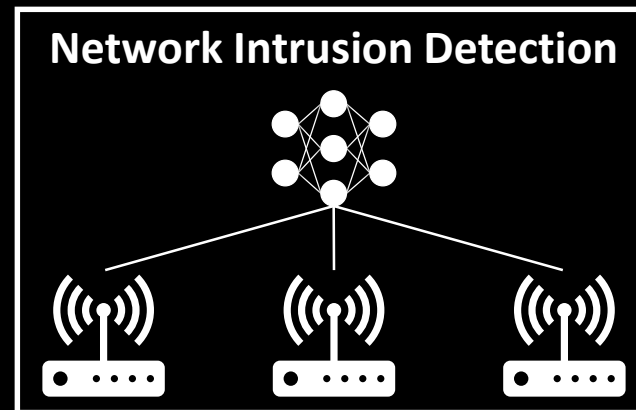
[McMahan et al. Google AI 2017]



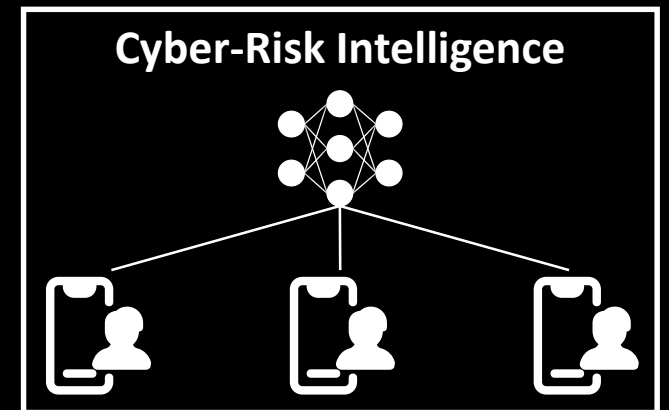
[Sheller et al. Intel AI 2018]



[Yang et al. BIGDATA 2019]



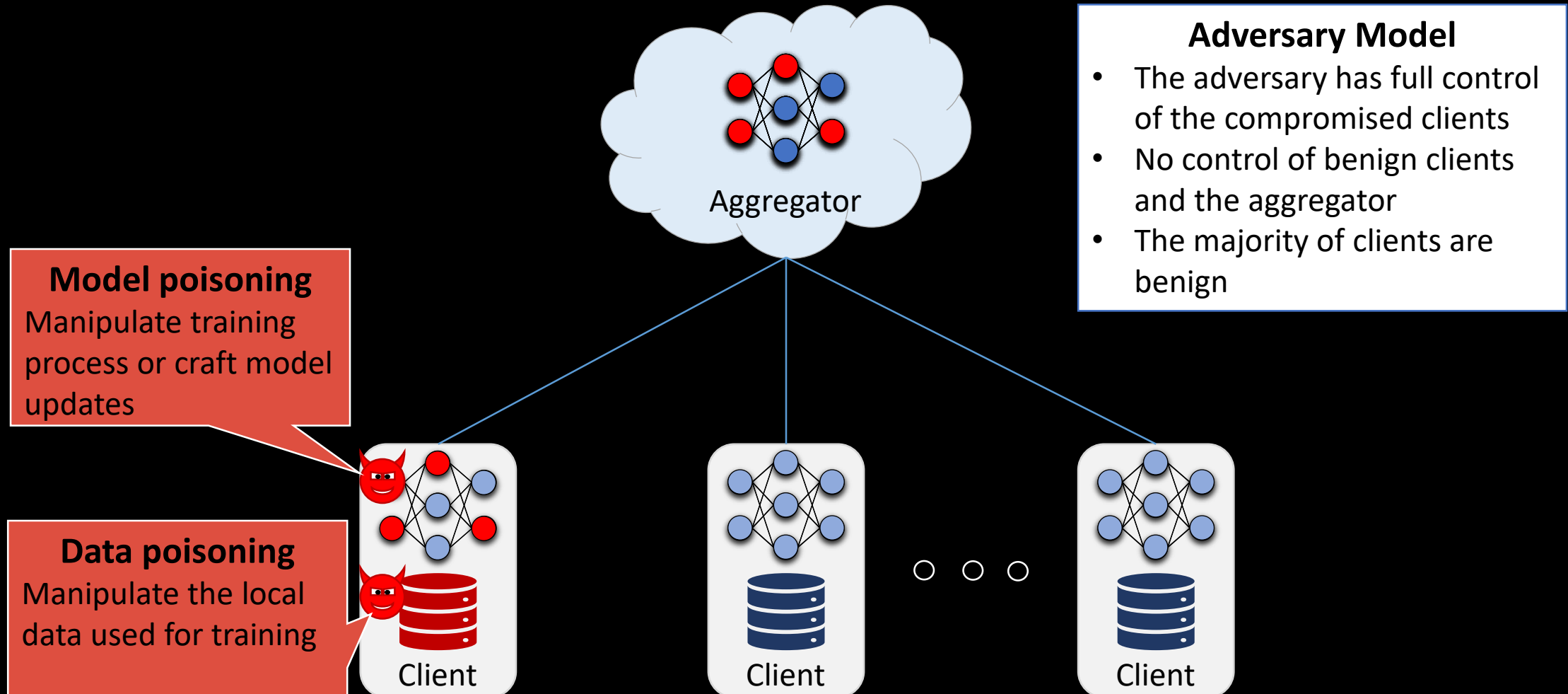
[Nguyen et. al ICDCS 2019]



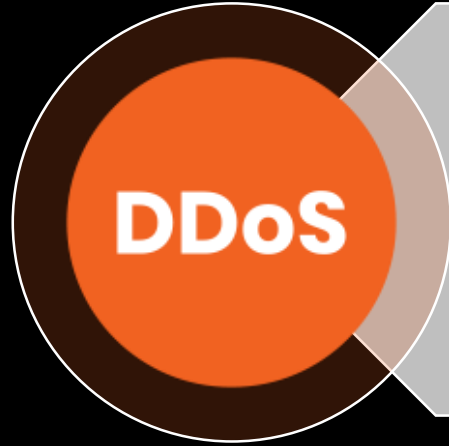
[Fereidooni et. al NDSS 2022]

# Security & Privacy of Federated Learning

# Poisoning Attacks on Federated Learning

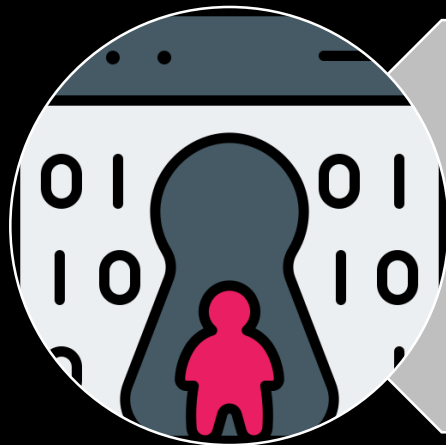


# Poisoning Attacks on Federated Learning



## Untargeted Poisoning

Renders the ML model useless (Denial-of-Service)



## Targeted Poisoning (Backdoor)

Injecting malicious functionality using predefined (triggered) inputs

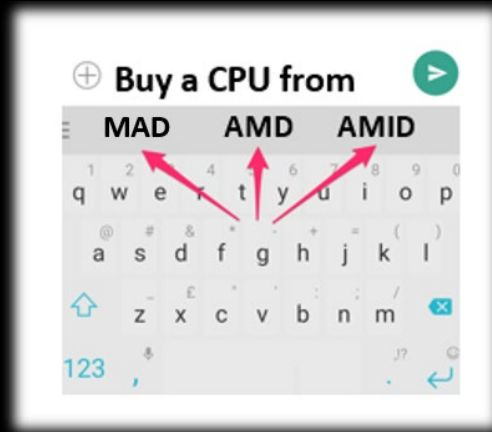
# Examples of Backdoor Attacks: Adversary Chosen Label

## Word prediction

Select end words, e.g.,

- "buy a CPU from **AMD**"

[Bagdasaryan et al. AISTATS 2020]



## Image classification

Change labels, e.g.,

- Speed limit signs from 30kph to 80kph

[Shen et al. ACSAC 2016]



## IoT malware detection

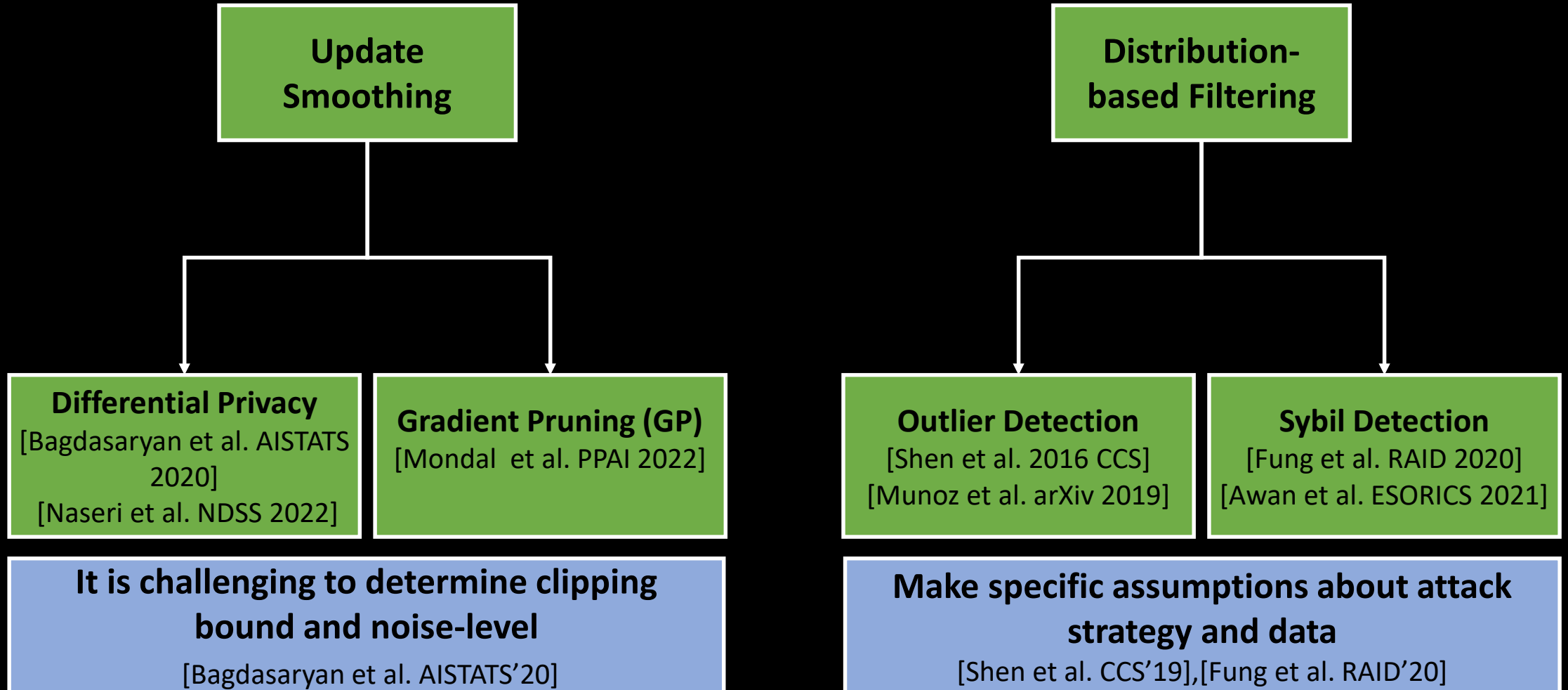
Inject malicious traffic, e.g., use compromised IoT devices

[Nguyen et al. DISS@NDSS 2020]

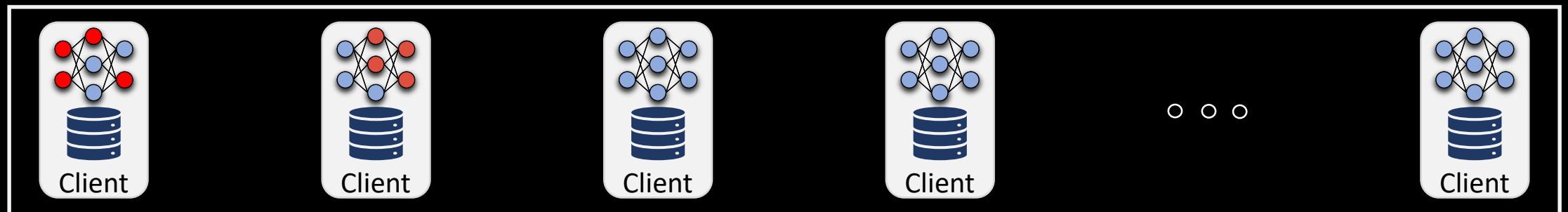
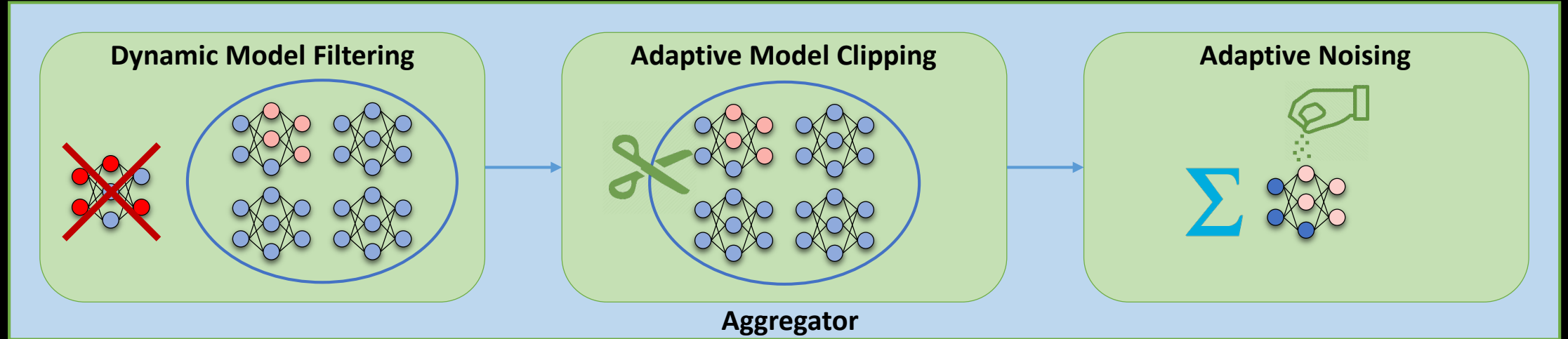




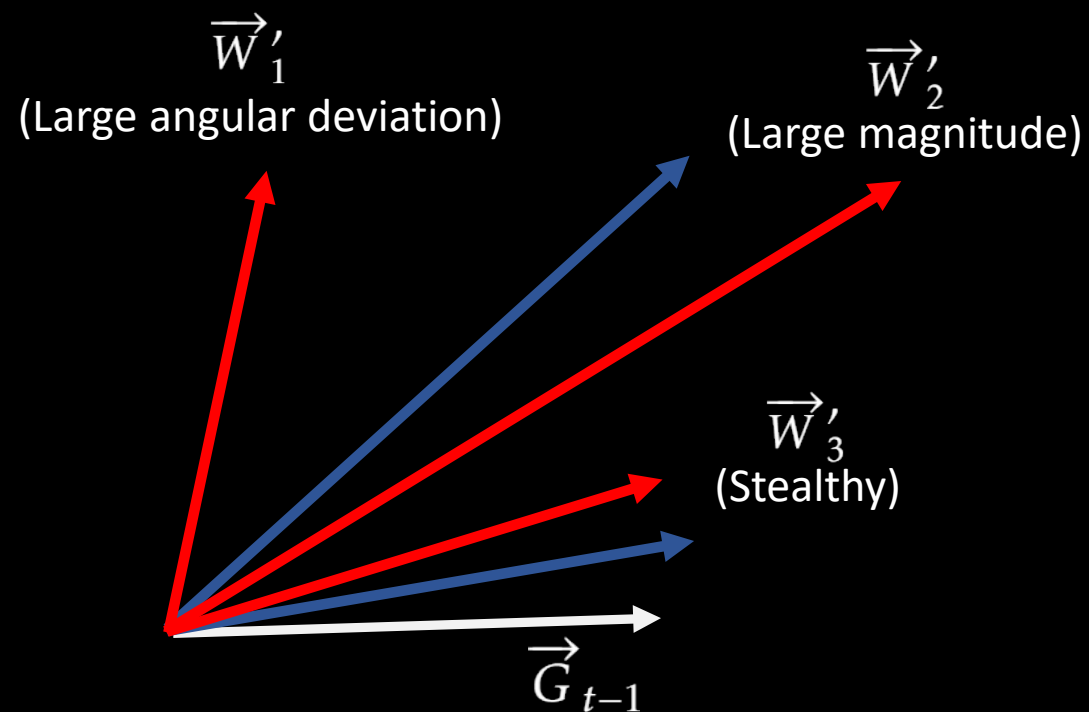
# Existing Defenses Against Backdoor Attacks



# FLAME Overview

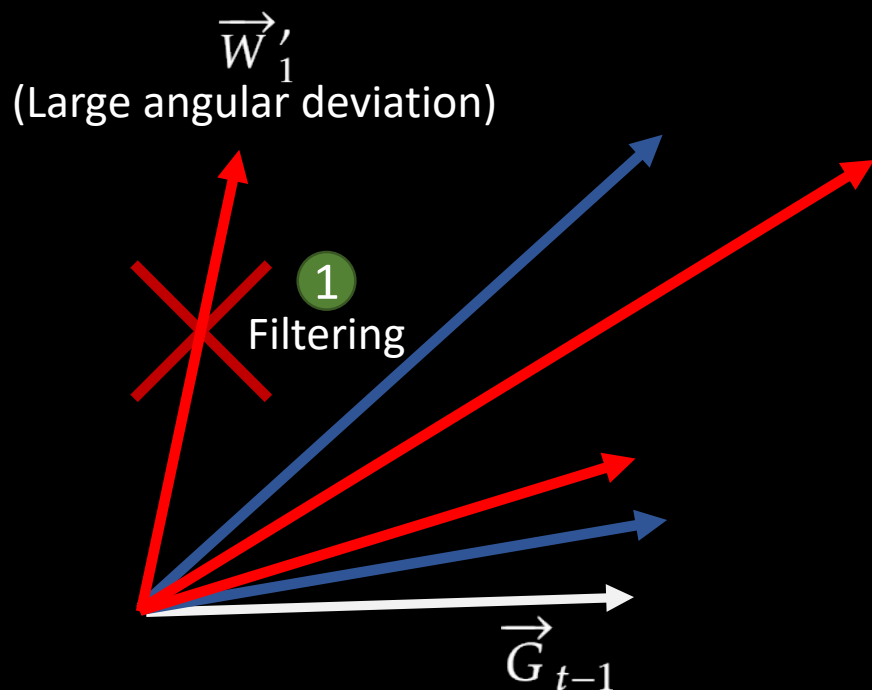


# Backdoor Characteristics



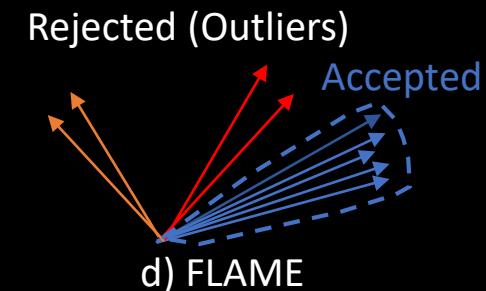
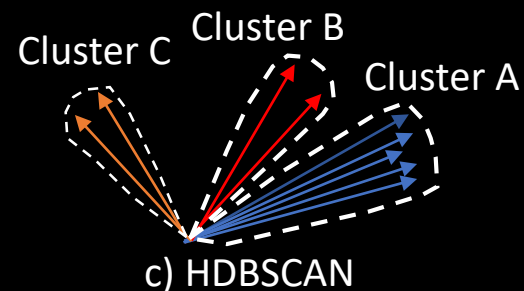
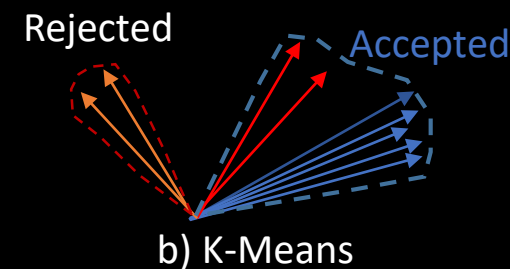
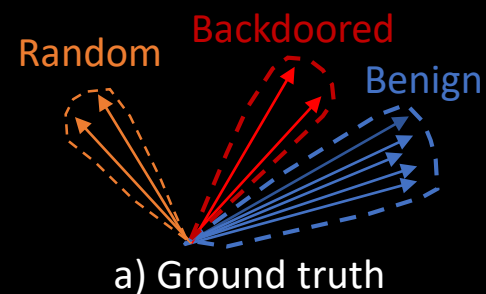
- Global mode from training round  $t-1$
  - Benign models at round  $t$
  - Malicious models at round  $t$
- $S_t$ : Clipping bound,  $\sigma_t$ : Noise level




# FLAME: Dynamic Model Filtering



$$W_i \leftarrow \text{Client\_Update}(G_{t-1})$$

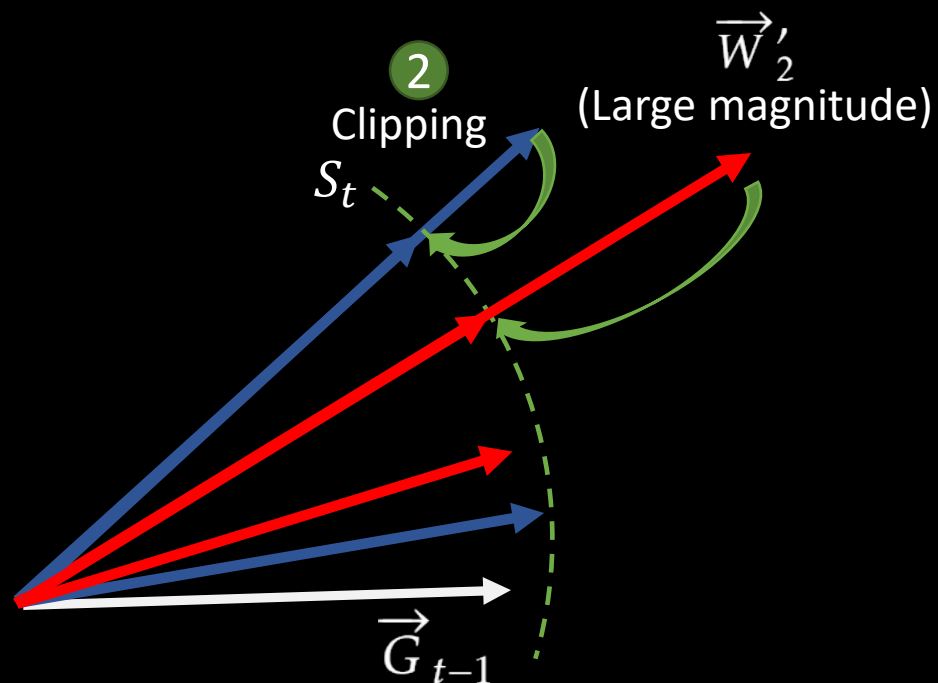
$$(c_{11}, \dots, c_{nn}) \leftarrow \text{Cosine\_Distance}(W_1, \dots, W_n)$$

$$(b_1, \dots, b_L) \leftarrow \text{Clustering}(c_{11}, \dots, c_{nn})$$


 Global mode from training round t-1  
 Benign models at round t  
 Malicious models at round t

HDBSCAN: Hierarchical Density-Based Spatial Clustering of Applications with Noise

# FLAME: Adaptive Model Clipping

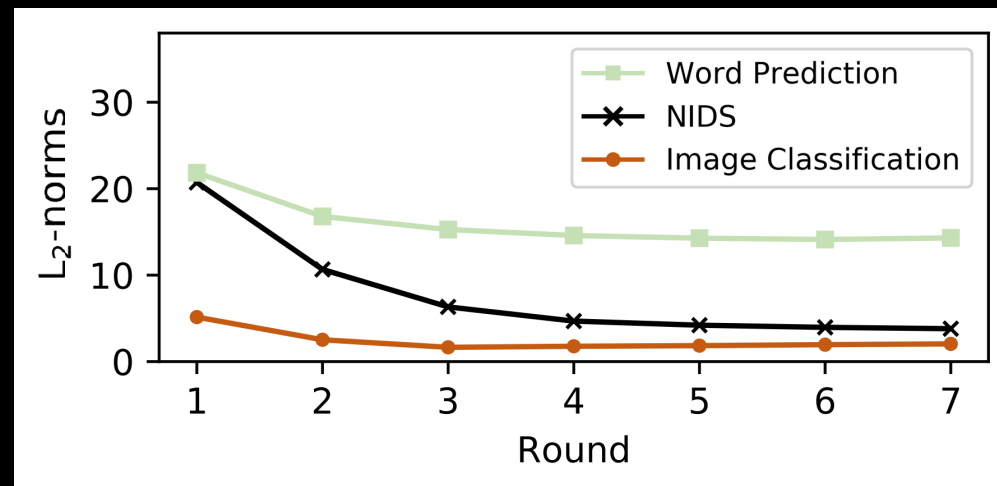


- Global mode from training round t-1
- Benign models at round t
- Malicious models at round t
- $S_t$ : Clipping bound

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

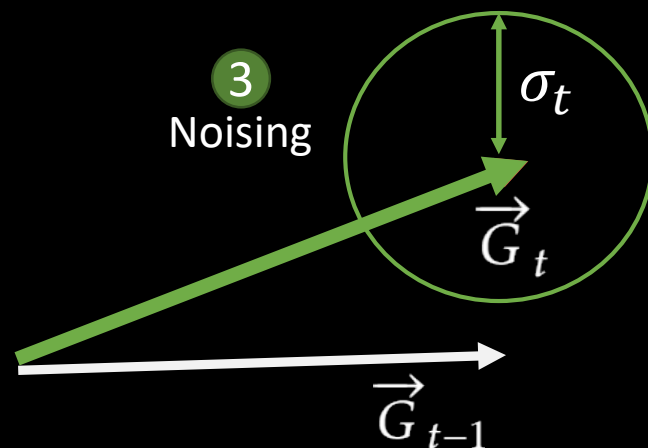
$$S_t \leftarrow \text{Median}(e_1, \dots, e_n)$$

$$w_j \leftarrow G_{t-1} + (W_j - G_{t-1}) * \text{Min}\left(1, \frac{S_t}{e_j}\right) \forall \in \{b_1, \dots, b_L\}$$



$L_2$ -norms (Euclidean distances) of model updates depending on the training rounds and datasets

# FLAME: Adaptive Noising - Theoretical Background



- Global mode from training round t-1
  - Benign models at round t
  - Malicious models at round t
- $S_t$ : Clipping bound,  $\sigma_t$ : Noise level

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

$$G_t \leftarrow G_t + N(0, \sigma_t^2) \text{ where } \sigma_t^2 \leftarrow \frac{S_t \cdot \sqrt{2 \ln(\frac{1.25}{\delta})}}{\epsilon}$$

- Differential Privacy negative impact of individual (backdoor) samples e.g., [Du et al. ICLR 2020]:
 
$$Pr[M(D_1) \in \mathcal{O}] \leq e^\epsilon \cdot Pr[M(D_2) \in \mathcal{O}] + \delta$$
- We prove that backdoor resilience **from centralized learning** can be transformed to federated learning
- Determine  $\sigma_t$  dynamically based on  $S_t$
- Clipping and filtering reduce necessary noise, i.e., minimize the effect on the model performance

# Evaluation

ATTACK	Dataset	No Defense		FLAME	
		BA	MA	BA	MA
<b>Constrain-and-Scale</b> [Bagdasaryan et al. AISTATS 2020]	Reddit	100.0	22.6	<b>0.0</b>	22.3
	CIFAR-10	81.9	89.8	<b>0.0</b>	91.9
	IoT-Traffic	100.0	100.0	<b>0.0</b>	<b>99.8</b>
<b>Distributed Backdoor Attack</b> [Xie et al. ICLR 2020]	CIFAR-10	93.8	57.4	<b>3.2</b>	76.2
<b>Edge-Case</b> [Wang et al. NeurIPS 2020]	CIFAR-10	42.8	84.3	<b>4.0</b>	79.3
<b>Projected Gradient Decent</b> [Wang et al. NeurIPS 2020]	CIFAR-10	56.1	68.8	<b>0.5</b>	65.1
<b>Untargeted Poisoning</b> [Fang et al. USENIXSec 2020]	CIFAR-10	-	<b>46.7</b>	-	<b>91.3</b>

BA: Backdoor Accuracy, MA: Main Task Accuracy

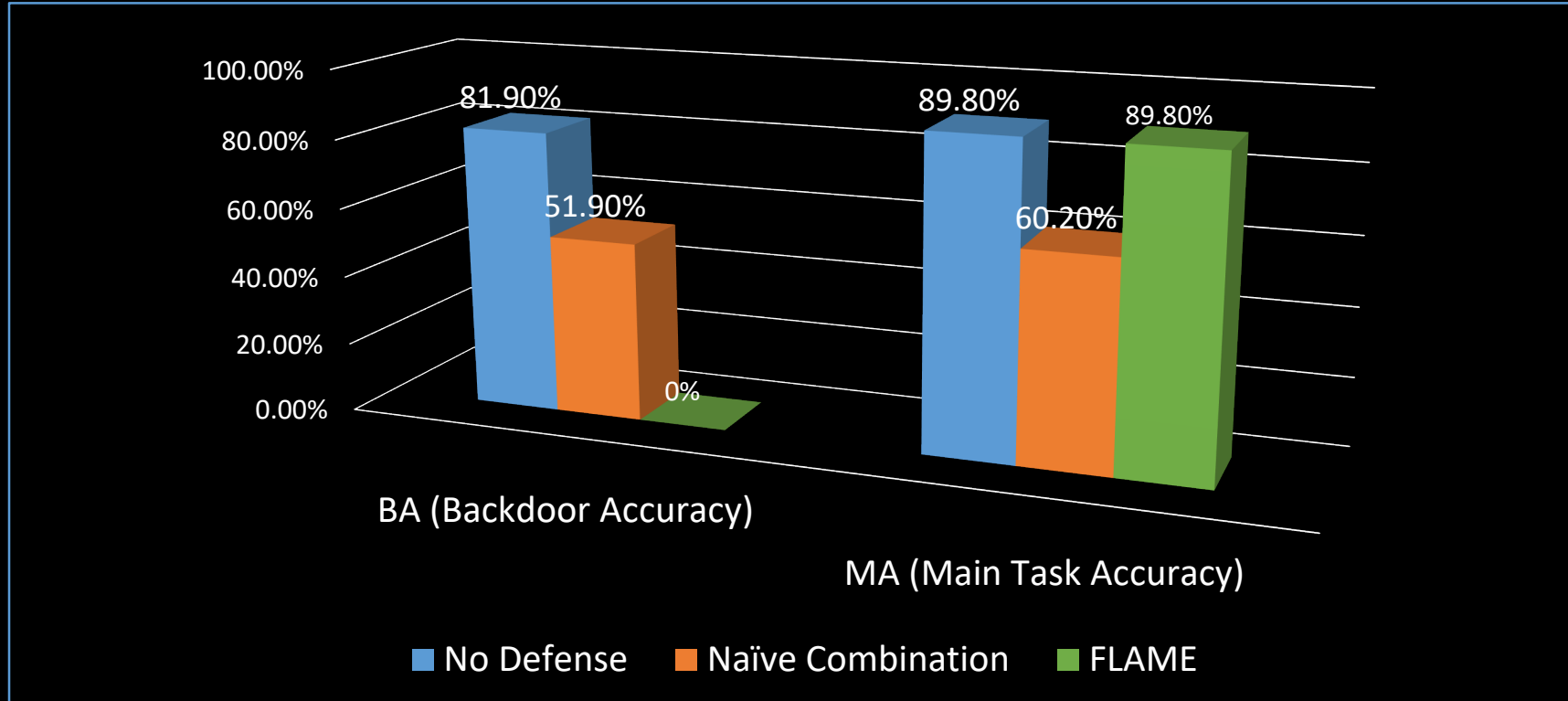
# FLAME vs. Existing Defenses

Defenses	Reddit		CIFAR-10		IoT-Traffic	
	BA	MA	BA	MA	BA	MA
Benign Setting	-	22.7	-	92.2	-	100.0
No defense	100.0	22.6	81.9	89.8	100.0	100.0
<b>Krum</b> [Blanchard et al. NIPS 2017]	100.0	9.6	100.0	56.7	100.0	84.0
<b>FoolsGold</b> [Fung et al. RAID 2020]	0.0	22.5	100.0	52.3	100.0	99.2
<b>Auror</b> [Shen et al. ACSAC 2016]	100.0	22.5	100.0	26.1	100.0	96.6
<b>AFA</b> [Muñoz-González et al. arXiv 2019]	100.0	22.4	0.0	91.7	100.0	87.4
<b>DP</b> [Sun et al. NeurIPS 2019]	14.0	18.9	0.0	78.9	14.8	82.3
<b>Median</b> [Yin et al ICML 2018]	0.0	22.0	0.0	50.1	0.0	87.7
<b>FLAME</b>	<b>0.0</b>	22.3	<b>0.0</b>	<b>91.9</b>	<b>0.0</b>	<b>99.8</b>

BA: Backdoor Accuracy, MA: Main Task Accuracy



# FLAME vs. Naïve Combination



Comparison between FLAME and a combination of existing defenses against constrain-and-scale attack [Bagdasaryan et al. AISTATS 2020] on the CIFAR-10 dataset

# Private FLAME: Privacy Preserving Aggregation

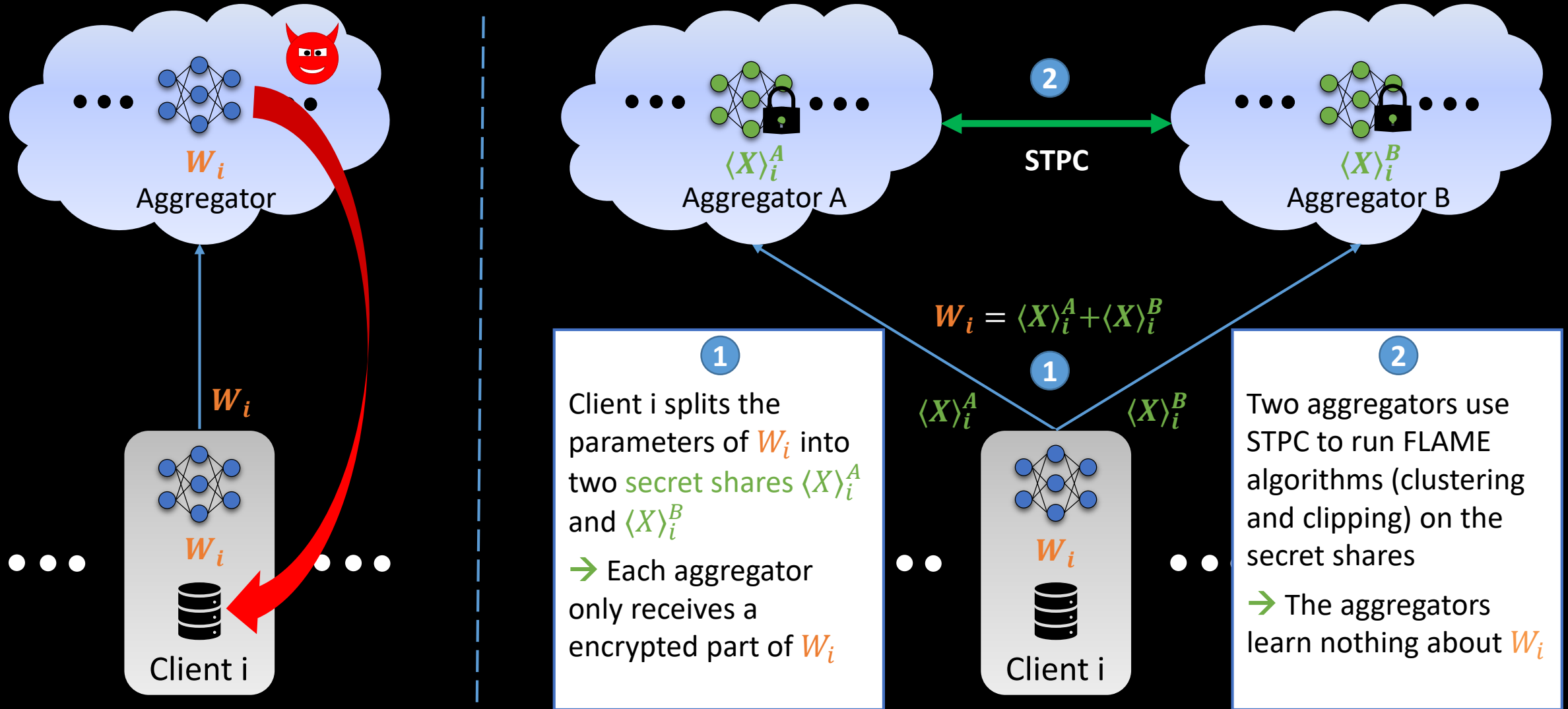
Using secure Multi-party Computation

# Private FLAME: Motivation

- Privacy attack: A **curious** aggregator can learn information from the training data by inference attacks
  - E.g., [Pyrgelis et al. NDSS 2018, Shokri et al. S&P 2017]
- Existing defenses **prohibit access** to the model updates to investigate backdoors
  - E.g., [Bonawitz et al. CCS 2017, Kairouz et al. PMLR 2021]
- Our goal: Introduce **private** FLAME such that FLAME algorithms are computed **under encryption**

# Private FALME: Solution

## Utilizing Secure Two-Party Computation (STPC)



# Private FLAME: Evaluation

- The runtime of Private FALME is significantly higher than standard FLAME
- However, such runtime overhead would be **acceptable** to maintain privacy
- Private FLAME provides **similar results w.r.t accuracy** in comparison to standard FLAME

#client	Reddit		CIFAR-10		IoT-Traffic	
	S	P	S	P	S	P
10	2.35	519.92	0.32	134.93	0.07	108.16
50	62.55	5895.70	2.62	766.12	0.69	269.35
100	252.13	<b>22081.65</b>	8.56	<b>2568.23</b>	2.11	<b>876.96</b>

Runtime in sec. of standard FALME (S) compared to private FALME (P) using secure two-party computation

	Reddit		CIFAR-10		IoT-Traffic	
	S	P	S	P	S	P
BA	0.0	0.0	0.0	0.0	0.0	0.0
MA	22.3	22.2	91.9	91.7	99.8	99.7
TPR	22.2	20.4	23.8	40.8	59.5	51.0
TNR	100.0	100.0	86.2	100.0	100.0	100.0

Effectiveness of standard FALME (S) in comparison to private FALME (P) using secure two-party computation

BA: Backdoor Accuracy, MA: Main Task Accuracy  
TPR: True Positive Rate, TNR: True Negative Rate

# Conclusion & Future Work

- FLAME, a novel backdoor defense for FL:
  - Mitigates state-of-the-art backdoor attacks effectively
  - Negligible impact on the benign performance of the models
  - Preserves privacy of clients' data
- Working on privacy-preserving poisoning defenses
  - Improving computational efficiency