

# Improving Intrusion Detection in Distributed Systems with Federated Learning

Defense replay at the SuperviZ Workshop

---

**Léo Lavaur**

Rennes, December 17<sup>th</sup>, 2024

Interdisciplinary Centre for Cybersecurity and Trust (SnT)  
University of Luxembourg

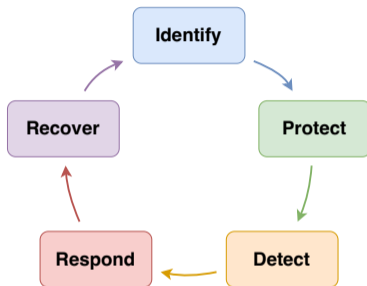
**Reviewers:** Anne-Marie Kermarrec · EPFL  
Eric Totel · Télécom SudParis

**Examiners:** Sonia Ben Mokhtar · CNRS  
Pierre-François Gimenez · Inria  
Vincent Nicomette · INSA Toulouse

**Supervisors:** Fabien Autrel · IMT Atlantique  
Marc-Oliver Pahl · IMT Atlantique

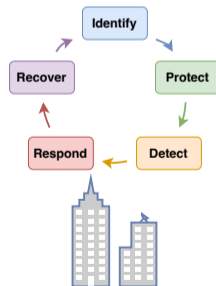
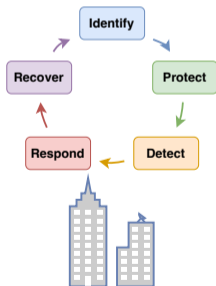
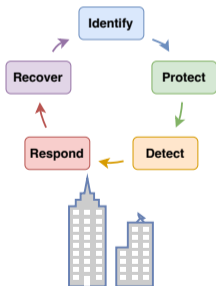
**Director:** Yann Busnel · IMT

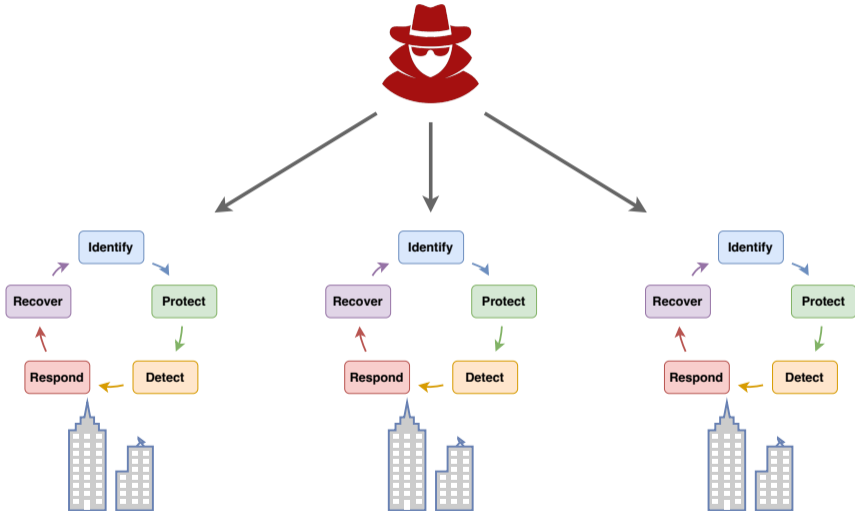


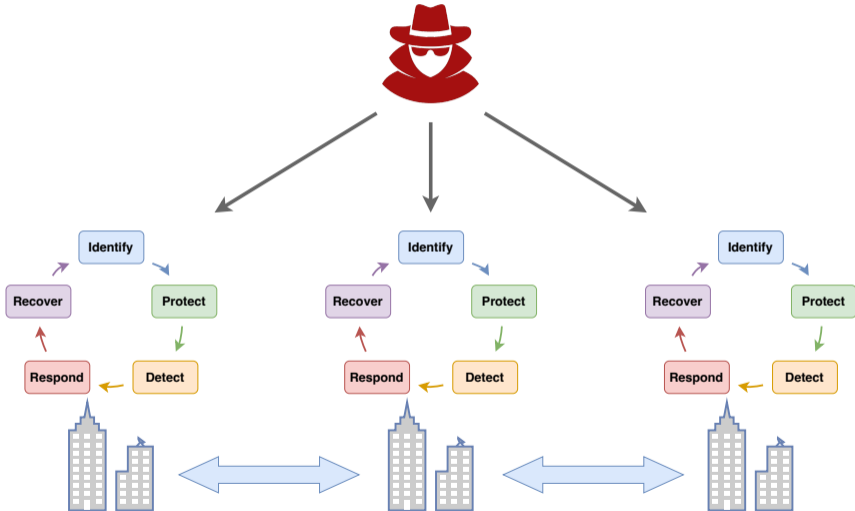


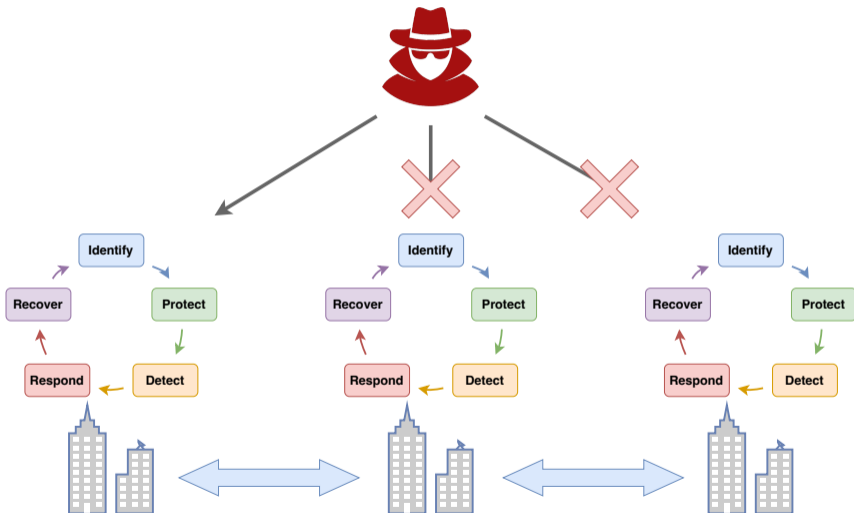
The security life-cycle [1].

[1] National Institute of Standards and Technology. *The NIST Cybersecurity Framework (CSF) 2.0*. 2024









- ▶ Collaboration pushed by:
  - common interest (*e.g.*, inter-SOCs<sup>1</sup>);

<sup>1</sup>Security Operational Center.

- ▶ Collaboration pushed by:
  - common interest (*e.g.*, inter-SOCs<sup>1</sup>);
  - national agencies (*e.g.*, NIST, ENISA, ANSSI);

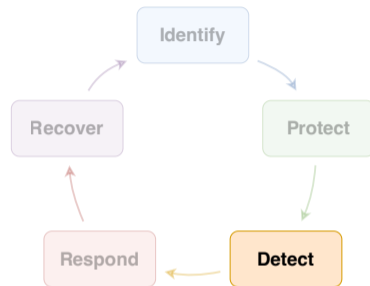
<sup>1</sup>Security Operational Center.



- ▶ Collaboration pushed by:
  - common interest (*e.g.*, inter-SOCs<sup>1</sup>);
  - national agencies (*e.g.*, NIST, ENISA, ANSSI);
  - regulation (*e.g.*, private-public information sharing in NIS2).

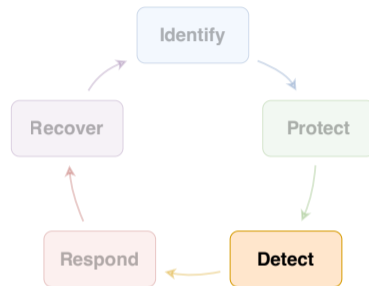
<sup>1</sup>Security Operational Center.

- ▶ Collaboration pushed by:
  - common interest (e.g., inter-SOCs<sup>1</sup>);
  - national agencies (e.g., NIST, ENISA, ANSSI);
  - regulation (e.g., private-public information sharing in NIS2).



<sup>1</sup>Security Operational Center.

- ▶ Collaboration pushed by:
  - common interest (e.g., inter-SOCs<sup>1</sup>);
  - national agencies (e.g., NIST, ENISA, ANSSI);
  - regulation (e.g., private-public information sharing in NIS2).



## Intrusion Detection System (IDS)

IDSs monitor the behavior of a system to detect malicious activities.

<sup>1</sup>Security Operational Center.

- ▶ Various types of algorithms: supervised, unsupervised, semi-supervised, reinforcement learning, *etc.*
- ▶ Great performance with Deep Learning (DL)...

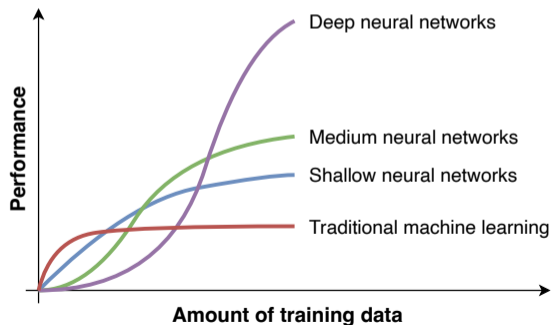
- ▶ Various types of algorithms: **supervised**, unsupervised, semi-supervised, reinforcement learning, *etc.*
- ▶ Great performance with Deep Learning (DL)...

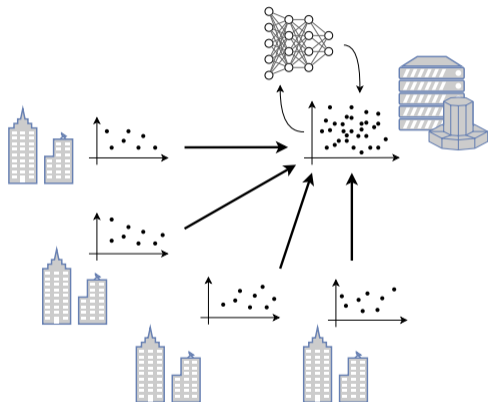
- ▶ Various types of algorithms: **supervised**, unsupervised, semi-supervised, reinforcement learning, *etc.*
- ▶ Great performance with Deep Learning (DL)... *on public datasets at least.*

- ▶ Various types of algorithms: **supervised**, unsupervised, semi-supervised, reinforcement learning, *etc.*
- ▶ Great performance with Deep Learning (DL)... *on public datasets at least.*

## Challenges of local training:

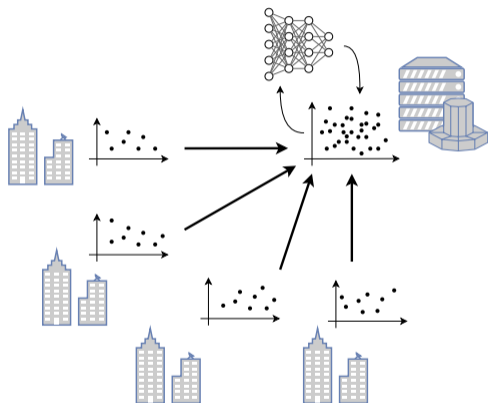
- ▶ not enough labelled data;
- ▶ risk of local bias or skewed data distribution.





Let's pool our data!





Let's pool our data! Although...

- ▶ Privacy concerns.
- ▶ Lack of trust in the data holder.
- ▶ Lack of trust in the learning process.
- ▶ ...

## Federated Learning (FL)

- ▶ Novel-*ish* distributed ML paradigm (Google) [2].

[2] McMahan et al. "Communication-Efficient Learning of Deep Networks from Decentralized Data". Proceedings of Machine Learning Research. 2017

## Federated Learning (FL)

- ▶ Novel-*ish* distributed ML paradigm (Google) [2].
- ▶ Distributed clients can train a common model without sharing their training data.

[2] McMahan et al. "Communication-Efficient Learning of Deep Networks from Decentralized Data". Proceedings of Machine Learning Research. 2017

## Federated Learning (FL)

- ▶ Novel-*ish* distributed ML paradigm (Google) [2].
- ▶ Distributed clients can train a common model without sharing their training data.
- ▶ **Privacy-preserving**: high level of abstraction for the shared models preventing data leakage.

[2] McMahan et al. "Communication-Efficient Learning of Deep Networks from Decentralized Data". Proceedings of Machine Learning Research. 2017

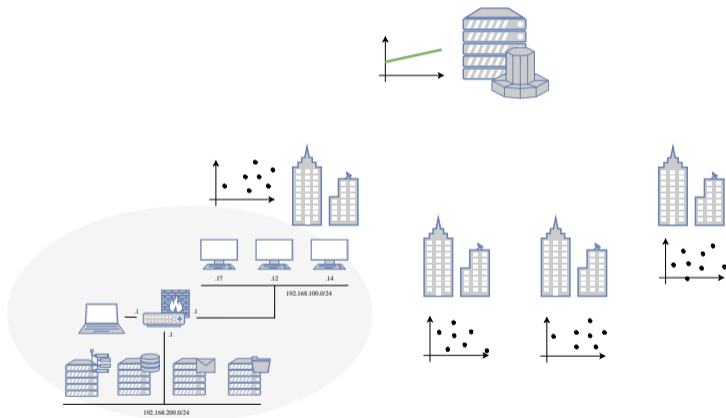


Figure: Typical FL workflow, applied to NIDSs.

1 Distribute the initial model

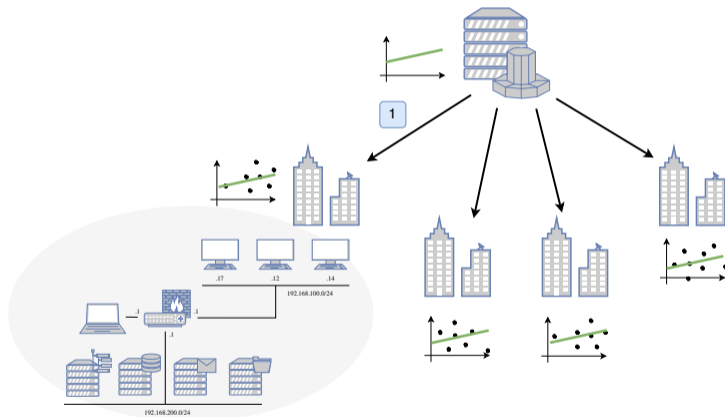


Figure: Typical FL workflow, applied to NIDSs.

- 1 Distribute the initial model
- 2 Train the model on local data

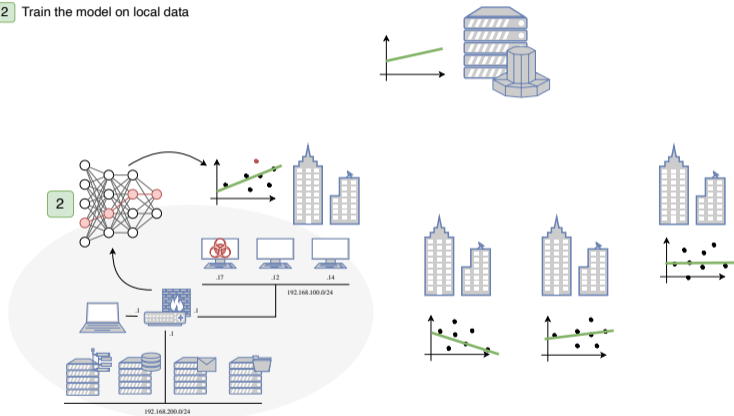


Figure: Typical FL workflow, applied to NIDSs.

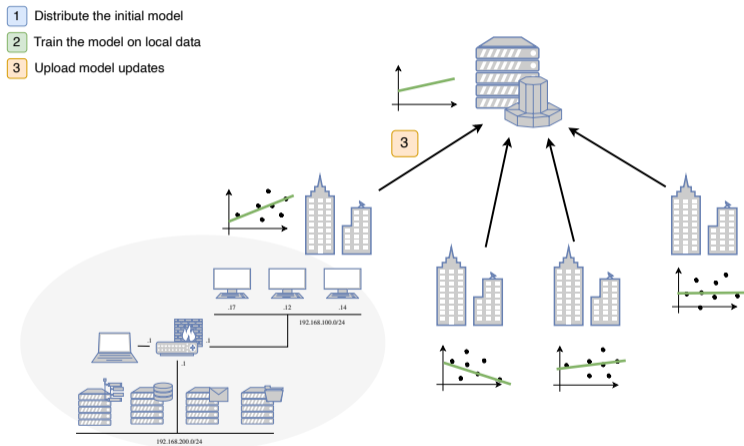
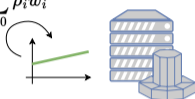


Figure: Typical FL workflow, applied to NIDSs.



- 1 Distribute the initial model
- 2 Train the model on local data
- 3 Upload model updates
- 4 Aggregate updates

$$4 \quad \frac{1}{n} \sum_{i=0}^n \rho_i w_i$$


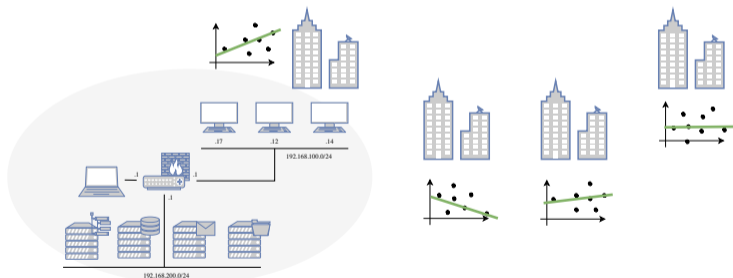


Figure: Typical FL workflow, applied to NIDSs.

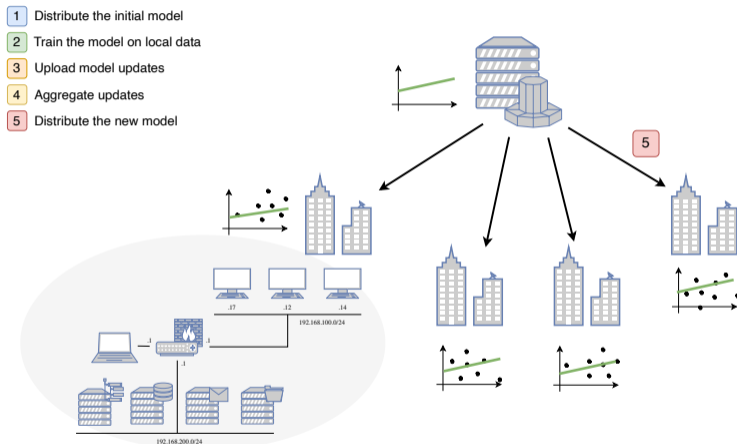


Figure: Typical FL workflow, applied to NIDSs.

- 1 Distribute the initial model
- 2 Train the model on local data
- 3 Upload model updates
- 4 Aggregate updates
- 5 Distribute the new model

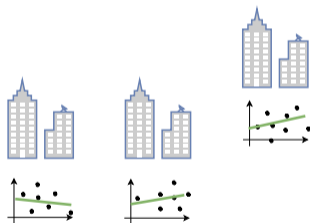
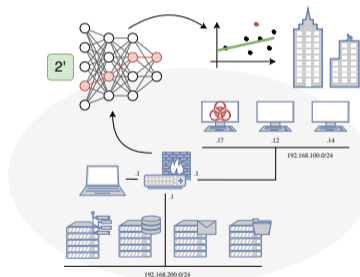
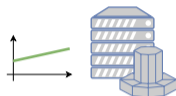


Figure: Typical FL workflow, applied to NIDSs.

## Collaborative Intrusion Detection between Distributed Organizations

- ▶ Each organization has its own NIDS<sup>2</sup> and monitors an information system.

<sup>2</sup>Network-based IDS

## Collaborative Intrusion Detection between Distributed Organizations

- ▶ Each organization has its own NIDS<sup>2</sup> and monitors an information system.
- ▶ Objective: improve their local detection performance.

## Collaborative Intrusion Detection between Distributed Organizations

- ▶ Each organization has its own NIDS<sup>2</sup> and monitors an information system.
- ▶ Objective: improve their local detection performance.
- ▶ Means: **expert knowledge** (*i.e.*, datasets) and **computing resources** (*i.e.*, model training).

<sup>2</sup>Network-based IDS

## Collaborative Intrusion Detection between Distributed Organizations

- ▶ Each organization has its own NIDS<sup>2</sup> and monitors an information system.
- ▶ Objective: improve their local detection performance.
- ▶ Means: **expert knowledge** (*i.e.*, datasets) and **computing resources** (*i.e.*, model training).



Figure: Typical workflow for ML-based NIDSs.

A cross-silo use case [3]:

- ▶ few clients (*i.e.*, 10–100);
- ▶ substantial amount of data, high heterogeneity;
- ▶ high availability, significant computing resources.

[3] Kairouz et al. “Advances and Open Problems in Federated Learning”. 2021



## Challenges from the Literature [4]

- ▶ *Functionality*: performance, heterogeneity, transferability, self-defense, and self-healing.

[4] Lavour et al. "The Evolution of Federated Learning-based Intrusion Detection and Mitigation: a Survey". *IEEE Transactions on Network and Service Management*. 2022

## Challenges from the Literature [4]

- ▶ *Functionality*: performance, heterogeneity, transferability, self-defense, and self-healing.
- ▶ *Deployment*: adaptability and scalability.

[4] Lavaur et al. "The Evolution of Federated Learning-based Intrusion Detection and Mitigation: a Survey". *IEEE Transactions on Network and Service Management*. 2022

## Challenges from the Literature [4]

- ▶ *Functionality*: performance, heterogeneity, transferability, self-defense, and self-healing.
- ▶ *Deployment*: adaptability and scalability.
- ▶ *Security and reliability*: security, privacy, trust, and reputation.

[4] Lavour et al. "The Evolution of Federated Learning-based Intrusion Detection and Mitigation: a Survey". *IEEE Transactions on Network and Service Management*. 2022

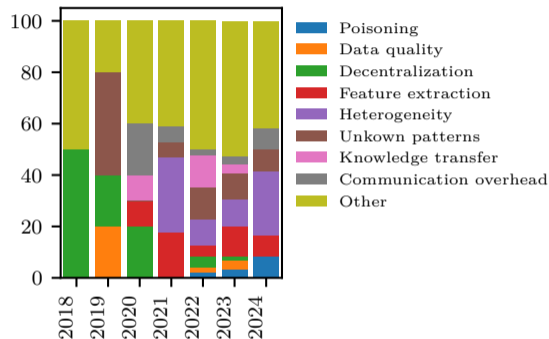
## Challenges from the Literature [4]

- ▶ *Functionality*: performance, heterogeneity, transferability, self-defense, and self-healing.
- ▶ *Deployment*: adaptability and scalability.
- ▶ *Security and reliability*: security, privacy, trust, and reputation.
- ▶ *Experimentation*: evaluation.

[4] Lavour et al. "The Evolution of Federated Learning-based Intrusion Detection and Mitigation: a Survey". *IEEE Transactions on Network and Service Management*. 2022

### Challenges from the Literature [4]

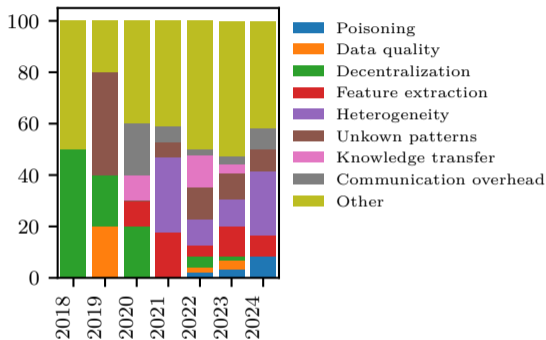
- ▶ *Functionality*: performance, heterogeneity, transferability, self-defense, and self-healing.
- ▶ *Deployment*: adaptability and scalability.
- ▶ *Security and reliability*: security, privacy, trust, and reputation.
- ▶ *Experimentation*: evaluation.



**Figure:** Challenges addressed by the literature (until 2024-04).

## Challenges from the Literature [4]

- ▶ *Functionality*: performance, heterogeneity, transferability, self-defense, and self-healing.
- ▶ *Deployment*: adaptability and scalability.
- ▶ *Security and reliability*: security, privacy, trust, and reputation.
- ▶ *Experimentation*: evaluation.



**Figure:** Challenges addressed by the literature (until 2024-04).

[4] Lavaur et al. "The Evolution of Federated Learning-based Intrusion Detection and Mitigation: a Survey". *IEEE Transactions on Network and Service Management*. 2022

## Challenges from the Literature [4]

- ▶ *Functionality*: performance, heterogeneity, transferability, self-defense, and self-healing.
- ▶ *Deployment*: adaptability and scalability.
- ▶ *Security and reliability*: security, privacy, trust, and reputation.
- ▶ *Experimentation*: evaluation.

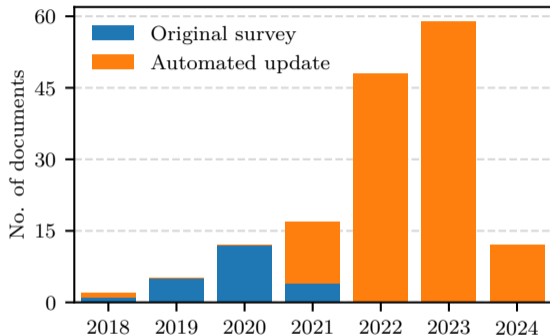


Figure: Publications on FL & IDS (until 2024-04).

[4] Lavour et al. "The Evolution of Federated Learning-based Intrusion Detection and Mitigation: a Survey". *IEEE Transactions on Network and Service Management*. 2022

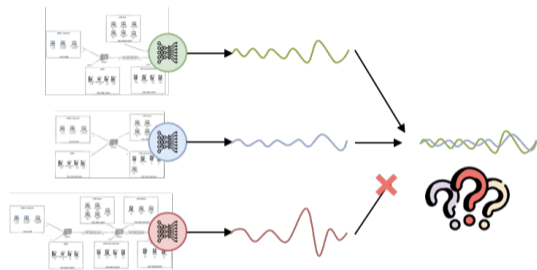


Figure: Heterogeneity headaches.

**Challenge I:** *Too much heterogeneity leads to poor performance...* [5]

**Challenge II:** *Difficult to identify malicious contributions when models are different...*

**Challenge III:** *No representative dataset of heterogeneous distributed intrusion detection...*

[5] Lavaur, Busnel, and Autrel. "Demo: Highlighting the Limits of Federated Learning in Intrusion Detection". *Proceedings of the 44th International Conference on Distributed Computing Systems (ICDCS)*. 2024



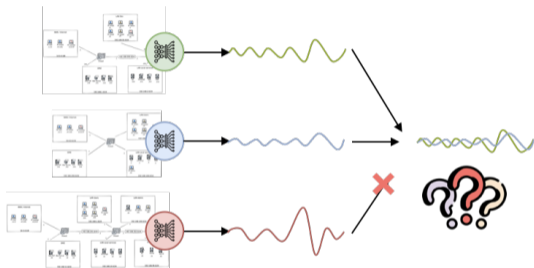


Figure: Heterogeneity headaches.

*Challenge I: Too much heterogeneity leads to poor performance...*

*Challenge II: Difficult to identify malicious contributions when models are different...*

*Challenge III: No representative dataset of heterogeneous distributed intrusion detection...*

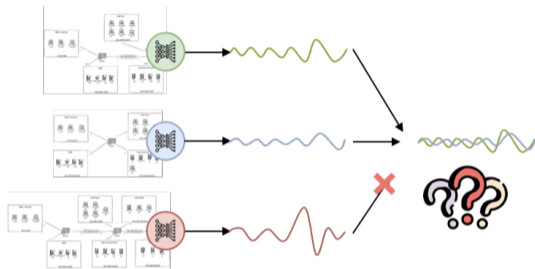


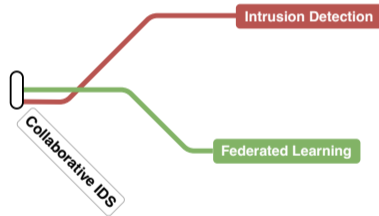
Figure: Heterogeneity headaches.

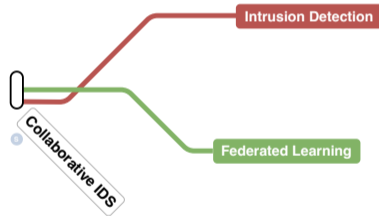
*Challenge I: Too much heterogeneity leads to poor performance...*

*Challenge II: Difficult to identify malicious contributions when models are different...*

**Challenge III: No representative dataset of heterogeneous distributed intrusion detection... [6]**

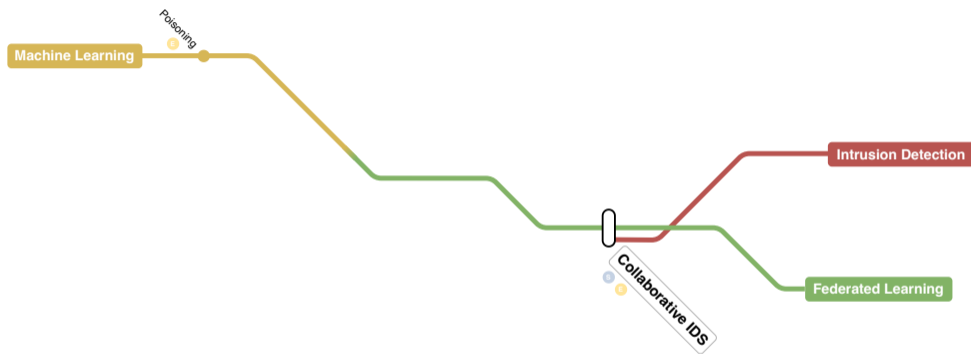
[6] Lavour et al. "Federated Learning as Enabler for Collaborative Security between Not Fully-Trusting Distributed Parties". *Proceedings of the 29th Computer & Electronics Security Application Rendezvous (C&ESAR)*. 2022





## CONTRIBUTIONS

**S** Systematic Literature Review





## CONTRIBUTIONS

- S Systematic Literature Review
- E Assessment & eiffel



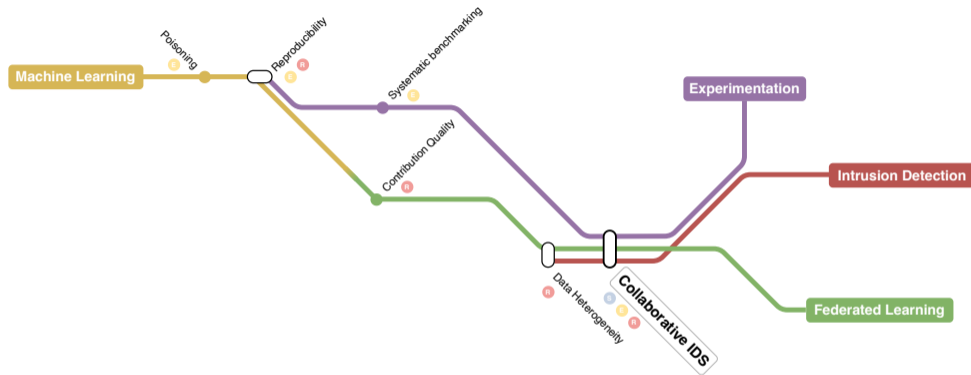
## CONTRIBUTIONS

-  Systematic Literature Review
-  Assessment & eiffel



## CONTRIBUTIONS

- S** Systematic Literature Review
- E** Assessment & eiffel
- R** RADAR



## CONTRIBUTIONS

- S Systematic Literature Review
- E Assessment & eiffel
- R RADAR





## CONTRIBUTIONS

- S Systematic Literature Review
- E Assessment & eiffel
- R RADAR



## CONTRIBUTIONS

- S Systematic Literature Review
- E Assessment & eiffel
- R RADAR
- F FedITN



## CONTRIBUTIONS

- S Systematic Literature Review
- E Assessment & eiffel
- R RADAR
- F FedITN



## CONTRIBUTIONS

- Systematic Literature Review
- Assessment & eiffel
- RADAR
- FedITN



|          |  |    |
|----------|--|----|
| <b>E</b> | Assessing the Impact of Label-Flipping Attacks .....             | 13 |
| <b>R</b> | Fighting Byzantine Contributions in Heterogeneous Settings ..... | 19 |

**CONTRIBUTIONS**

- Systematic Literature Review
- Assessment & eiffel
- RADAR
- FedITN

## E Assessing the Impact of Label-Flipping Attacks

---

[7] Lavaur, Busnel, and Autrel. "Systematic Analysis of Label-flipping Attacks against Federated Learning in Collaborative Intrusion Detection Systems". *Proceedings of the 19th International Conference on Availability, Reliability and Security (ARES)*. 2024

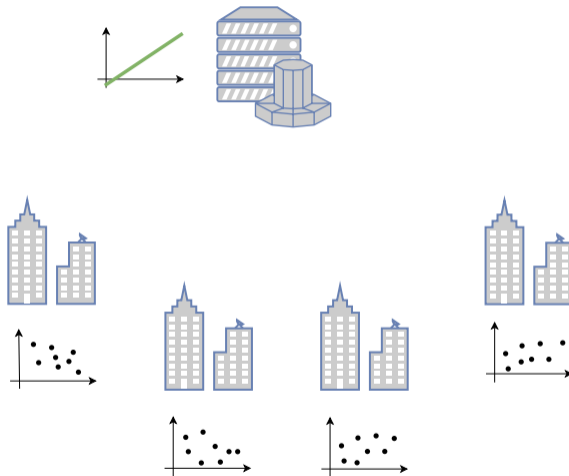


Figure: Poisoning attacks on FL.

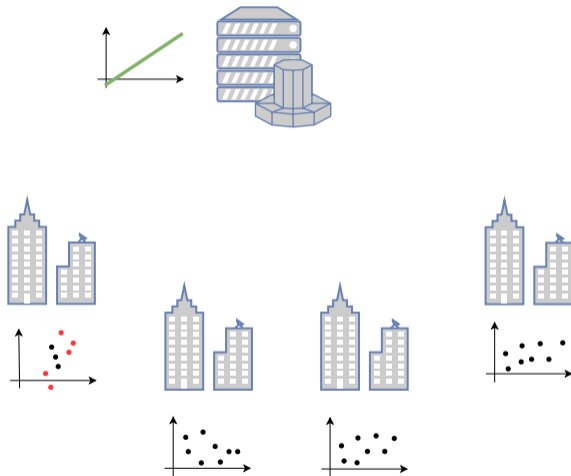


Figure: Poisoning attacks on FL.



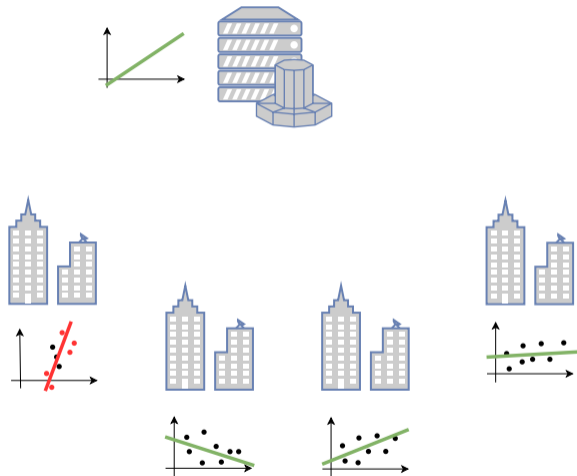


Figure: Poisoning attacks on FL.

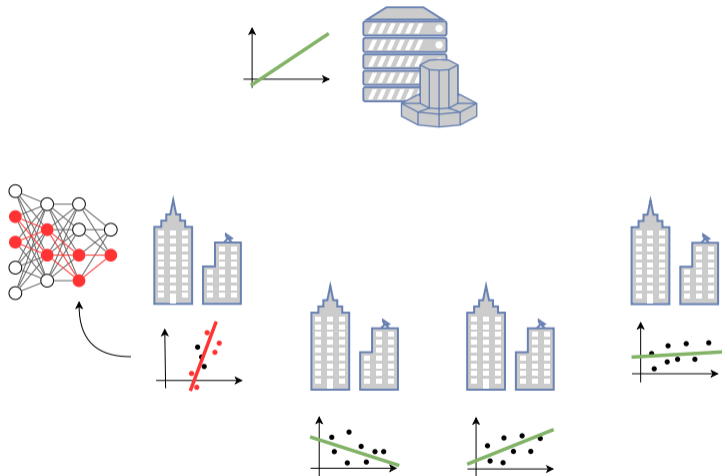


Figure: Poisoning attacks on FL.

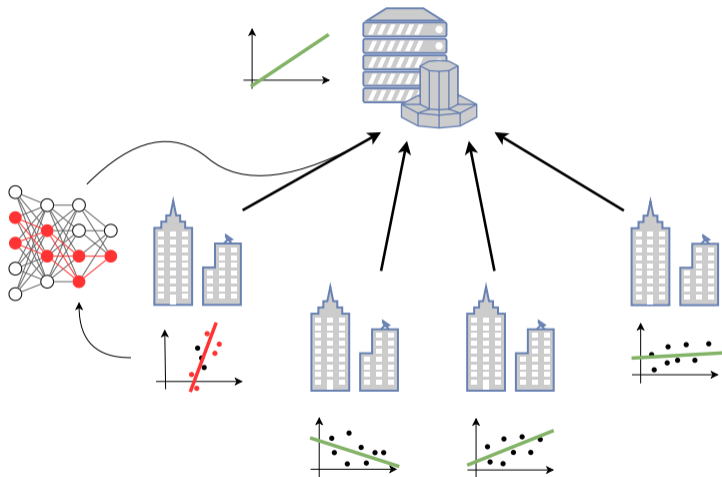


Figure: Poisoning attacks on FL.

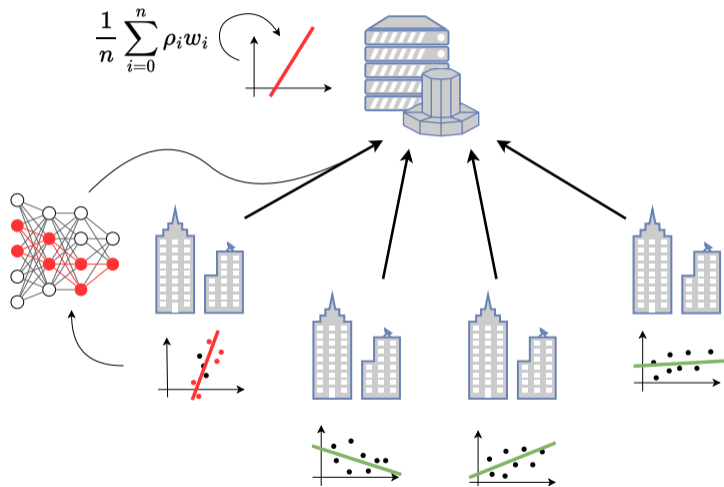


Figure: Poisoning attacks on FL.

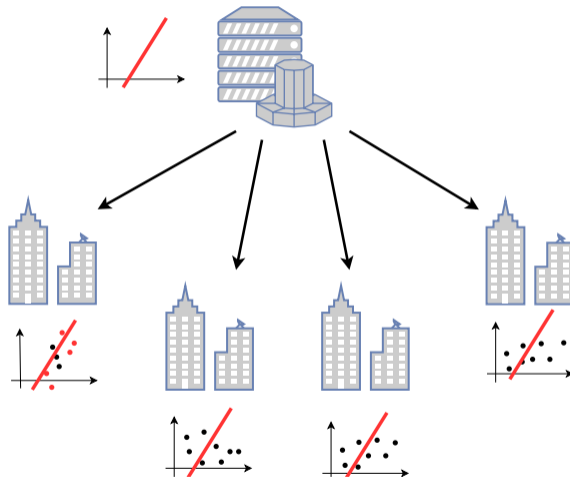


Figure: Poisoning attacks on FL.

Poisoning attacks

## Poisoning attacks

### COMPONENT

- ▶ Data poisoning (e.g., label-flipping, clean-label)
- ▶ Model poisoning (e.g., gradient boosting)

## Poisoning attacks

### COMPONENT

- ▶ Data poisoning (e.g., label-flipping, clean-label)
- ▶ Model poisoning (e.g., gradient boosting)

### OBJECTIVE

- ▶ Untargeted: impact model performance
- ▶ Targeted: modify behavior for specific samples



## Poisoning attacks

### COMPONENT

- ▶ Data poisoning (e.g., label-flipping, clean-label)
- ▶ Model poisoning (e.g., gradient boosting)

### OBJECTIVE

- ▶ Untargeted: impact model performance
- ▶ Targeted: modify behavior for specific samples

### PROPORTION

- ▶ Single attacker
- ▶ Colluding attackers: multiple coordinated adversaries

## Poisoning attacks

### COMPONENT

- ▶ Data poisoning (e.g., **label-flipping**, clean-label)
- ▶ Model poisoning (e.g., gradient boosting)

### OBJECTIVE

- ▶ Untargeted: impact model performance
- ▶ Targeted: modify behavior for specific samples

### PROPORTION

- ▶ Single attacker
- ▶ Colluding attackers: multiple coordinated adversaries

## Existing studies

- ▶ Often partial, focusing on challenging a specific defense mechanism.
- ▶ Lack of reproducibility and comparability (different datasets, models, and attacks).
- ▶ No targeted attacks binary classification.

## Existing studies

- ▶ Often partial, focusing on challenging a specific defense mechanism.
- ▶ Lack of reproducibility and comparability (different datasets, models, and attacks).
- ▶ No targeted attacks binary classification.

## Research Questions

1. Is the behavior of poisoning attacks predictable?
2. Do hyperparameters influence the impact of poisoning attacks?
3. Are IDS backdoors realistic using label-flipping attacks?
4. Is there a critical threshold where label-flipping attacks begin to impact performance?
5. Is gradient similarity enough to detect label-flipping attacks?

## Existing studies

- ▶ Often partial, focusing on challenging a specific defense mechanism.
- ▶ Lack of reproducibility and comparability (different datasets, models, and attacks).
- ▶ No targeted attacks binary classification.

## Research Questions

1. Is the behavior of poisoning attacks predictable?
2. Do hyperparameters influence the impact of poisoning attacks?
3. Are IDS backdoors realistic using label-flipping attacks?
4. Is there a critical threshold where label-flipping attacks begin to impact performance?
5. **Is gradient similarity enough to detect label-flipping attacks?**

## RQ5: IS GRADIENT SIMILARITY ENOUGH TO DETECT LABEL-FLIPPING ATTACKS?

- ▶ Known technique to detect poisoning attacks [8].

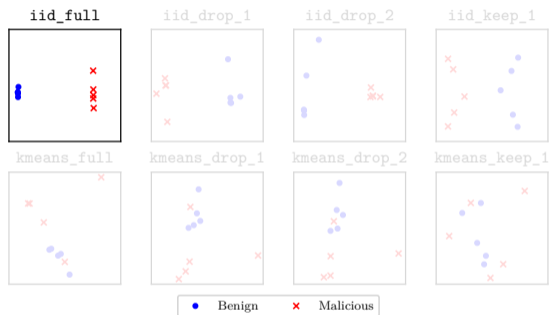


Figure: PCA projection of the gradients in 2D (CICIDS).

[8] Tolpegin et al. "Data Poisoning Attacks Against Federated Learning Systems". Lecture Notes in Computer Science. 2020

## RQ5: IS GRADIENT SIMILARITY ENOUGH TO DETECT LABEL-FLIPPING ATTACKS?

- ▶ Known technique to detect poisoning attacks [8].
- ▶ High heterogeneity makes it harder to detect attackers.

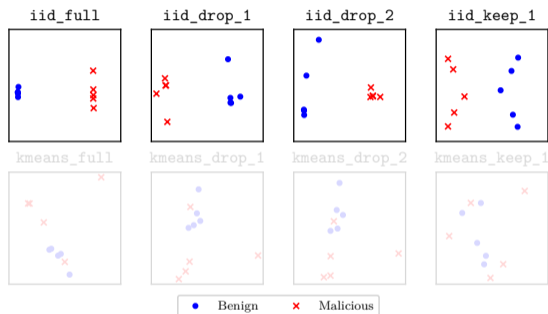


Figure: PCA projection of the gradients in 2D (CICIDS).

## RQ5: IS GRADIENT SIMILARITY ENOUGH TO DETECT LABEL-FLIPPING ATTACKS?

- ▶ Known technique to detect poisoning attacks [8].
- ▶ High heterogeneity makes it harder to detect attackers.

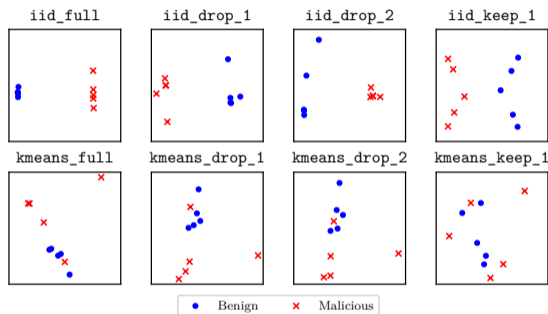


Figure: PCA projection of the gradients in 2D (CICIDS).

[8] Tolpegin et al. "Data Poisoning Attacks Against Federated Learning Systems". Lecture Notes in Computer Science. 2020



1. A *deeper* understanding of the behavior of label-flipping attacks in FL-based CIDSs.
  - Similarity-based detection techniques show limitations in detecting poisoning attacks.
  - Limited by the models' generalization capabilities and the characteristic overlap between classes.
  - Hyperparameter dependencies, but not on the average performance impact.

1. A *deeper* understanding of the behavior of label-flipping attacks in FL-based CIDSs.
  - Similarity-based detection techniques show limitations in detecting poisoning attacks.
  - Limited by the models' generalization capabilities and the characteristic overlap between classes.
  - Hyperparameter dependencies, but not on the average performance impact.
2. A **reproducible** evaluation framework to study the impact of label-flipping attacks in FIDS using FL.
  - Reproducible, extendable, and available in open-access<sup>3</sup>.
  - Calls to be extended to other poisoning attacks, datasets, and partitioning strategies.

<sup>3</sup><https://github.com/leolavaur/eiffel>

## R Fighting Byzantine Contributions in Heterogeneous Settings

---

## Case study reminder

- ▶ Multiple organizations collaborating on a federated Intrusion Detection System.
- ▶ Partial heterogeneity in the datasets: different data distributions but existing similarities.

## Case study reminder

- ▶ Multiple organizations collaborating on a federated Intrusion Detection System.
- ▶ Partial heterogeneity in the datasets: different data distributions but existing similarities.

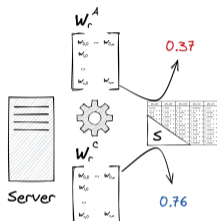
## Byzantine contributions:

- ▶ data quality issues (*e.g.*, labelling, noise);
- ▶ distribution mismatches; and
- ▶ adversaries, *possibly colluding*.

### Quality Assessment in Heterogeneous Settings

For  $n$  participants  $p_i$  and their local datasets  $d_i$  of unknown similarity, each participant uploads a model update  $w_i^r$  at each round  $r$ . Given  $P = \{p_1, p_2, \dots, p_n\}$  and  $W = \{w_1^r, w_2^r, \dots, w_n^r\}$ , how can one assess the quality of each participant's contribution without making assumptions on the data distribution across the datasets  $d_i$ ?

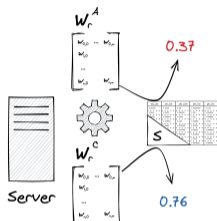
## Server-side evaluation [10]



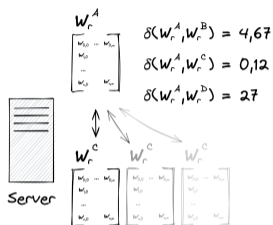
- ▶ Only applicable in IID settings.
- ▶ Single source of truth.

[10] Zhou et al. "A Differentially Private Federated Learning Model against Poisoning Attacks in Edge Computing". *IEEE Transactions on Dependable and Secure Computing*. 2022

## Server-side evaluation [10]



## Server-side comparison [11]



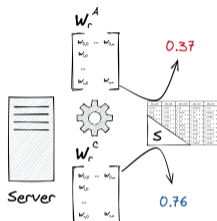
- ▶ Only applicable in IID settings.
- ▶ Single source of truth.
- ▶ Less related to client data.

[10] Zhou et al. "A Differentially Private Federated Learning Model against Poisoning Attacks in Edge Computing". *IEEE Transactions on Dependable and Secure Computing*. 2022

[11] Briggs, Fan, and Andras. "Federated Learning with Hierarchical Clustering of Local Updates to Improve Training on Non-IID Data". *2020 International Joint Conference on Neural Networks (IJCNN)*. 2020

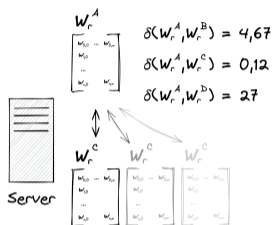


## Server-side evaluation [10]



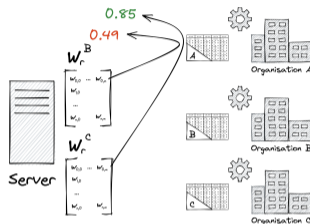
- ▶ Only applicable in IID settings.
- ▶ Single source of truth.

## Server-side comparison [11]



- ▶ Less related to client data.

## Client-side evaluation [12]

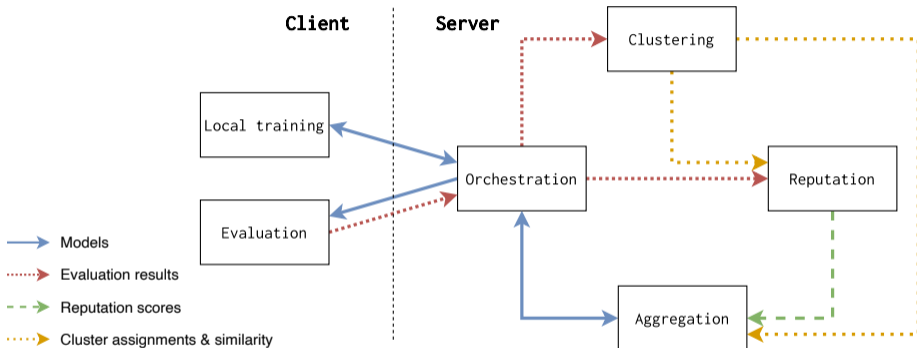


- ▶ High cost in cross-device.
- ▶ More susceptible to badmouthing.

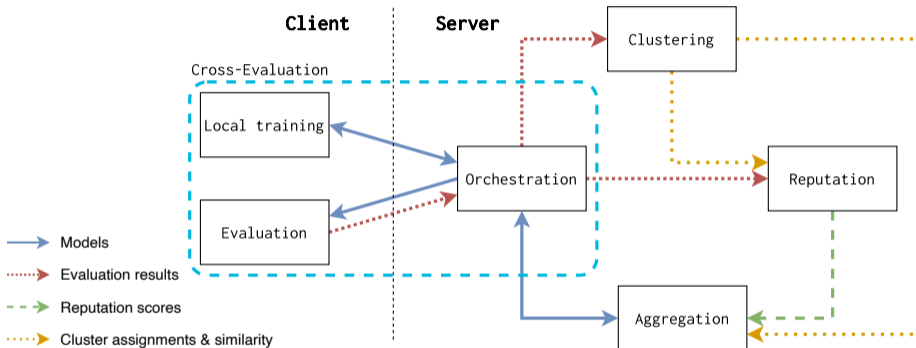
[10] Zhou et al. "A Differentially Private Federated Learning Model against Poisoning Attacks in Edge Computing". *IEEE Transactions on Dependable and Secure Computing*. 2022

[11] Briggs, Fan, and Andras. "Federated Learning with Hierarchical Clustering of Local Updates to Improve Training on Non-IID Data". *2020 International Joint Conference on Neural Networks (IJCNN)*. 2020

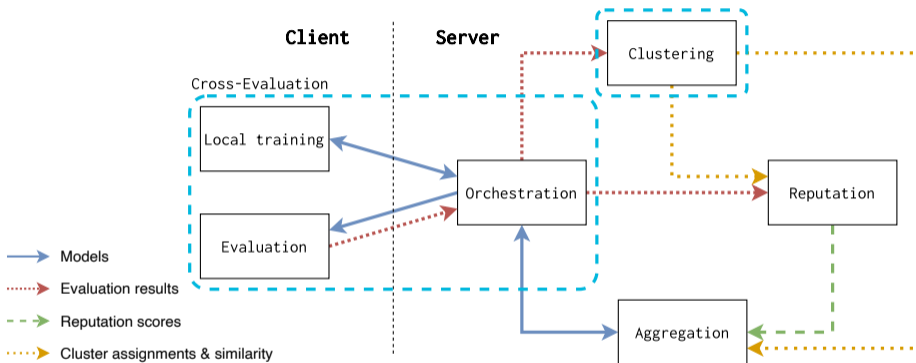
[12] Zhao et al. *Shielding Collaborative Learning: Mitigating Poisoning Attacks through Client-Side Detection*. 2020



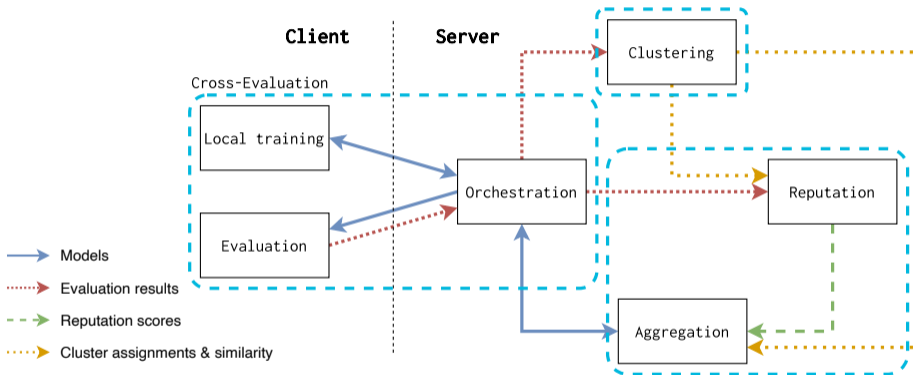
RADAR architecture.



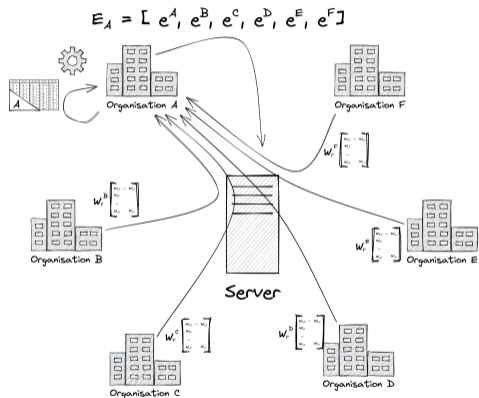
RADAR architecture.



RADAR architecture.

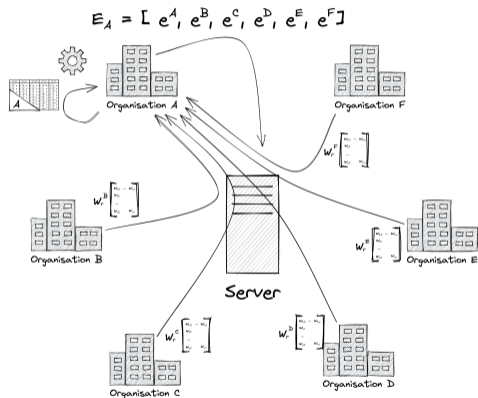


RADAR architecture.



## Advantages

- ▶ Exhaustive overview of the entire system at each round  $r$ . **No need of prior knowledge!**
- ▶ Evaluations (e.g., accuracy, F1 score) representative of participants' data.

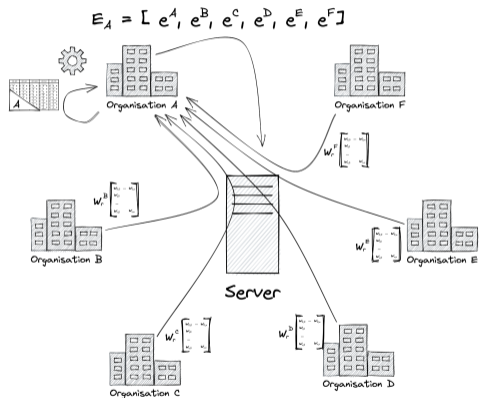


## Advantages

- ▶ Exhaustive overview of the entire system at each round  $r$ . **No need of prior knowledge!**
- ▶ Evaluations (e.g., accuracy, F1 score) representative of participants' data.

## Drawbacks

- ▶ High communication and computation costs.
- ▶ Does not scale well.



## Advantages

- ▶ Exhaustive overview of the entire system at each round  $r$ . **No need of prior knowledge!**
- ▶ Evaluations (e.g., accuracy, F1 score) representative of participants' data.

## Drawbacks

- ▶ High communication and computation costs.
- ▶ Does not scale well.

## But...

- ▶ Cross-silo use case: few clients, with reasonable computing capacity.
- ▶ Slow workflow: long time between rounds.



## Objective

- ▶ Build *more* homogeneous communities of participants to facilitate model aggregation.

## Objective

- ▶ Build *more* homogeneous communities of participants to facilitate model aggregation.
  
- ▶ Distance metric
  - Based on **cross-evaluation** results.
  - **Cosine similarity** [11].

## Objective

- ▶ Build *more* homogeneous communities of participants to facilitate model aggregation.

### ▶ Distance metric

- Based on **cross-evaluation** results.
- **Cosine similarity** [11].

### ▶ Algorithm

- **Hierarchical clustering.** [11]
- Dynamic aggregation threshold.

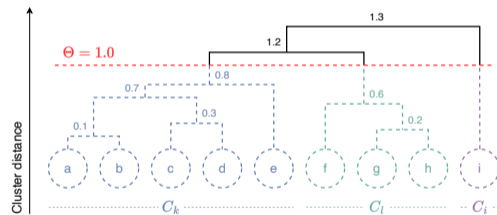


Figure: Hierarchical clustering.

[11] Briggs, Fan, and Andras. "Federated Learning with Hierarchical Clustering of Local Updates to Improve Training on Non-IID Data". 2020 International Joint Conference on Neural Networks (IJCNN). 2020

## Definition: Reputation Systems [13]

- ▶ Long-lived entities expecting future interaction.
- ▶ Capture and distribution of feedback about current interactions (such information must be visible in the future).
- ▶ Use of feedback to guide trust decisions.

[13] Resnick et al. "Reputation Systems". *Communications of the ACM*. 2000

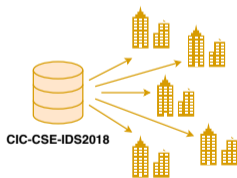
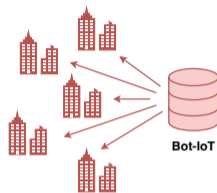
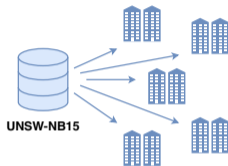
## Definition: Reputation Systems [13]

- ▶ Long-lived entities expecting future interaction.
  - ▶ Capture and distribution of feedback about current interactions (such information must be visible in the future).
  - ▶ Use of feedback to guide trust decisions.
- 
- ▶ Votes weighted by the similarity inside each cluster.
  - ▶ Exponential decay for potential redemption.

[13] Resnick et al. "Reputation Systems". *Communications of the ACM*. 2000

## Datasets

- ▶ Heterogeneous datasets, but some participants can share similarities.
- ▶ 4 datasets: CIC-CSE-IDS2018, UNSW-NB15, Bot-IoT, ToN\_IoT.
- ▶ NF-V2 [14] feature set (*i.e.*, NetFlow V9).



[14] Sarhan, Layeghy, and Portmann. *Towards a Standard Feature Set for Network Intrusion Detection System Datasets*. 2021

## Parameters

- ▶ *Target*: Affected classes.
- ▶ *Data Poisoning Rate (DPR)*: proportion of targeted data with flipped labels.
- ▶ *Model Poisoning Rate (MPR)*: number of attackers in the cluster.



**colluding minority 100T**  
*(i.e., 2 attackers, 100% DPR on Reconnaissance class).*

**Table:** *Effect of different attack configurations (untargeted) on all baselines.* RA is RADAR, FG is FoolsGold, FA is FedAvg (on all participants), and FC is FedAvg ideally clustered per dataset.

| Scenario                 | ASR (%)     |             |       |             |
|--------------------------|-------------|-------------|-------|-------------|
|                          | RA          | FG          | FA    | FC          |
| <b>Targeted (100T)</b>   |             |             |       |             |
| Benign                   | <b>0.00</b> | 5.17        | 5.10  | 0.09        |
| Lone                     | <b>0.00</b> | 93.82       | 6.73  | 0.45        |
| Collud. min.             | <b>0.00</b> | 2.97        | 9.99  | 53.40       |
| Collud. maj.             | 73.39       | <b>8.10</b> | 17.65 | 59.36       |
| <b>Untargeted (100U)</b> |             |             |       |             |
| Benign                   | 0.09        | 0.39        | 33.30 | <b>0.06</b> |
| Lone                     | <b>0.08</b> | 99.89       | 54.70 | 0.12        |
| Collud. min.             | 0.10        | <b>0.04</b> | 44.53 | 6.26        |
| Collud. maj.             | <b>0.08</b> | 38.98       | 59.49 | 94.36       |

*lower is better*



**Table:** *Effect of different attack configurations (untargeted) on all baselines.* RA is RADAR, FG is FoolsGold, FA is FedAvg (on all participants), and FC is FedAvg ideally clustered per dataset.

| Scenario                 | ASR (%)     |             |       |             |
|--------------------------|-------------|-------------|-------|-------------|
|                          | RA          | FG          | FA    | FC          |
| <b>Targeted (100T)</b>   |             |             |       |             |
| Benign                   | <b>0.00</b> | 5.17        | 5.10  | 0.09        |
| Lone                     | <b>0.00</b> | 93.82       | 6.73  | 0.45        |
| Collud. min.             | <b>0.00</b> | 2.97        | 9.99  | 53.40       |
| Collud. maj.             | 73.39       | <b>8.10</b> | 17.65 | 59.36       |
| <b>Untargeted (100U)</b> |             |             |       |             |
| Benign                   | 0.09        | 0.39        | 33.30 | <b>0.06</b> |
| Lone                     | <b>0.08</b> | 99.89       | 54.70 | 0.12        |
| Collud. min.             | 0.10        | <b>0.04</b> | 44.53 | 6.26        |
| Collud. maj.             | <b>0.08</b> | 38.98       | 59.49 | 94.36       |

*lower is better*

**Table:** *Effect of different attack configurations (untargeted) on all baselines.* RA is RADAR, FG is FoolsGold, FA is FedAvg (on all participants), and FC is FedAvg ideally clustered per dataset.

| Scenario                 | ASR (%)      |             |       |             |
|--------------------------|--------------|-------------|-------|-------------|
|                          | RA           | FG          | FA    | FC          |
| <b>Targeted (100T)</b>   |              |             |       |             |
| Benign                   | <b>0.00</b>  | 5.17        | 5.10  | 0.09        |
| Lone                     | <b>0.00</b>  | 93.82       | 6.73  | 0.45        |
| Collud. min.             | <b>0.00</b>  | 2.97        | 9.99  | 53.40       |
| Collud. maj.             | <b>73.39</b> | <b>8.10</b> | 17.65 | 59.36       |
| <b>Untargeted (100U)</b> |              |             |       |             |
| Benign                   | 0.09         | 0.39        | 33.30 | <b>0.06</b> |
| Lone                     | <b>0.08</b>  | 99.89       | 54.70 | 0.12        |
| Collud. min.             | 0.10         | <b>0.04</b> | 44.53 | 6.26        |
| Collud. maj.             | <b>0.08</b>  | 38.98       | 59.49 | 94.36       |

*lower is better*

## 1. RADAR can:

- leverage **cross-evaluation**, **clustering** and **reputation** to address heterogeneity and Byzantine contributions;
- adjust rapidly to changes in behavior; and
- mitigate most tested scenarios (limiting case handled up to 80% of poisoned data).

## 1. RADAR can:

- leverage **cross-evaluation**, **clustering** and **reputation** to address heterogeneity and Byzantine contributions;
- adjust rapidly to changes in behavior; and
- mitigate most tested scenarios (limiting case handled up to 80% of poisoned data).

## 2. How generic?

- Only few conditions: parametric models, locally owned evaluation set, a small-scale use case, and a trusted central server.

## 1. RADAR can:

- leverage **cross-evaluation**, **clustering** and **reputation** to address heterogeneity and Byzantine contributions;
- adjust rapidly to changes in behavior; and
- mitigate most tested scenarios (limiting case handled up to 80% of poisoned data).

## 2. How generic?

- Only few conditions: parametric models, locally owned evaluation set, a **small-scale use case**, and a **trusted central server**.

### 1. RADAR can:

- leverage **cross-evaluation**, **clustering** and **reputation** to address heterogeneity and Byzantine contributions;
- adjust rapidly to changes in behavior; and
- mitigate most tested scenarios (limiting case handled up to 80% of poisoned data).

### 2. How generic?

- Only few conditions: parametric models, locally owned evaluation set, a **small-scale use case**, and a **trusted central server**.

### 3. Future works:

- Remove the central server dependency for **increased trust and scalability**.
- Test the approach in more realistic heterogeneous settings.

## Conclusion

---



## CONTRIBUTIONS

- Systematic Literature Review
- E Assessment & eiffel
- R RADAR
- FedITN





### CONTRIBUTIONS

- S Systematic Literature Review
- E Assessment & eiffel
- R RADAR
- F FedITN



- CONTRIBUTIONS**
- S Systematic Literature Review
  - E Assessment & eiffel
  - R RADAR
  - F FedITN



## CONTRIBUTIONS

- S Systematic Literature Review
- E Assessment & eiffel
- R RADAR
- F FedITN



## CONTRIBUTIONS

- S Systematic Literature Review
- E Assessment & eiffel
- R RADAR
- F FedITN



## CONTRIBUTIONS

- S Systematic Literature Review
- E Assessment & eiffel
- R RADAR
- F FedITN



### CONTRIBUTIONS

- S Systematic Literature Review
- E Assessment & eiffel
- R RADAR
- F FedITN



### CONTRIBUTIONS

- S Systematic Literature Review
- E Assessment & eiffel
- R RADAR
- F FedITN



### CONTRIBUTIONS

- S Systematic Literature Review
- E Assessment & eiffel
- R RADAR
- F FedITN













# THANK YOU FOR YOUR ATTENTION!

## Improving Intrusion Detection in Distributed Systems with Federated Learning

- ▶ Three publications in international **conferences**: ICDCS 2024, ARES (BASS) 2023, and SRDS 2024.
- ▶ One article in an international **journal**: IEEE TNSM.
- ▶ National and international **tutorials** on Federated Learning for Intrusion Detection: EUR CyberSchool's Spring Research School 2023, NoF 2023 and ICDCS 2024.

## REFERENCES I

- [1] National Institute of Standards and Technology. ***The NIST Cybersecurity Framework (CSF) 2.0***. NIST CSWP 29. Gaithersburg, MD: National Institute of Standards and Technology, Feb. 26, 2024, NIST CSWP 29. DOI: [10.6028/NIST.CSWP.29](https://doi.org/10.6028/NIST.CSWP.29). URL: <https://nvlpubs.nist.gov/nistpubs/CSWP/NIST.CSWP.29.pdf> (visited on 05/23/2024).
- [2] Brendan McMahan et al. **“Communication-Efficient Learning of Deep Networks from Decentralized Data”**. In: *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*. Ed. by Aarti Singh and Jerry Zhu. Vol. 54. Proceedings of Machine Learning Research. PMLR, Apr. 20–22, 2017, pp. 1273–1282. URL: <https://proceedings.mlr.press/v54/mcmahan17a.html>.
- [3] Peter Kairouz et al. **“Advances and Open Problems in Federated Learning”**. Mar. 8, 2021. arXiv: [1912.04977](https://arxiv.org/abs/1912.04977) [cs, stat]. URL: <http://arxiv.org/abs/1912.04977> (visited on 04/01/2022).
- [4] Léo Lavaur et al. **“The Evolution of Federated Learning-based Intrusion Detection and Mitigation: a Survey”**. In: *IEEE Transactions on Network and Service Management*. Special Issue on Network Security Management (June 2022).

## REFERENCES II

- [5] Léo Lavaur, Yann Busnel, and Fabien Autrel. **“Demo: Highlighting the Limits of Federated Learning in Intrusion Detection”**. In: *Proceedings of the 44th International Conference on Distributed Computing Systems (ICDCS)*. Jersey City, NJ, USA, July 2024.
- [6] Léo Lavaur et al. **“Federated Learning as Enabler for Collaborative Security between Not Fully-Trusting Distributed Parties”**. In: *Proceedings of the 29th Computer & Electronics Security Application Rendezvous (C&ESAR)*. Rennes, France, Oct. 2022.
- [7] Léo Lavaur, Yann Busnel, and Fabien Autrel. **“Systematic Analysis of Label-flipping Attacks against Federated Learning in Collaborative Intrusion Detection Systems”**. In: *Proceedings of the 19th International Conference on Availability, Reliability and Security (ARES). Workshop on Behavioral Authentication for System Security (BASS)*. Vienna, Austria, Aug. 2024.
- [8] Vale Tolpegin et al. **“Data Poisoning Attacks Against Federated Learning Systems”**. In: *Computer Security – ESORICS 2020*. Ed. by Liqun Chen et al. Lecture Notes in Computer Science. Cham: Springer International Publishing, 2020, pp. 480–501. ISBN: 978-3-030-58951-6. DOI: [10.1007/978-3-030-58951-6\\_24](https://doi.org/10.1007/978-3-030-58951-6_24).
- [9] Léo Lavaur et al. **“RADAR: Model Quality Assessment for Reputation-aware Collaborative Federated Learning”**. In: *Proceedings of the 43rd International Symposium on Reliable Distributed Systems (SRDS)*. Charlotte, NC, USA, Sept. 2024.



## REFERENCES III

- [10] Jun Zhou et al. **“A Differentially Private Federated Learning Model against Poisoning Attacks in Edge Computing”**. In: *IEEE Transactions on Dependable and Secure Computing* (2022), pp. 1–1. ISSN: 1941-0018. DOI: [10.1109/TDSC.2022.3168556](https://doi.org/10.1109/TDSC.2022.3168556).
- [11] Christopher Briggs, Zhong Fan, and Peter Andras. **“Federated Learning with Hierarchical Clustering of Local Updates to Improve Training on Non-IID Data”**. In: *2020 International Joint Conference on Neural Networks (IJCNN)*. 2020 International Joint Conference on Neural Networks (IJCNN). July 2020, pp. 1–9. DOI: [10.1109/IJCNN48605.2020.9207469](https://doi.org/10.1109/IJCNN48605.2020.9207469).
- [12] Lingchen Zhao et al. ***Shielding Collaborative Learning: Mitigating Poisoning Attacks through Client-Side Detection***. Mar. 9, 2020. arXiv: [1910.13111 \[cs\]](https://arxiv.org/abs/1910.13111). URL: <http://arxiv.org/abs/1910.13111> (visited on 08/28/2022). Pre-published.
- [13] Paul Resnick et al. **“Reputation Systems”**. In: *Communications of the ACM* 43.12 (Dec. 1, 2000), pp. 45–48. ISSN: 0001-0782. DOI: [10.1145/355112.355122](https://doi.org/10.1145/355112.355122). URL: <https://doi.org/10.1145/355112.355122> (visited on 02/01/2023).

## REFERENCES IV

- [14] Mohanad Sarhan, Siamak Layeghy, and Marius Portmann. *Towards a Standard Feature Set for Network Intrusion Detection System Datasets*. May 14, 2021. arXiv: [2101.11315 \[cs\]](https://arxiv.org/abs/2101.11315). URL: <http://arxiv.org/abs/2101.11315> (visited on 09/12/2022). Pre-published.
- [15] Rafael Uetz et al. “Reproducible and Adaptable Log Data Generation for Sound Cybersecurity Experiments”. In: *Annual Computer Security Applications Conference. ACSAC '21: Annual Computer Security Applications Conference. Virtual Event USA*: ACM, Dec. 6, 2021, pp. 690–705. ISBN: 978-1-4503-8579-4. DOI: [10.1145/3485832.3488020](https://doi.org/10.1145/3485832.3488020). URL: <https://dl.acm.org/doi/10.1145/3485832.3488020> (visited on 08/08/2022).
- [16] ACM. *Artifact Review and Badging v1.1*. Aug. 24, 2020. URL: <https://www.acm.org/publications/policies/artifact-review-and-badging-current> (visited on 08/17/2022).
- [17] Daniel J Beutel et al. “Flower: A Friendly Federated Learning Research Framework”. 2020. arXiv: [2007.14390](https://arxiv.org/abs/2007.14390).
- [18] Eelco Dolstra. “The Purely Functional Software Deployment Model”. S.l.: s.n., 2006.

## Extra Slides

---

Assessment

Sound experiments [15]; [16]:

- ▶ *valid* (i.e., well-defined and unrefutable);
- ▶ *controllable* (e.g., parameterized); and
- ▶ *reproducible* (i.e., the same results can be obtained by another group using the author's artefact).

[15] Uetz et al. "Reproducible and Adaptable Log Data Generation for Sound Cybersecurity Experiments". *Annual Computer Security Applications Conference*. 2021

[16] ACM. *Artifact Review and Badging v1.1*. 2020

## EXPERIMENTAL SETUP

Experiment orchestration using **Eiffel** [5].

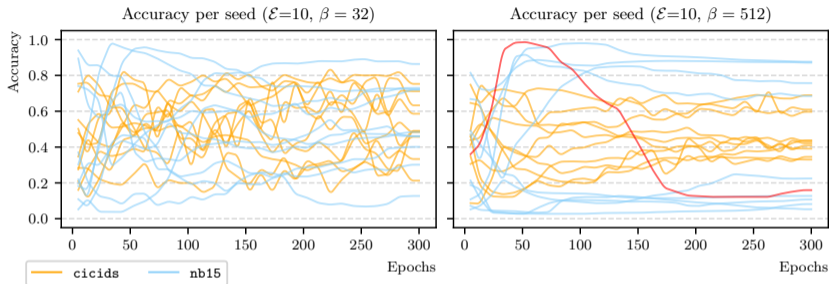
- ▶ **Flower** simulation framework [17] for Federated Learning (FL).
- ▶ **Hydra** for experiment generation and configuration.
- ▶ Custom-made poisoning engine with different attack strategies.
- ▶ Nix [18] and Poetry to fix system and Python dependencies, enabling reproducibility.

1,067 experiments  $\times$  10 seeds (1,613 hours of computation.)

[17] Beutel et al. “Flower: A Friendly Federated Learning Research Framework”. 2020

[18] Dolstra. “The Purely Functional Software Deployment Model”. 2006

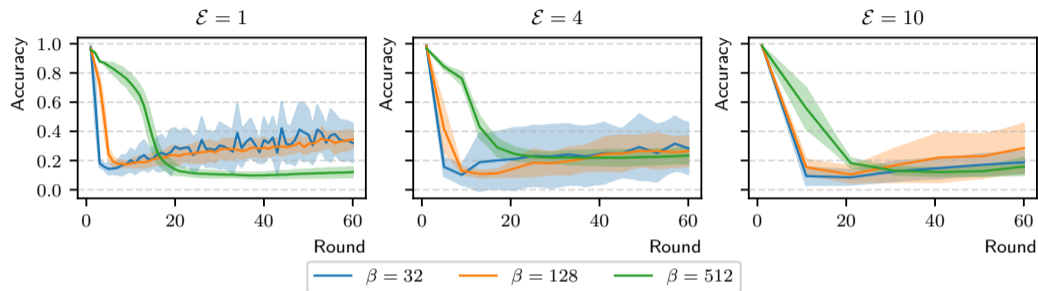
## RQ1: ARE POISONING ATTACKS PREDICTABLE?



**Figure:** Predictability of label-flipping attacks.

- ▶ Very high variance in the results, but tends to stabilize (on different values) after a few rounds.
- ▶ The impact of the attack is highly dependent on the seed.
  - Initial parameters, data shuffling, partitioning, ...

## RQ2: DO HYPERPARAMETERS INFLUENCE THE IMPACT OF POISONING ATTACKS?



**Figure:** Effect of the hyperparameters on the accuracy of the poisoned model in the late scenario (50% attackers, CICIDS).

- ▶ **late-3** scenario: attackers start poisoning after 3 rounds
- ▶ High batch size leads to more inertia, less instantaneous impact  
→ More impactful in constrained environments

## Extra Slides

---

RADAR



# RESULTS

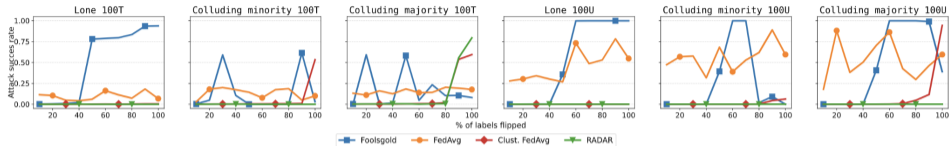


Figure: Baseline comparison.

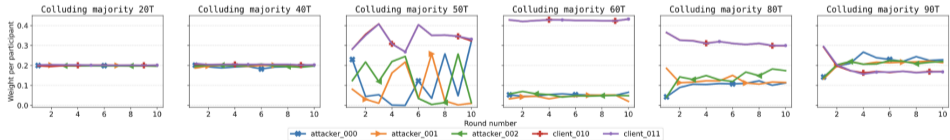


Figure: RADAR's limiting scenario.