# RepuNet: A Reputation System for Mitigating Malicious Clients in DFL

Isaac Marroqui Penalva, Enrique Tomás Martínez Beltrán, Manuel Gil Pérez, Alberto Huertas Celdrán

*Abstract*—Decentralized Federated Learning (DFL) enables nodes to collaboratively train models without a central server, introducing new vulnerabilities since each node independently selects peers for model aggregation. Malicious nodes may exploit this autonomy by sending corrupted models (model poisoning), delaying model submissions (delay attack), or flooding the network with excessive messages, negatively affecting system performance. Existing solutions often depend on rigid configurations or additional infrastructures such as blockchain, leading to computational overhead, scalability issues, or limited adaptability. To overcome these limitations, this paper proposes RepuNet, a decentralized reputation system that categorizes threats in DFL and dynamically evaluates node behavior using metrics like model similarity, parameter changes, message latency, and communication volume. Nodes' influence in model aggregation is adjusted based on their reputation scores. RepuNet was integrated into the Nebula DFL platform and experimentally evaluated with MNIST and CIFAR-10 datasets under non-IID distributions, using federations of up to 25 nodes in both fully connected and random topologies. Different attack intensities, frequencies, and activation intervals were tested. Results demonstrated that RepuNet effectively detects and mitigates malicious behavior, achieving F1 scores above 95% for MNIST scenarios and approximately 76% for CIFAR-10 cases. These outcomes highlight RepuNet's adaptability, robustness, and practical potential for mitigating threats in decentralized federated learning environments.

*Index Terms*—Federated Learning, reputation systems, malicious node detection, anomaly detection, reputation-based aggregation, attack mitigation, adversarial behavior in DFL

## I. INTRODUCTION

FEDERATED Learning (FL) [1] has emerged as a key solution in the era of Artificial Intelligence (AI) to train models collaboratively without compromising data privacy. Traditionally, centralized FL systems have dominated this field, but their reliance on a single server introduces security vulnerabilities, as this central point can become a target for malicious attacks. To mitigate these risks, DFL has gained relevance [2], eliminating the need for a central server, removing the single point of failure, and enhancing system robustness. This is crucial in sensitive domains such as healthcare or finance. However, DFL introduces unique security challenges due to the lack of central supervision during training. In the first phase, nodes train local models independently, which makes them vulnerable to *model poisoning* attacks, where adversaries inject distorted updates to degrade global convergence [3]. In the second phase, nodes exchange parameters with neighbors, a stage exposed to *communication*

Isaac Marroqui Penalva, Enrique Tomás Martínez Beltrán, Manuel Gil Perez, and Alberto Huertas Celdrán are with the University of Murcia, Spain.

*delay* attacks. These can cause desynchronization and unstable learning dynamics, as outdated models propagate through the network [4]. Finally, during aggregation, malicious participants can launch *flooding or DoS* attacks by overwhelming the network with excessive messages, disrupting communication and slowing convergence [5]. These threats call for robust mechanisms to mitigate malicious behavior. There are various approaches to mitigating these attacks in DFL, such as the use of blockchain [6], cluster-based aggregations [7], suspicious model filtering [8], and cryptographic defense techniques [9]. However, these approaches have certain limitations. The effectiveness of cluster-based systems depends on proper configuration, making them vulnerable to manipulation if a group contains a high proportion of attackers. Furthermore, nodes may coordinate to influence reputation within their own groups, avoiding detection. If hierarchical aggregation is later performed, additional computational overhead is introduced, which may not be well-suited for decentralized environments. In the case of blockchain usage, while it provides traceability and resistance to manipulation, it introduces high computational and storage costs, along with potential scalability issues as the number of nodes in the network grows. Model filtering can be too restrictive or generate false positives, affecting learning diversity. On the other hand, cryptographic solutions, although effective in some contexts, add computational overhead, making them unfeasible in resource-limited scenarios. To address these limitations, this paper proposes RepuNet, a decentralized reputation system specifically designed for DFL. RepuNet serves as a dynamic mechanism for detecting and mitigating attacks by continuously evaluating the behavior of nodes over time. It penalizes those that exhibit anomalous patterns—such as communication delays, transmission of inconsistent parameters, or excessive traffic generation. By integrating reputation scores into each training round, the system can exclude malicious node contributions without significantly affecting network stability, thereby enhancing both security and global model convergence.

**The main contributions of this work are as follows:**

- **Design of a threat-aware and adaptive reputation system for DFL.** This work defines a threat model that identifies and categorizes key attacks in DFL—such as model poisoning, intentional communication delays, and flooding—and uses it as the basis for RepuNet, a decentralized reputation mechanism. The system evaluates the behavior of neighboring nodes using dynamic metrics (e.g., model similarity, parameter change ratio, latency, and message volume), enabling adaptive weighting and

detection of malicious patterns.

- **Progressive penalization and controlled exclusion mechanism.** RepuNet gradually reduces the influence of nodes with low reputation in the aggregation process without disconnecting them. This supports the detection of persistent adversarial behavior while allowing reintegration if behavior improves in later rounds.
- **Deployment of RepuNet on the Nebula platform under realistic DFL conditions.** The system has been implemented and evaluated using the Nebula platform [10], which supports configurable topologies, attack scenarios, and real-time metric visualization. Experiments were performed in emulated environments with up to 25 nodes, fully connected and random topologies, and three types of attacks. The evaluation used MNIST and CIFAR-10 datasets under non-IID distributions, demonstrating significant improvements in anomaly detection, robust behavior under different levels of adversarial activity, and stable model performance across rounds.

The structure of this article is as follows: Section II reviews related work on attack mitigation in DFL. Section III introduces the threat model and explains the rationale for focusing on specific attacks. Section IV details the design and operation of the proposed reputation system, RepuNet. Section V describes the experimental setup and attack scenarios. Section VI presents and analyzes the results. Finally, Section VII summarizes the main contributions and outlines future work.

## II. RELATED WORK

This section reviews reputation systems in FL and defense mechanisms against malicious behavior in decentralized networks. It covers blockchain-based FL, hierarchical aggregation, and attack detection strategies—such as data and model poisoning or DDoS—as well as both centralized and decentralized approaches that promote honest participation and penalize adversarial nodes.

### A. Reputation Systems in Federated Learning

A centralized approach [11] updates reputation scores using behavioral history and a beta distribution. Nodes must exceed an initial reputation threshold to participate, ensuring reliability from the outset. Decentralized systems offer transparency and resilience. [12] stores contributions on blockchain; reputation uses direct and indirect feedback. [13] applies blockchain reputation in IoT-FL to promote cooperation. Cooperative incentives have emerged. [14] proposes a co-utility rule system to reward compliant behavior without central coordination. [15] explores hierarchical aggregation using accuracy and participation to select models. [16] proposes FGFL, a blockchain system using reputation and contribution metrics to adjust incentives. Together, these approaches reflect a shift from static, centralized schemes toward adaptive, incentive-aware, and decentralized reputation frameworks suitable for heterogeneous federated learning environments.

### B. Attack Mitigation Strategies

DFL is vulnerable to various attacks that compromise the integrity, efficiency, or convergence of the learning process. These include model manipulation, flooding, and message delays. Several studies address model poisoning and data manipulation. [17] studies poisoning in intrusion detection and proposes model filtering. In smart grid contexts, [3] employs deep learning to detect falsified data without compromising privacy. To mitigate network saturation, [5] proposes a decentralized FL-based architecture for anomaly detection in industrial IoT. Likewise, [18] demonstrates how distributed traffic analysis in IoT-Edge environments can identify malicious patterns while preserving data locality. Delay attacks—where nodes transmit outdated updates to desynchronize the network—are addressed in [4], which uses hierarchical FL to reduce the influence of slower participants. [19] surveys anomaly detection, cryptography, and reputation mechanisms. Blockchain has also been proposed to enhance trust and traceability. For example, [20] uses sharding and layered validation to secure aggregation, while [21] targets passive threats like "lazybone" nodes through multidimensional evaluation of contribution quality. These strategies highlight the need for hybrid defenses combining anomaly detection, contribution filtering, decentralized coordination, and adaptive reputation, to ensure resilience in diverse DFL scenarios.

### C. Comparison with Existing Approaches

RepuNet allows decentralized reputation where nodes assess neighbors without central authority or blockchain. Reputation evolves dynamically from local metrics such as model similarity, parameter variation, latency, and message volume, using a neutral initial value to avoid bias. Unlike static or irreversible schemes, RepuNet allows nodes to recover reputation through consistent behavior, balancing penalization and reintegration. Compared to blockchain-based systems, RepuNet avoids the heavy computational and communication overhead. It eliminates the need for immutable global records and relies on lightweight, adaptive updates performed locally at each node. Additionally, peer feedback can be optionally incorporated to strengthen individual assessments, especially in scenarios with partial observations, newly joined nodes, or sparse topologies. This combination enhances system resilience without compromising decentralization. As shown in Table I, RepuNet integrates more metrics—including communication anomalies—and avoids overhead from blockchain-based systems.

## III. THREAT MODEL

DFL faces several vulnerabilities due to the autonomy of participating nodes, which may act maliciously and compromise the quality, efficiency, or stability of the collaborative process.

The most relevant threats in DFL can be grouped into four categories:

- **Model-based attacks:** manipulation of local parameters (e.g., *model poisoning*, *backdoor attacks*).
- **Communication-based attacks:** such as *delayer attacks*, *flooding*, or DoS.

TABLE I
COMPARISON OF ANOMALY DETECTION CRITERIA, DETECTED ATTACKS, AND MITIGATION STRATEGIES IN FL SYSTEMS.

| Work (Year) | FL Type | Detection / Penalization Criteria | Feedback | Detected Attack(s) | Anomaly Detection | Action Taken |
|---|---|---|---|---|---|---|
| [13] (2019) | DFL (Blockchain) | Crowdsourced contribution history on blockchain | X | – | X | Trustworthy nodes rewarded; reputation affects participation |
| [12] (2019) | DFL (Blockchain) | Trust via contribution quality and peer reputation (direct/indirect) | ✓ | – | X | Adjusts reputation; low scores discourage future inclusion |
| [17] (2019) | IoT | Model deviation from expected behavior (filtering) | X | Poisoning | ✓ | Suspicious updates discarded before aggregation |
| [14] (2021) | DFL (Blockchain) | Violation of cooperation rules (co-utility model) | X | – | X | Rewards or penalties applied via rule compliance |
| [11] (2022) | CFL | Deviation in behavioral history (beta-based reputation) | X | – | X | Reputation score adjusted (centralized); no exclusion |
| [16] (2022) | DFL (Blockchain) | Quality and frequency of valid contributions | X | Lazy / low-effort participation | ✓ | Contributors receive dynamic rewards or are excluded if passive |
| [3] (2022) | Smart Grid | Anomalous data pattern detection using deep learning | X | FDIA (false data injection) | ✓ | Classifier-based filtering of malicious updates |
| [5] (2022) | Industrial IoT | Traffic anomalies detected via federated classifiers | X | DDoS | ✓ | Detection used to alert and reconfigure network paths |
| [18] (2023) | IoT Edge | Distributed traffic anomaly pattern detection in IoT | X | DDoS, poisoning | ✓ | Local anomaly response; contributes to distributed decision-making |
| [15] (2023) | Hierarchical CFL | Accuracy, consistency, and past participation history | X | – | X | Contribution filtered before aggregation based on reputation |
| [19] (2024) | DFL (Blockchain) | Overview of multiple metrics: gradients, consistency, etc. (survey) | X | Poisoning, backdoor, Sybil, free-riding | ✓ | Describes possible responses: filtering, exclusion, weight adjustment |
| [4] (2024) | FL (Hierarchical) | Model submission latency (ranking-based classification) | X | Delay | X | Low-ranked nodes deprioritized or excluded in hierarchy |
| [20] (2025) | DFL (Blockchain) | Distributed model validation through consensus (sharded blockchain) | X | Sybil, poisoning | ✓ | Invalid models rejected during consensus verification |
| [21] (2025) | CFL | Contribution scoring via loss, gradient norm, and activity level | ✓ | Lazybone | ✓ | Low-contributing nodes excluded from aggregation |
| **RepuNet** (2025) | DFL | Similarity, parameter change, arrival latency, message volume | ✓ | Poisoning, delay, flooding | ✓ | Dynamic penalization of reputation, exclusion from aggregation, and adaptive reintegration |

- **Data-based attacks:** noise injection, false data submission, or inference of private information.
- **Identity/participation attacks:** including *Sybil*, *free-riding*, or *lazybone* behaviors.

This work targets three threats due to their impact and feasibility of local detection:

- **Model poisoning:** injection of malicious updates to degrade the global model.
- **Delayer attack:** intentional delay in sending models to disrupt synchronization.
- **Flooding attack:** excessive message generation to saturate the network.

These attacks degrade model quality, synchronization, or communication, and can be detected using local metrics such as model similarity, latency, or message count. They are common in real-world deployments. Other threats (e.g., inference, Sybil, or passive participation) require identity validation or secure exchanges and are beyond the scope of this work. By focusing on locally observable threats, the proposed system remains practical and suitable for controlled experimental validation, as described in Section V.

## IV. REPUNET: AN ADAPTIVE REPUTATION SYSTEM TO MITIGATE ATTACKS IN DFL

This section presents RepuNet, a decentralized and adaptive reputation system specifically designed to detect and mitigate attacks in DFL. RepuNet continuously evaluates the behavior of participating nodes using a set of observable metrics and adjusts their influence in model aggregation according to their reputation. This section details the reputation scoring mechanism, the metrics used, the step-by-step operational flow, the reputation update process, exclusion and recovery criteria, and its practical integration in the Nebula experimental platform.

### A. Reputation System Foundations

RepuNet is designed for fully DFL environments with the goal of evaluating node behavior to detect and mitigate malicious activity while preserving contribution quality and minimizing overhead. Each node autonomously decides to accept, weight, or exclude incoming models. Its design emphasizes training efficiency and system stability, avoiding bottlenecks or mechanisms that could hinder global convergence. This enables a robust and resilient collaborative environment, even in the presence of anomalous behavior.
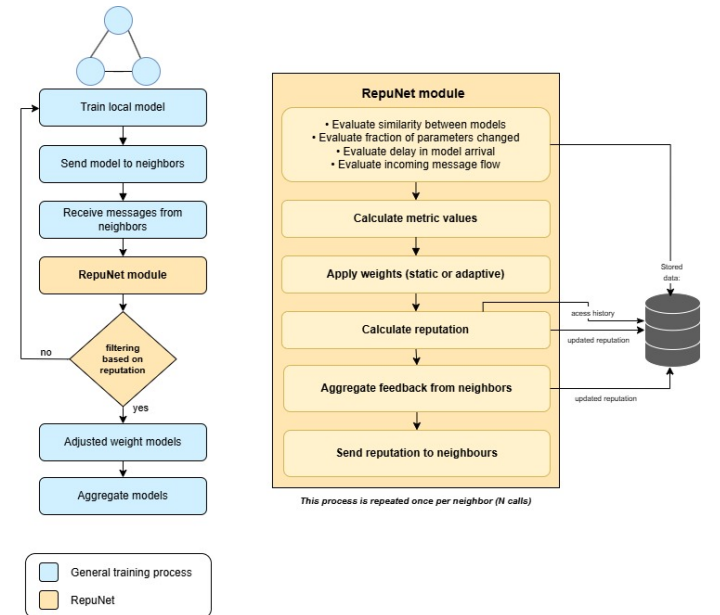
### B. Decentralized Architecture and Operation



Fig. 1. Interaction between a federated node and the RepuNet module, showing the flow of training, reputation evaluation, and aggregation.

RepuNet is modular and decentralized. Each node in the federated network includes an autonomous reputation component that evaluates the behavior of its neighbors during each training round. This evaluation determines whether received models should be filtered or adjusted before aggregation,

depending on the sender's reputation. Figure 1 shows two main components. On the left is the standard local training flow: local model training, model broadcasting, reception of external models, filtering based on reputation, and aggregation. On the right, the internal RepuNet module calculates a reputation score for each neighbor using observable behavioral metrics. The outcome of this evaluation guides the filtering and adjustment decisions. Upon receiving external models, the node calls the RepuNet module. This module evaluates each neighbor based on metrics such as model similarity, fraction of parameters changed, delay in model arrival, and incoming message flow. These metrics are processed and weighted—either statically or adaptively—to produce a reputation score. Feedback from other nodes is also aggregated to enhance the evaluation. Reputations are stored locally and then disseminated to neighbors in each round as part of the protocol. The updated reputation influences the filtering decision: if the reputation is below a given threshold, the model is discarded; otherwise, it is incorporated with a weight proportional to its reputation. The evaluation is performed once per neighbor per round.

The complete operational flow is summarized below:

1) **Local training:** The node updates its model using local private data.
2) **Model broadcast:** The model is shared with neighboring nodes.
3) **Message reception:** Models are received from neighbors.
4) **Reputation evaluation:** RepuNet evaluates each neighbor using local metrics and aggregated feedback.
5) **Filtering decision:** If a neighbor's reputation is too low, the model is discarded.
6) **Adjusted weighting:** Accepted models are weighted based on the sender's reputation.
7) **Aggregation:** Reputation-weighted models are aggregated with the local model.
8) **Reputation dissemination:** Updated reputation scores are stored and sent to all neighbors.

Unlike binary filters, RepuNet adjusts each model's aggregation weight. This adaptive integration allows the contribution of each node to evolve dynamically based on behavioral history. The next sections detail the computation of reputation and the metrics involved. The operational logic can be formalized through Algorithm 1, which summarizes the reputation-guided federated training cycle.

### C. Metrics Used in RepuNet

RepuNet relies on four primary metrics to evaluate the behavior of neighbors in a DFL network. These metrics are divided into two categories: *model quality* and *communication behavior*. Each metric returns a normalized value between 0 and 1, where 1 indicates optimal behavior. The system applies soft penalties to deviations, and the metrics are combined with dynamic weights to generate the reputation score.

*1) Model Similarity:* This metric assesses the alignment between the local model and the one received. Four classic

---

**Algorithm 1** Federated round with reputation-guided adaptive aggregation

```
1:  COMMUNICATION_ROUND()
2:      TrainLocalModel()
3:      SendModelToNeighbors()
4:      messages ← ReceiveMessagesFromNeighbors()
5:  for each message in messages do
6:      rep ← GetReputation(message.sender)
7:      if IsReputationSufficient(rep) then
8:          weight ← ComputeWeightFromReputation(rep)
9:          StoreAcceptedModel(message.model, weight)
10:     end if
11: end for
12: AggregateAcceptedModels()
13:
14: UPDATE_REPUTATION()
15: for each neighbor do
16:     metrics ← CalculateMetrics(neighbor)
17:     score ← ∑_j weights[j] · NormalizeMetric(metrics[j])
18:     history ← GetReputationHistory(neighbor)
19:     rep ← ∑_{r∈history} ω^{(r)} · rep[r] + ω^{(t)} · score
20:     if feedback available then
21:         avg_fb ← average reputation received from neighbors
22:         rep ← η · rep + (1 − η) · avg_fb
23:     end if
24:     StoreReputation(neighbor, rep)
25:     SendReputationToNeighbor(neighbor, rep)
26: end for
```

---

distance measures are used: cosine, Euclidean, Manhattan, and Pearson correlation. The final similarity is a weighted average:

$$\text{similarity} = \sum_k \gamma_k \cdot S_k \tag{1}$$

where $S_k$ represents each individual measure and $\gamma_k$ its weight. Low values indicate anomalous deviations in model parameters.

*2) Fraction of Parameters Changed:* RepuNet computes this metric by maintaining a history of change fraction values ($f$) and their associated thresholds ($t$), from which the mean and standard deviation are derived:

$$\mu_f = \frac{1}{N} \sum_{i=0}^{N} f_i, \qquad \sigma_f = \sqrt{\frac{1}{N} \sum_{i=0}^{N} (f_i - \mu_f)^2} \tag{2}$$

$$\mu_t = \frac{1}{N} \sum_{i=0}^{N} t_i, \qquad \sigma_t = \sqrt{\frac{1}{N} \sum_{i=0}^{N} (t_i - \mu_t)^2} \tag{3}$$

From these statistics, acceptable upper limits are defined:

$$\text{limit}_f = (\mu_f + \sigma_f) \times 1.05, \quad \text{limit}_t = (\mu_t + \sigma_t) \times 1.10 \tag{4}$$

If the current value exceeds the corresponding threshold, it is considered anomalous. A penalty proportional to the deviation is then calculated:

$$P = \frac{|f_{\text{current}} - \mu_f|}{\mu_f} \tag{5}$$

This penalty is transformed into a smoothed score using a sigmoid function:

$$S = 1 - \left( \frac{1}{1 + e^{-P}} \right) \tag{6}$$

If no anomaly is detected, the score is perfect: $S = 1.0$. The scores for the change fraction ($S_f$) and threshold ($S_t$) are combined with equal weights:

$$f_{\text{final}} = 0.5 \cdot S_f + 0.5 \cdot S_t \tag{7}$$

To avoid overpenalizing isolated variations, a temporal smoothing mechanism is applied:

$$f_{\text{final}}^{(t)} = \lambda^{(t)} \cdot f_{\text{actual}}^{(t)} + (1 - \lambda^{(t)}) \cdot f_{\text{final}}^{(t-1)} \tag{8}$$

where $\lambda^{(t)} \in [0,1]$ is a sensitivity coefficient for recent changes.

Finally, if a node does not receive a model from a neighbor during an iteration, this metric is automatically penalized by reducing the score by 50% from its previous value.

*3) Model Arrival Latency:* This metric evaluates how promptly a neighboring node delivers its model in each communication round. Its goal is to detect **anomalous delays** that may affect node synchronization and aggregation stability.

The computation begins by logging the model arrival time $t_{\text{arrival}}$, compared to the start time of the round. Depending on the temporal origin of the model, latency is defined as:

$$
\text{latency} = \begin{cases} t_{\text{arrival}} - t_{\text{start\_current\_round}} & \text{current round model} \\ t_{\text{arrival}} - t_{\text{start\_model\_round}} & \text{previous round model} \end{cases} \quad (9)
$$

Based on past observations, the historical latency average is computed as:

$$
\overline{\text{latency}} = \frac{1}{N} \sum_{i=1}^{N} \text{latency}_i \quad (10)
$$

When the attack begins from the first iteration, round 0 is used as the baseline, and the tolerance threshold is set at 150% of the average latency. The deviation is calculated as the difference between the current latency and the historical average latency. The final score is assigned as follows:

$$
\text{score} = \begin{cases} 1.0 & \text{if within threshold} \\ \dfrac{1}{1 + \exp\left(\frac{|\Delta t|}{\tau^{(t)}}\right)} & \text{if above} \end{cases} \quad (11)
$$

where $\Delta t$ is the observed deviation from the historical mean and $\tau^{(t)}$ is a tunable temporal tolerance parameter. To smooth inter-round variation, an exponential update is applied:

$$
\bar{s}_{\text{lat}}^{(t)} = \mu^{(t)} \cdot s_{\text{lat}}^{(t)} + (1 - \mu^{(t)}) \cdot \bar{s}_{\text{lat}}^{(t-1)} \quad (12)
$$

In early rounds, where history is limited, a bootstrapping penalty is introduced:

$$
\bar{s}_{\text{lat}}^{(t)} = s_{\text{lat}}^{(t)} \cdot (1 - \delta^{(t)}), \quad \delta^{(t)} \in [0,1], \quad \text{typically } \delta^{(t)} = 0.05 \quad (13)
$$

If no model is received from a node during a round, its score is reduced by 50 % from the previous round.

This metric helps detect both intentional delays and *lack of participation*, contributing to the preservation of temporal coherence in the federation.

*4) Incoming Message Flow:* This metric evaluates the communicative activity of a node by comparing the number of messages sent to a specific neighbor during the current round with the collective behavior of the previous round. Its goal is to detect anomalies like excessive or persistent overcommunication typical of flooding attacks.

For each node pair $(a,b)$, the number of messages sent in round $r$ is denoted $m_r^{(a,b)}$, and the previous round $(r-1)$ message counts across all pairs $(i,j)$ are collected. From this, the 25th percentile $P_{25}$ is computed as a reference, along with standard deviation $\sigma_{r-1}$ and mean $\mu_{r-1}$.

The relative increase is defined as:

$$
\text{rel\_incr} = \max\left(\frac{m_r^{(a,b)} - P_{25}}{P_{25}}, 0\right) \quad (14)
$$

A dynamic tolerance margin is calculated based on variability:

$$
\text{dynamic\_margin} = \frac{\sigma_{r-1} + 1}{\log(1 + P_{25}) + 1} \quad (15)
$$

If the relative increase exceeds the margin, a smoothed logarithmic exponential penalty is applied:

$$
s_{\text{msg}}^{(r)} \leftarrow s_{\text{msg}}^{(r)} \cdot \exp\left(-\left(\frac{\log\left(1 + \text{rel\_incr} - \text{dynamic\_margin}\right)}{\log\left(1 + \text{dynamic\_margin}\right) + \varepsilon}\right)^2\right) \quad (16)
$$

If the node also exceeds the global mean:

$$
\text{amplification} = 1 + \frac{\text{increase\_mean}}{\mu_{r-1} + \varepsilon} \quad (17)
$$

$$
s_{\text{msg}}^{(r)} \leftarrow s_{\text{msg}}^{(r)} \cdot \exp\left(-(\text{extra\_penalty} \cdot \text{amplification})^2\right) \quad (18)
$$

where increase_mean adjusts automatically based on the system phase.

A per-node and per-neighbor temporal history is maintained. Repeated low scores across rounds trigger additional multiplicative penalties. To avoid permanent exclusion, a minimum adaptive bound is enforced:

$$
s_{\text{msg}}^{(r)} \leftarrow \max\left(s_{\text{msg}}^{(r)}, f_{\min}(r, a, b)\right) \quad (19)
$$

Values are smoothed with a weighted average of up to three prior rounds. This formulation detects and penalizes meaningful communication deviations proportionally and dynamically, without requiring fixed thresholds.

These four metrics form the foundation of the RepuNet reputation system and are summarized in Table II.

### D. Weight Assignment and Reputation Computation

Once the metrics for each neighbor have been computed and normalized, RepuNet assigns dynamic weights that reflect the discriminative power of each metric in the current round. This assignment is performed individually per node, based on deviation from a reference historical behavior. If the system has a stable enough history, deviations are computed with respect to the accumulated historical mean. In attack scenarios activated in later rounds (e.g., after round 7), the prior history is used as the reference pattern. In contrast, if the attack begins in the first iteration, only the data from round 0 is used as a baseline. Let $m_j^{(t)}$ be the observed value of metric $j$ for a node in round $t$, and $\mu_j$ its historical mean. The absolute deviation is calculated as:

$$
d_j^{(t)} = \left| m_j^{(t)} - \mu_j \right| \quad (20)
$$

These deviations are aggregated into a total value:

$$
D^{(t)} = \sum_k d_k^{(t)} \quad (21)
$$

and used to compute the normalized weight assigned to each metric:

$$
w_j^{(t)} = \frac{d_j^{(t)}}{D^{(t)}} \quad (22)
$$

If all deviations are zero (i.e., the node behaves exactly as expected), randomly normalized weights are assigned to prevent stagnation in the reputation process. Additionally, a dynamic lower bound is applied to ensure no relevant metric is ignored. Once the weights are set, an intermediate reputation

TABLE II
SUMMARY OF METRICS USED IN REPUTATION COMPUTATION

| Metric | Type | Purpose | Symbol | Range | Penalty trigger | Normalization |
|---|---|---|---|---|---|---|
| Model similarity | Model quality | Align local and received models | $S$ | $[0,1]$ | Low similarity | Weighted average (cosine, Euclidean, Manhattan, Pearson) |
| Fraction of parameters changed | Model quality | Detect abrupt model changes | $F$ | $[0,1]$ | High deviation from past | Sigmoid penalty on deviation |
| Model arrival latency | Communication | Penalize delayed models | $L$ | $[0,1]$ | Above historical latency | Sigmoid based on delay deviation |
| Incoming message flow | Communication | Detect flooding behavior | $M$ | $[0,1]$ | High message volume | Exponential penalty vs. dynamic margin |

score $S_i^{(t)}$ is computed for each neighbor $i$, using a weighted sum of the metrics:

$$S_i^{(t)} = \sum_j w_j^{(t)} \cdot m_j^{(t)} \tag{23}$$

This score is used to update the node's overall reputation, taking into account the history of past rounds $\mathcal{H}_i^{(t)}$. The reputation is updated as a weighted combination of past values and the current score:

$$R_i^{(t)} = \sum_{r \in \mathcal{H}_i^{(t)}} \omega^{(r)} \cdot R_i^{(r)} + \omega^{(t)} \cdot S_i^{(t)} \tag{24}$$

where $\omega^{(r)}$ and $\omega^{(t)}$ are weights satisfying $\sum \omega = 1$, allowing balance between historical behavior and current observation. Finally, if reputation feedback is enabled, RepuNet can incorporate the opinions of neighboring nodes. Given a set of received values $R_{j \to i}^{(t)}$, an adjusted reputation is computed using a convex combination controlled by a parameter $\eta \in [0, 1]$:

$$\tilde{R}_i^{(t)} = \eta \cdot R_i^{(t)} + (1 - \eta) \cdot \text{mean}\left(\left\{R_{j \to i}^{(t)}\right\}\right) \tag{25}$$

The resulting value $\tilde{R}_i^{(t)}$ represents the final reputation of the node for that round and will be used to determine its influence in the next aggregation iteration.

### E. Aggregation, Exclusion, and Recovery Control

Once the final reputation of each node has been computed, RepuNet uses this value to make local decisions about including or excluding their models in the aggregation process. If a neighbor's reputation falls below a predefined trust threshold, their model is automatically discarded for that round. This exclusion prevents anomalous, malicious, or inconsistent behaviors from influencing the global model. It is important to note that exclusion does not imply disconnection or blocking of the node in the system: the penalty only applies to the use of its model, allowing continued participation in subsequent rounds and re-evaluation. This ensures the architecture remains open and tolerant toward nodes that may recover their behavior after temporary penalties. The exclusion threshold is configurable and can be adapted to the nature of the environment or the attack. In typical configurations, reputation values below 0.6 lead to model rejection, while higher reputations allow the node's contribution to be weighted accordingly. Furthermore, reputation is evaluated and updated in every round. If a node improves its behavior and its metrics begin to align with the expected pattern, its reputation gradually increases, allowing for natural and progressive reintegration into aggregation. This recovery property avoids permanent penalties for transient issues or isolated failures. Overall, this control mechanism provides a flexible, adaptive, and resilient system that acts as a dynamic filter against untrustworthy or malicious contributions while encouraging the reintegration of rehabilitated nodes.

### F. Distributed Feedback and Robustness

In addition to local metrics, RepuNet allows each node to incorporate reputation feedback from trusted neighbors. This complementary information reinforces the individual assessment, especially in scenarios where a node may lack sufficient observations or its perception may be distorted due to specific network conditions or atypical behavior. Each node may receive reputation scores issued by its neighbors in the current round. This feedback is incorporated through a weighted combination of the local evaluation $R_i^{(t)}$ and the average of the received evaluations $\text{mean}(\{R_{j \to i}^{(t)}\})$, regulated by a coefficient $\eta \in [0, 1]$ that controls the level of trust in the internal assessment:

$$\tilde{R}_i^{(t)} = \eta \cdot R_i^{(t)} + (1 - \eta) \cdot \text{mean}\left(\left\{R_{j \to i}^{(t)}\right\}\right) \tag{26}$$

Since only nodes with an active reputation module are allowed to issue evaluations, it is guaranteed that all received scores originate from valid and trustworthy participants, thereby preventing the propagation of manipulated information. This feedback mechanism serves as a decentralized consensus approach, enabling the smoothing of isolated evaluation errors, mitigating local fluctuations, and reinforcing system robustness against targeted attacks or biased perceptions. This preserves the stability and adaptability of the federated process. In subsequent rounds, the system repeats the same evaluation cycle: metric computation, dynamic weight assignment, reputation update, neighbor feedback incorporation, and dissemination of the final reputation. This continuous and adaptive mechanism allows the federated system to remain resilient against malicious nodes while supporting the reintegration of nodes that recover reliable behavior.

### G. Integration of RepuNet into the Nebula Platform

RepuNet has been implemented as a modular, autonomous, and platform-independent component. For experimental evaluation, it was integrated into the *Nebula* platform, which supports decentralized FL scenarios and includes real-time visualization via TensorBoard [10]. The integration enables testing across topologies and attack types.

Nebula's architecture consists of three main layers:

- **Frontend:** a graphical interface where users can:
  - Activate or deactivate RepuNet.
  - Select evaluation metrics and their weighting scheme.
  - Configure parameters such as the exclusion threshold, feedback weight $\eta$, and temporal smoothing.
- **Controller:** an intermediate layer that propagates configurations to all nodes, ensuring consistency across rounds.
- **Core:** the execution layer for each node, which handles:

– Local training on private data.
– Model exchange with neighbors.
– Metric evaluation and reputation updates.
– Weighted aggregation based on reputation.

Each node maintains its own history of metrics and reputations, enabling smoothing, recurrence detection, and adaptive behavior. The integrated visualization tools help track the evolution of each metric, monitor node reputation across rounds, and observe the influence of distributed feedback.
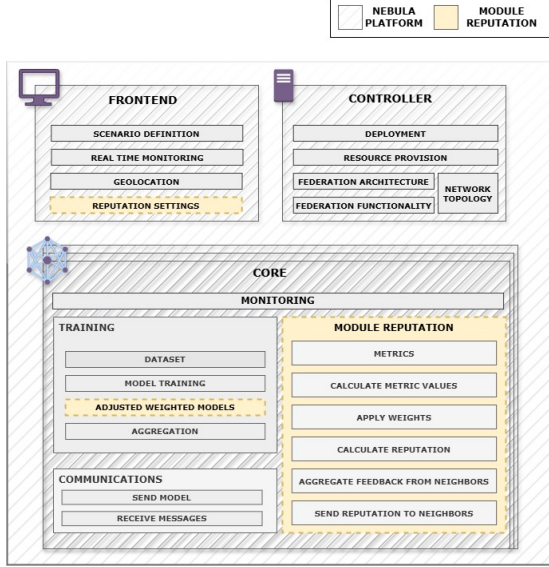


Fig. 2. Integration of the RepuNet reputation system into the Nebula platform architecture.

## V. REPUNET EVALUATION

This section evaluates RepuNet in adversarial DFL environments. The objective is to assess its effectiveness in detecting and penalizing malicious nodes, preserving model performance and training stability. The experiments vary the proportion of attackers and the timing of activation, analyzing the evolution of reputation scores and their effect on aggregation, with and without defense.

### A. Experimental Environment

Experiments were run on Nebula, which emulates DFL via virtual scenarios. Nodes train, evaluate, and aggregate locally without a central server. The entire federation runs as independent processes on a single machine, enabling full control and real-time monitoring. Table III summarizes the general setup parameters, including datasets, topologies, hardware configuration, and training schedule used across all experiments.

We emulated model poisoning, delayer, and flooding attacks to test robustness. For each, the experimental task and measurement objective are clearly defined in Table IV.

### B. Model Poisoning Attack

This attack degrades the global model by injecting malicious updates that differ significantly from honest contributions and

TABLE III
SUMMARY OF EXPERIMENTAL SETUP PARAMETERS

| Parameter | Description |
|---|---|
| Platform | Nebula (Decentralized FL emulation platform) |
| Datasets | MNIST, CIFAR-10 |
| Nodes per federation | 10–25 nodes |
| Topology types | Fully connected, Random |
| Data partition | Dirichlet (Non-IID), $\alpha \in \{0.1, 0.3, 0.5, 0.9\}$ |
| Local epochs | MNIST: 1 epoch/round; CIFAR-10: 10 epochs/round |
| Aggregation algorithm | FedAvg |
| Hardware specs | Intel Xeon CPU E5-2697, 62 GB RAM, Ubuntu 22.04 |
| Number of rounds | 20 communication rounds |
| Timeout for aggregation | 30–60 seconds |
| Metrics used | Model Similarity, Fraction of Parameters Changed, Arrival Latency, Incoming Message Flow |

TABLE IV
ATTACK TYPES, TASKS, AND MEASUREMENT OBJECTIVES

| Attack Type | Task | Measurement Objective |
|---|---|---|
| Model Poisoning | Image classification (MNIST, CIFAR-10) | Evaluate degradation of global model accuracy (F1-score) due to manipulated updates. |
| Delayer Attack | Image classification (MNIST) | Assess synchronization issues and increased aggregation latency caused by delayed submissions. |
| Flooding Attack | Image classification (MNIST) | Measure communication efficiency degradation and CPU overhead caused by excessive messaging. |

may hinder convergence if undetected. RepuNet evaluates each neighbor before aggregation based on behavioral metrics. Experiments were conducted using MNIST and CIFAR-10 datasets under different topologies, attacker proportions, and data heterogeneity levels. Table V summarizes the configuration of all poisoning scenarios, including early and intermittent variants. Each scenario is labeled with a short identifier (e.g., *Base*, *2.1*, *6.1*) used consistently across results and discussion.

TABLE V
MODEL POISONING ATTACK SCENARIOS: CONFIGURATION SUMMARY.

| Scenario Group | Topology | Malicious Nodes | Dirichlet. $\alpha$ |
|---|---|---|---|
| Base / 1.x / 2.x | Fully | 30–60% | 0.5 |
| 3.1 | Ring | 30% | 0.5 |
| 3.2 | Random | 30% | 0.5 |
| 4.1 / 4.2 | Fully | 30% | 0.3 / 0.9 |
| 5.1 / 5.2 / 5.3 | Fully | 30% | CIFAR10: 0.5 / 0.1 |
| 6.1 | Fully | 30% | 0.5 (early) |
| 7.1 / 7.2 | Fully | 30% | 0.5 (interval 2/3) |

Table VI reports the F1-score for each scenario with and without RepuNet. Scores are measured at round 8 and 11, respectively. The difference $\Delta F_1$ quantifies the improvement, and the last column categorizes the observed impact.

Figure 3 provides a visual summary of the impact levels across all scenario groups. Subfigures organize experiments by activation type or configuration and highlight thresholds for low, medium, and high impact levels. Three representative scenarios illustrate RepuNet's response to poisoning attacks. These experiments use the MNIST dataset and simulate attacks in federations with 10 nodes. Each case includes two plots: one showing the average reputation of all nodes, and another focused on malicious nodes only. The results are shown in Figure 4, which illustrates how reputation evolves in response to adversarial behavior under different topologies and attacker

(a) General scenarios

(b) Early attacks (round 1)

(c) Intermittent attacks

(d) Dirichlet $\alpha$ variation (4.1 vs 4.2)

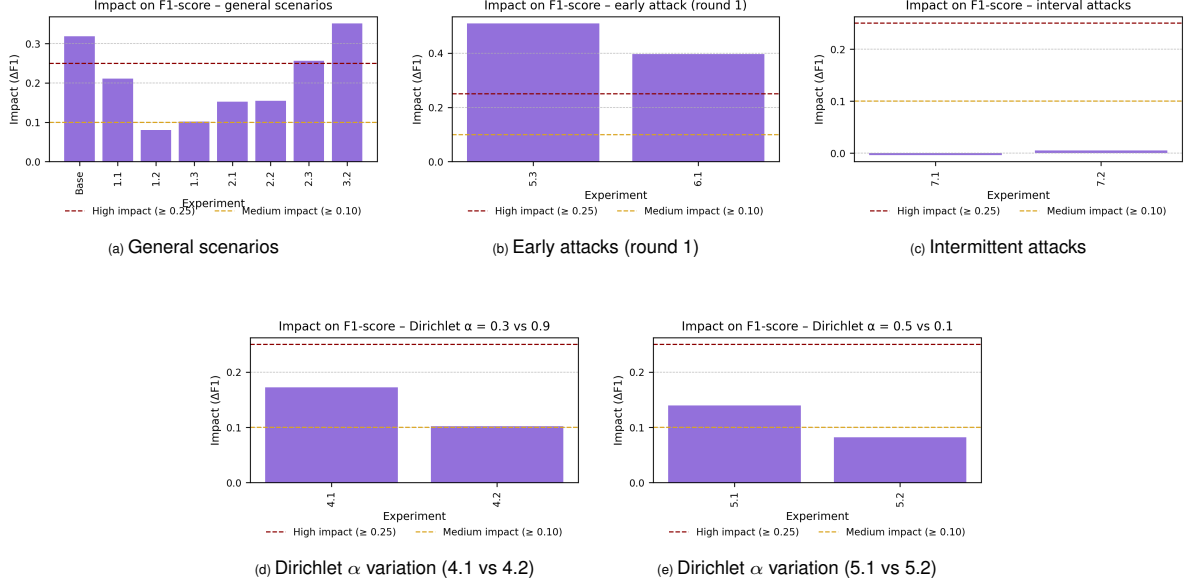(e) Dirichlet $\alpha$ variation (5.1 vs 5.2)

Fig. 3. Impact of the RepuNet reputation system in various *model poisoning* attack scenarios. (a) General scenarios, (b) attacks from round 1, (c) intermittent attacks, (d–e) sensitivity to the heterogeneity parameter $\alpha$ of the Dirichlet distribution. Horizontal lines mark the thresholds for high impact ($\geq 0.25$) and medium impact ($\geq 0.10$).

TABLE VI
F1-SCORE COMPARISON ACROSS MODEL POISONING SCENARIOS WITH AND WITHOUT REPUNET.

| Scenario | F1 w/rep (r8) | F1 w/o rep (r11) | $\Delta F_1$ | Impact |
|---|---|---|---|---|
| Base | 0.6800 | 0.3610 | 0.3190 | High |
| 1.1 | 0.6340 | 0.4226 | 0.2114 | Medium |
| 1.2 | 0.6670 | 0.5868 | 0.0802 | Low |
| 1.3 | 0.6447 | 0.5430 | 0.1017 | Medium |
| 2.1 | 0.5972 | 0.4444 | 0.1528 | Medium |
| 2.2 | 0.5051 | 0.3506 | 0.1545 | Medium |
| 2.3 | 0.4143 | 0.1575 | 0.2568 | High |
| 3.1 | 0.6821 | 0.4017 | 0.2804 | High |
| 3.2 | 0.6744 | 0.3231 | 0.3513 | High |
| 4.1 | 0.6736 | 0.5012 | 0.1724 | Medium |
| 4.2 | 0.6787 | 0.5763 | 0.1024 | Medium |
| 5.1 | 0.5360 | 0.3965 | 0.1395 | Medium |
| 5.2 | 0.5122 | 0.4301 | 0.0821 | Low |
| 5.3 | 0.5166 | 0.0720 | 0.4446 | High |
| 6.1 | 0.6191 | 0.2225 | 0.3966 | High |
| 7.1 | 0.6796 | 0.6833 | -0.0037 | Low |
| 7.2 | 0.6879 | 0.6825 | 0.0054 | Low |



(a) Base – All nodes

(b) 3.2 – All nodes

(c) 2.3 – All nodes

(d) Base – Malicious nodes

(e) 3.2 – Malicious nodes
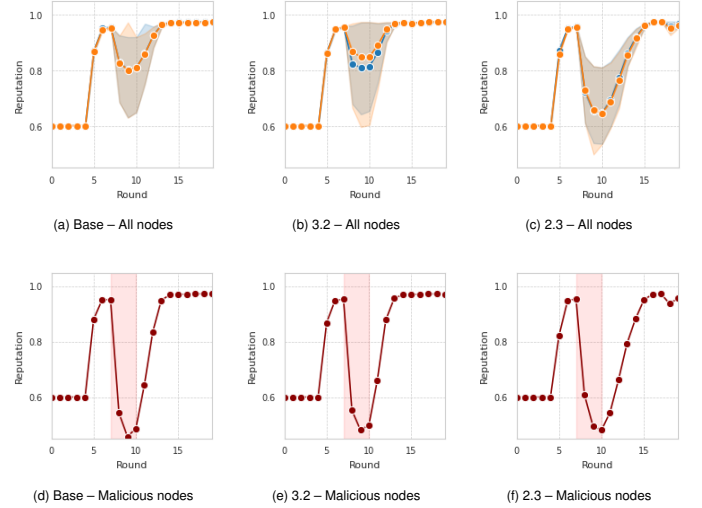
(f) 2.3 – Malicious nodes

Fig. 4. Reputation evolution in three scenarios with *model poisoning* attacks. Top row: average reputation of all federation nodes. Bottom row: average reputation of malicious nodes.

proportions.

Scenario **Base** represents a fully connected topology with 30% malicious nodes. RepuNet responds quickly after the attack is triggered in round 7, penalizing adversarial behavior. The reputation of honest nodes remains high, while that of malicious nodes drops noticeably. Nodes 192.168.51.10:45009 and 192.168.51.11:45010 were selected for illustration. Scenario **3.2** uses a random topology and maintains the same proportion of attackers. Greater variability in reputation is observed due to the network's sparse connectivity, which reduces the influence of malicious nodes. The drop in reputation begins in round 7, confirming RepuNet's ability to adapt in less structured environments. Nodes 192.168.51.2:45001 and 192.168.51.9:45008 are shown. Scenario **2.3** corresponds to a more adversarial environment, with a fully connected topology

and 60% malicious nodes. Despite the increased pressure, RepuNet still penalizes malicious behavior consistently. The Dirichlet parameter $\alpha = 0.5$ used for data partitioning has limited effect in this setting. Nodes 192.168.51.4:45003 and 192.168.51.9:45008 were selected for observation. To illustrate the internal adjustment process, Table VII shows how node 10 tracks a malicious neighbor across rounds. In rounds 7 and 8, after the attack starts, the model similarity and fraction of parameters changed drop significantly. This leads to an automatic reweighting, increasing the importance of these more discriminative metrics. The system adapts to contextual anomalies, penalizing malicious behavior without manual in-

tervention.

TABLE VII
METRIC AND WEIGHT EVOLUTION FOR A MALICIOUS NEIGHBOR (NODE 10).

| Metric | Round 5 | Round 6 | Round 7 | Round 8 |
|---|---|---|---|---|
| Model similarity | 0.84 | 0.82 | 0.42 | 0.36 |
| Fraction of parameters changed | 0.79 | 0.75 | 0.33 | 0.29 |
| Arrival latency | 0.91 | 0.93 | 0.89 | 0.94 |
| Messages exchanged | 0.86 | 0.88 | 0.85 | 0.87 |
| **Weight: similarity** | 0.25 | 0.24 | 0.41 | 0.43 |
| **Weight: fraction** | 0.22 | 0.22 | 0.38 | 0.39 |
| **Weight: latency** | 0.28 | 0.28 | 0.11 | 0.10 |
| **Weight: messages** | 0.25 | 0.26 | 0.10 | 0.08 |

### C. Delayer Attack

The *delayer* attack disrupts the federated process by intentionally delaying the transmission of model updates. While it does not alter model content, it exploits the temporal dimension—latency—to desynchronize nodes, degrade convergence, and stall training. Experiments on MNIST were conducted using the Nebula platform, exploring different delay durations, topologies, and attacker proportions. **Table VIII** summarizes the configurations evaluated.

TABLE VIII
SUMMARY OF DELAY ATTACK CONFIGURATIONS TESTED.

| Scenario Group | Topology | Delay (s) | Malicious Nodes | Interval |
|---|---|---|---|---|
| Base / 2.x | Fully | 20 | 30–60% | 1 |
| 3.2 | Random | 30 | 30% | 1 |
| 4.x | Fully | 40–60 | 30% | 1 |
| 5.x (intermittent) | Fully | 80 | 30% | 2–4 |
| 6.x | Fully | 100 | 30–60% | 1 |

Unlike poisoning, this attack targets efficiency. Figure 7 presents the total training time by node type, showing how lag propagates and increases desynchronization when no mitigation is applied. For reference, the baseline training duration without attack was 5 minutes; with reputation enabled, it increased slightly to 5.7 minutes. To assess how reputation adapts, Figure 5 shows average reputation evolution under different attack timings. Figures (a) and (b) correspond to the global reputation for attacks starting in round 7 and round 1, respectively. Figures (c) and (d) show the reputation of malicious nodes in the same scenarios. Finally, (e) and (f) present global and malicious reputations under intermittent attacks with intervals 2, 3, and 4.

Figure 6 shows round gaps between benign and malicious nodes. Subfigure (a) refers to scenarios where the attack begins in round 1; (b) shows scenarios with attack starting from round 7. In both cases, average reputation of malicious nodes is included for correlation.

Lastly, Figure 8 reports the average number of aggregated models depending on the percentage of malicious nodes. Panel (a) refers to round 9 (attack from round 7), while (b) shows round 2 (attack from round 1). These results demonstrate how model aggregation is reduced in the presence of delays.



(a) Global reputation – attack from round 7

(b) Malicious reputation – attack from round 7

(c) Global reputation – attack from round 1

(d) Malicious reputation – attack from round 1

(e) Global reputation – intermittent attack
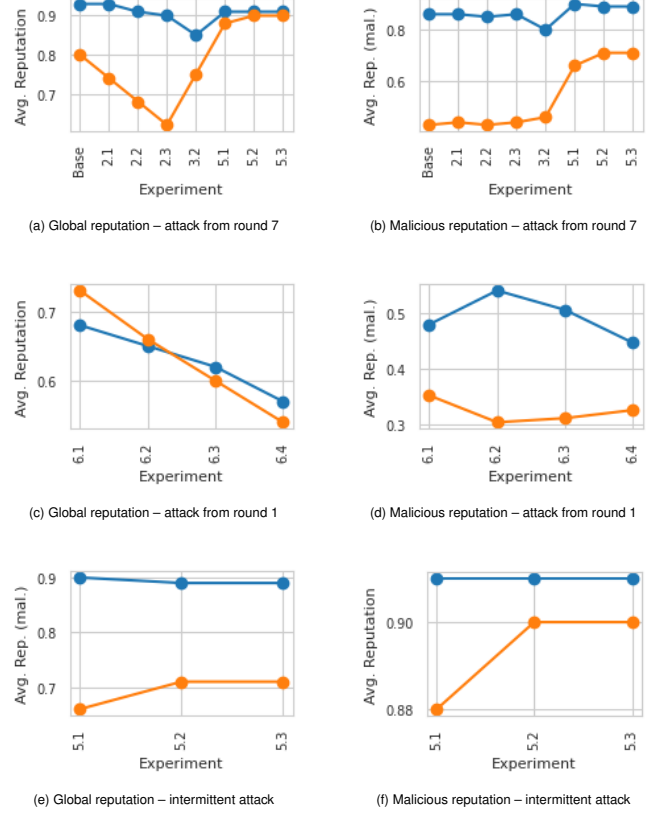
(f) Malicious reputation – intermittent attack

Fig. 5. Evolution of average reputation across different *delayer* attack configurations. Each row compares global reputation (left) and malicious node reputation (right) under a specific attack timing: from round 7, from round 1, and with intermittent activation.
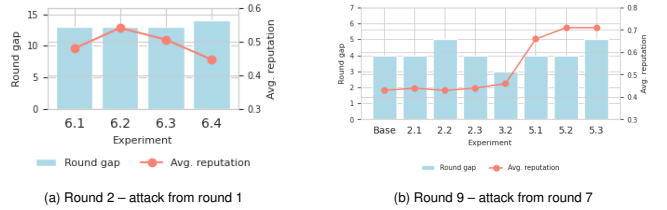


(a) Round 2 – attack from round 1

(b) Round 9 – attack from round 7

Fig. 6. Relationship between the round gap and average reputation of malicious nodes in two scenarios.

### D. Flooding Attack

The *flooding* attack aims to saturate the communication layer by overwhelming the network with a high volume of messages from malicious nodes. This behavior degrades communication efficiency, increases CPU overhead, and can ultimately impact model convergence. RepuNet addresses this threat by penalizing nodes that exhibit abnormal message volume. Table IX summarizes the configurations evaluated for this attack, including variations in topology, proportion of malicious nodes, activation round, and message injection interval.

Figure 9 shows the evolution of the average reputation under flooding scenarios starting at round 7. Malicious nodes exhibit a sharp and consistent decline in reputation, while benign participants maintain high values, highlighting RepuNet's ability
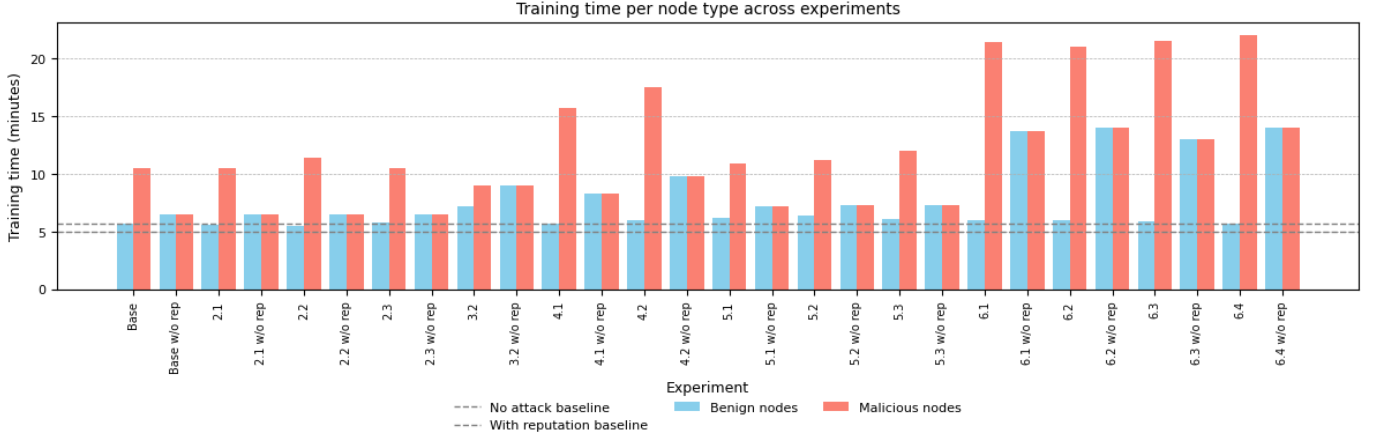
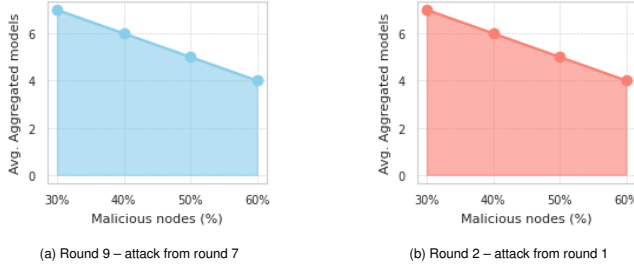Fig. 7. Total training time per node type in the different scenarios.



Fig. 8. Average number of models aggregated by percentage of malicious nodes.

TABLE IX
FLOODING ATTACK SCENARIOS: CONFIGURATION SUMMARY (GROUPED).

| Scenario Group | Topology | Malicious Nodes | Start Round |
|---|---|---|---|
| Base / 2.x | Fully | 30–60% | 7 |
| 3.1 | Random | 30% | 7 |
| 4.1 | Fully | 30% | 7 (interval 2) |
| 5.x | Fully | 30–60% | 1 |

to penalize anomalies without affecting honest contributions.



Fig. 9. Average reputation evolution of all and malicious nodes during the flooding attack (activated in round 7).

Additional scenarios were tested where the flooding attack was activated from the first round and ran until round 10. In these early-activation cases, round 0 served as the only observation phase, with all nodes initialized to a default reputation of 0.6. Table X presents the reputation scores observed in these scenarios.

TABLE X
BENIGN AND MALICIOUS REPUTATION UNDER *flooding* ATTACKS
STARTING AT ROUND 1.

| Scenario | Benign Rep. (r2) | Benign Rep. (r6) | Malicious Rep. (r2) | Malicious Rep. (r6) |
|---|---|---|---|---|
| 5.1 | 0.9387 | 0.9568 | 0.3256 | 0.2718 |
| 5.2 | 0.9385 | 0.9347 | 0.2936 | 0.2970 |
| 5.3 | 0.9050 | 0.9335 | 0.3025 | 0.3239 |
| 5.4 | 0.9168 | 0.9313 | 0.2950 | 0.3914 |

The impact on the aggregation process was also assessed by measuring the average number of models accepted in round 9. As shown in Figure 10, the number of aggregated models decreases proportionally with the increase in malicious nodes. This behavior holds for both early and delayed attack activation, confirming RepuNet's robustness even with limited prior information.
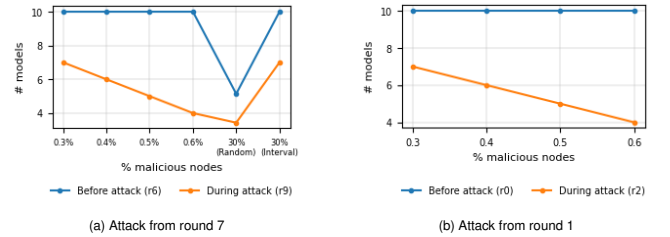


Fig. 10. Average number of models aggregated in round 9 under different flooding attack activation points. The Y-axis represents the average number of models aggregated by the federation nodes. The system consistently reduces the inclusion of malicious contributions as their proportion increases.

Finally, to evaluate the indirect cost of flooding attacks, Table XI reports the variation in average CPU usage observed in benign nodes. This variation reflects the computational overhead caused by processing excessive messages. Scenarios with higher percentages of malicious nodes tend to exhibit increased CPU load, underscoring the importance of mitigating such attacks not only to preserve aggregation quality but also to reduce resource consumption.

## VI. DISCUSSION

The experimental evaluation demonstrates that RepuNet provides effective protection against diverse adversarial be-

TABLE XI
CPU USAGE VARIATION IN BENIGN NODES UNDER *flooding* ATTACKS.

| Scenario | Attack round | Avg. CPU Variation (%) |
|---|---|---|
| Base | 7–12 | 5.23 |
| 2.1 | 7–12 | -1.03 |
| 2.2 | 7–12 | 15.14 |
| 2.3 | 7–12 | 3.70 |
| 3.1 | 7–12 | 11.21 |
| 4.1 | 7–12 | -4.49 |
| 5.1 | 1–10 | -8.14 |
| 5.2 | 1–10 | 5.70 |
| 5.3 | 1–10 | -4.80 |
| 5.4 | 1–10 | -9.90 |

haviors in DFL scenarios. Rather than reiterating the internal mechanics of the reputation engine, this section focuses on analyzing its observed behavior, highlighting strengths, limitations, and areas for improvement. Across all experiments, RepuNet consistently reduced the influence of malicious nodes during aggregation without compromising the participation structure of the network. Reputation trajectories showed rapid penalization after attack activation, confirming the system's responsiveness. However, certain patterns—such as intermittent or oscillating attacks—exposed scenarios where nodes could recover influence too quickly, suggesting that stricter memory-based penalization strategies may be needed. The ability to adapt dynamically to contextual changes, without centralized control or static heuristics, positioned RepuNet as a flexible defense strategy. Nonetheless, further refinement of metric weighting and exclusion thresholds could enhance robustness under more stealthy or coordinated threats. The following subsections delve into specific attack types and their mitigation outcomes.

### A. Model Poisoning

Table VI presents the F1-scores obtained with and without RepuNet across several poisoning scenarios. The most significant improvements were observed in scenarios *Base*, *3.2*, and *6.1*, where RepuNet increased performance by over 30 points in F1-score, reaching differences above 0.35. In early-activation cases such as *5.3*, where attacks began at round 1, the gap was even greater: F1-score improved from 0.0720 to 0.5166. This demonstrates RepuNet's capacity to recover learning even under high pressure and without a prior observation phase. Scenarios with skewed data distributions (e.g., *5.2* with $\alpha = 0.1$) showed lower gains, as expected under higher heterogeneity, but still maintained improvements over the baseline. In intermittent scenarios (*7.1*, *7.2*), RepuNet showed limited effect, suggesting the system could benefit from memory-based metrics to counter low-frequency adversarial patterns. Overall, the results confirm that RepuNet substantially mitigates poisoning attacks across various conditions, adapting both to topology and distribution challenges.

### B. Delayer Attack

As shown in Table VIII, the delayer attack was tested under various delays (20–100 seconds), attacker proportions,

and topologies. RepuNet successfully penalized delayed nodes across all configurations. In early rounds (r2), scenarios such as *2.3* and *6.1* already showed reputation gaps of more than 0.3 between benign and malicious participants. Reputation remained low for attackers in both fully connected and random topologies, as shown in Figure 5, confirming the system's capacity to suppress desynchronizing behavior. Intermittent delay scenarios revealed partial reintegration between inactive phases, which may be addressed by increasing the exclusion threshold or applying temporal penalties. Despite this, overall training time remained stable for honest nodes, while malicious participants completed fewer rounds (Figure 7). The reduction in model aggregation under delay conditions was also consistent. Figure 8 shows that as the delay severity or attacker proportion increased, the number of models selected decreased accordingly—particularly in scenarios *4.1* and *6.3*.

### C. Flooding Attack

Table IX summarizes the configurations evaluated, including topology, proportion of malicious nodes, activation round, and interval. When the flooding attack was activated at round 7 (scenarios *Base* to *4.1*), RepuNet responded effectively: the reputation of malicious nodes dropped from initial values near 0.98 to values between 0.38 and 0.48, while honest nodes maintained values around 0.80. In early-activation scenarios (round 1, *5.1–5.4*), all nodes started from a neutral reputation of 0.6. Table X shows the evolution of reputation for benign and malicious nodes. Benign nodes rapidly recovered to values above 0.93 by round 6, while adversarial nodes remained low or exhibited only limited recovery. The number of aggregated models in round 9 (Figure 10) decreased progressively with the proportion of malicious nodes, both for early and delayed attacks, indicating consistent exclusion of unreliable contributions. Intermittent attacks such as *4.1* revealed partial reintegration of malicious nodes between flooding bursts, suggesting a limitation of short-term metrics. Despite this, RepuNet re-penalized such behavior in subsequent rounds, preserving overall robustness. Table XI presents CPU usage variations for benign nodes. In most cases, RepuNet reduced or stabilized overhead, with the highest efficiency gains observed in scenarios *5.1* and *5.3*. These results confirm that reputation-based filtering improves not only model quality but also resource efficiency under flooding conditions.

## VII. CONCLUSIONS

This work presented a reputation mechanism for DFL that detects and mitigates malicious behavior during training. The system evaluates each node based on locally observable metrics—model similarity, parameter changes, latency, and message volume—and updates reputation scores dynamically. Nodes with low reputation are excluded from aggregation, while recovery is possible upon improved behavior. Experiments confirm the system's effectiveness against model poisoning, delay, and flooding attacks. In all cases, malicious nodes suffered a sharp drop in reputation after initiating the attack, enabling their exclusion without degrading the aggregated model.

RepuNet could be extended to deal with threats like data poisoning, which opens possibilities for integrating new metrics that assess model coherence with local data. Future improvements could also include a dynamic exclusion threshold, instead of a fixed one, and improved weight assignment by analyzing Z-score normalization, Bayesian schemes, or attention mechanisms to adjust metric relevance based on context.

## REFERENCES

[1] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, 2017, pp. 1273–1282.

[2] E. T. Martínez Beltrán, M. Q. Pérez, P. M. S. Sánchez, S. L. Bernal, G. Bovet, M. G. Pérez, G. M. Pérez, and A. H. Celdrán, "Decentralized federated learning: Fundamentals, state of the art, frameworks, trends, and challenges," *IEEE Communications Surveys & Tutorials*, vol. 25, no. 4, pp. 2983–3013, 2023.

[3] Y. Li, X. Wei, Y. Li, Z. Dong, and M. Shahidehpour, "Detection of false data injection attacks in smart grid: A secure federated deep learning approach," *IEEE Transactions on Smart Grid*, vol. 13, no. 6, pp. 4862–4872, 2022.

[4] F. P.-C. Lin, S. Hosseinalipour, N. Michelusi, and C. G. Brinton, "Delay-aware hierarchical federated learning," *IEEE Transactions on Cognitive Communications and Networking*, vol. 10, no. 2, pp. 674–688, April 2024.

[5] J. Li, L. Lyu, X. Liu, X. Zhang, and X. Lyu, "FLEAM: A federated learning empowered architecture to mitigate DDoS in industrial IoT," *IEEE Transactions on Industrial Informatics*, vol. 18, no. 6, pp. 4059–4068, 2022.

[6] A. Taherpour and X. Wang, "HybridChain: Fast, accurate, and secure transaction processing with distributed learning," *IEEE Transactions on Parallel and Distributed Systems*, vol. 35, no. 6, pp. 968–982, 2024.

[7] M. Panigrahi, S. Bharti, and A. Sharma, "A reputation-aware hierarchical aggregation framework for federated learning," *Computers and Electrical Engineering*, vol. 111, p. 108900, 2023.

[8] E. Gabrielli, D. Belli, Z. Matrullo, V. Miori, and G. Tolomei, "Protecting federated learning from extreme model poisoning attacks via multidimensional time series anomaly detection," 2024. [Online]. Available: https://doi.org/10.48550/arXiv.2303.16668

[9] M. Mohamad, M. Önen, W. Ben Jaballah, and M. Conti, "Sok: Secure aggregation based on cryptographic schemes for federated learning," in *Proceedings on Privacy Enhancing Technologies*, 2023, pp. 140–157.

[10] CyberDataLab, "Nebula: A platform for decentralized federated learning," https://github.com/CyberDataLab/nebula, 2025, accessed: 2025-06-13.

[11] Z. Song, H. Sun, H. H. Yang, X. Wang, Y. Zhang, and T. Q. S. Quek, "Reputation-based federated learning for secure wireless networks," *IEEE Internet of Things Journal*, vol. 9, no. 2, pp. 1212–1226, 2022.

[12] J. Kang, Z. Xiong, D. Niyato, S. Xie, and J. Zhang, "Incentive mechanism for reliable federated learning: A joint optimization approach to combining reputation and contract theory," *IEEE Internet of Things Journal*, vol. 6, no. 6, pp. 10 700–10 714, 2019.

[13] Y. Zhao, J. Zhao, L. Jiang, R. Tan, D. Niyato, Z. Li, L. Lyu, and Y. Liu, "Privacy-preserving blockchain-based federated learning for IoT devices," 2021. [Online]. Available: https://doi.org/10.48550/arXiv.1906.10893

[14] J. Domingo-Ferrer, A. Blanco-Justicia, J. Manjón, and D. Sánchez, "Secure and privacy-preserving federated learning via co-utility," *IEEE Internet of Things Journal*, vol. 9, no. 5, pp. 3988–4000, 2021.

[15] M. Panigrahi, S. Bharti, and A. Sharma, "A reputation-aware hierarchical aggregation framework for federated learning," *Computers and Electrical Engineering*, vol. 111, p. 108900, 2023.

[16] L. Gao, L. Li, Y. Chen, C. Xu, and M. Xu, "FGFL: A blockchain-based fair incentive governor for federated learning," *Journal of Parallel and Distributed Computing*, vol. 163, pp. 283–299, 2022.

[17] T. Nguyen, P. Rieger, M. Miettinen, and A.-R. Sadeghi, "Poisoning attacks on federated learning-based IoT intrusion detection system," in *Workshop on Decentralized IoT Systems and Security*, 2020, pp. 1–7.

[18] M. Zang, C. Zheng, T. Koziak, N. Zilberman, and L. Dittmann, "Federated learning-based in-network traffic analysis on IoT edge," in *2023 IFIP Networking Conference (IFIP Networking)*, 2023, pp. 1–6.

[19] E. Hallaji, R. Razavi-Far, M. Saif, and Q. Yang, "Decentralized federated learning: A survey on security and privacy," *IEEE Transactions on Big Data*, vol. 10, pp. 194–213, 2024.

[20] S. Yuan, B. Cao, Y. Sun, Z. Wan, and M. Peng, "Secure and efficient federated learning through layering and sharding blockchain," *IEEE Transactions on Network Science and Engineering*, vol. 11, no. 3, pp. 3120–3134, 2024.

[21] H. Wang, H. Zhang, L. Wang, S. Xuan, and Q. Zhang, "Fedeval: Defending against lazybone attack via multi-dimension evaluation in federated learning," *ACM Transactions on Sensor Networks*, vol. 21, no. 1, pp. 1–23, 2025.

**Isaac Marroquí Penalva** received his M.Sc. degree in Software Engineering from the University of Murcia, Spain, in 2023. His master's thesis focused on the design and evaluation of reputation mechanisms for decentralized federated learning systems, with applications in cybersecurity and the Internet of Things. His current research interests include trust management, collaborative learning, and secure distributed systems.

**Enrique Tomás Martínez Beltrán** is working towards a Ph.D. in Computer Science at the University of Murcia, Spain. His research interests include cybersecurity, IoT, and collaborative learning using FL.

**Manuel Gil Perez** is an Associate Professor in the Department of Information and Communication Engineering of the University of Murcia, Murcia, Spain. His scientific activity is mainly devoted to cyber security, including intrusion detection systems, trust management, privacy-preserving data sharing, and security operations in highly dynamic scenarios.

**Alberto Huertas Celdrán** is an assistant professor at the University of Murcia and a guest researcher at the Communication Systems Group CSG, Department of Informatics IfI, University of Zürich. He received his PhD in Computer Science from the University of Murcia, Spain. His scientific interests include cybersecurity, FL, and computer networks.