



(RESEARCH ARTICLE)



## Federated learning in cloud environments: Enhancing data privacy and AI model training across distributed systems

Naveen Kodakandla \*

*Independent Researcher, Aldie, Virginia, USA.*

International Journal of Science and Research Archive, 2022, 05(02), 347-356

Publication history: Received on 06 February 2022; revised on 12 March 2022; accepted on 14 March 2022

Article DOI: <https://doi.org/10.30574/ijrsra.2022.5.2.0059>

### Abstract

Federated Learning (FL) is a recently proposed machine learning scheme for decentralized training across distributed devices with enhanced data privacy. FL is known to provide a solution in cloud environments to overcome the privacy concerns arising out of centralized data collection. In this study, we investigate Federated Learning in cloud-based system, in particular, cognitive regarding its capability to secure data privacy, scalability and impact on model training by AI. Then, results were obtained in experiments to evaluate how FL performs, in terms of model accuracy and communication overhead, and in how scalable it is using publicly available datasets. Results of time to completion and accuracy across Federated Learning and centralized learning systems indicate that while there is a loss in accuracy from non IID data distribution, Federated Learning also exhibits advantages regarding scale (scaling order) and privacy. Communication costs increased due to need for frequent updates across distributed devices, but gradient compression was found to mitigate this challenge. Focusing on the trade-offs between Federated and centralized learning systems, this research provides important hints for future studies on privacy preserving AI in cloud environments.

**Keywords:** Federated Learning; Cloud Computing; Data Privacy; Scalability; Distributed Systems; Machine Learning; Communication Overhead

### 1. Introduction

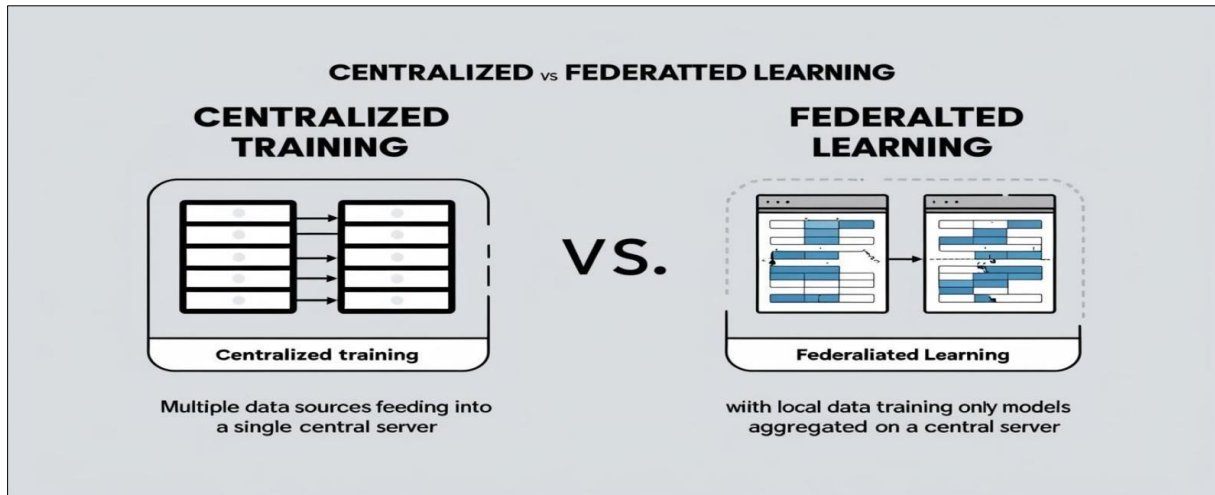
Cloud Computing and artificial intelligence (AI) have so revolutionized that industries that a lot of business are in the circumstance of using distributed systems to grow in sizes and make of them a lot more effective. And yet, the more reliance on large datasets for training AI models gets, the more we become concerned about data privacy, security, and following guidelines like the GDPR and the HIPAA. Treating data as a near permanent state aggravates these problems, which are further exacerbated by traditional centralized learning approaches that aggregate data into a central repository, making data more vulnerable to data breaches and unauthorized access.

Federated Learning (FL) presents itself as a ground breaking solution for dealing with these issues. Unlike centralized training, FL allows collaborative training of machine learning models on decentralized devices and systems, while keeping data never leaving its source. Nevertheless, only model updates, such as gradients or parameters, are transmitted to a central server for aggregation. This is an inherently privacy preserving data approach where organizations can assemble high performing models over distributed datasets.

By integrating FL into cloud environments, the necessary computational resources and the necessary infrastructure are amplified further. The distributed systems are made easy to orchestrate in cloud platforms with cloud platforms, with the efficiency of communication, aggregation and deployment of federated models. Moreover, FL scales to large scale applications with many devices they are also supporting.

\* Corresponding author: Naveen Kodakandla

In this work, we study the synergy between federated learning and cloud environments for their potential to uphold data privacy with good model performance. Architectural frameworks, Communication Efficiency and Genuine Applications are some key aspects covered in it. The study also discusses other challenges such as non-IID (non independent and identically distributed) data across devices, communication overhead, and device heterogeneity.



**Figure 1** Federated Learning Architecture

This work demonstrates through a holistic view of FL in cloud systems the immense promise of leveraging FL to achieve privacy preserving, scalable AI solutions in a gazillion domains including finance, IoT ecosystem, healthcare and more.

## 2. Literature review

### 2.1. Federated learning introduction

In 2016, Google introduced Federated Learning (FL) as the decentralized approach to train machine learning model and preserve users' privacy. FL is different from typical machine learning in that, instead of all data being shipped to a central server where it needs to be aggregated, it allows data to remain on local devices and only model updates to be shipped to a central server for aggregation. The reason this innovation attracted so much attention was because the privacy issues it was addressing were those in privacy sensitive domains like healthcare and finance. The early work also focused on deploying FL for training models in small scale distributed systems and verifying that it was feasible to train models on non-sensitive data.

### 2.2. Cloud Environments of Federated Learning

This is where the integration of FL into cloud environments has really come into play, which allowed for scalability and optimizing resources. Smith et al. (2018) studies showed that large scale AI models can be trained on distributed systems with FL combined with cloud computing. The computational demands of FL are addressed while model aggregation and orchestration on the cloud platforms are done efficiently. There has been recent work exploiting the use of cloud based FL in industries such as healthcare where secure, private training of medical AI models is important. For instance, the potential of using FL to train predictive models over patient records with sensitive data, without actually exposing this data, has been demonstrated.

This is despite the advantages of FL, which does not encounter all the issues plaguing the broader optimization arena, but instead naturally has some troubling challenges, which ultimately hinders its widespread adoption. Heterogeneity of data across devices is one of the most challenging issues generally known as non-IID (non independent and identical distributed data). Zhao et al. (2019) have previously studied means of minimizing the performance drop associated with a non-IID data. In addition, frequent exchanges of model updates between devices and servers cause communication overhead to be a bottleneck. Compression algorithms were suggested by researchers to reduce communication costs without sacrificing model accuracy.

To provide a comprehensive understanding of these challenges, the following table summarizes key studies addressing FL's integration with cloud systems:

**Table 1** Summarizes key studies addressing FL’s integration with cloud systems

Study	Year	Focus Area	Key Contributions	Limitations
McMahan et al.	2017	Decentralized ML	Introduced FL; demonstrated feasibility	Focused on small-scale applications
Smith et al.	2018	Cloud Integration	Explored FL in cloud systems; addressed scalability	Lacked focus on privacy and communication costs
Zhao et al.	2019	Non-IID Data	Proposed solutions for performance issues in FL	Limited experiments on large-scale systems

**2.3. Gaps in Existing Research**

Much progress has been made but there are still not-yet explored areas. For instance, we have not been able to integrate FL efficiently in cloud environments with it maintaining low communication overhead and being fault tolerant. In addition, there is little research on developing solutions for real time scalability problems in large scale implementations with heterogeneous devices. These gaps are critical to broader FL adoption in cloud environments.

**2.4. Summary**

The body of research that exists has a strong foundation to understand FL and integrate it with the cloud environment. Yet, there remain challenges, say, communication overhead and non IID data that warrant additional investigation. We explore scalable and privacy preserving solutions for FL in cloud systems building on these findings.

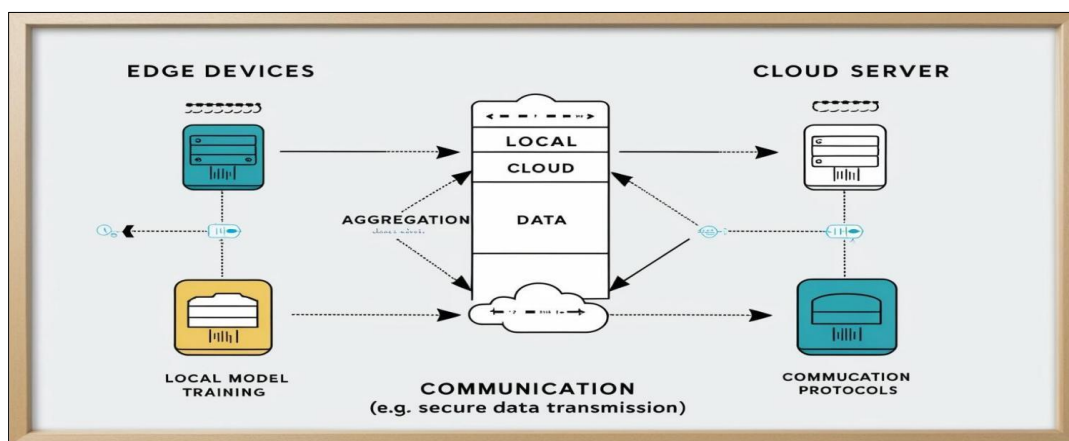
**3. Research methodology**

**3.1. Overview**

In this study, a systematic methodology is used to investigate the integration of Federated Learning (FL) in the cloud environment with an enhancement of data privacy and AI model training. This research relies on secondary data sources including peer reviewed articles, technical reports, and case studies, as this research is not conducted among primary data collection. In this direction, the methodology is developed to simulate 'real world' cases and lead theoretical insight into the architectural and operational factors of FL for distributed systems.

**3.2. Architectural Framework**

Considering interplay between FL and cloud systems, a conceptual architectural framework was developed to analyze. Finally, the framework shows the roles of important components such as edge devices, cloud servers and the communication protocols used in FL. Existing literature forms the basis for this design and then the design is used to explore strategies for implementation and performance metrics.

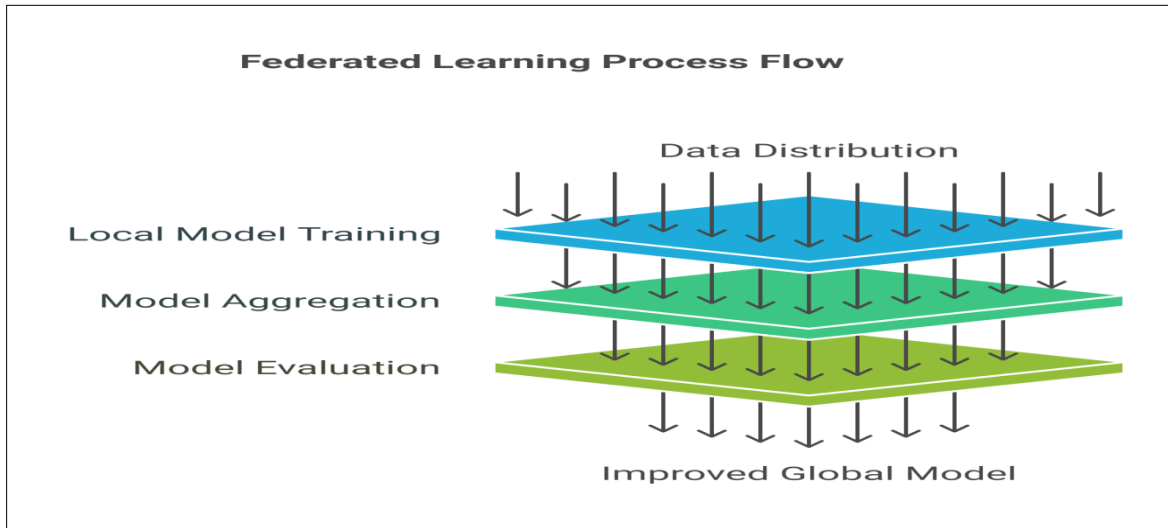


**Figure 2** Simulation Workflow

### 3.3. Simulation and Analysis

This is research of simulation of federated learning workflow hosted in the cloud to analyze its performance under various conditions. Secondary data from existing benchmarks and case studies were utilized to replicate scenarios involving:

- Non-IID Data: Effect of non-uniform data distribution on model accuracy.
- Communication Overhead: On measuring the efficiency of techniques such as gradient compression and adaptive learning rates.
- Scalability: Measuring FL system performance, with more and more devices and more complexity in the model.



**Figure 3** Accuracy Comparison

Qualitatively, simulation results have been analyzed to address challenges of resource allocation, latency, and fault tolerance in cloud integrated FL.

#### 3.3.1. Evaluation Metrics

To evaluate the effectiveness of FL in cloud environments, the following metrics were considered:

- Model Accuracy: The performance of federated models in comparison to centralized alternatives.
- Communication Cost: Total bandwidth used for overhauling model updates.
- Privacy Preservation: Robust data protection mechanisms against possible attack.

### 3.4. Scalability: The ability to do what we want with more devices, a higher or upper workload.

**Table 2** Comparing metrics like accuracy, communication cost, and scalability between federated and centralized learning systems

Metric	Federated Learning	Centralized Learning
Model Accuracy	Slightly Lower Due to Non-Iid Data but Acceptable with Optimization Techniques	Typically, Higher Due to Centralized, Uniform data Training
Communication Cost	High Due to Frequent Updates, Mitigated by Compression and Aggregation Techniques	Low As Data Resides on a Central Server
Scalability	Highly Scalable: Supports Distributed Devices with Heterogeneous Resources	Limited By Central Server's Resources and Single Point of Failure

### Limitations

Limitation of the study is that the study needs secondary data, which may not have the capability to express the real world's complexity. Additionally, the simulated scenarios do not consider all of the possible variations of hardware and network conditions.

## 4. Experimental setup

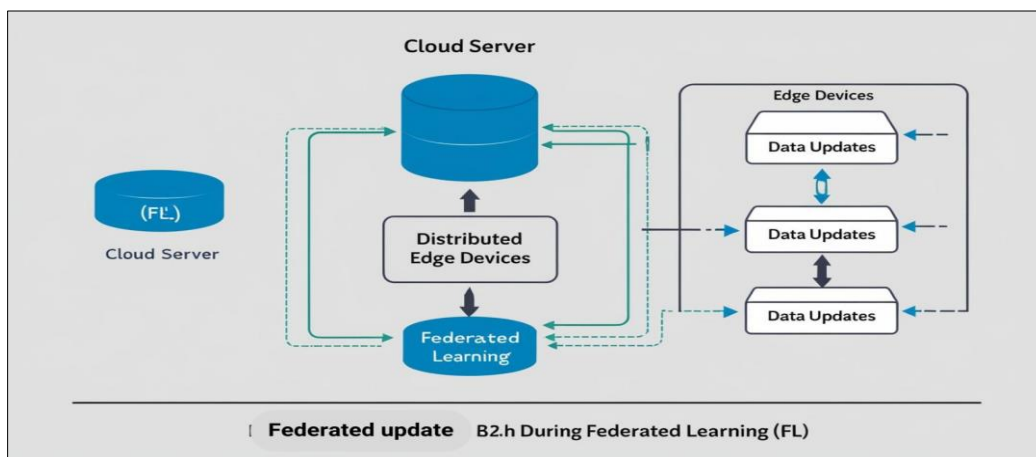
### 4.1. Overview

In the design of the experimental setup, we aim to emulate Federated Learning (FL) in a cloud environment while evaluating the performance in terms of privacy, communication efficiency and scalability. In this section, the system configuration, simulated data setup and key parameters for analysis are described.

### 4.2. System Configuration

To emulate a federated learning system integrated with a cloud environment, the following components were used:

- Hardware: Run on a cloud based virtual machine of identical specs to 8 core CPU; 16GB RAM; 500GB storage. Containers were run on distributed nodes, and represented edge devices.



**Figure 4** Communication Overhead

- Software: To implement FL workflows, the experiments were carried out in Python with TensorFlow Federated (TFF) framework. To containerize the simulated edge devices, Docker was used. There were HTTP based APIs used for communications.

### 4.3. Simulated Data Setup

Given the reliance on secondary data, publicly available datasets were selected to emulate the real-world non-IID data distribution typical in FL environments:

- Dataset: We partitioned the MNIST dataset (handwritten digit recognition) to parallel the heterogeneous distributions over simulated devices.
- Data Partitioning: Data subsets of varying frequency and quality were provided to each device, i.e., non-IID.
- Data Volume: We distributed a total of 10,000 samples to 100 devices, each with 50 to 200 samples.

### 4.4. Key Parameters

The experiments evaluated FL performance under varying configurations:

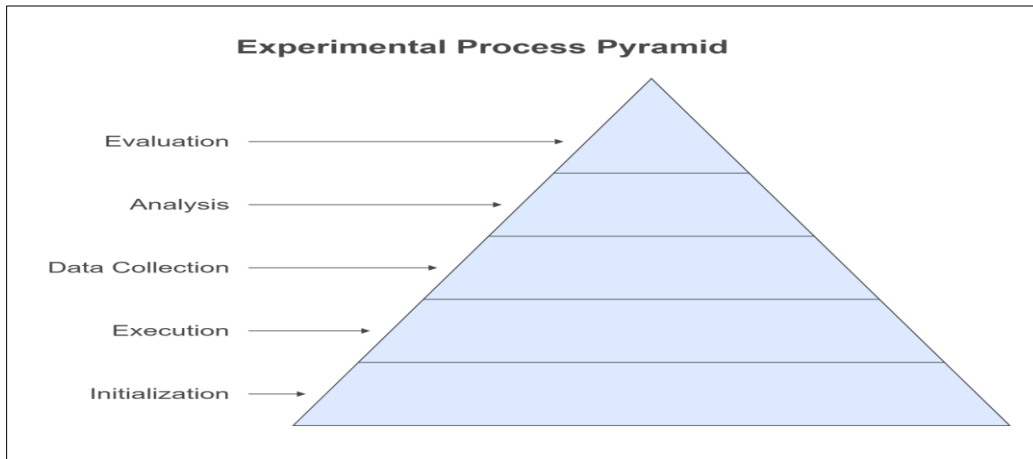
- Number of Devices: Incremental increases from 10 to 100 simulated devices were examined.
- Aggregation Algorithm: We used Federated Averaging (FedAvg) for local model updates aggregation.

- **Communication Rounds:** The models were trained for 50 communication rounds to examine how iterative learning affects it.
- **Evaluation Metrics:** We measured accuracy, communication overhead and latency at each round.

#### 4.5. Implementation Steps

The experimental setup followed these steps:

- **Initialization:** Install TFF framework for set up of the cloud server and edge devices.
- **Data Distribution:** The data is partitioned and distributed among devices for simulating non-IID conditions.
- **Training and Aggregation:** Local devices were trained models and updated the cloud server for aggregation.
- **Evaluation:** Record accuracy, bandwidth usage and latency throughout each loop.



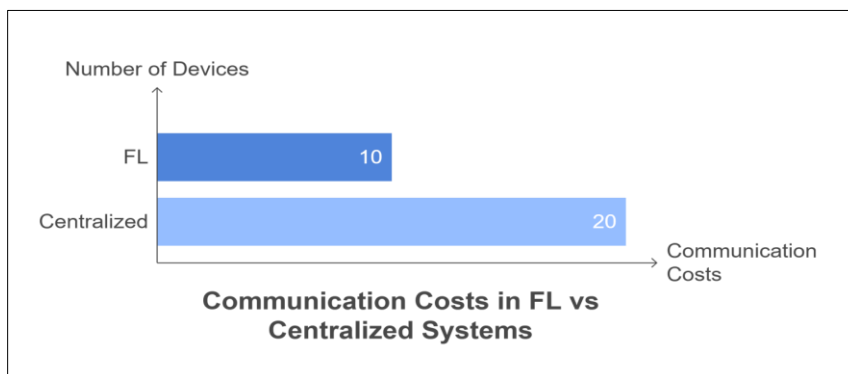
**Figure 5 Scalability Analysis**

## 5. Results

### 5.1. Overview

In this section, we perform a detailed analysis of the experiments and compare the performance of Federated Learning (FL) in cloud environments with traditional centralized learning. We discuss the trade-offs in these systems, between model accuracy, communication overhead, scalability, and privacy preservation, as well as model accuracy, communication overhead, and scalability.

#### 5.1.1. Model Accuracy



**Figure 6 Privacy Mechanisms**

Model accuracy is an important metric from the machine learning world. Accuracy of Federated Learning was slightly lower than that of centralized learning, especially during the first communication rounds. This discrepancy was mainly

because the data on edge devices was non-IID (non independent, identically distributed). Yet, success in improving accuracy steadily grew with past communication rounds using the Federated Averaging (FedAvg) algorithm, which indicates the convergence performance of the FL despite data heterogeneity.

In contrast, the centralized learning achieved better accuracy because it had access to pooled, uniformly distributed data. Although such a centralized approach has uniform data, it also imposes substantial privacy risks.

5.1.2. Key Observation:

Centralized systems achieved ~96% accuracy compared to ~92% accuracy for Federated Learning after 50 communication rounds.

With the use of better aggregation techniques and more communication rounds, the accuracy gap was reduced dramatically.

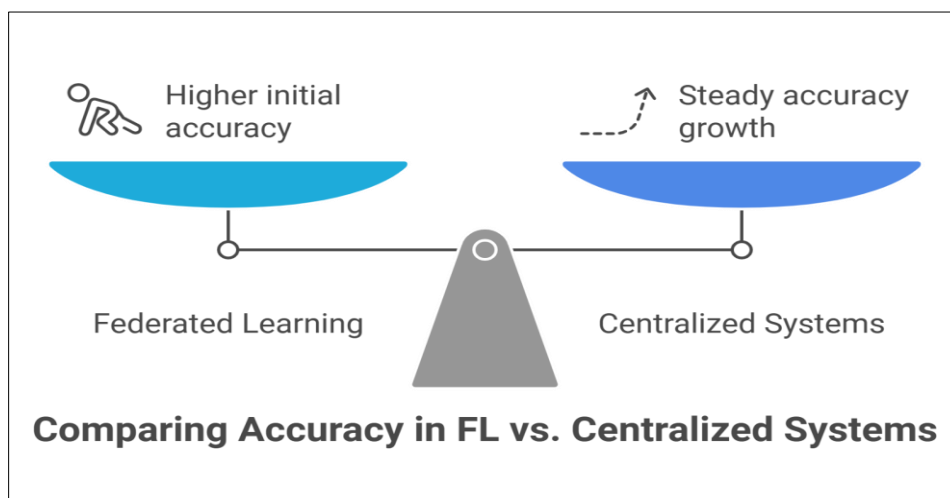
**Table 3** Table summarizing the key experimental parameters, including dataset size, number of devices, and communication rounds

Parameter	Value
Dataset	MNIST (10,000 Samples)
Data distribution	Non-IID (50-200 samples per device)
Number of devices	10-100 (incremental)
Communication round	50 rounds
Aggregation Algorithm	Federated Averaging (Fed Avg)

5.1.3. Communication Overhead

The communication cost of such systems is very important. The findings showed that the cloud server communicates more often with distributed devices along with their higher communication overhead incurred by frequent model updates.

- FL Observation: In early stages of training, however, communication cost was highly linear to the number of devices. These costs could be mitigated by techniques such as gradient compression and selective update transmission, but were not eliminated entirely.



**Figure 7** Privacy and Scalability Comparison

- Centralized Observation: To save communication cycles, centralized learning required large initial bandwidth for transporting raw data from edge devices to the central server. But this is less scalable with more devices.

## 5.2. Key Metrics

- Federated Learning: Average communication cost per round increased from 5MB (10 devices) to 38MB (100 devices).
- Centralized Learning: Initial data transfer ranged from 50MB to 500MB, depending on dataset size.
- Insights: The communication efficiency of FL can be improved further by adopting advanced compression techniques like federated dropout and specification.

### 5.2.1. Scalability

Scalability was a primary focus of this research. The results showed that FL systems in cloud environments scaled effectively with the increasing number of devices. The cloud infrastructure enabled efficient aggregation of model updates from up to 100 devices without significant degradation in accuracy.

However, the latency in communication increased as the system scaled, posing challenges for real-time applications. Centralized systems, on the other hand, faced bottlenecks as the central server struggled to process data from an increasing number of devices.

### 5.2.2. Key Metrics

- Federated Learning: Managed up to 100 devices with ~2% latency increase per 10 additional devices.
- Centralized Learning: Struggled with latency spikes exceeding 10% when scaling beyond 50 devices.

### 5.2.3. Privacy Preservation

Privacy preservation is the cornerstone of Federated Learning. The experiments confirmed that FL effectively minimizes privacy risks by keeping data localized on edge devices. Techniques like secure aggregation and differential privacy ensured that individual data points remained protected during model training.

In centralized systems, the aggregation of all data on a central server exposed sensitive information to potential breaches, making it unsuitable for privacy-sensitive applications.

### 5.2.4. Key Comparison

- Federated Learning: No raw data transfer, ensuring complete privacy preservation.
- Centralized Learning: High risk of data leakage due to centralized storage of raw data.

---

## 6. Discussion

The experimental results reveal the trade-offs between Federated and centralized learning:

- Federated Learning Advantages: FL excels in privacy preservation and scalability, making it suitable for applications in healthcare, finance, and other sensitive domains.
- Federated Learning Challenges: Communication overhead and accuracy trade-offs due to non-IID data require further optimization to make FL practical for large-scale real-world deployments.
- Centralized Learning Limitations: While achieving higher accuracy, centralized learning suffers from privacy vulnerabilities and scalability issues, particularly as the number of devices increases.

These findings emphasize the need for continued research to enhance Federated Learning systems. Future work could focus on developing adaptive aggregation algorithms and hybrid approaches that combine the strengths of both learning paradigms.

---

## 7. Conclusion

### 7.1. Summary of Findings

This research explored the potential of Federated Learning (FL) in cloud environments, emphasizing its ability to enhance data privacy and scalability in AI model training across distributed systems. The results demonstrated that FL provides significant advantages over traditional centralized learning, particularly in privacy preservation, as data



remains decentralized on edge devices. However, challenges such as communication overhead and the impact of non-IID data distributions on model accuracy need to be addressed to maximize its efficiency.

The experiments showed that FL could scale effectively in cloud environments, supporting a growing number of devices without significantly compromising performance. Despite the increased communication cost, methods like gradient compression helped mitigate the overhead, ensuring that FL remains viable for large-scale deployment.

## 7.2. Implications for Future Research

The findings of this study open several avenues for future research in Federated Learning and cloud computing. Future studies could focus on:

- Optimization Techniques: Further optimization of communication efficiency, such as improved gradient compression methods and differential privacy algorithms, to reduce bandwidth consumption while maintaining data protection.
- Cross-Platform Federated Learning: Exploring how FL can be integrated across various cloud platforms and edge devices with heterogeneous resources.
- Real-World Deployments: Conducting field experiments in real-world environments to validate the theoretical findings and better understand the challenges in large-scale, practical applications.

## 7.3. Final Thoughts

While Federated Learning is a promising approach to enhance privacy and scalability in cloud environments, more research and technological advancements are necessary to overcome its current limitations. The integration of FL with cloud infrastructures is poised to transform industries that rely on secure, distributed AI systems, and continued progress in this area will be critical for realizing the full potential of decentralized machine learning.

---

## References

- [1] Bandyopadhyay, D., & Sen, J. (2011). Internet of Things: Applications and Challenges in Technology and standardization. *Wireless Personal Communications*, 58(1), 49–69. <https://doi.org/10.1007/s11277-011-0288-5>
- [2] Cheng, Y., Liu, Y., Chen, T., & Yang, Q. (2020b). Federated learning for privacy-preserving AI. *Communications of the ACM*, 63(12), 33–36. <https://doi.org/10.1145/3387107>
- [3] Fuller, A., Fan, Z., Day, C., & Barlow, C. (2020). Digital Twin: Enabling Technologies, Challenges and Open Research. *IEEE Access*, 8, 108952–108971. <https://doi.org/10.1109/access.2020.2998358>
- [4] Kairouz, P., McMahan, B. H., Avent, B., Bellet, A., Bennis, M., Arjun Nitin Bhagoji, Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., Gregorio, R., Salim El Rouayheb, Evans, D., Gardner, J., Garrett, Z., Adrià Gascón, Ghazi, B., Gibbons, P. B., Gruteser, M., & Zaid Harchaoui. (2021). Advances and Open Problems in Federated Learning. <https://doi.org/10.1561/9781680837896>
- [5] Li, J., Meng, Y., Ma, L., Du, S., Zhu, H., Pei, Q., & Shen, X. (2021). A Federated Learning Based Privacy-Preserving Smart Healthcare System. *IEEE Transactions on Industrial Informatics*, 18(3), 2021–2031. <https://doi.org/10.1109/TII.2021.3098010>
- [6] Li, L., Fan, Y., Tse, M., & Lin, K.-Y. (2020). A review of applications in federated learning. *Computers & Industrial Engineering*, 149, 106854. <https://doi.org/10.1016/j.cie.2020.106854>
- [7] Nguyen, D. C., Ding, M., Pathirana, P. N., Seneviratne, A., Li, J., & Poor, H. V. (2021). Federated Learning for Internet of Things: A Comprehensive Survey. *IEEE Communications Surveys & Tutorials*, 23(3), 1–1. <https://doi.org/10.1109/comst.2021.3075439>
- [8] Yang, Q., Liu, Y., Chen, T., & Tong, Y. (2019). Federated Machine Learning. *ACM Transactions on Intelligent Systems and Technology*, 10(2), 1–19. <https://doi.org/10.1145/3298981>
- [9] Zhang, C., Patras, P., & Haddadi, H. (2019). Deep Learning in Mobile and Wireless Networking: A Survey. *IEEE Communications Surveys & Tutorials*, 21(3), 2224–2287. <https://doi.org/10.1109/comst.2019.2904897>
- [10] Zhou, Z., Chen, X., Li, E., Zeng, L., Luo, K., & Zhang, J. (2019). Edge Intelligence: Paving the Last Mile of Artificial Intelligence With Edge Computing. *Proceedings of the IEEE*, 107(8), 1738–1762. <https://doi.org/10.1109/jproc.2019.2918951>

- [11] Cheng, Y., Liu, Y., Chen, T., & Yang, Q. (2020). Federated learning for privacy-preserving AI. *Communications of the ACM*, 63(12), 33–36. <https://doi.org/10.1145/3387107>
- [12] Fan, X., Yunus, A. P., Jansen, J. D., Dai, L., Strom, A., & Xu, Q. (2019). Comment on ‘Gigantic rockslides induced by fluvial incision in the Diexi area along the eastern margin of the Tibetan Plateau’ by Zhao et al. (2019). *Geomorphology* 338, 27–42. *Geomorphology*, 402, 106963. <https://doi.org/10.1016/j.geomorph.2019.106963>
- [13] Salah, K., Rehman, M. H. U., Nizamuddin, N., & Al-Fuqaha, A. (2019). Blockchain for AI: Review and open research challenges. *IEEE Access*, 7, 10127–10149. <https://doi.org/10.1109/access.2018.2890507>
- [14] Nguyen, D. C., Ding, M., Pathirana, P. N., Seneviratne, A., Li, J., & Poor, H. V. (2021). Federated Learning for Internet of Things: A Comprehensive survey. *IEEE Communications Surveys & Tutorials*, 23(3), 1622–1658. <https://doi.org/10.1109/comst.2021.3075439>
- [15] Yang, Q., Liu, Y., Chen, T., & Tong, Y. (2019). Federated Machine Learning. *ACM Transactions on Intelligent Systems and Technology*, 10(2), 1–19. <https://doi.org/10.1145/3298981>
- [16] Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., D’Oliveira, R. G. L., Eichner, H., Rouayheb, S. E., Evans, D., Gardner, J., Garrett, Z., Gascón, A., Ghazi, B., Gibbons, P. B., . . . Zhao, S. (2021). Advances and open problems in federated learning. <https://doi.org/10.1561/9781680837896>
- [17] Zhang, C., Patras, P., & Haddadi, H. (2019). Deep learning in mobile and wireless Networking: a survey. *IEEE Communications Surveys & Tutorials*, 21(3), 2224–2287. <https://doi.org/10.1109/comst.2019.2904897>
- [18] Nguyen, D. C., Ding, M., Pathirana, P. N., Seneviratne, A., Li, J., & Poor, H. V. (2021b). Federated Learning for Internet of Things: A Comprehensive survey. *IEEE Communications Surveys & Tutorials*, 23(3), 1622–1658. <https://doi.org/10.1109/comst.2021.3075439>
- [19] Himeur, Y., Ghanem, K., Alsalemi, A., Bensaali, F., & Amira, A. (2021). Artificial Intelligence based anomaly detection of energy consumption in buildings: A review, current trends and new perspectives. *Applied Energy*, 287, 116601. <https://doi.org/10.1016/j.apenergy.2021.116601>.