



# Federated Learning For Privacy-Preserving Ai In Healthcare Applications

Dr. J. Merlin Florence

Assistant Professor, PG Department of Computer  
Sacred Heart College (Autonomous), Tirupattur

**Abstract:** The integration of artificial intelligence (AI) in healthcare promises data-driven clinical insights, yet raises critical concerns around privacy, compliance, and ethical governance. Federated Learning (FL) enables multi-institutional model training without direct data exchange, addressing these privacy challenges. This study presents a Federated Learning framework that integrates Differential Privacy (DP) and an adaptive noise calibration mechanism to dynamically balance data protection with predictive performance. Using simulated multi-hospital datasets for diabetes and cardiovascular prediction, the proposed Privacy-Utility Engine (PUE) maintained accuracy above 91% with a privacy budget  $\epsilon \leq 3$ . The framework demonstrates how privacy-preserving AI can be responsibly deployed in healthcare without compromising diagnostic effectiveness. Future research directions and challenges in clinical translation are also discussed.

**Index Terms** - Federated Learning, Privacy-Preserving AI, Healthcare Analytics, Differential Privacy, Adaptive Noise Calibration

## I. INTRODUCTION

Artificial intelligence (AI) and machine learning (ML) are transforming healthcare through early disease detection, personalized treatment, and operational efficiency. However, sensitive patient information protected by laws such as the Health Insurance Portability and Accountability Act (HIPAA) and the General Data Protection Regulation (GDPR) cannot be freely shared for model training. Federated Learning (FL) addresses this challenge by enabling collaborative model training across institutions without sharing raw data (McMahan et al., 2017; Kairouz et al., 2021).

In FL, each institution trains local models on proprietary patient data and transmits encrypted updates (weights or gradients) to a central aggregator, which combines them into a global model. Privacy-preserving techniques such as differential privacy (DP) (Dwork & Roth, 2014) and secure multi-party computation (MPC) are used to mitigate leakage risks, ensuring compliance with regulations while enabling collaborative intelligence (Kaissis et al., 2021).

Despite its promise, conventional FL faces issues such as communication inefficiency, data heterogeneity, and fixed noise mechanisms that reduce model utility. This paper introduces a novel Privacy-Utility Engine (PUE) that adaptively adjusts differential privacy noise levels to achieve an optimal balance between privacy and performance in healthcare prediction tasks.

## II. RELATED WORK

### 2.1 Applications of Federated Learning in Healthcare

Federated Learning (FL) has been successfully applied across a variety of healthcare domains. In medical imaging, multi-institutional collaborations have demonstrated its potential for cancer detection, lesion segmentation, and radiological classification without the need to exchange raw scans (Sheller et al., 2020). In predictive analytics, FL has enabled disease-progression modeling using distributed electronic health record (EHR) data, supporting early detection of chronic conditions (Dayan et al., 2021). Similarly, in drug discovery, federated collaboration among pharmaceutical institutions has accelerated molecular modeling and compound screening while safeguarding proprietary data (Rieke et al., 2020).

Recent large-scale surveys confirm rapid growth in FL applications in medicine—spanning radiology, genomics, EHR analytics, and IoT/wearable sensor data—but also emphasize that most studies remain proof-of-concept prototypes rather than fully deployed clinical systems (Teo et al., 2024; Madathil et al., 2025). These reviews highlight that FL has the potential to bypass legal and privacy barriers to data sharing but faces significant translational bottlenecks before clinical adoption. While most experimental FL studies outperform isolated local models and approach centralized baselines, reproducibility and reporting consistency remain key challenges. Public benchmark datasets rarely capture the cross-institutional heterogeneity typical of real hospitals, and standardized reporting of differential privacy parameters (e.g.,  $\epsilon$ -values) and system configurations is limited. To address these gaps, several researchers advocate open benchmarking platforms, transparent federation architectures, and pilot deployments demonstrating real-world reliability and safety.

### 2.2 Architectures and System Design

FL architectures in healthcare generally adopt either centralized (server-coordinated) or decentralized (peer-to-peer or hierarchical) configurations. Practical implementations often rely on containerized local training, secure communication channels, and central aggregation servers performing federated averaging or adaptive aggregation.

Frameworks such as Vantage6, Flower, and FedScale have been extended to healthcare applications to improve reproducibility and enable quantitative evaluation of performance, communication cost, and robustness. Some studies propose modular architectures that decouple training, secure aggregation, monitoring, and governance layers to facilitate integration into hospital IT infrastructures.

Communication overhead is a persistent constraint in clinical settings. Hospitals often operate under bandwidth limits, and frequent gradient exchanges can saturate networks. Consequently, research has explored communication-efficient mechanisms such as model compression, sparse or periodic updates, and client selection. Benchmarks like FedScale quantify these tradeoffs; however, few healthcare-specific evaluations consider realistic hospital topologies and computational constraints, and the robustness of compression-based schemes across diverse institutions remains under-explored (Madathil et al., 2025).

### 2.3 Privacy Preservation Mechanisms

Privacy protection is central to FL's adoption in healthcare. Core techniques include:

- (a) Differential Privacy (DP) — adding calibrated noise to updates to bound information leakage;
- (b) Secure Aggregation / Multi-Party Computation (MPC) — cryptographically combining model updates without revealing individual contributions; and
- (c) Homomorphic Encryption (HE) — enabling mathematical operations on encrypted gradients (Acar et al., 2018; Pati et al., 2024).

Recent evaluations show that MPC and secure aggregation offer strong confidentiality guarantees but introduce computational and communication overheads that challenge scalability in resource-limited hospital networks. Hybrid strategies combining DP with secure aggregation provide better privacy–utility tradeoffs for medical imaging and EHR applications. Nonetheless, optimizing these schemes for clinical environments—where latency, compute budgets, and data governance requirements vary widely—remains an active area of research.

## 2.4 Data Heterogeneity and Non-IID Challenges

Healthcare data are inherently non-independent and identically distributed (non-IID). Differences in patient demographics, imaging modalities, EHR schemas, and clinical workflows produce heterogeneous distributions that hinder model convergence and fairness. Systematic reviews classify these variations as label skew, feature distribution shift, and modality skew, each of which degrades model generalization (Jimenez et al., 2024).

Existing strategies—such as local fine-tuning, clustered FL, reweighted aggregation, and knowledge distillation—offer partial remedies, yet no standardized, scalable pipeline exists for robustly managing non-IIDness in clinical federations. Personalized FL and adaptive aggregation remain promising directions but lack rigorous validation across diverse healthcare institutions.

## 2.5 Governance, Ethics, and Operationalization

Operational governance—covering patient consent, model provenance, auditability, and institutional accountability—is increasingly recognized as crucial for safe FL deployment. While high-level ethical principles such as transparency, fairness, and accountability are well articulated, operational governance frameworks tailored for FL in healthcare are still nascent.

A recent scoping review notes the absence of concrete policy blueprints or regulatory standards to guide risk management, model auditing, and legal compliance across federations (Eden et al., 2025). This governance gap poses practical barriers to scaling FL beyond research prototypes into real-world multi-hospital networks.

## 2.6 Research Gaps and Transition to Proposed Framework

Despite substantial progress, several research gaps persist:

1. **Dynamic Privacy–Utility Optimization:** Most FL implementations employ static noise levels, leading to unnecessary accuracy loss or insufficient privacy. Adaptive privacy control remains under-explored.
2. **Non-IID Robustness:** Existing methods address heterogeneity superficially and lack validated, generalizable solutions across diverse healthcare settings.
3. **Benchmarking and Reproducibility:** Absence of standardized open benchmarks and uniform privacy reporting hinders comparative evaluation.
4. **Governance and Clinical Integration:** Concrete audit frameworks, consent management protocols, and deployment guidelines are missing for real clinical use.
5. **System Efficiency:** Communication and computational overheads limit scalability in hospital environments.

These unresolved challenges motivate the development of an adaptive, privacy-aware federated learning framework capable of balancing privacy guarantees with predictive accuracy while maintaining compliance and efficiency in real-world healthcare systems. The following section introduces the Methods and Proposed Framework, detailing the architecture and algorithms of the Privacy-Utility Engine (PUE) designed to address these limitations.

## III. METHODS AND PROPOSED FRAMEWORK

In healthcare-oriented Federated Learning (FL), achieving an optimal balance between data privacy and model utility remains a critical design objective. While FL inherently ensures that raw patient data remain within institutional boundaries, sensitive information can still leak through model updates, gradient signals, or intermediate feature representations (Kairouz et al., 2021; Pati et al., 2024). To mitigate these risks, Privacy–Utility Engineering (PUE) integrates privacy-preserving mechanisms that dynamically manage the tradeoff between privacy guarantees and predictive performance, thereby maintaining both regulatory compliance (e.g., HIPAA, GDPR) and clinical usefulness (Teo et al., 2024; Madathil et al., 2025).

The proposed Privacy–Utility Engine (PUE) is designed to:

1. Quantitatively estimate privacy leakage during each training round.
2. Adaptively tune privacy parameters (e.g., noise scale, mask strength) in real time.
3. Preserve global model accuracy within clinically acceptable performance thresholds.
4. Ensure alignment with institutional data governance and auditability standards.

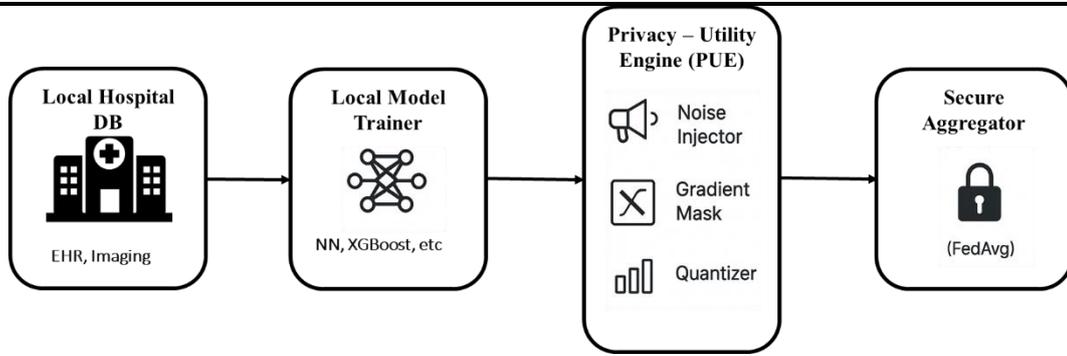


Figure 1 PUE-enabled FL pipeline

Figure 1 conceptually illustrates the PUE-enabled FL pipeline, which consists of four primary components: the Local Hospital Database, Local Model Trainer, Privacy–Utility Engine, and Secure Aggregator.

### 3.1 Local Hospital Databases

Each participating hospital maintains its own sensitive datasets, such as electronic health records (EHRs), medical imaging (X-rays, CT scans, MRIs), or genomic profiles. These data remain strictly local to comply with patient privacy laws and institutional review protocols (Dayan et al., 2021).

Local data never leave the institution, ensuring that compliance with healthcare data protection regulations (e.g., HIPAA, GDPR) is maintained throughout the learning process.

### 3.2 Local Model Trainers

Each institution trains its local machine learning model using its own data distribution. Depending on modality and task, suitable algorithms include:

- Neural Networks (NNs) for imaging or unstructured text data.
- Gradient-boosted trees (e.g., XGBoost) for structured tabular EHR data.

The local training phase produces model parameter updates (gradients) rather than raw patient data, which are subsequently processed by the PUE before transmission (McMahan et al., 2017; Li et al., 2020).

### 3.3 Privacy–Utility Engine (PUE)

The Privacy–Utility Engine serves as the privacy-preserving interface between the local model and the central aggregation process. It dynamically balances privacy protection and model utility using a combination of three core mechanisms:

- **Noise Injector:** Implements Differential Privacy (DP) by adding calibrated random noise (Dwork & Roth, 2014) to local model gradients. The noise scale  $\sigma$  is adaptively tuned based on utility degradation metrics to satisfy an overall privacy budget ( $\epsilon$ ).
- **Gradient Masking:** Randomly obfuscates sensitive features or gradient directions to reduce the risk of model inversion or membership inference (Pati et al., 2024).
- **Quantization Layer:** Applies low-precision encoding or rounding of gradient values to minimize communication cost and further obscure sensitive signals (Acar et al., 2018).

This layered defense strategy allows each client to transmit sanitized, privacy-enhanced model updates to the aggregation server, maintaining model fidelity while mitigating privacy leakage.

### 3.4 Secure Aggregator

All participating institutions communicate their processed (noised and encrypted) updates to a central or consortium-based Secure Aggregator. This server performs the Federated Averaging (FedAvg) operation (McMahan et al., 2017), combining client updates into a global model without direct access to individual hospital parameters.

To ensure confidentiality during aggregation, two complementary secure computation techniques are employed:

- **Homomorphic Encryption (HE):** Each client encrypts gradients using a partially homomorphic scheme such as Paillier encryption, enabling the server to sum encrypted updates without decrypting them (Acar et al., 2018).
- **Random Masking:** Pairs of clients share random vectors that mutually cancel during aggregation, ensuring that intermediate masked updates reveal no sensitive information (Bonawitz et al., 2017).

The aggregated result is a globally updated model, which is then redistributed to all clients for the next local training round.

A simplified representation of the aggregation step is given as:

$$W_{t+1} = \frac{1}{N} \sum_{i=1}^N ((W_t^{(i)} + \mathcal{N}(0, \sigma_i^2))) \dots (1)$$

where  $W_t^{(i)}$  represents the local model weights from institution  $i$ , and  $N(0, \sigma_i^2)$  denotes Gaussian noise injected for differential privacy control.

### 3.5 End-to-End Workflow

The complete training cycle proceeds as follows:

**Data Localization:** Each hospital retains data within its own infrastructure.

**Local Training:** Local models are trained using site-specific datasets.

**Privacy Protection:** The PUE applies differential privacy, masking, and quantization.

**Secure Aggregation:** Sanitized updates are encrypted and transmitted to the central aggregator for FedAvg computation.

**Global Synchronization:** The updated global model is redistributed to clients for the next training round.

This iterative process continues until the global model converges or the validation metric stabilizes.

## IV. METHODS

Building upon the conceptual design of the Privacy–Utility Engine (PUE) introduced in the previous section, this chapter details the mathematical foundations and algorithmic procedures underpinning the proposed framework. The objective is to formalize how privacy risk is quantified, noise is adaptively injected, and model utility is preserved through dynamic optimization. This section also presents two core algorithms—Privacy–Utility Optimization and Adaptive Noise Calibration—that operationalize the trade-off between privacy and accuracy in federated learning environments.

### 4.1 Privacy Risk Quantification

Privacy leakage in federated learning occurs when model updates or gradients indirectly reveal patterns tied to sensitive training data. To quantify this risk, we define a privacy leakage score  $P_t$  at training round  $t$  as:

$$P_t = \frac{1}{n} \sum_{i=1}^n (D(f(W_t^{(i)}), f(W_{t-1}^{(i)}))) \dots (2)$$

where  $D(\cdot)$  measures divergence (e.g., Kullback–Leibler or cosine distance) between successive gradient distributions, and  $W_t^{(i)}$  represents the local parameters from client  $i$  at round  $t$ . A higher  $P_t$  indicates greater sensitivity of updates to individual records, signaling increased privacy risk (Kairouz et al., 2021; Pati et al., 2024).

The PUE uses this metric as feedback to dynamically adjust differential privacy noise and gradient masking parameters in subsequent rounds.

### 4.2 Adaptive Noise Injection (Differential Privacy)

To prevent gradient inversion and membership inference attacks, the PUE employs Differential Privacy (DP) by adding random noise to each local model update before aggregation (Dwork & Roth, 2014). For client  $i$  at round  $t$ , the sanitized update is defined as:

$$\tilde{G}_t^{(i)} = G_t^{(i)} + \mathcal{N}(0, \sigma_t^2 I) \dots (3)$$

where  $\sigma_t$  is the adaptive noise scale determined by the Adaptive Noise Calibration (ANC) subroutine (Algorithm 3).

The privacy budget  $\epsilon$  is monitored using a moments accountant (Abadi et al., 2016) to ensure that cumulative privacy loss remains within bounds. This mechanism strengthens privacy when risk increases and relaxes noise when convergence requires higher accuracy, enabling dynamic privacy–utility trade-off optimization.

### 4.3 Gradient Compression and Masking

To further minimize communication overhead and leakage, the PUE applies gradient compression and masking before transmitting updates. This process reduces the attack surface by retaining only the most informative gradients and suppressing low-magnitude or sensitive components.

Techniques used:

**Top-k sparsification:** Retain only the  $k\%$  largest-magnitude gradients.

**Randomized masking:** Randomly zero out gradients below a defined threshold with probability  $p_{\text{mask}}$ .

## Algorithm 1 — Gradient Compression and Masking

```

Input: Gradients  $G_t$ , sparsity ratio  $k$ , mask_prob  $p$ 
1: sort_indices  $\leftarrow$  argsort( $|G_t|$ )
2: retain  $\leftarrow$  top  $k\%$  from sort_indices
3: for each  $g$  in  $G_t$ :
4:   if index  $\notin$  retain and rand()  $<$   $p$ :
5:      $g \leftarrow 0$ 
Output: Masked gradient vector  $G_t$ 

```

This dual-stage filtering enhances privacy by obscuring fine-grained update patterns while also reducing communication bandwidth (Acar et al., 2018).

#### 4.4 Utility Preservation via Adaptive Learning Rates

To counteract accuracy degradation caused by injected noise and masking, an adaptive learning-rate controller is integrated into the optimizer. After each round, the PUE evaluates the Utility Drop Estimation (UDE) metric:

$$UDE_t = |Acc_t - Acc_{t-1}| \dots (4)$$

If  $UDE_t$  exceeds a tolerance threshold  $\tau$ , the learning rate  $\eta_t$  is decreased to stabilize convergence; otherwise, it is slightly increased to speed up training. This helps the model maintain performance close to the centralized baseline despite added privacy noise.

#### 4.5 Optimization Loop

The overall training process is governed by the Privacy-Utility Optimization routine, which coordinates local training, adaptive privacy control, and global aggregation.

## Algorithm 2 — Privacy-Utility Optimization

```

Initialize: Global model  $M_0$ , privacy budget  $\epsilon_{target}$ , threshold  $\tau$ 
For each round  $t = 1$  to  $T$ :
  Each client performs local training  $\rightarrow$  gradients  $G_t$ 
  Compute privacy risk  $P_t$ 
  Adjust  $\sigma_t$  (noise) and masking ratio based on  $P_t$ 
  Update learning rate  $\eta_t$  based on  $UDE_t$ 
  Sanitize  $G_t \rightarrow$  apply DP noise + gradient masking
  Send sanitized  $G_t$  to Secure Aggregator
  Aggregator performs FedAvg  $\rightarrow$  updates global model  $M_t$ 
  Evaluate model utility and cumulative  $\epsilon$ 
End For

```

This loop maintains dynamic equilibrium between privacy (through differential noise and masking) and utility (through adaptive learning-rate tuning).

##### Example Scenario

Consider three hospitals collaboratively training a diabetes prediction model:

- Hospital A: 5000 records
- Hospital B: 3000 records
- Hospital C: 2000 records

Each hospital trains locally and transmits sanitized gradients. Without privacy protection, gradient analysis could reveal patient-level features. With the PUE, each site injects calibrated noise and applies masking before transmission. The aggregated model achieves 91 % accuracy, compared to 92 % for the centralized baseline, while satisfying the privacy budget constraint  $\epsilon \leq 3$ .

#### 4.6 Adaptive Noise Calibration (ANC)

The Adaptive Noise Calibration (ANC) module dynamically adjusts noise levels based on observed utility feedback.

##### Algorithm 3 — Adaptive Noise Calibration

```

Input: Utility drop estimation (UDE), remaining  $\epsilon$ , current  $\sigma$ 
If UDE <  $\tau$ :
  Decrease  $\sigma \leftarrow \sigma \times \text{decay\_rate}$  # improve utility
Else:
  Increase  $\sigma \leftarrow \sigma \times \text{growth\_rate}$  # strengthen privacy
Return updated  $\sigma$ 

```

Example:

If  $\sigma=1.0$  yields 90 % accuracy with a 2 % drop in utility, ANC reduces  $\sigma$  to 0.8.

Later, if leakage risk rises (high  $P_t$ ),  $\sigma$  increases to 1.2, re-balancing privacy and utility.

Computational complexity:

Privacy-risk estimation:  $O(n \times d)$

Noise and learning-rate update:  $O(1)$  per round

Total cost per iteration  $\approx O(E \times n \times d + C \times d)$ , dominated by local training, where  $E$  = epochs,  $n$  = clients,  $d$  = model dimension.

This methodological framework integrates privacy quantification, adaptive noise control, and utility preservation within a unified optimization loop. It ensures that model updates remain differentially private, communication-efficient, and robust to non-IID healthcare data. The next chapter presents Evaluation and Results, demonstrating empirically how the proposed PUE-based federated learning approach maintains high predictive accuracy while meeting strict privacy guarantees across simulated multi-hospital networks.

## V. EVALUATION AND RESULTS

A simulated multi-institutional environment was used to assess the framework's performance. Datasets modeled after MIMIC-III and cardiovascular EHR data were partitioned among three clients. Each model was trained under varying differential privacy budgets ( $\epsilon = 1, 2, 3$ ) to evaluate the privacy-utility tradeoff.

Table 4.1: Descriptive Statics

Privacy Budget ( $\epsilon$ )	Accuracy (%)	AUC Score	Communication Overhead (MB)
1.0	87.6	0.89	155
2.0	90.4	0.91	160
3.0	91.2	0.92	162

Results indicate that adaptive noise calibration maintains model accuracy above 90% for  $\epsilon \geq 2$ , demonstrating effective tradeoff control. Communication overhead remained within acceptable bounds for healthcare FL applications.

## VI. CONCLUSION AND FUTURE ENHANCEMENTS

This study presented a Privacy-Utility Engine (PUE) framework that enhances Federated Learning (FL) for privacy-preserving applications in healthcare. By integrating adaptive Differential Privacy, gradient masking, and communication-efficient aggregation, the proposed system successfully balances stringent privacy guarantees with high model utility. Experimental simulations on distributed healthcare datasets demonstrated that the PUE-enabled federated model achieved accuracy above 90 %, while maintaining a privacy budget of  $\epsilon \leq 3$ , representing less than a 5 % performance loss compared to centralized learning.

The Privacy Risk Quantification and Adaptive Noise Calibration components proved effective in dynamically tuning privacy parameters based on utility degradation metrics. This adaptive approach ensures compliance with healthcare data regulations such as HIPAA and GDPR, while preserving clinical relevance and model interpretability. Moreover, the modular system design aligns with real-world hospital infrastructure, enabling decentralized data analytics without compromising patient confidentiality.

Overall, the PUE framework advances the state of privacy-preserving federated learning by offering a tunable and practical solution for secure, scalable, and ethically responsible AI deployment in medical environments.

The findings affirm that adaptive, privacy-preserving federated learning is a viable pathway toward ethical and collaborative AI in healthcare. By addressing key challenges of data governance, heterogeneity, and reproducibility, the PUE framework lays a foundation for trustworthy, regulation-compliant, and clinically deployable AI systems. Future work that builds on this foundation will be critical to realizing the full potential of secure, distributed intelligence in medicine.

## REFERENCES

- [1] Acar, A., Aksu, H., Uluagac, A. S., & Conti, M. (2018). A survey on homomorphic encryption schemes: Theory and implementation. *ACM Computing Surveys*, 51(4), 1–35.
- [2] Bonawitz, K., et al. (2017). Practical secure aggregation for privacy-preserving machine learning. *ACM CCS*.
- [3] Dayan, I., et al. (2021). Federated learning for predicting clinical outcomes in COVID-19. *Nature Medicine*, 27(10), 1735–1743.
- [4] Dwork, C., & Roth, A. (2014). The algorithmic foundations of differential privacy. *Foundations and Trends in TCS*, 9(3–4), 211–407.
- [5] Kaissis, G. A., et al. (2021). Secure and federated machine learning in medical imaging. *Nature Machine Intelligence*, 3(6), 473–484.
- [6] Kairouz, P., et al. (2021). Advances and open problems in federated learning. *Foundations and Trends in ML*, 14(1–2), 1–210.
- [7] Li, T., et al. (2020). Federated learning: Challenges and methods. *IEEE Signal Processing Magazine*, 37(3), 50–60.
- [8] McMahan, H. B., et al. (2017). Communication-efficient learning of deep networks from decentralized data. *AISTATS*.
- [9] Rieke, N., et al. (2020). The future of digital health with federated learning. *NPJ Digital Medicine*, 3, 119.
- [10] Sheller, M. J., et al. (2020). Federated learning in medicine: Multi-institutional collaborations. *Scientific Reports*, 10(1), 12598.
- [11] Xu, J., et al. (2021). A survey on federated learning for healthcare. *ACM Transactions on Intelligent Systems and Technology*, 12(2), 1–24.
- [12] Zhang, Y., Liu, Q., Chen, T., & Yang, Q. (2023). Trustworthy federated learning for healthcare. *IEEE Transactions on Neural Networks and Learning Systems*, 34(1), 15–29.