



Cyprus
University of
Technology



Securing Federated Sensitive Topic Classification against Poisoning Attacks

Tianyue Chu¹

Alvaro Garcia-Recuero¹

Costas Iordanou²

Georgios Smaragdakis³

Nikolaos Laoutaris¹

¹ IMDEA Networks Institute

² Cyprus University of Technology

³ TU Delft

[Developing the
Science of Networks]



(a) Zone: European Union



(b) Zone: California (USA)

□ Legislations (GDPR)

- define sensitive data : related to health, political opinions, religious beliefs, sexual orientation and racial ethnic origin
- ensure data protection and safeguard online content that contains sensitive data



(a) Zone: European Union



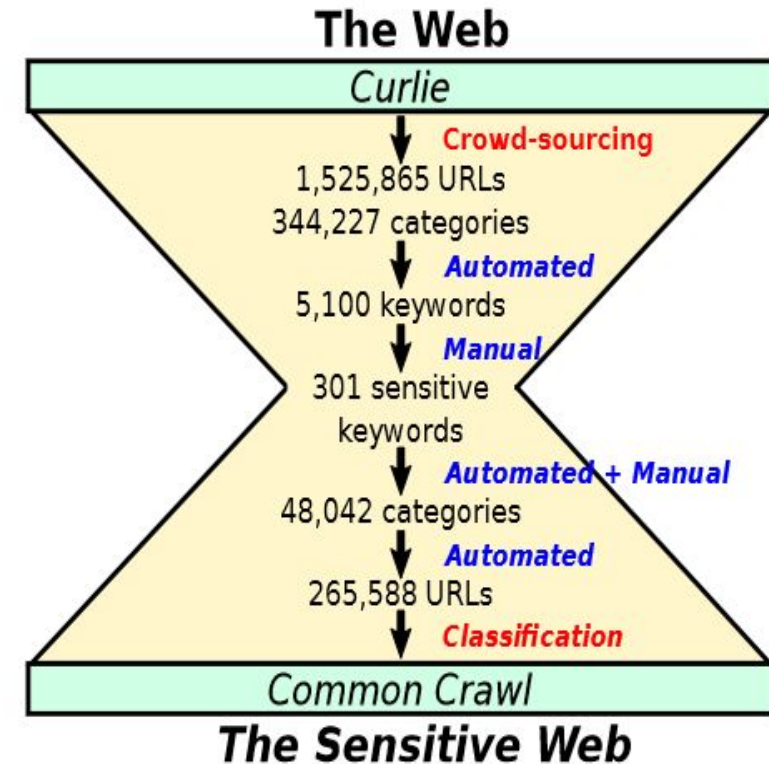
(b) Zone: California (USA)

□ Legislations (GDPR)

- define sensitive data : related to health, political opinions, religious beliefs, sexual orientation and racial ethnic origin
- ensure data protection and safeguard online content that contains sensitive data

How to detect whether the content of a URL relates to any of the sensitive data?

- A classifier for detecting GDPR sensitive data [2]
 - Defined on GDPR sensitive categories
 - train a centralized classifier to detect whether the content of a URL relates to sensitive categories
 - trained using 156k sensitive URLs from Curlie with over 90% accuracy



[1] Matic, Srdjan, Costas Iordanou, Georgios Smaragdakis, and Nikolaos Laoutaris. "Identifying sensitive urls at web-scale." In *Proceedings of the ACM Internet Measurement Conference*, pp. 619-633. 2020.

- Why do we need to identify sensitive URLs
 - 30% of sensitive URLs are hosted in domains that fail to use HTTPS.
 - In sensitive web pages with third-party cookies, 87% of the third-parties set at least one persistent cookie.

- Why do we need to identify sensitive URLs

- 30% of sensitive URLs are hosted in domains that fail to use HTTPS.



- In sensitive web pages with third-party cookies, 87% of the third-parties set at least one persistent cookie.



Being tracked when visiting web pages that contain sensitive content

- Why do we need to identify sensitive URLs

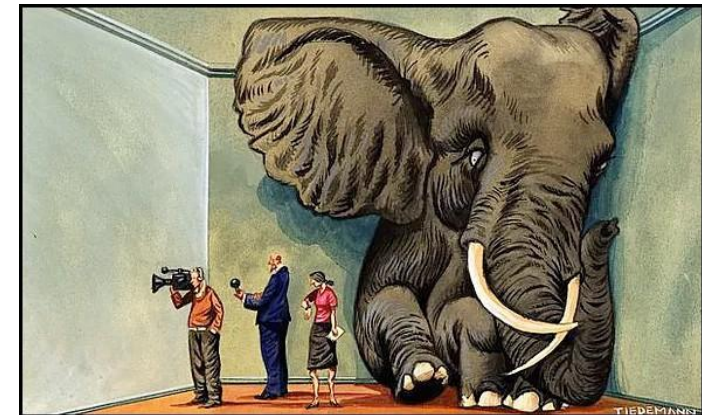
- 30% of sensitive URLs are hosted in domains that fail to use HTTPS.



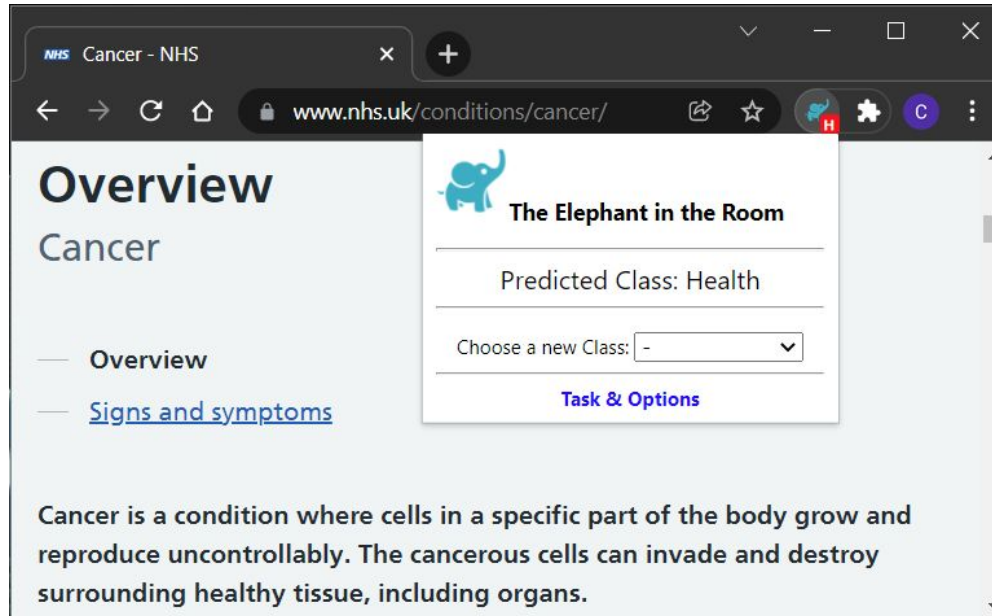
- In sensitive web pages with third-party cookies, 87% of the third-parties set at least one persistent cookie.



the “Elephant in the Room” of privacy

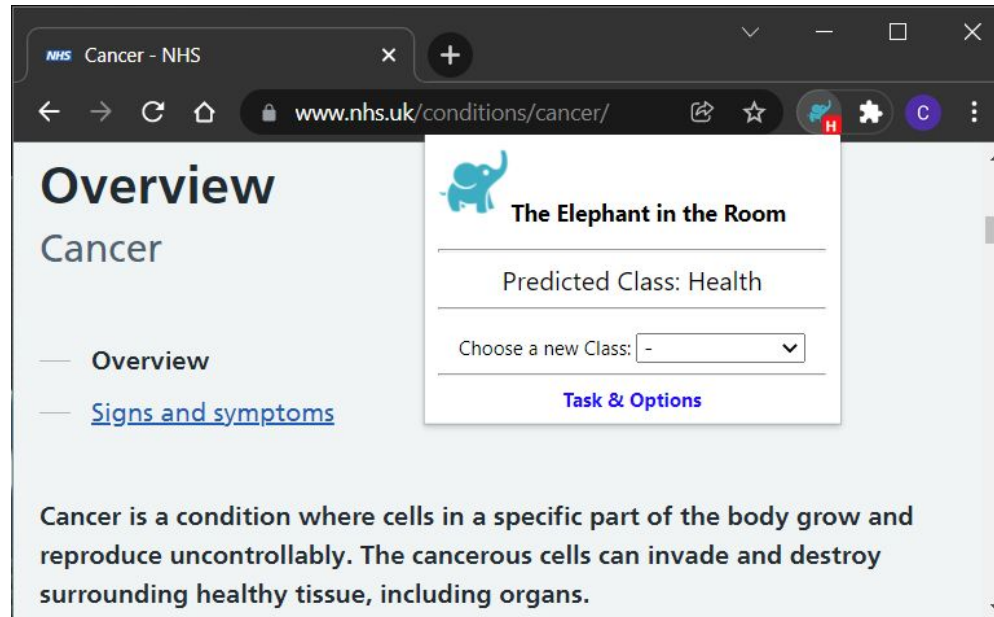


- What do we want to do?



A classifier capable of detecting URLs containing sensitive content in real-time

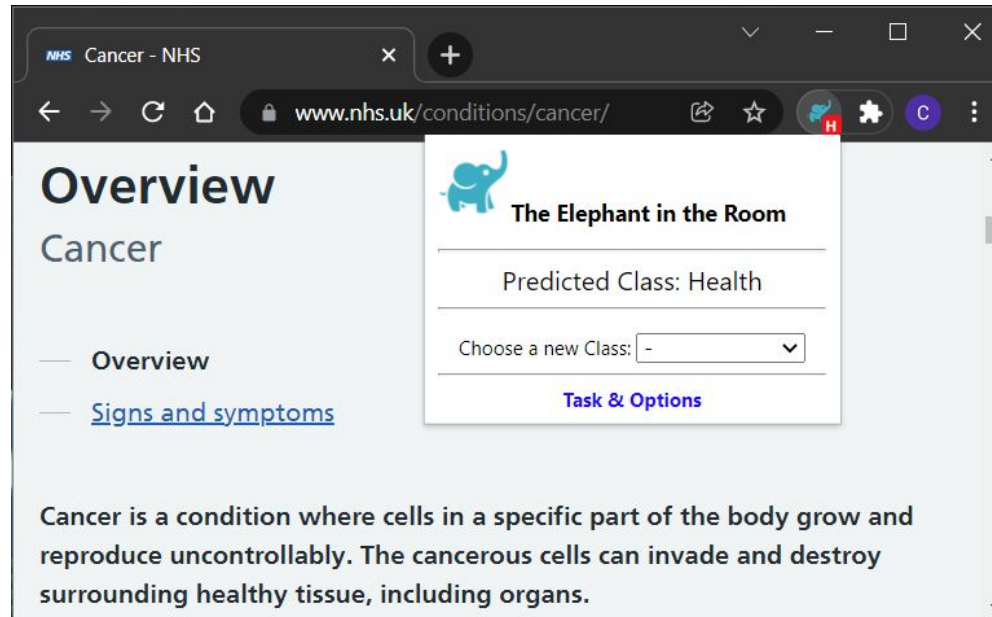
- What do we want to do?



A classifier capable of detecting URLs containing sensitive content in real-time

when you visit a webpage, it pops up which category this webpage falls in

- What do we want to do?

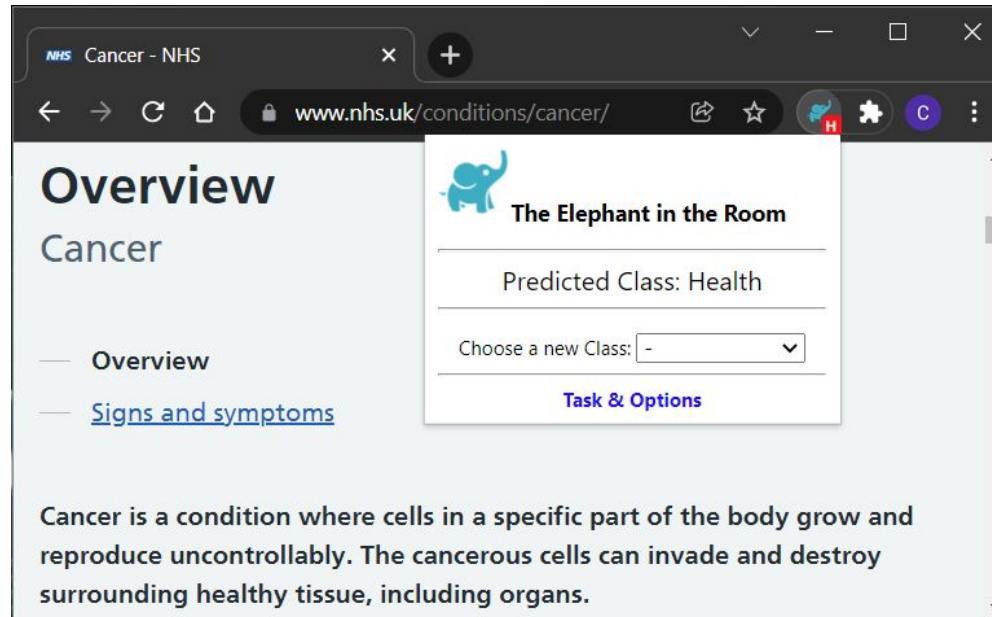


A classifier capable of detecting URLs containing sensitive content in real-time



when you visit a webpage, it pops up which category this webpage falls in

- What do we want to do?

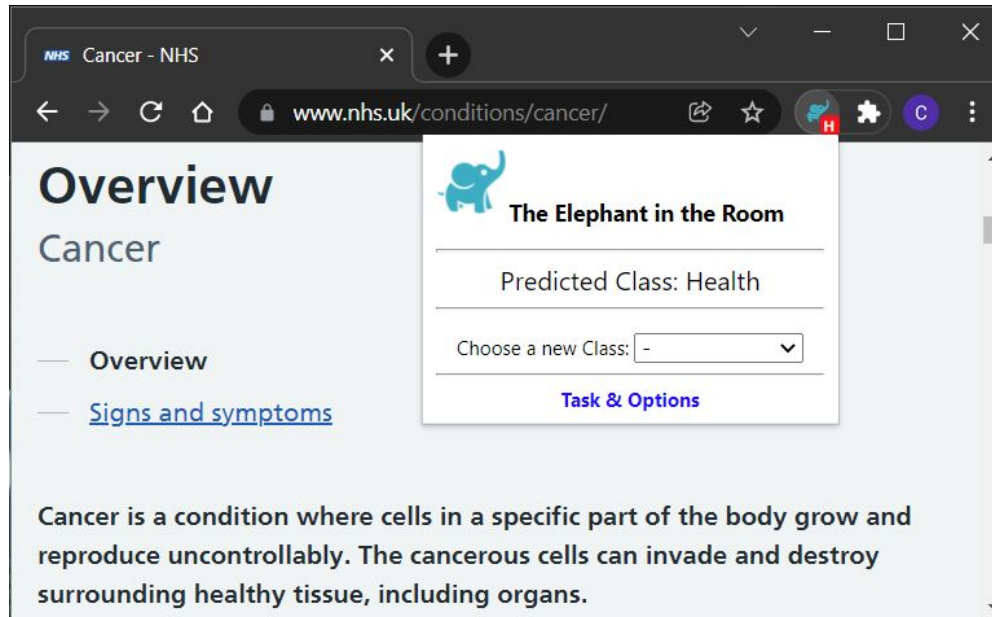


A classifier capable of detecting URLs containing sensitive content in real-time



when you visit a webpage, it pops up which category this webpage falls in

- What do we want to do?



A classifier capable of detecting URLs containing sensitive content in real-time



when you visit a webpage, it pops up which category this webpage falls in

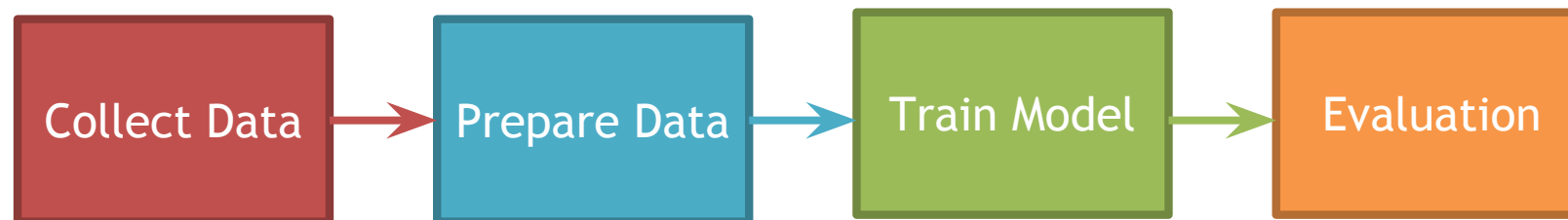


- **Limitations of centralised classifier**
 - tied to a fixed training set:
difficult to quickly cover labels related to **yet unseen** sensitive content
 - being centralized: cannot be used to drive a privacy-preserving distributed classification system

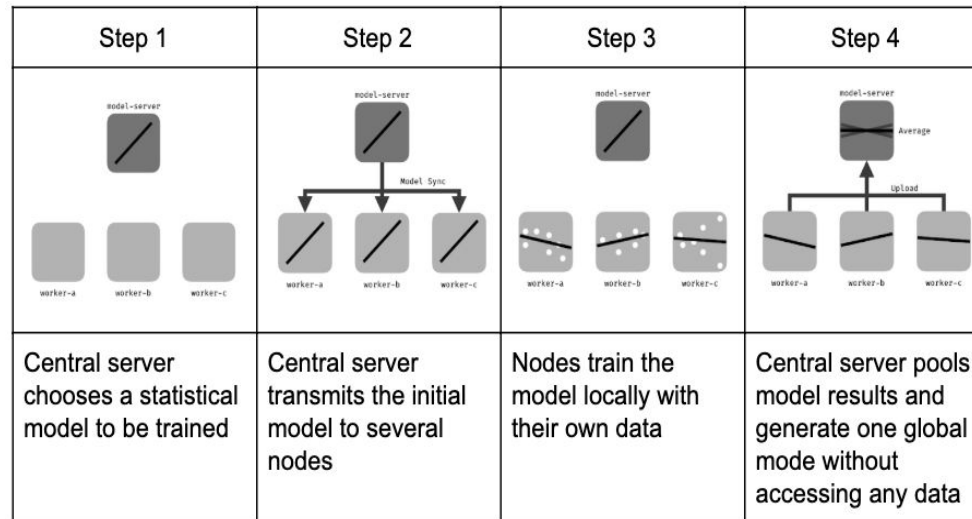
- Limitations of centralised classifier
 - tied to a fixed training set:
difficult to quickly cover labels related to yet unseen sensitive content
 - being centralized: cannot be used to drive a privacy-preserving distributed classification system



- Limitations of centralised classifier
 - tied to a fixed training set:
difficult to quickly cover labels related to yet unseen sensitive content
 - being centralized: cannot be used to drive a privacy-preserving distributed classification system

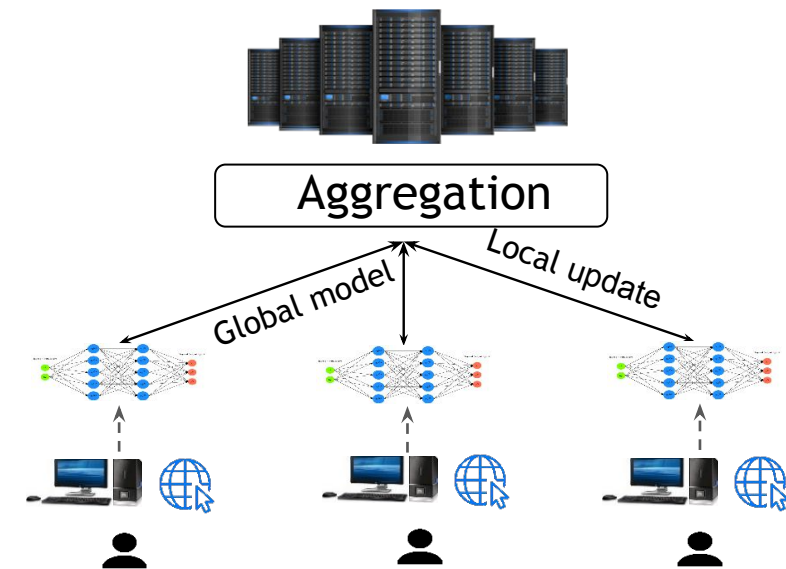
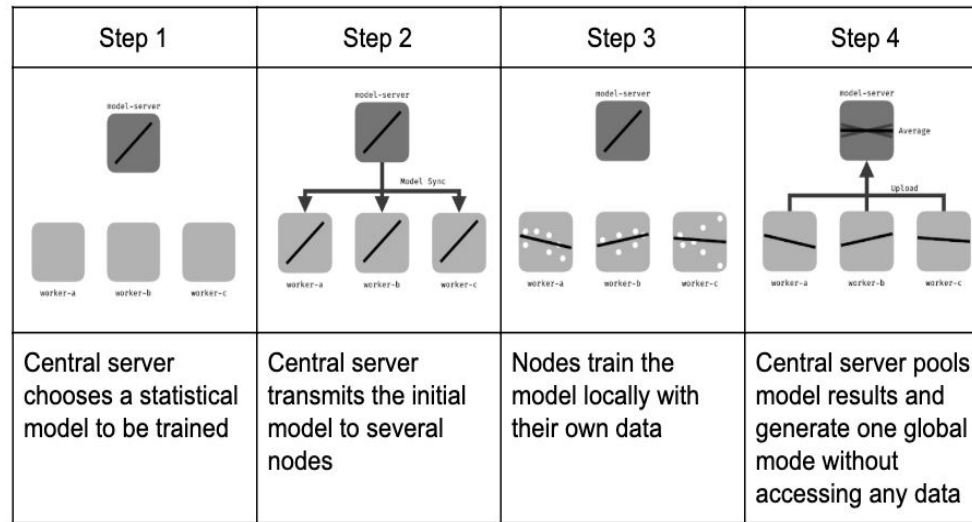


Solution: Federated learning (FL) [3]



[3] McMahan, Brendan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. "Communication-efficient learning of deep networks from decentralized data." In *Artificial intelligence and statistics*, pp. 1273-1282. PMLR, 2017.

Solution: Federated learning (FL) [3]

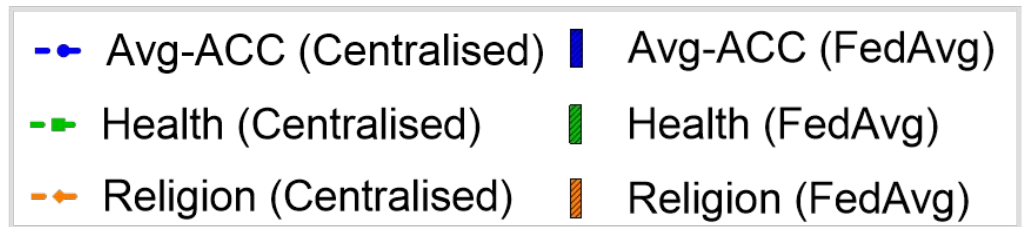
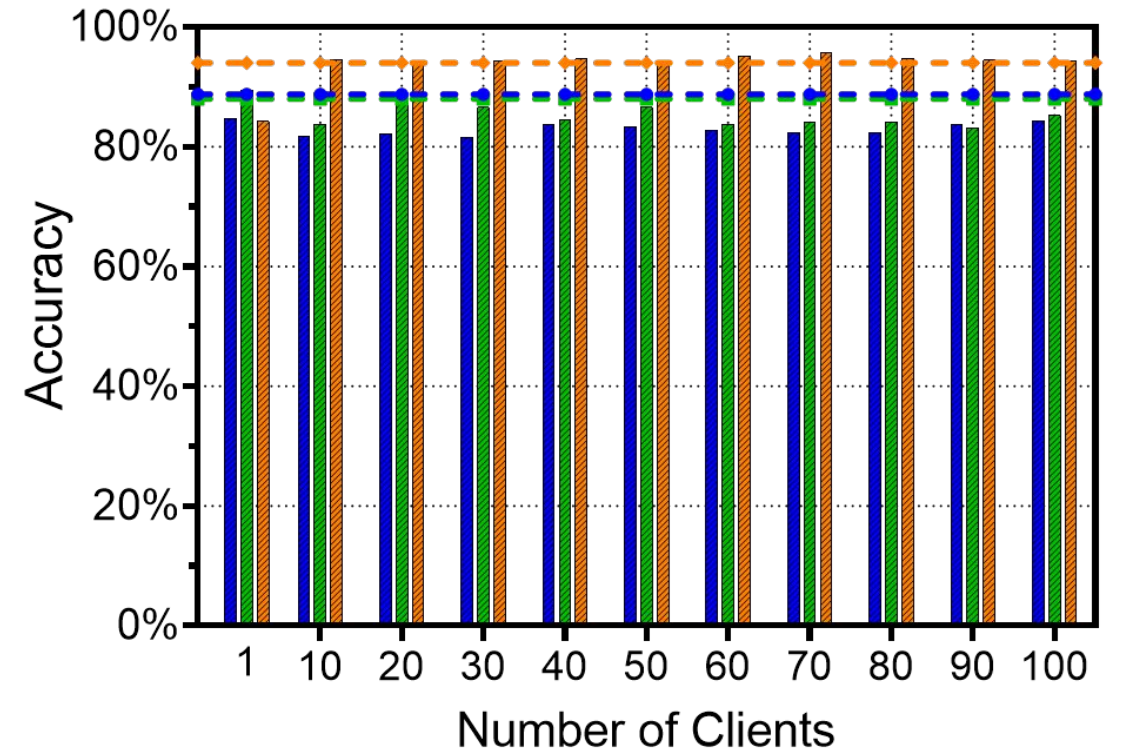


- **Up-to-date:** continuously learn from real-time web data gathered by users
- **Being distributed with privacy:** users train the classifier locally using personal data

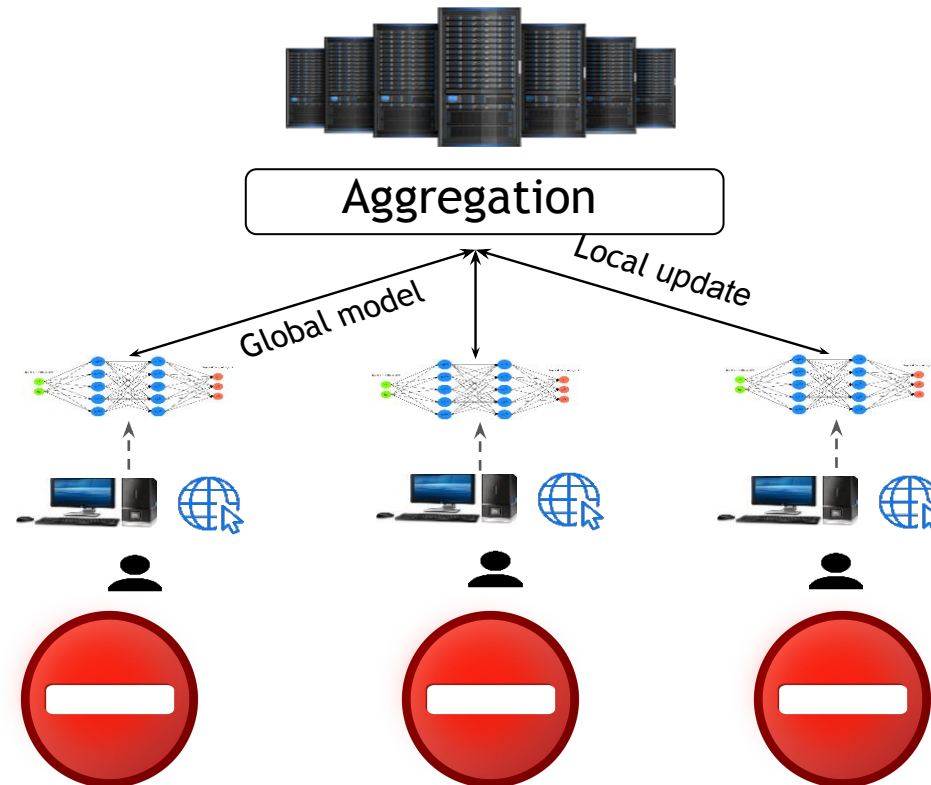
[3] McMahan, Brendan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. "Communication-efficient learning of deep networks from decentralized data." In *Artificial intelligence and statistics*, pp. 1273-1282. PMLR, 2017.

- Centralized classifier
 - Naïve-Bayes
- FL-based classifier
 - A simple neural network
 - FedAvg [3]
 - Performance

Accuracy	FL	Centralized
Health	85%	88%
Religion	93%	94%

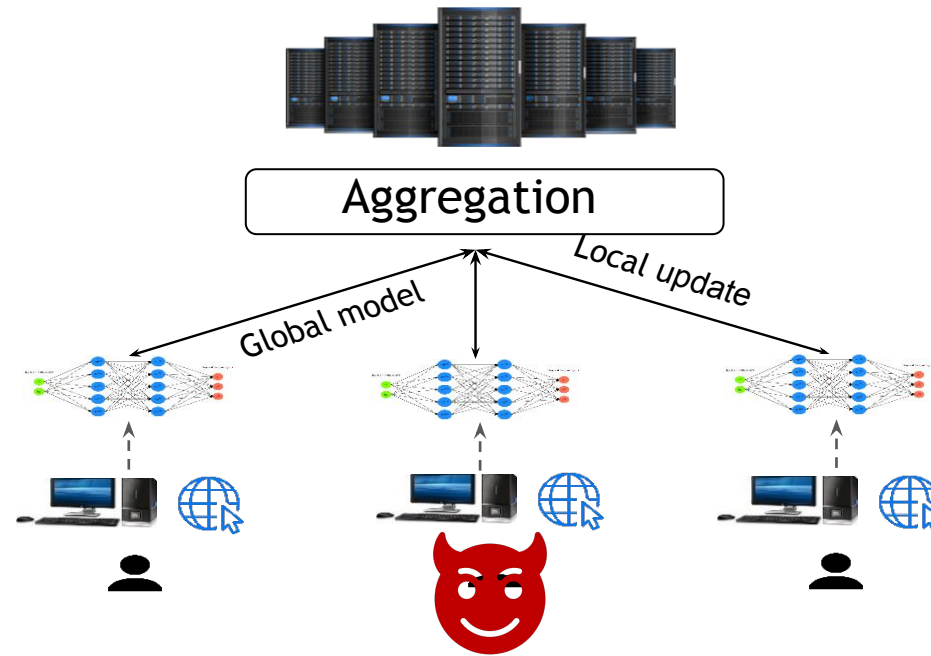


- Federated Learning (FL) is vulnerable



Federated learning

- Federated Learning (FL) is vulnerable



FL is vulnerable to so-called poisoning attacks

Byzantine-robust FL methods

- Several Byzantine-robust defense methods have been developed
- Recent studies [4], [5], [6] have shown some methods are forgetful by not tracking information from previous aggregation rounds

[4] Bagdasaryan, Eugene, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. "How to backdoor federated learning." In *International Conference on Artificial Intelligence and Statistics*, pp. 2938-2948. PMLR, 2020.

[5] Bhagoji, Arjun Nitin, Supriyo Chakraborty, Prateek Mittal, and Seraphin Calo. "Analyzing federated learning through an adversarial lens." In *International Conference on Machine Learning*, pp. 634-643. PMLR, 2019.

[6] Karimireddy, Sai Praneeth, Lie He, and Martin Jaggi. "Learning from history for byzantine robust optimization." In *International Conference on Machine Learning*, pp. 5311-5319. PMLR, 2021.

Byzantine-robust FL methods

- Several Byzantine-robust defense methods have been developed
- Recent studies [4], [5], [6] have shown some methods are forgetful by not tracking information from previous aggregation rounds



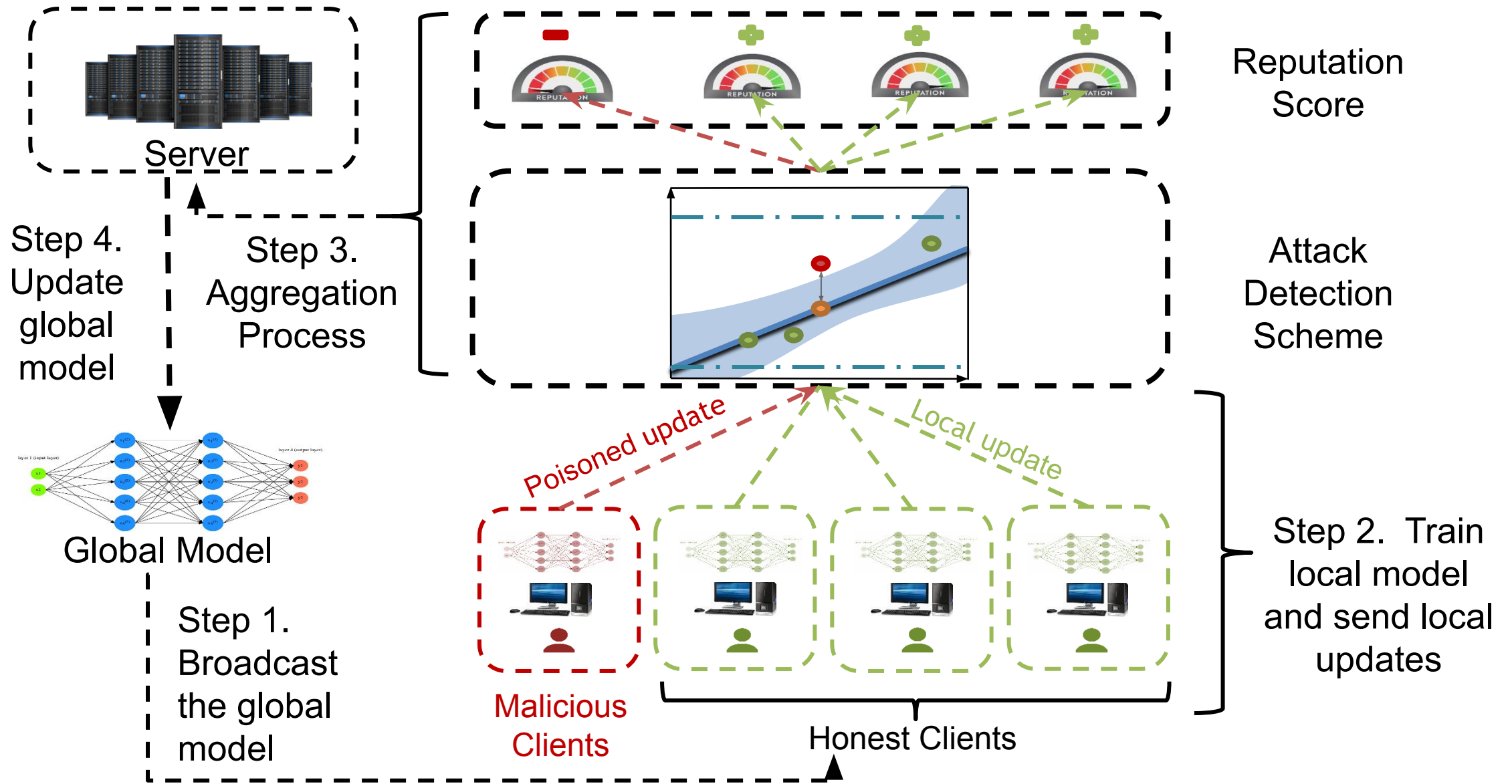
We design a robust FL aggregation method to generate reputation scores of clients based on their historical behaviors

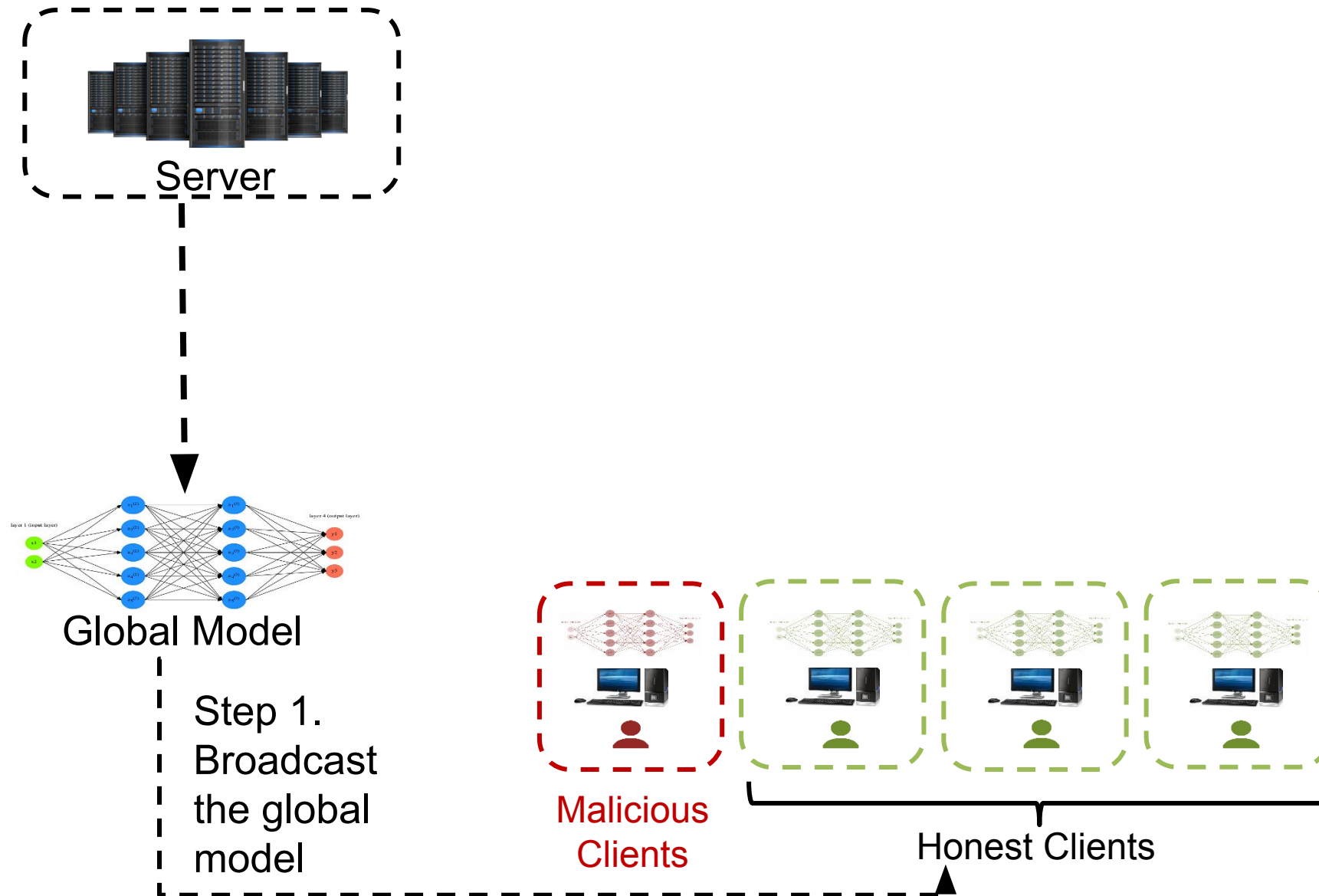
[4] Bagdasaryan, Eugene, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. "How to backdoor federated learning." In *International Conference on Artificial Intelligence and Statistics*, pp. 2938-2948. PMLR, 2020.

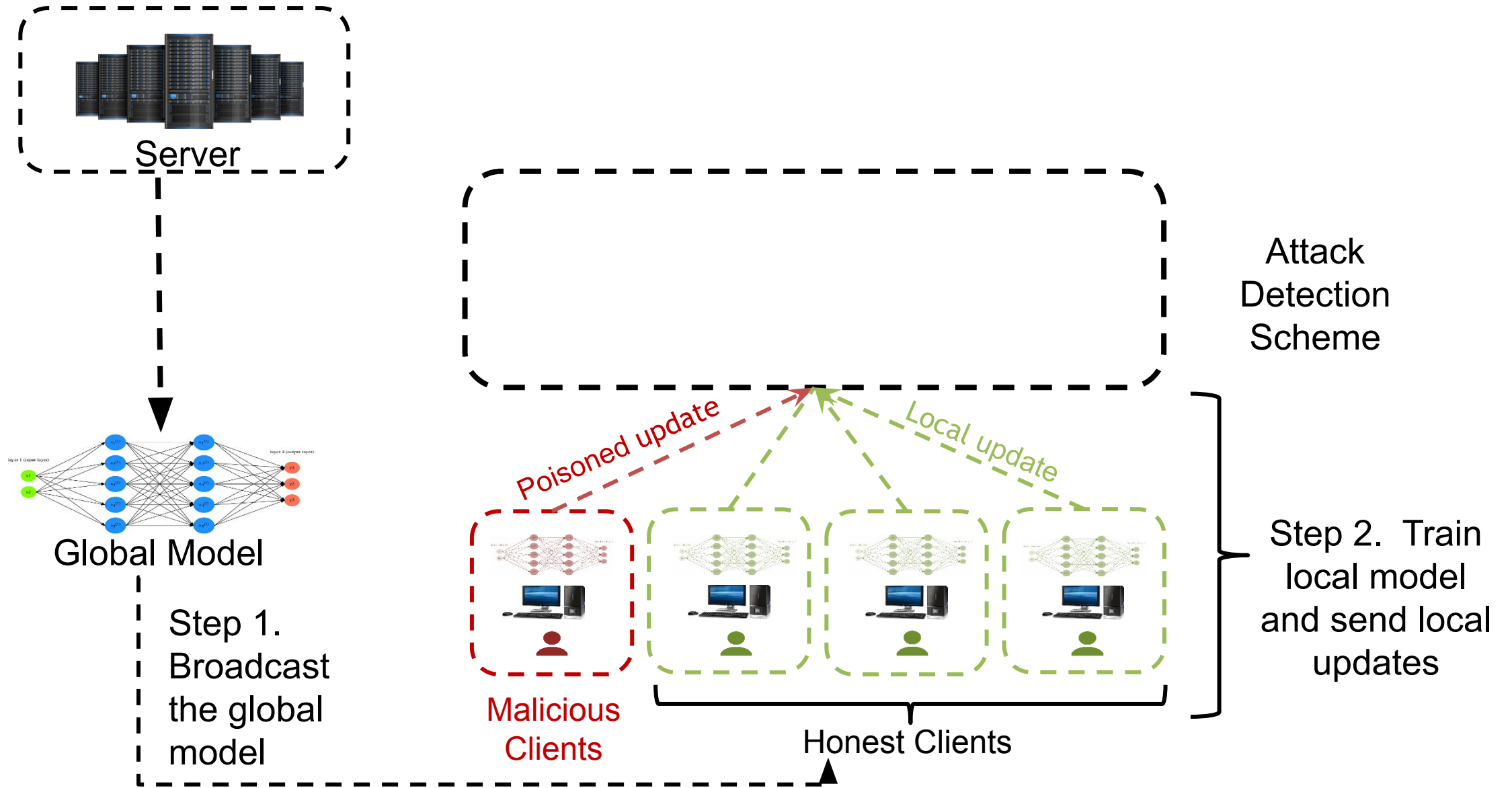
[5] Bhagoji, Arjun Nitin, Supriyo Chakraborty, Prateek Mittal, and Seraphin Calo. "Analyzing federated learning through an adversarial lens." In *International Conference on Machine Learning*, pp. 634-643. PMLR, 2019.

[6] Karimireddy, Sai Praneeth, Lie He, and Martin Jaggi. "Learning from history for byzantine robust optimization." In *International Conference on Machine Learning*, pp. 5311-5319. PMLR, 2021.

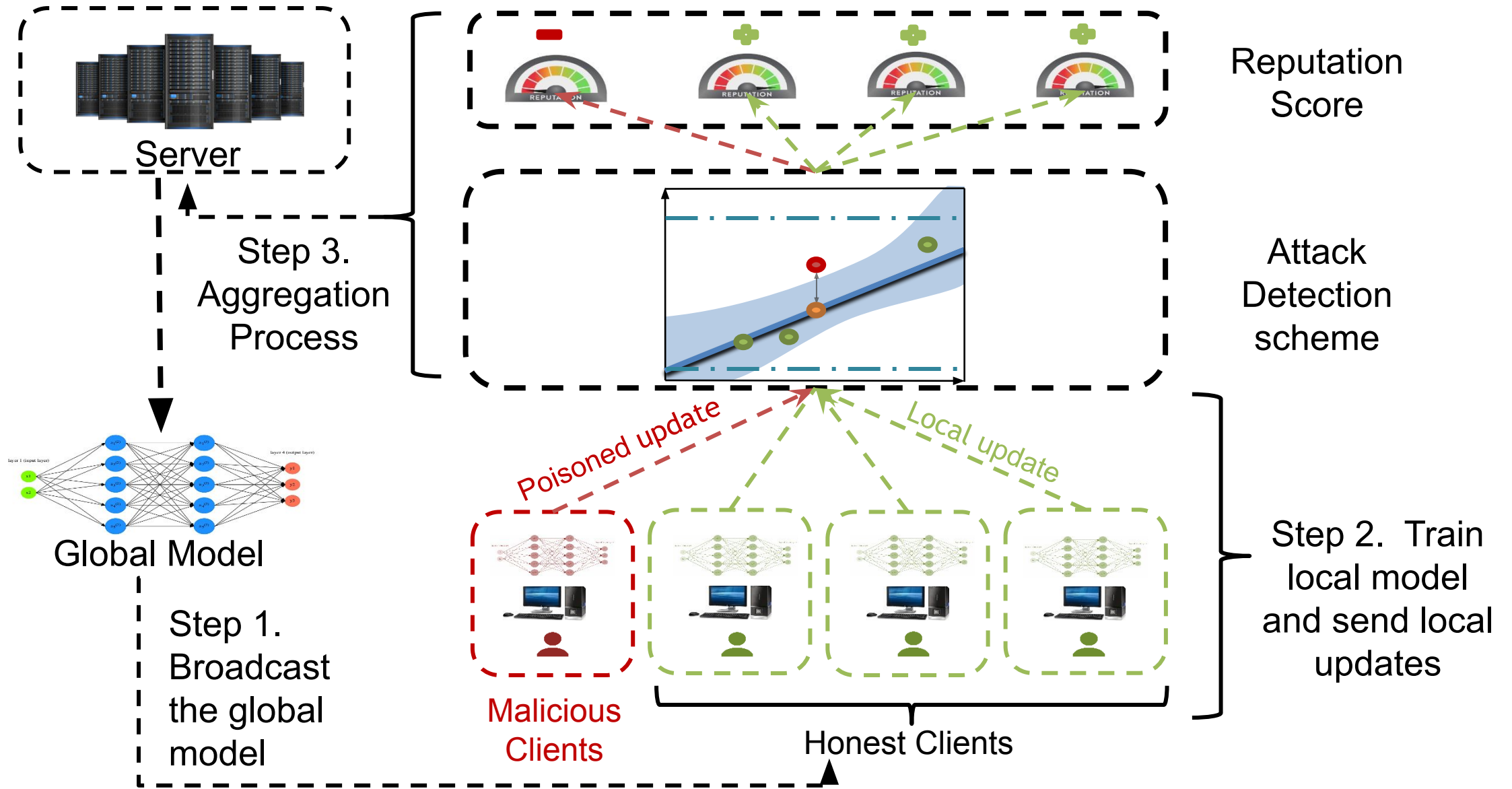
Algorithm



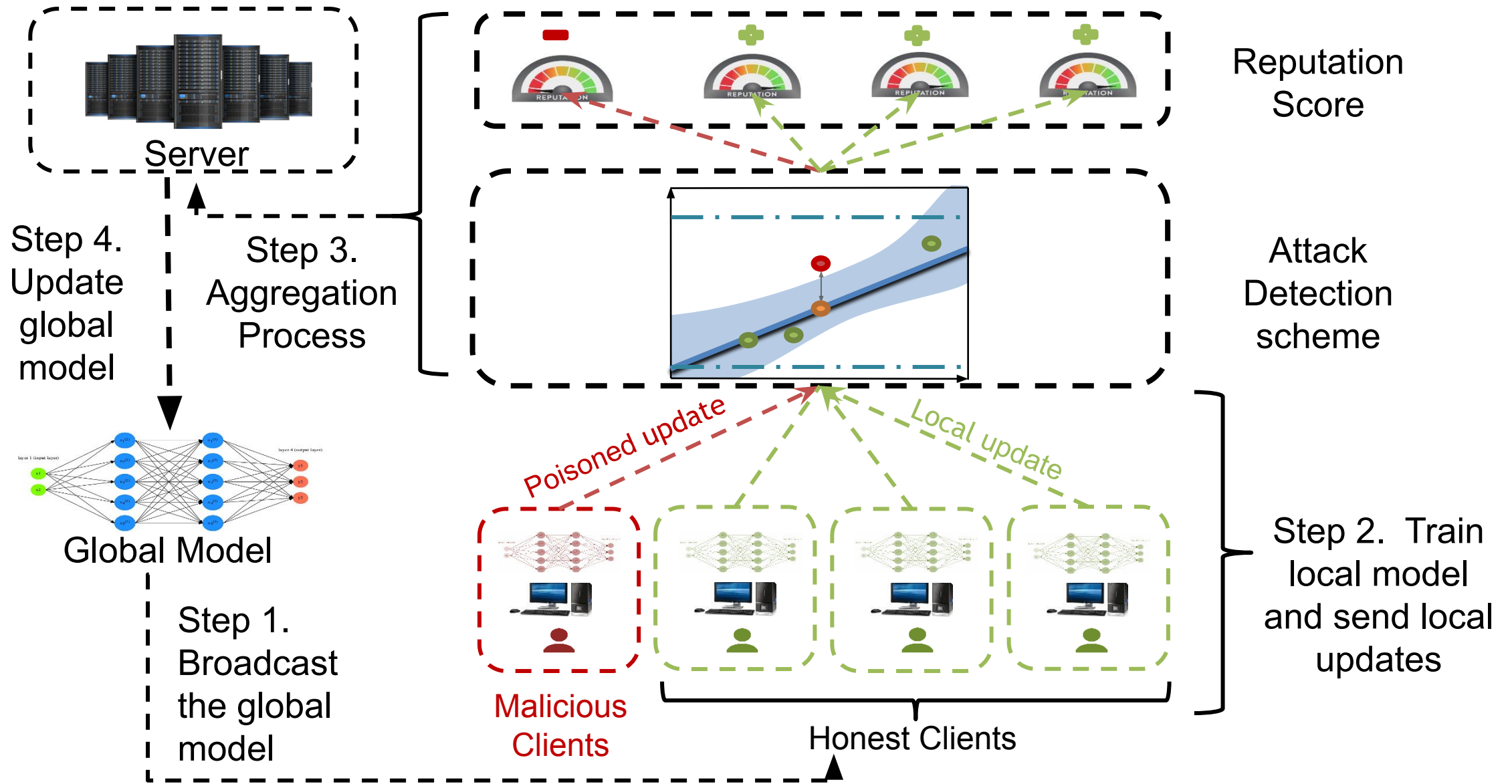




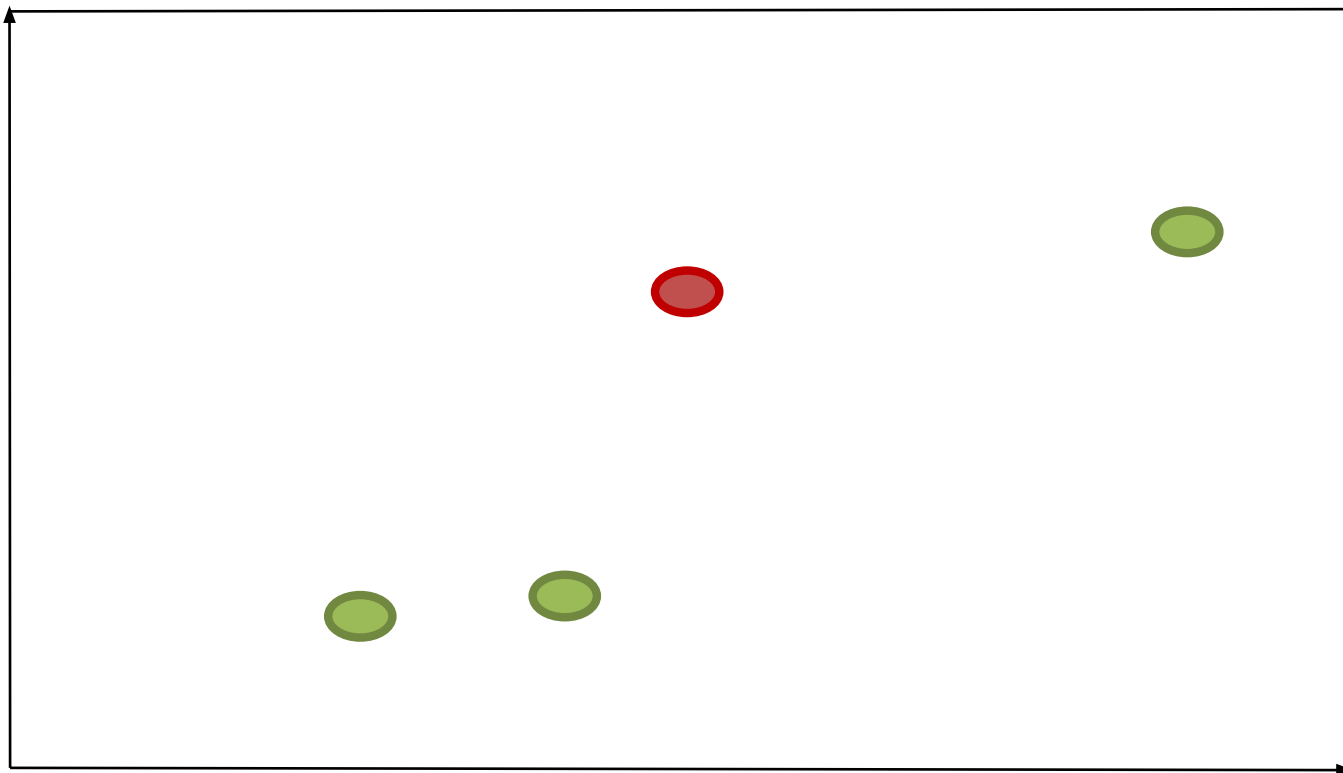
Algorithm



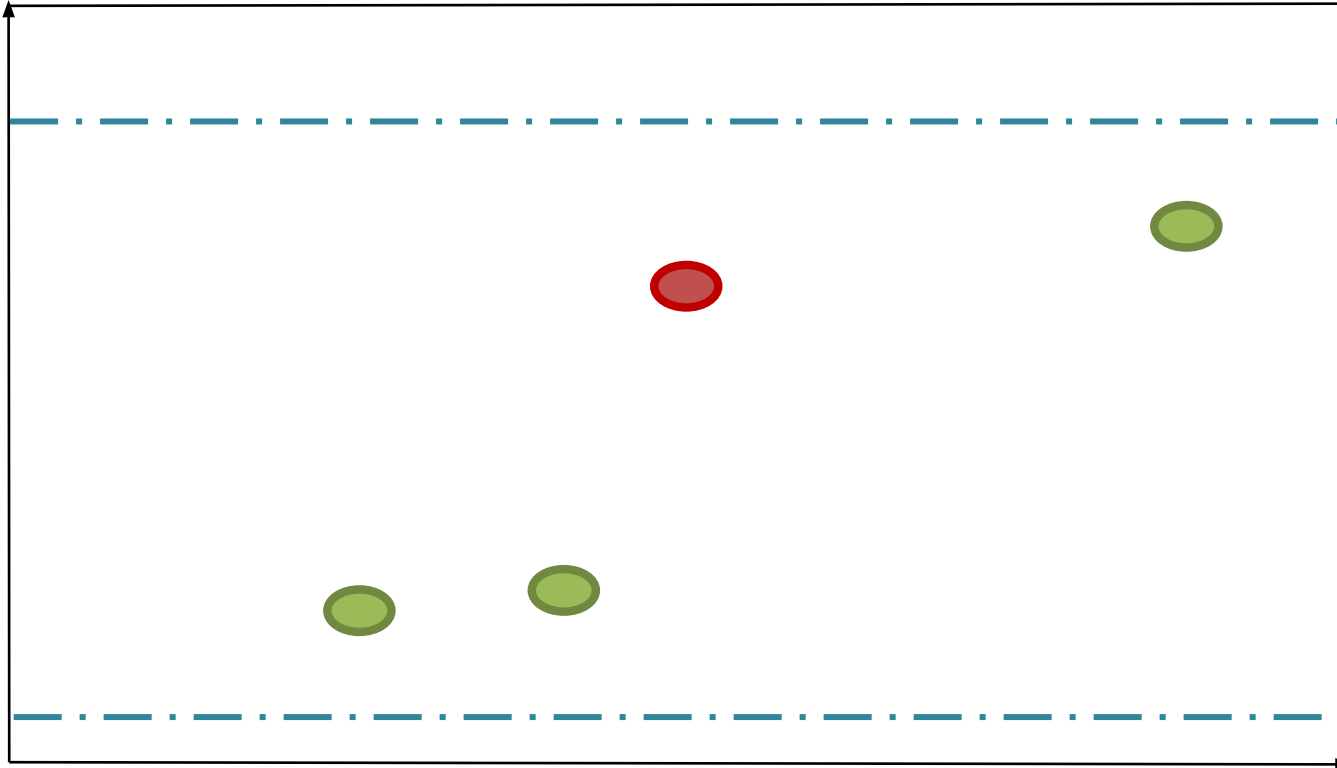
Algorithm



Attack Detection Scheme



Attack Detection Scheme



Algorithm 1: Rescale(w)

Input : $\{w_{i,n}^t\}$ \leftarrow Local Model parameters in round t

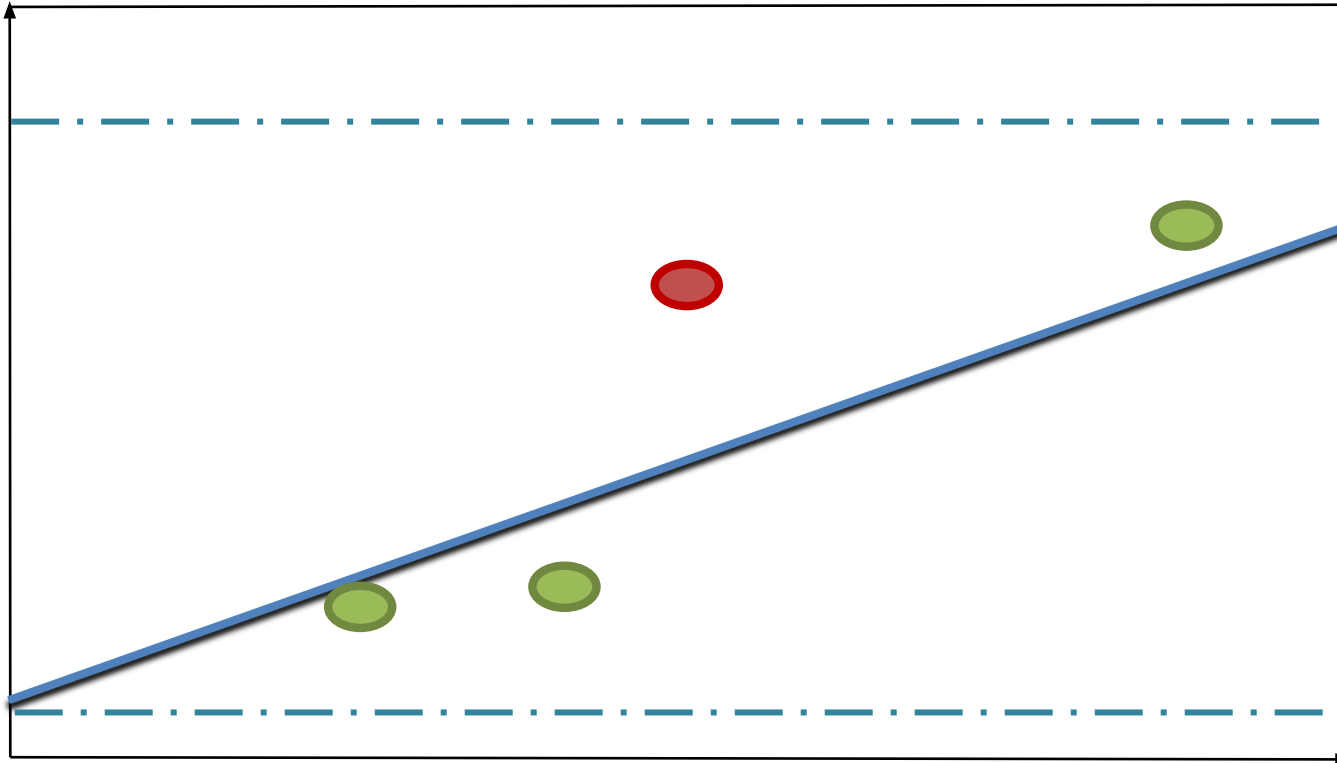
Output: $\{w_{i,n}^t\}$ with range of value less than ϖ

```

1 for  $n \leftarrow 1$  to  $N$  do
2   // Determine the maximum range
3    $Rm = \max w_{i,n}^t - \min w_{i,n}^t = w_{i,n}^{t,(Max)} - w_{i,n}^{t,(Min)}$ 
4   while  $Rm > \varpi$  do
5     // Rescale range based on standard
      deviation.
6      $w_{i,n}^{t,(Max)} := w_{i,n}^{t,(Max)} - \sigma(w_{i,n}^t);$ 
7      $w_{i,n}^{t,(Min)} := w_{i,n}^{t,(Min)} + \sigma(w_{i,n}^t);$ 
8     // Updated  $Rm$ .
9      $Rm = \max w_{i,n}^t - \min w_{i,n}^t =$ 
       $w_{i,n}^{t,(Max)} - w_{i,n}^{t,(Min)}$ 
10  end while
11 end for

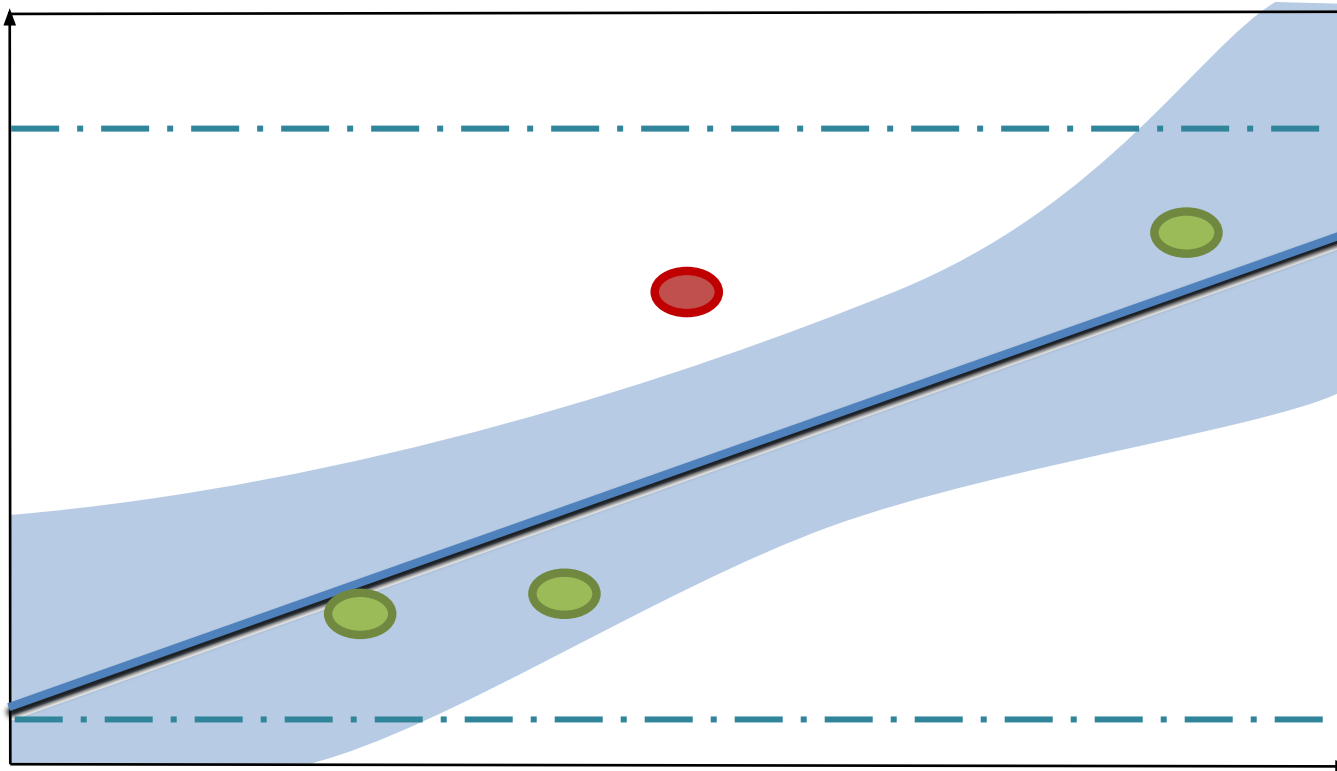
```

Attack Detection Scheme



$$\hat{B}_n = \operatorname{median}_i \left\{ \operatorname{median}_{i \neq j} \{B_n(i, j)\} \right\}$$
$$\hat{A}_n = \operatorname{median}_i \left\{ w_{i,n} - \hat{B}_n x_{i,n} \right\}$$

Attack Detection Scheme

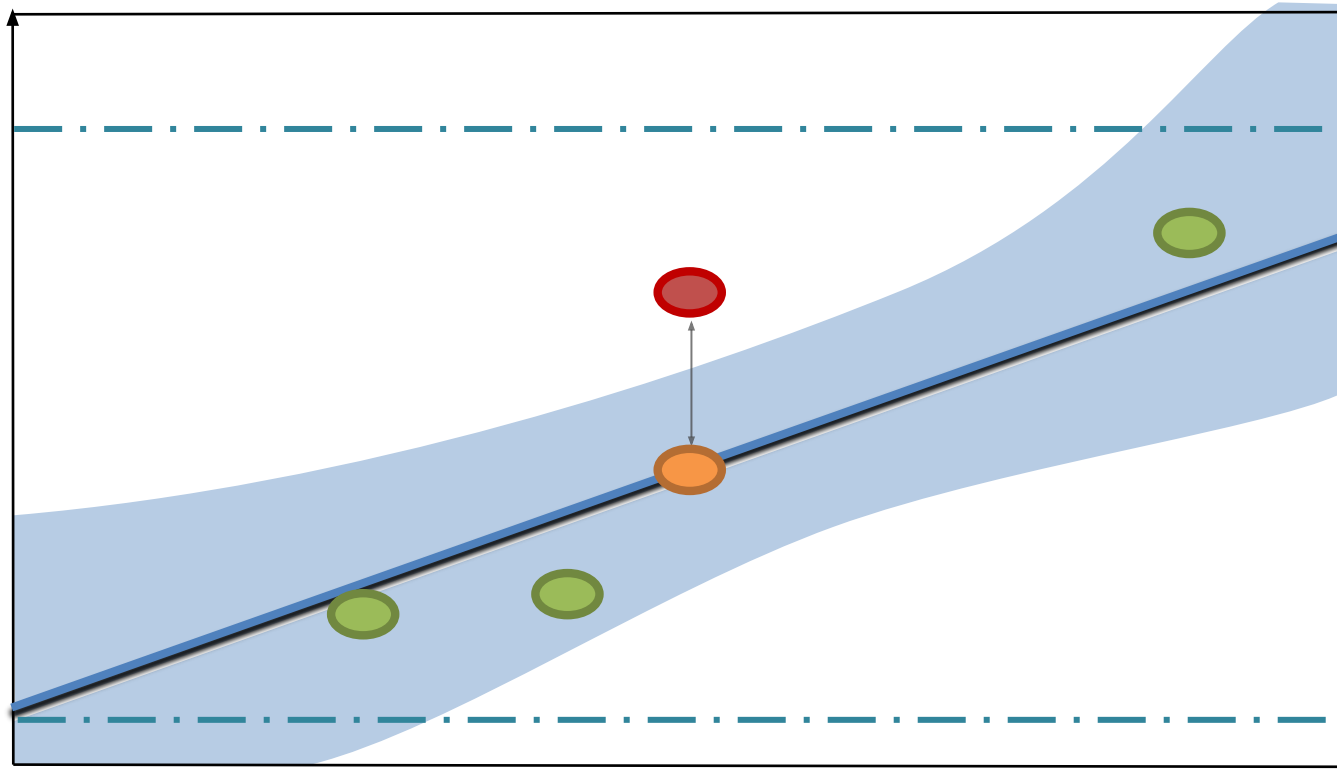


$$s_{i,n}^t = \frac{\sqrt{1 - \text{diag}(H_n^t)}}{e_{i,n}^t} \Psi \left(\frac{e_{i,n}^t}{\sqrt{1 - \text{diag}(H_n^t)}} \right)$$

where confidence interval $\Psi(x)$:

$$\Psi(x) = \max\{-\lambda\sqrt{2/M}, \min(\lambda\sqrt{2/M}, x)\}$$

Attack Detection Scheme



An update with confidence value less than threshold δ , it replaces it with the median

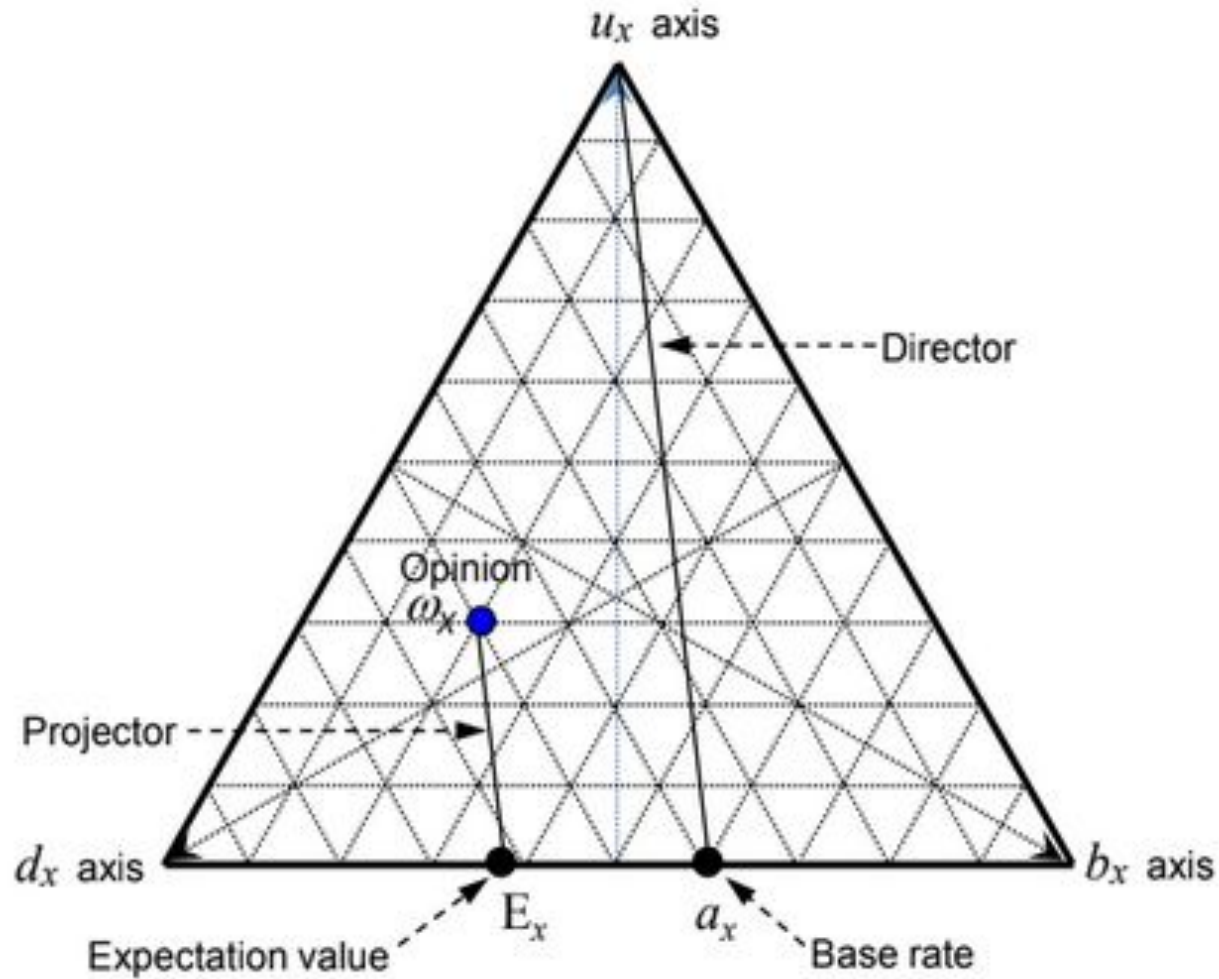
$$\mathbf{w}_{i,n}^t = \begin{cases} \mathbf{w}_{i,n}^t & \text{if } s_{i,n}^t > \delta \\ \text{median}_i \{ \mathbf{w}_{i,n}^t \} & \text{if } s_{i,n}^t \leq \delta \end{cases}$$

$$s_{i,n}^t = \frac{\sqrt{1 - \text{diag}(H_n^t)}}{e_{i,n}^t} \Psi \left(\frac{e_{i,n}^t}{\sqrt{1 - \text{diag}(H_n^t)}} \right)$$

where confidence interval $\Psi(x)$:

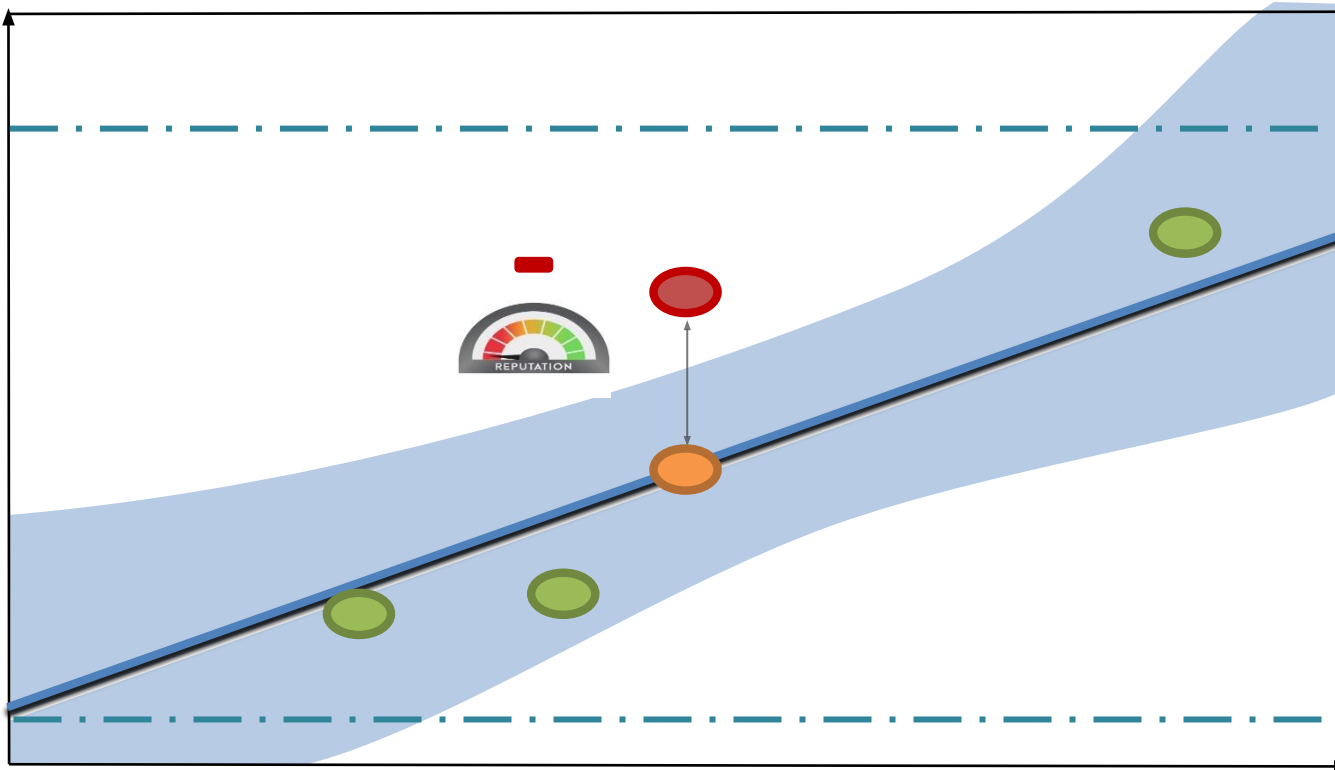
$$\Psi(x) = \max\{-\lambda\sqrt{2/M}, \min(\lambda\sqrt{2/M}, x)\}$$

Subjective logic model



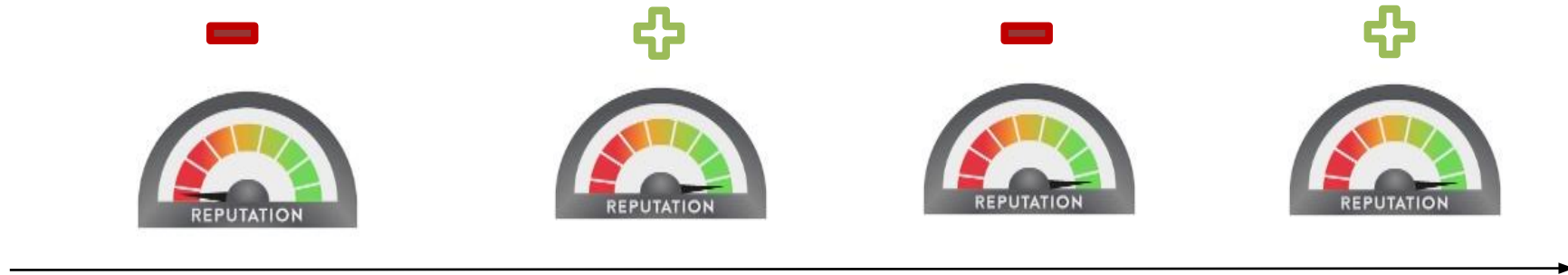
$$\begin{cases} b_i^t = \frac{\kappa P_i^t}{\kappa P_i^t + \eta N_i^t + W} \\ d_i^t = \frac{\eta N_i^t}{\kappa P_i^t + \eta N_i^t + W} \\ u_i^t = \frac{W}{\kappa P_i^t + \eta N_i^t + W} \end{cases}$$

Subject logic model

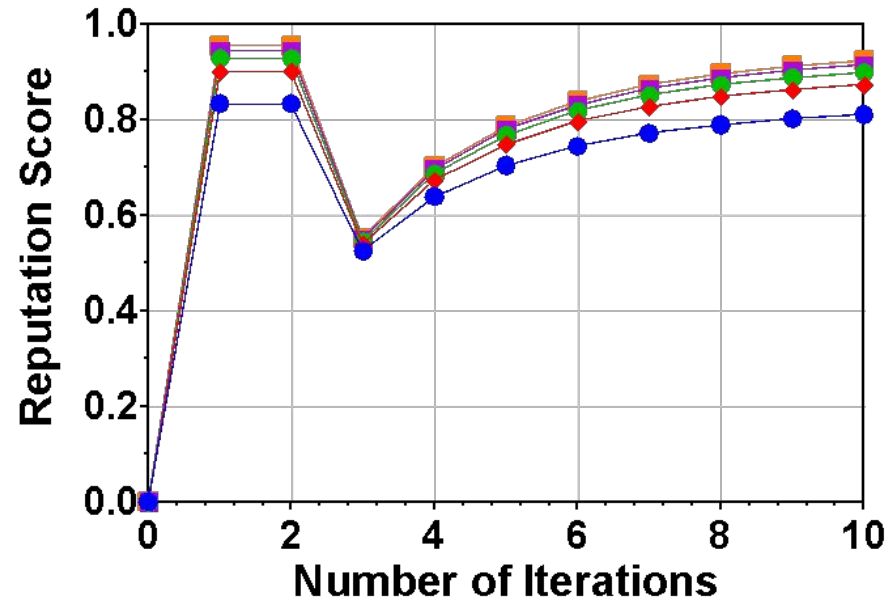


$$\begin{cases} b_i^t = \frac{\kappa P_i^t}{\kappa P_i^t + \eta N_i^t + W} \\ d_i^t = \frac{\eta N_i^t}{\kappa P_i^t + \eta N_i^t + W} \\ u_i^t = \frac{W}{\kappa P_i^t + \eta N_i^t + W} \end{cases}$$

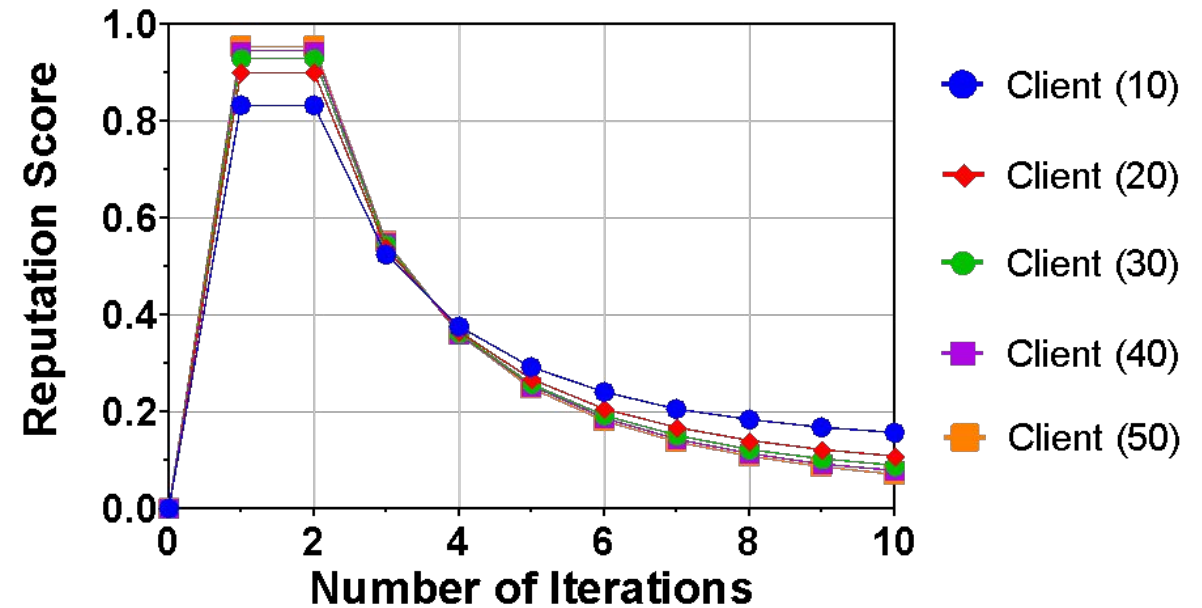
Subject logic model



$$\theta_{j,t} = \exp(-c(t - j)) \longrightarrow \tilde{R}_i^t = \frac{\sum_{j=\tilde{s}}^t \theta_{j,t} R_i^j}{\sum_{j=\tilde{s}}^t \theta_{j,t}}$$



(a) single attack

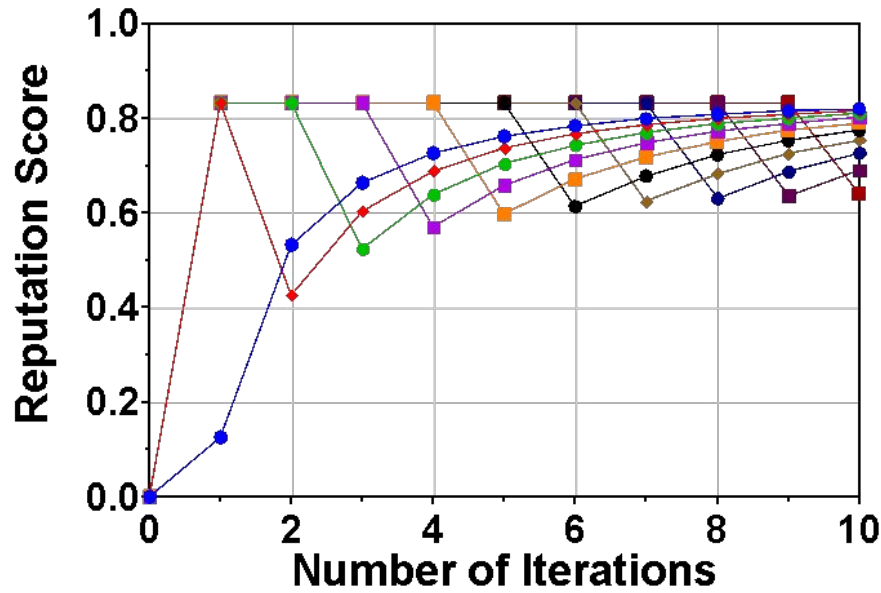


(b) continuous attack

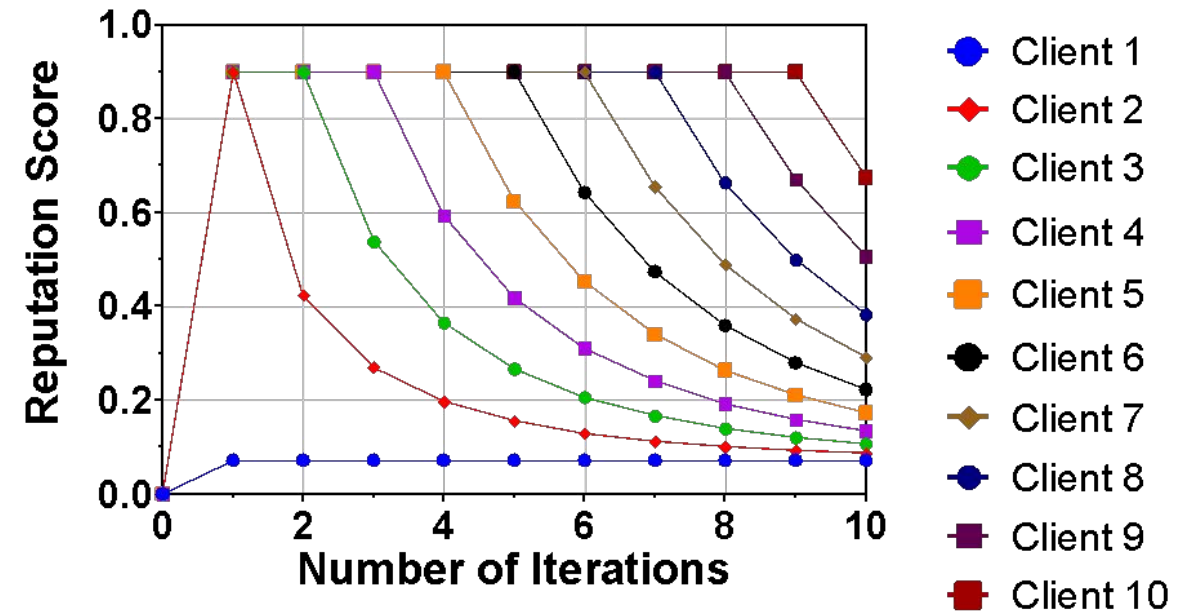
The decay of reputation score in Client (X) with X model parameters when they

(a) attack once at 3rd iteration

(b) attack continuously at and after the 3rd iteration

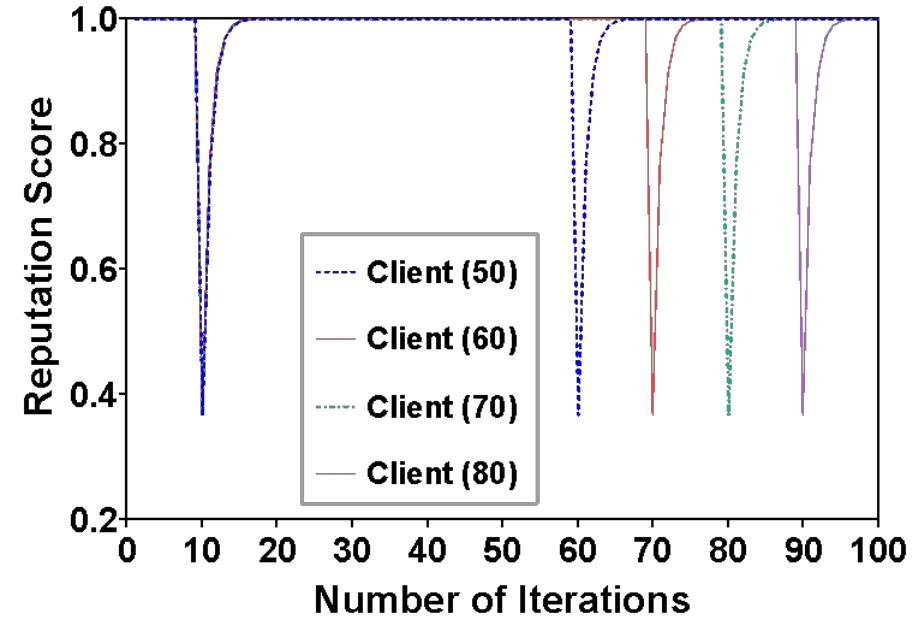
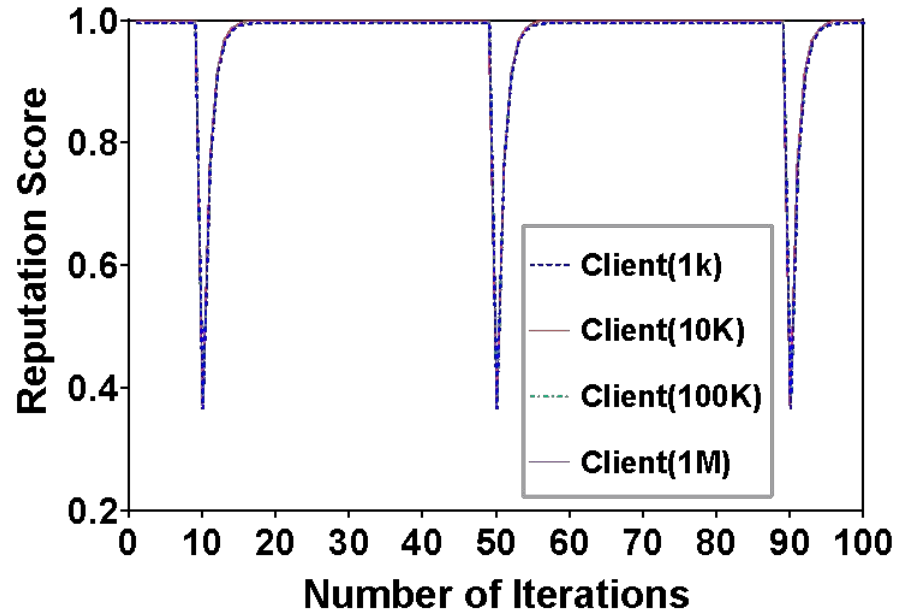


(a) single attack



(b) continuous attack

- The decay of reputation score in Client X with same model parameters when they
 - attack once at X iteration and
 - attack continuously after starting to attack at X iteration



- The decay of reputation score in
 - (a) client with X model parameters when they attack at 10, 50 and 90 iteration;
 - (b) client X with 1M parameters when they attack at 10 and $10 + X$ iteration

Theorem shows the convergence is guaranteed

Theorem

Under Assumptions, $\exists \epsilon > 0$ that:

$$\sqrt{\frac{d \log(1 + \hat{Q}MLD)}{M(1 - p)}} + C \frac{\mathcal{G}_w}{\sqrt{\hat{Q}}} + p \leq \frac{1}{2} - \epsilon \quad (1)$$

After t rounds, Our Algorithm converges with probability at least

$$1 - \xi \in \left[1 - \frac{4d}{(1 + \hat{Q}MLv)^d}, 1 \right) \text{ as}$$

$$\|w^t - w^*\|_2 \leq (1 - Lr)^t \|w^0 - w^*\|_2 + \frac{\sqrt{N}}{L} \Delta_1 + \frac{1}{L} \Delta_2 \quad (2)$$

The Corollary establishes the converge rate and error rate

Corollary

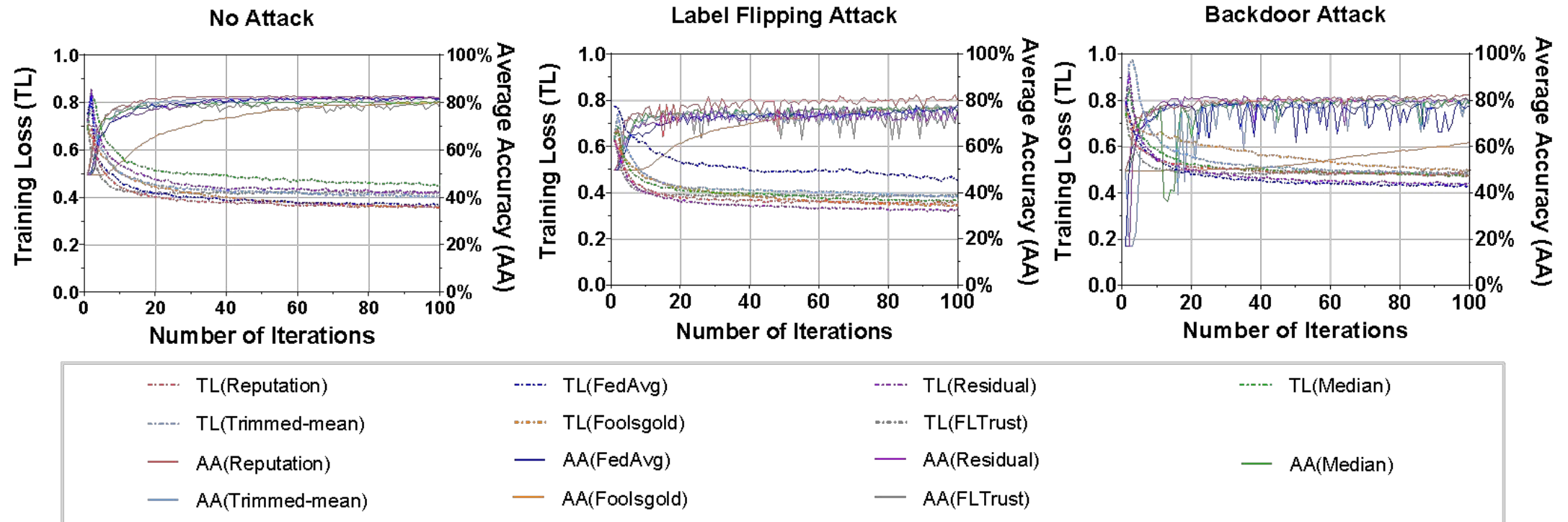
Continuing with Theorem 1, when the iterations satisfy

$t \geq \frac{1}{Lr} \log \left(\frac{L}{\sqrt{N}\Delta_1 + \Delta_2} \|w^0 - w^*\|_2 \right), \exists \xi \in \left(0, \frac{4d}{(1 + \hat{Q}MLv)^d} \right],$ we have:

$$\mathbb{P} \left(\|w^t - w^*\|_2 \leq \frac{2\sqrt{N}}{L} \Delta_1 + \frac{2}{L} \Delta_2 \right) \geq 1 - \xi$$

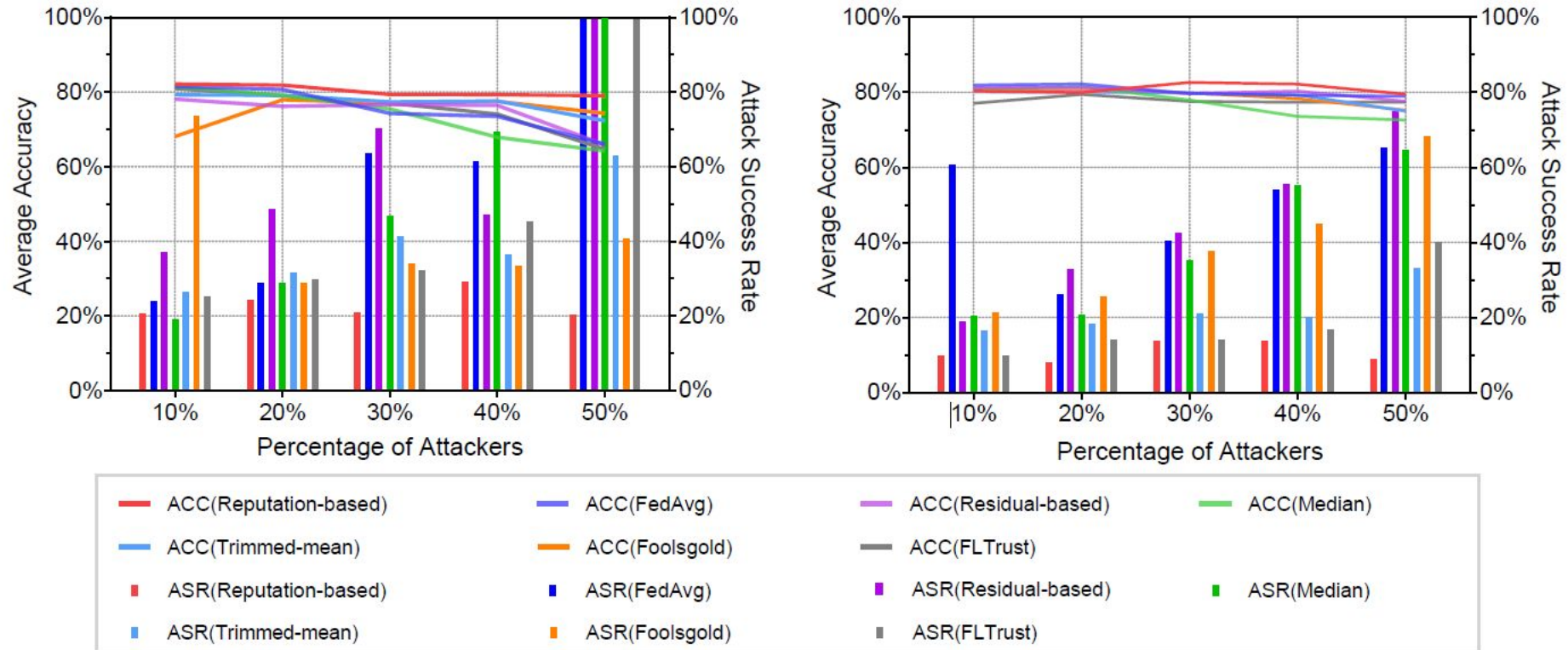
- The convergence is guaranteed in bounded time
- The trade-off between converge rate and error rate
- Guidance for hyper-parameters tuning

Convergence and Accuracy



- Converges 1.6× to at least 4.2× faster than all competing state-of-the-art methods
- Provides the same or better accuracy than competing methods

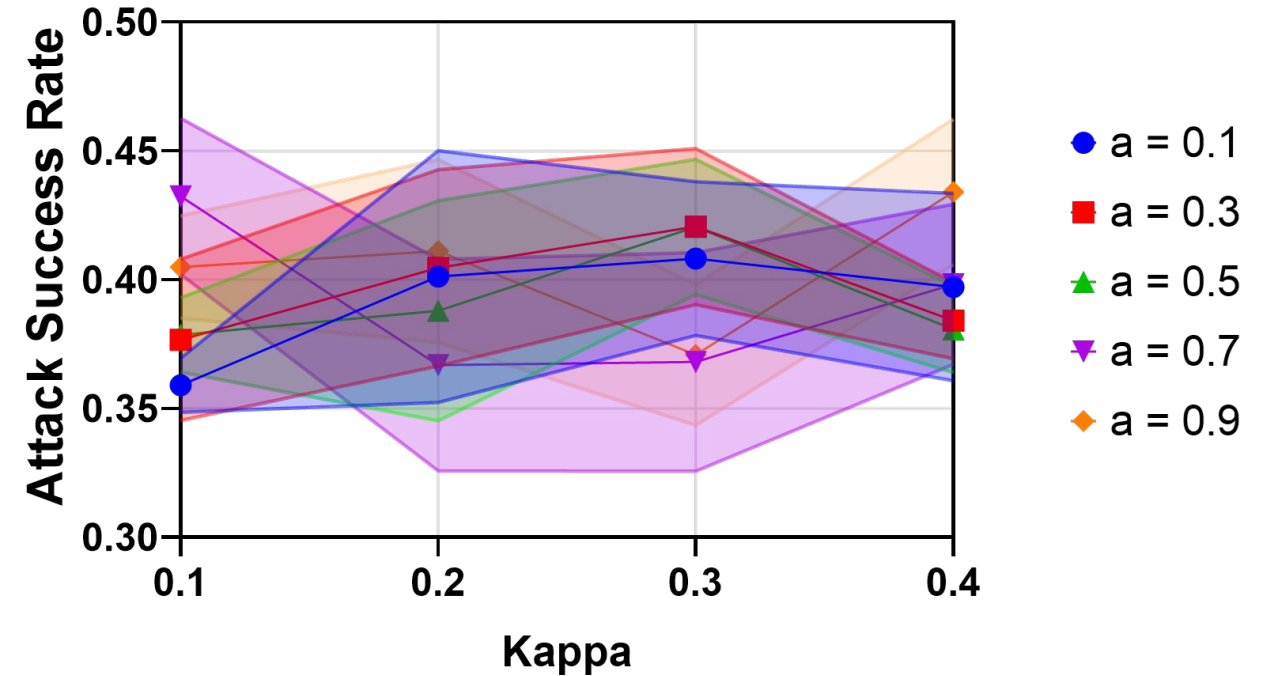
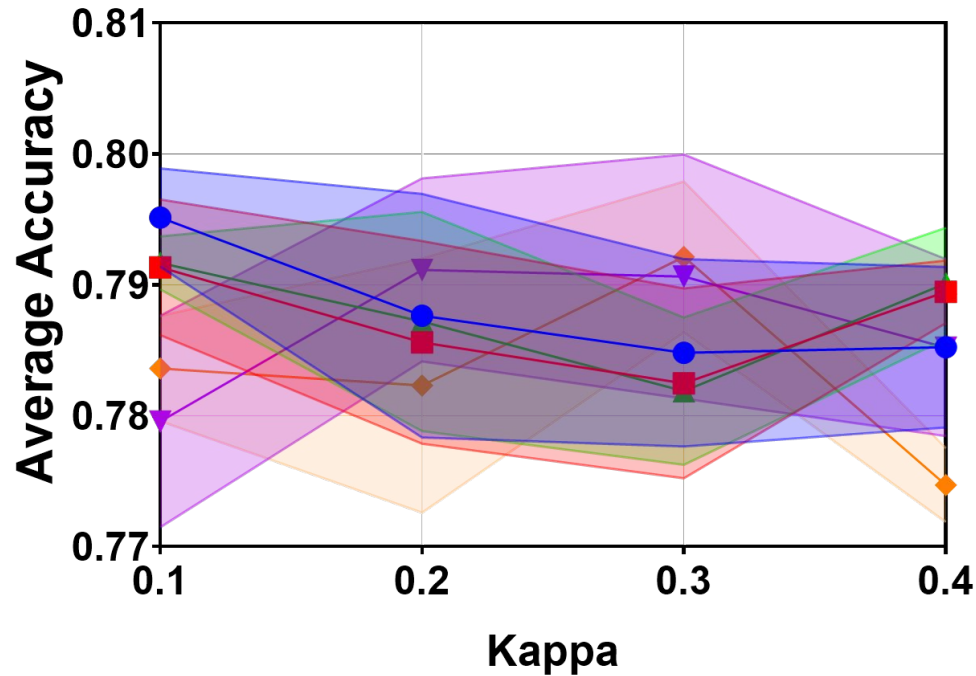
Attack Success rate



Varying percentage of attackers from 10% to 50%

- Yields the **lowest ASR** compared to all other methods, with the average ASR of them being at least 72% higher than ours

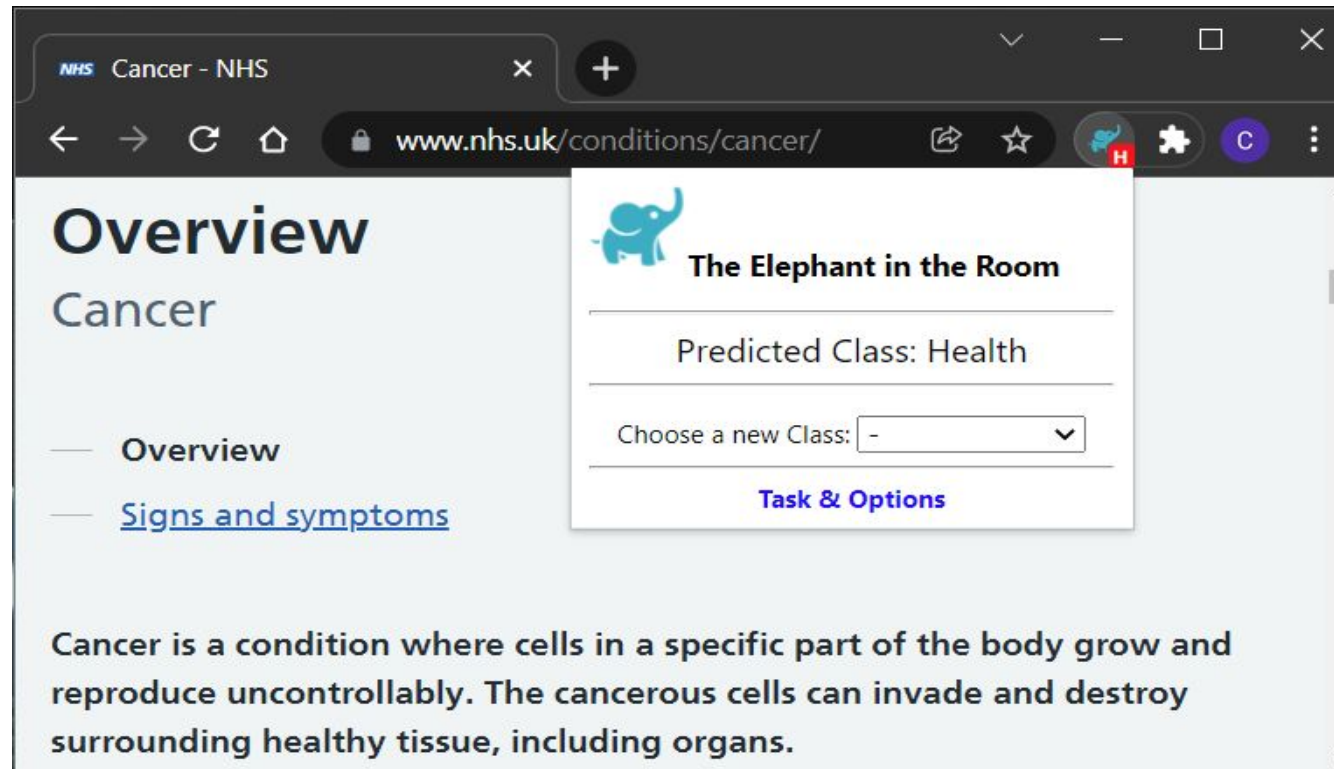
Hyper-parameters Searching



The result demonstrates that our method is

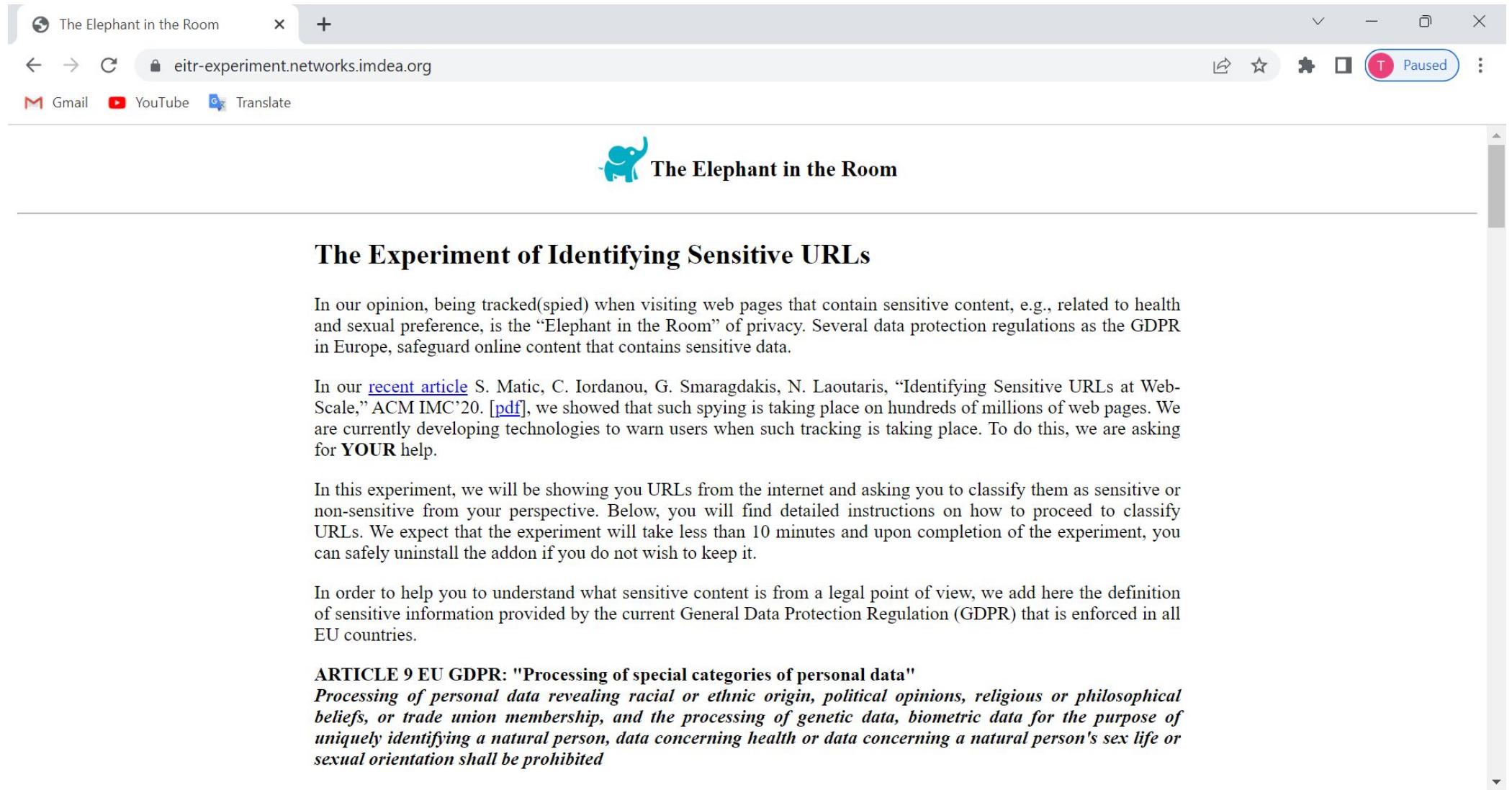
- very stable and efficient in terms of hyper-parameter selection
- and it achieves a high degree of precision
- the result is compatible with the theoretical analysis

- A research prototype to evaluate the robustness of our algorithm in a simple real-world setting with real users



<https://eit-experiment.networks.imdea.org/>





The Elephant in the Room

The Experiment of Identifying Sensitive URLs

In our opinion, being tracked(spied) when visiting web pages that contain sensitive content, e.g., related to health and sexual preference, is the “Elephant in the Room” of privacy. Several data protection regulations as the GDPR in Europe, safeguard online content that contains sensitive data.

In our [recent article](#) S. Matic, C. Iordanou, G. Smaragdakis, N. Laoutaris, “Identifying Sensitive URLs at Web-Scale,” ACM IMC’20. [[pdf](#)], we showed that such spying is taking place on hundreds of millions of web pages. We are currently developing technologies to warn users when such tracking is taking place. To do this, we are asking for **YOUR** help.

In this experiment, we will be showing you URLs from the internet and asking you to classify them as sensitive or non-sensitive from your perspective. Below, you will find detailed instructions on how to proceed to classify URLs. We expect that the experiment will take less than 10 minutes and upon completion of the experiment, you can safely uninstall the addon if you do not wish to keep it.

In order to help you to understand what sensitive content is from a legal point of view, we add here the definition of sensitive information provided by the current General Data Protection Regulation (GDPR) that is enforced in all EU countries.


ARTICLE 9 EU GDPR: "Processing of special categories of personal data"
Processing of personal data revealing racial or ethnic origin, political opinions, religious or philosophical beliefs, or trade union membership, and the processing of genetic data, biometric data for the purpose of uniquely identifying a natural person, data concerning health or data concerning a natural person's sex life or sexual orientation shall be prohibited

Real-user Experiment

New Tab x The Elephant in the Room x +

The Elephant in the Room | chrome-extension://bgneigpdlakopojkadjpkajcipcjcbpc/src/internal_pages/experiment.html

Gmail YouTube Translate

 **The Elephant in the Room**

Your email:

Your Task

URL	Visited
https://www.bbc.co.uk/news/world-asia-china-51466362	true
https://www.breitbart.com/national-security/2020/02/11/famed-dissident-china-welding-people-shut-in-their-homes-to-fight-coronavirus/	false
https://www.scientiststudy.com/2020/02/research-wuhan-coronavirus-also-spreads.html	true
http://www.munising.org/	false
https://foreignpolicy.com/2020/01/31/coronavirus-china-trump-united-states-public-health-emergency-response/	false
https://www.theguardian.com/world/2020/feb/11/coronavirus-expert-warns-infection-could-reach-60-of-worlds-population	false
https://yalliban.com/2020/02/10/it-could-take-years-to-make-a-vaccine-for-the-wuhan-coronavirus/	false
https://www.statista.com/chart/20785/coronavirus-recoveries/	false
https://www.technologyreview.com/2020/02/11/china-has-launched-an-app-so-people-can-check-their-risk-of-catching-the-coronavirus/	false
http://foxcitiesregionalpartnership.com/	false
https://www.nytimes.com/2020/02/11/world/asia/coronavirus-china.html	false
https://www.nytimes.com/2020/02/11/briefing/coronavirus-new-hampshire-t-mobile.html	false
http://www.exit170.ca/	false
http://outbreaknewstoday.com/caribbean-princess-outbreak-case-count-tops-350-causative-agent-still-not-known-41474/	false
https://dcdirtylaundry.com/rigged-china-changes-the-definition-of-infected-to-ignore-coronavirus-patients-who-test-positive-but-show-no-symptoms/	false
https://str.sg/1jyV/	false
https://www.theguardian.com/animals-farmed/2020/feb/04/animals-farmed-live-exports-risk-of-disease-china-goes-big-on-pork-and-eu-meat-tax	false
https://arynews.tv/en/health-safety-oakistanis-china-coronavirus/	false
http://bbc.in/37blve	false
https://www.grahamcluley.com/coronavirus-phishing-attack-cdc/	false

Extensions | The Elephant in the Room | Research: Wuhan coronavirus also spreads | scientiststudy.com/2020/02/research-wuhan-coronavirus-also-spreads.html

1k Shares

2019 Novel Coronavirus

This illustration, created by the Centers for Disease Control and Prevention (CDC), reveals the ultrastructural morphology exposed by the new coronavirus (2019-nCoV). Note the peaks that adorn the outer surface of the virus, which give the appearance of a crown surrounding the virion when viewed under an electron microscope. The protein particles E, S, M and HE, also located on the outer surface of the particle, have all been indicated. Credits: Wikipedia

"It is important to note that 2019-nCoV has been reported elsewhere in the feces of patients with atypical abdominal symptoms, similar to SARS which has also been excreted in the urine, suggesting a route of faecal transmission that is highly transmissible", said William Keevil at the UK Science Media Center, professor of environmental health at the University of Southampton.

Computed tomography images on day 5 after symptom onset

Computed tomography images after treatment on day 19 after symptom onset

Tomographic images of the chest of a 52-year-old patient infected with the new coronavirus (2019-nCoV). A: chest tomography images obtained on January 7, 2020, showing opacity in both lungs, on the 5th day after the onset of symptoms. B: images taken on January 21, 2020, showing the absorption of bilateral opacity after extracorporeal membrane oxygenation treatment from January 7 to 12, in the intensive care unit. Credits: Zhiyong Peng / Zhongnan Hospital of Wuhan University

This possibility is not really surprising for scientists, given that the new virus belongs to the same family as SARS. In 2003, faecal transmission of SARS contributed to the infection of hundreds of people in the Amoy Gardens subdivision in Hong Kong. A plume of warm air from the bathroom had contaminated several apartments and had been blown by wind to

The Elephant in the Room

Predicted Class: Health

Choose a new Class: -

Task & C Non-sensitive

Religion

Health

Politics

Ethnicity

Sexual

PHYSICS

Planet and Environment

Plants & Animals

Robotics

Space & Astrophysics

Technology

Unusual

Main Tags

Archeology and Paleontology

Biology & Health

Chemistry

Computer Science

Computing

Electronics

Energy

Materials

Mathematics

Medical Science

Nanotechnology

News

Other Science

Physics

Planet and Environment

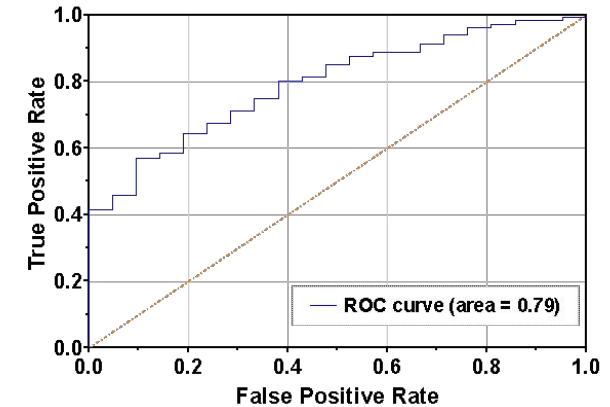
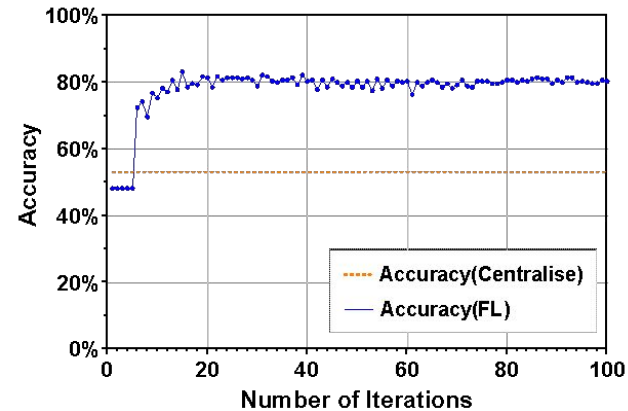
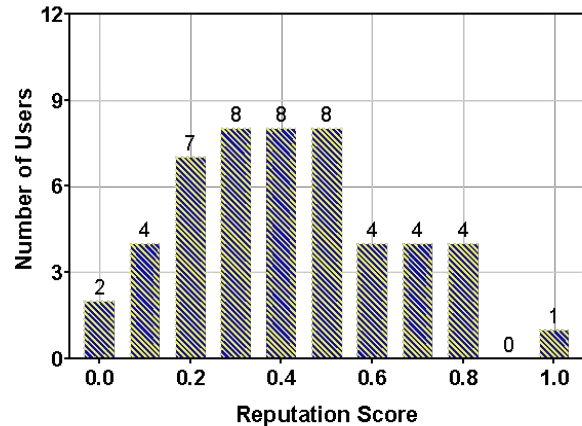
Plants & Animals

Robotics

Space & Astrophysics

Technology

Unusual



- Result with real users

- The divergence of the user's interpretation of the sensitive information
- Our method converges as rapidly as in simulation and achieves an average accuracy of 80%.

Our FL-based solution can quickly learn to classify sites about COVID even with inconsistent input provided by real users

Conclusion



Employ FL for sensitive content classification



Implement EITR browser extension



Design a robust FL aggregation with reputation

- Theoretical analysis
- Experimental evaluation



Implementation available online:
<https://github.com/FRM-Sec/FRM>



Questions:
Email: tianyue.chu@imdea.org

THANK YOU!

Attack detection scheme

Algorithm 1: Rescale(w)

Input : $\{w_{i,n}^t\} \leftarrow$ Local Model parameters in round t
Output: $\{w_{i,n}^t\}$ with range of value less than ϖ

```

1 for  $n \leftarrow 1$  to  $N$  do
2   // Determine the maximum range
3    $Rm = \max w_{i,n}^t - \min w_{i,n}^t = w_{i,n}^{t,(Max)} - w_{i,n}^{t,(Min)}$ 
4   while  $Rm > \varpi$  do
5     // Rescale range based on standard deviation.
6      $w_{i,n}^{t,(Max)} := w_{i,n}^{t,(Max)} - \sigma(w_{i,n}^t);$ 
7      $w_{i,n}^{t,(Min)} := w_{i,n}^{t,(Min)} + \sigma(w_{i,n}^t);$ 
8     // Updated  $Rm$ .
9      $Rm = \max w_{i,n}^t - \min w_{i,n}^t = w_{i,n}^{t,(Max)} - w_{i,n}^{t,(Min)}$ 
10    end while
11 end for
```

Reputation:

Time decay function

Algorithm 2: Aggregation Algorithm

Server :
Input : $w^0 \leftarrow$ Pretrained Model
 $\kappa, \eta, a, W, c, s \leftarrow$ Reputation parameters
Output: Global model M_{global} with w^T

```

1 for Iteration  $t \leftarrow 1$  to  $T$  do
2   // Broadcast global model to clients
3   send( $w^{t-1}$ );
4   // Wait until all updates arrive
5   receive( $w^t$ );
6   // Rescale parameters by Algorithm 1
7    $w^t \leftarrow$  Rescale( $w^t$ );
8   for  $n \leftarrow 1$  to  $N$  do
9     for  $i \leftarrow 1$  to  $M$  do
10      // Compute parameter confidence
11       $s_{i,n}^t = Eq_{\square}(w_{i,n}^t);$ 
12      // Rectify abnormal parameters
13       $w_{i,n}^t := Eq_{\boxtimes}(s_{i,n}^t, \delta);$ 
14      record ( $P_i^t, N_i^t$ );
15    end for
16  end for
17  for  $i \leftarrow 1$  to  $M$  do
18    // Calculate reputation score
19     $\tilde{R}_i^t = Eq_{\boxdot}(P_i^t, N_i^t, \kappa, \eta, a, W, c, s);$ 
20  end for
21  // Normalisation
22   $\bar{R}^t \leftarrow$  Norm( $\tilde{R}^t$ );
23  for  $n \leftarrow 1$  to  $N$  do
24    // Update the parameters
25     $w_n^t := \sum_{i=1}^M \frac{\bar{R}_i^t}{\sum_{i=1}^M \bar{R}_i^t} w_{i,n}^t;$ 
26  end for
27  // Obtain parameters for global model
28   $w^t := [w_1^t, \dots, w_n^t];$ 
29 end for
```

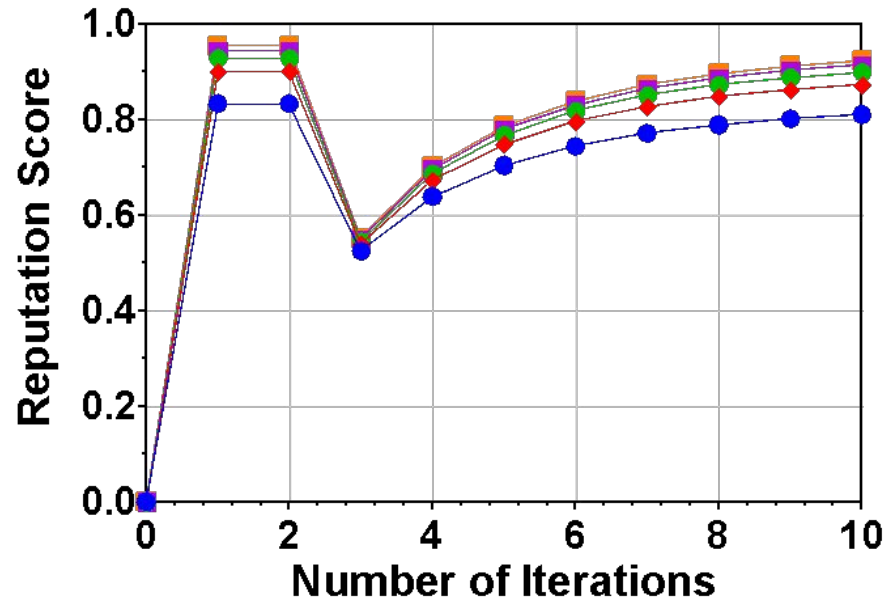
Client :

```

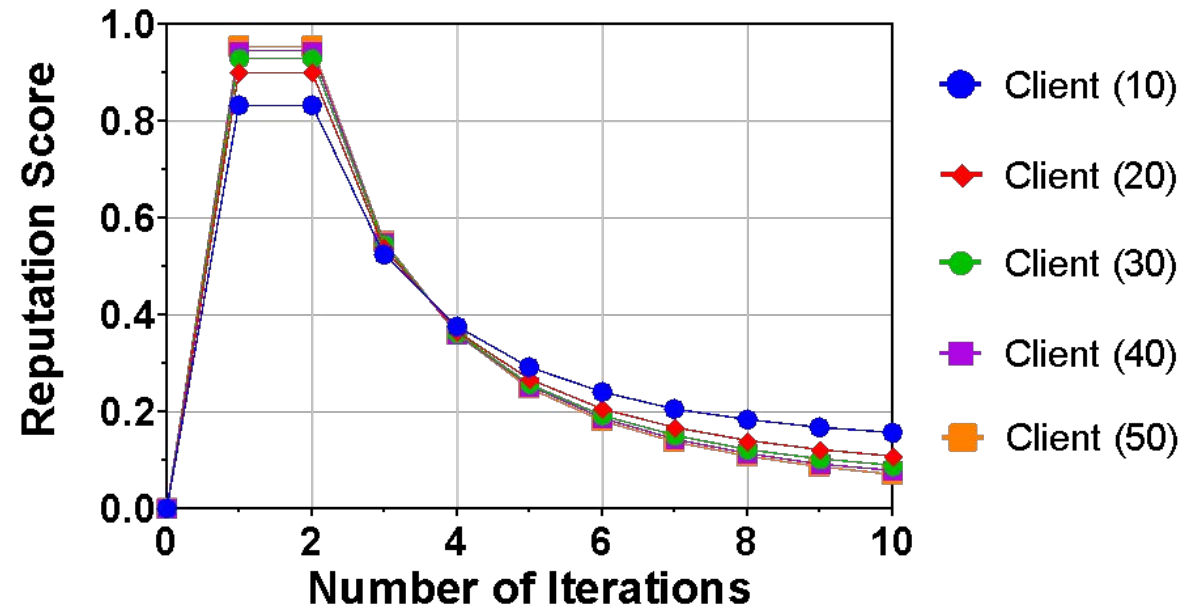
1 for Client  $i \leftarrow 1$  to  $M$  do in parallel
2   receive( $w^{t-1}$ );
3   // Train local model
4    $w_i^t \leftarrow w_i^{t-1} - r \frac{\partial \ell_i(w_i^{t-1})}{\partial w};$ 
5   send( $w_i^t$ );
6 end forpar
```

Repeated Median

IRLS scheme



(a) single attack



(b) continuous attack

The decay of reputation score in Client (X) with X model parameters when they

(a) attack once at 3rd iteration

(b) attack continuously at and after the 3rd iteration

Centralized vs Distributed

Pros: centralised training performs better for some tasks, but for the task of detecting sensitive URLs to protect personal online privacy, a distributed classifier is a preferable solution.

Cons:

- Privacy: In a centralized manner, even with semi-supervised learning, manual labelling of certain training data is still required. However, because users are labelling sensitive data, their privacy will be harmed if the server has access to their labelling data for centralised training. Therefore, FL is a natural privacy-preserving method for conducting distributed collaborative model training across clients that do not disclose their local data.
- Efficiency: for the centralized training, the server has to update the dataset for learn new content, maybe by paying people to do it. but for FL, it can voluntarily and continuously learn from real-time web data gathered by users, and will represent the user's interpretation of sensitive content.

The number of attackers

we acknowledge that constraining the number of attackers is not realistic, but

- as with other works [7][15][19][20][21] on this topic, we follow the same assumption that the percentage of attackers is lower than the percentage of benign users in the system for our evaluation.
- In Sec. 3C Theoretical Guarantees, we also demonstrate how the number of attackers influences the performance of our algorithm.
- Additionally, we did not restrict the number of attackers in the real-user experiment, the results show that our method still performs well.

In a label flipping attack, the attacker flips the labels of training samples to a targeted label and trains the model accordingly.

In our case, the attacker changes the label of “Health” to “Non-sensitive”.

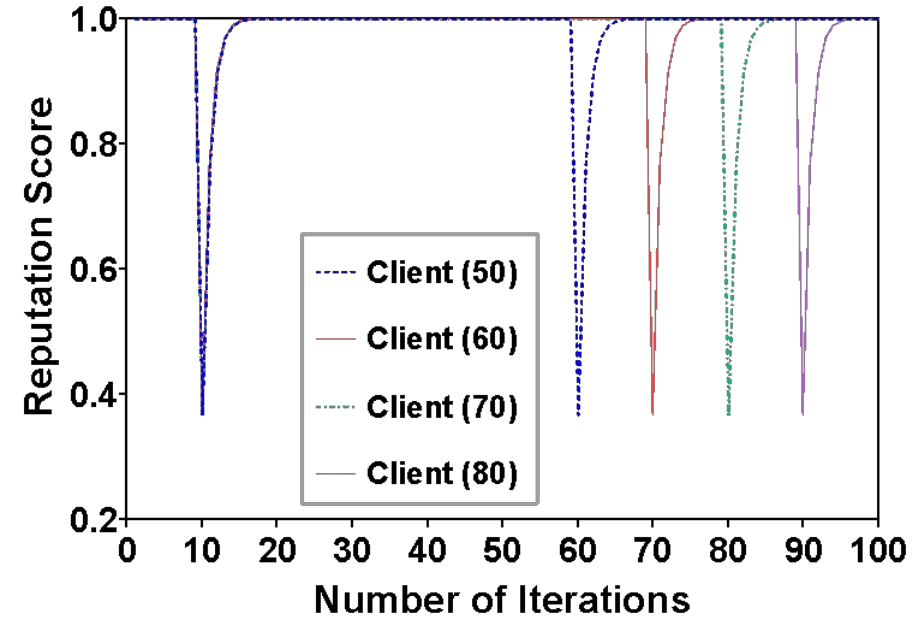
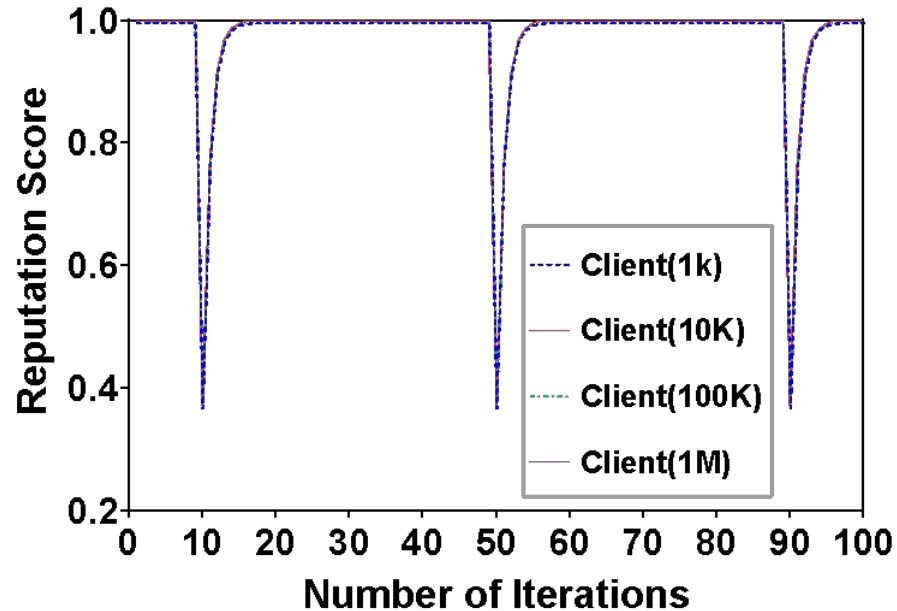
In a backdoor attack, attackers inject a designed pattern into their local data and train these manipulated data with clean data, in order to develop a local model that learns to recognise such pattern.

We realise backdoor attacks inserting the top 10 frequent words with their frequencies for the “Health” category. Therein the backdoor targets are the labels “non-sensitive”. A successful backdoor attack would acquire a global model that predicts the backdoor target label for data along with specific pattern

Consider two attacks:

We focus on two common attack strategies for sensitive context classification, namely, (i) label flipping attack [24] and (ii) backdoor attack [12].

Comparing to other attacks, for example model poisoning attack [29], [63], these two data poisoning attacks are more likely to be carried out by real users in the real world via our browser extension described in Section V, since polluting data is easier than manipulating model updates using the browser extension. Note that privacy attacks including membership inference attack [65] and property inference attack [64], are out of the scope of this paper, but form part of our ongoing and future work



These figures show that even if attackers spread our poisoning over multiple iterations and then try to recover their reputation score by acting benignly, our detection scheme can still identify them.

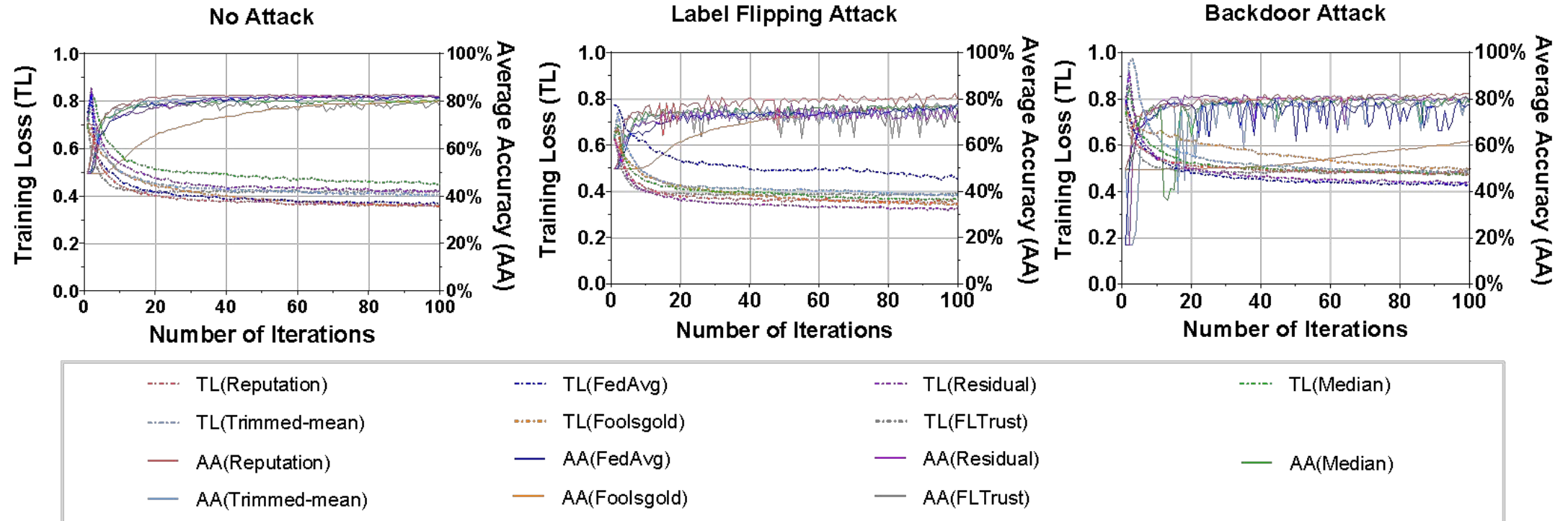
This is because our attack detection and reputation schemes work in sequence. The attack detection scheme detects malicious updates without considering any reputation scores and rectifies them to mitigate damage.

Then, the reputation scheme modifies the reputation scores based on the detection results. Also, attackers that employ a higher number of model parameters suffer a slightly higher reduction of reputation, which is consistent with Corollary

- **Objectives**

- Compare our method with other SOA
- Use a text based real-world dataset of sensitive categories
- Show our experimental result are consistent with theoretical analysis
- Under different scenarios: no attack, under attack

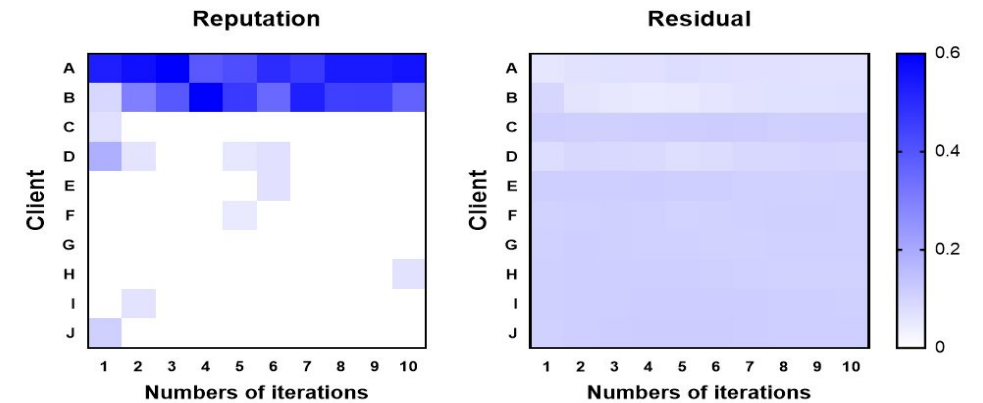
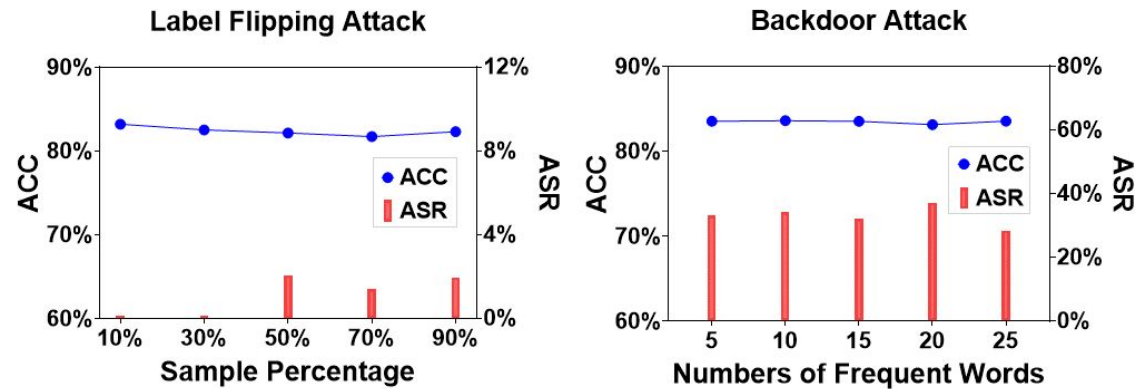
Convergence and Accuracy



why converge faster:

This is due to the fact that in our algorithm we give higher weights to the clients with high-quality updates, as illustrated in Figure~\ref{fig:comparision}[nextt slides], causing the model to converge rapidly and retain consistent accuracy. In addition, even under the two different attacks, our method:

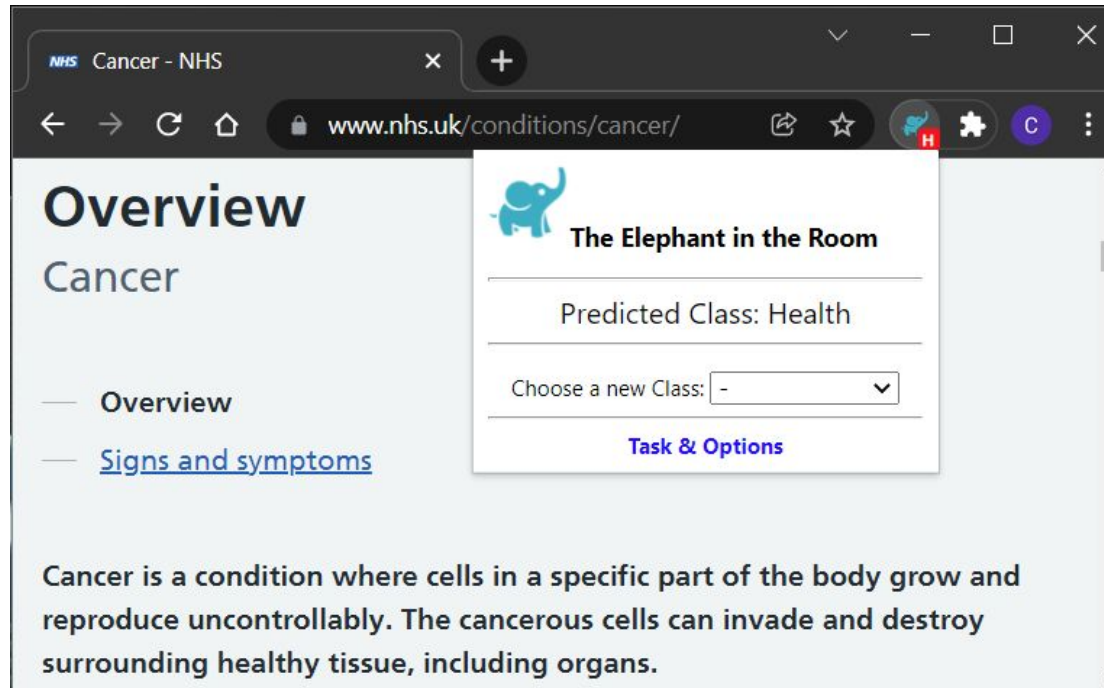
Other Evaluations



- The percentage of the poisoned sample is increased, it leads to the decrease of the accuracy of the model and to a slight increase of ASR

- our reputation method the aggregation weights of malicious clients, which are their reputations, are rectified near to 0, outperform the residual-based method

- A research prototype to evaluate the robustness of our algorithm in a simple real-world setting with real users



The back-end server is responsible to distribute the initial classification model and the consequently updated model(s) to the clients, and receive new annotations from the different clients of the system.

Fake account

We agree that in the real world, attackers can create a large number of new accounts and launch a single attack to damage the reputation mechanism. One solution to multiple account creation is to attach participation via certificates to the real-world identities of users. The method of authenticating and binding identities will be implemented in future work.

1. Commission, E. Data protection in the EU, The General Data Protection Regulation (GDPR); Regulation (EU) 2016/679. (2018), <https://ec.europa.eu/info/law/law-topic/data-protection>
2. Matic, Srdjan, Costas Iordanou, Georgios Smaragdakis, and Nikolaos Laoutaris. "Identifying sensitive urls at web-scale." In *Proceedings of the ACM Internet Measurement Conference*, pp. 619-633. 2020.
3. McMahan, Brendan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. "Communication-efficient learning of deep networks from decentralized data." In *Artificial intelligence and statistics*, pp. 1273-1282. PMLR, 2017.