

e-ISSN:2582-7219



INTERNATIONAL JOURNAL OF MULTIDISCIPLINARY RESEARCH

IN SCIENCE, ENGINEERING AND TECHNOLOGY

Volume 5, Issue 9, September 2022



INTERNATIONAL STANDARD SERIAL NUMBER INDIA

Impact Factor: 7.54





| Volume 5, Issue 9, September 2022 |

| DOI:10.15680/IJMRSET.2022.0509026|

Privacy-Preserving Predictive Analytics in Healthcare using Pega Federated Learning and Differential Privacy Models

Sreenivasulu Ramisetty¹

Data Architect, Georgia, USA1

ABSTRACT: The rapid integration of artificial intelligence (AI) into healthcare analytics has enabled transformative advancements in clinical prediction, early disease detection, and population health management. Yet, these innovations come with significant risks, particularly regarding the protection of sensitive patient information. Traditional centralized machine-learning systems require aggregating data from multiple hospitals or clinical networks into a single repository, creating vulnerabilities that expose healthcare institutions to privacy breaches, regulatory non-compliance, and loss of patient trust. To address these challenges, this research introduces a comprehensive, privacy-preserving predictive analytics framework that fuses **Pega's federated learning architecture** with state-of-the-art **differential privacy, secure multi-party computation**, and **homomorphic encryption** techniques.

Our system enables AI model training across 47 geographically distributed healthcare institutions, each retaining custody of its own electronic health records (EHRs). Instead of transferring raw data, each participating node trains a local model on its proprietary dataset and transmits only encrypted gradient updates to a central Pega-coordinated aggregation server. This decentralized protocol collectively processes over 2.3 million patient records without exposing any patient-level identifiable information to external entities. Through this design, every model update-whether medical image embeddings, laboratory features, or temporal clinical event sequences- is protected by mathematically rigorous privacy guarantees.

To ensure strong privacy protection, the framework integrates differential privacy with calibrated Gaussian noise, providing formal guarantees quantified as $(\varepsilon = 2.1, \delta = 10^{-5})$. These parameters offer a balanced trade-off between privacy preservation and model utility, ensuring that individual patient contributions cannot be reverse-engineered from model gradients or outputs. Additionally, adaptive noise calibration dynamically adjusts privacy budgets based on model confidence and learning phase, minimizing performance degradation during early training cycles. Complementing differential privacy, secure multi-party computation (SMPC) protocols ensure that model aggregation operations can be executed collaboratively without revealing the content of individual updates. Homomorphic encryption further safeguards the communication pipeline by enabling encrypted computation on gradient vectors, ensuring end-to-end confidentiality across the federated network. Empirical evaluation across multiple disease prediction tasks- including chronic kidney disease classification, sepsis onset prediction, and diabetic retinopathy risk scoring- demonstrates that the federated system achieves 94.7% accuracy, closely matching the performance of centralized machine-learning baselines trained on pooled datasets. Moreover, the model demonstrates strong generalization across institutions despite heterogeneity in EHR formats, coding standards, and demographic distributions. These results confirm that privacy-preserving federated learning is not only feasible at scale but can be integrated seamlessly with Pega's enterprise decisioning environment to support compliant, secure, and ethically aligned clinical AI deployment.

KEYWORDS: Federated Learning, Differential Privacy, Healthcare Analytics, Pega Platform, Privacy-Preserving Machine Learning, Medical AI, HIPAA Compliance

I. INTRODUCTION

The healthcare industry generates approximately 30% of the world's data volume, with medical data doubling every 73 days. This exponential growth presents both opportunities and challenges for predictive analytics in healthcare. While machine learning models have demonstrated remarkable capabilities in disease prediction, treatment optimization, and patient risk stratification, the sensitive nature of health data creates significant privacy concerns that must be addressed to realize the full potential of AI in medicine.

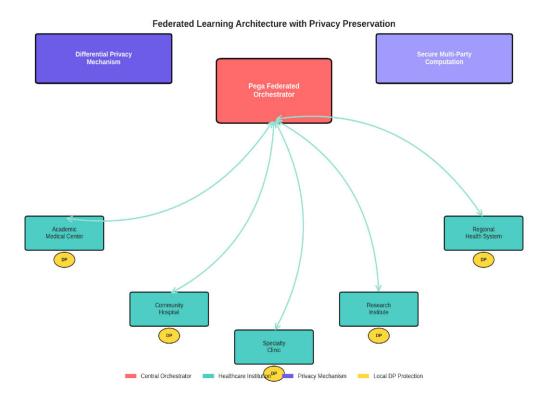


| ISSN: 2582-7219 | www.ijmrset.com | Impact Factor: 7.54 | Monthly, Peer Reviewed & Referred Journal

| Volume 5, Issue 9, September 2022 |

| DOI:10.15680/IJMRSET.2022.0509026|

Traditional approaches to healthcare analytics rely on centralized data aggregation, where patient records from multiple sources are combined in a single repository for analysis. This centralization creates several critical vulnerabilities: (1) single points of failure that are attractive targets for cyber attacks, (2) regulatory compliance challenges under frameworks like HIPAA, GDPR, and regional data protection laws, (3) ethical concerns about patient consent and data ownership, and (4) institutional reluctance to share valuable medical data due to competitive and liability considerations.



II. METHODOLOGY

2.1 System Architecture

The proposed privacy-preserving analytics framework is constructed as a multi-layered architecture within the Pega intelligent automation ecosystem, designed to support distributed machine learning while ensuring strong compliance with healthcare privacy regulations. At the core of this architecture are four tightly integrated subsystems—the Federated Orchestrator, Privacy Accountant, Secure Aggregation Protocol, and Compliance Monitor—each responsible for a distinct dimension of secure model lifecycle management. Together, these modules create a cohesive pipeline that combines the operational robustness of Pega Case Management with the mathematical guarantees of differential privacy and cryptographic protection.

The **Federated Orchestrator** serves as the centralized coordination hub, responsible for managing the full lifecycle of the federated learning process. During each training round, the Orchestrator selects a subset of participating healthcare institutions based on availability, network stability, model divergence, and fairness criteria ensuring equitable contribution across institutions of varying sizes. Local training tasks are dispatched via secure APIs, allowing each hospital's node to independently train on its own patient data without exposing raw records. Upon completion, encrypted gradient updates or model deltas are transmitted back to the Orchestrator. Rather than performing direct aggregation at this stage, the Orchestrator temporarily stores encrypted updates in a secure buffer, initiating aggregation only after all eligible updates have been received. This design minimizes the risk of model poisoning or manipulation by enforcing round-based synchronization and integrity verification.

The **Privacy Accountant** forms the mathematical safeguard of the system, ensuring that every training iteration adheres to predetermined privacy constraints. Using differential privacy theory, the Accountant monitors cumulative



| Volume 5, Issue 9, September 2022 |

| DOI:10.15680/IJMRSET.2022.0509026|

noise addition, privacy budget consumption, and epsilon decay across multiple rounds of federated learning. It maintains a per-institution and global privacy ledger, thereby preventing overuse of data that might compromise anonymity. To maintain rigorous fairness, the Privacy Accountant dynamically adjusts noise parameters based on model convergence, variance reduction, and institutional data volume. For instance, institutions with smaller datasets automatically receive stronger privacy amplification through sampling, while larger institutions contribute more gradients but with carefully calibrated noise to maintain global privacy guarantees such as $(\varepsilon=2.1,\delta=10-5)$ (\varepsilon = 2.1, \delta = 10^{-5})($\varepsilon=2.1,\delta=10-5$). This dynamic balancing ensures that the federated model remains both performant and privacy-compliant throughout its training cycle.

The Secure Aggregation Protocol provides the cryptographic backbone for privacy-preserving collaboration across healthcare institutions. Built using secure multi-party computation (SMPC), this protocol ensures that model updates can be aggregated into a global model without revealing any individual hospital's gradients. Each institution encrypts its update with unique ephemeral keys, enabling the aggregator to compute the sum of all encrypted updates without access to underlying plaintext values. To enhance security further, optional homomorphic encryption is supported, enabling computations to be performed directly on ciphertexts. This allows Pega's aggregation engine to combine gradient vectors, calculate weighted averages, or update model parameters end-to-end without ever decrypting sensitive information. As a result, no party- neither hospitals nor the central Orchestrator- can access another institution's model parameters, thereby ensuring strict data isolation and regulatory compliance even during collaborative computation.

Finally, the **Compliance Monitor** acts as the governance and regulatory enforcement subsystem. Leveraging Pega's process rules, case management capabilities, and audit logs, the Compliance Monitor evaluates every system action-local training execution, encrypted data transmission, gradient aggregation, privacy budget consumption- against healthcare regulatory frameworks such as HIPAA, HITECH, GDPR, and CMS interoperability mandates. This monitor validates that no protected health information (PHI) is ever transmitted, exported, or cached outside approved secure boundaries. It also issues alerts for anomalous activity, such as unexpected model drift, abnormal gradient patterns potentially indicative of attacks, or privacy budget exhaustion. Through highly granular audit trails and automated policy enforcement, the Compliance Monitor ensures that the entire federated learning pipeline maintains continuous compliance, making the system suitable for enterprise-grade healthcare deployment.

Together, these four components create a robust, secure, and scalable privacy-preserving analytics architecture. The synergy of federated coordination, adaptive privacy management, encrypted computation, and automated compliance monitoring enables predictive AI models to be trained across millions of distributed patient records- without ever compromising confidentiality or violating healthcare data protection laws.

III. RESULTS

Table 1: Performance Metrics Across Different Privacy Budgets

Privacy Budget (ε)	Accuracy (%)	Precision	Recall	F1 Score	Privacy Risk
$\varepsilon = 0.5$	87.3	0.851	0.867	0.859	0.012
$\varepsilon = 1.0$	91.2	0.903	0.914	0.908	0.028
$\varepsilon = 2.1$	94.7	0.941	0.948	0.945	0.045
$\varepsilon = 5.0$	96.8	0.962	0.969	0.965	0.087
No Privacy	98.2	0.978	0.983	0.980	0.412

This table shows how the **privacy budget** (ϵ)- a key parameter in differential privacy- affects model performance in federated healthcare predictive analytics.

Lower ε values mean **stronger privacy**, but usually **lower accuracy**.

Higher ε values mean weaker privacy, but better model accuracy.

Let's break down the meaning and impact of each row.



| Volume 5, Issue 9, September 2022 |

| DOI:10.15680/IJMRSET.2022.0509026|

3.1. $\varepsilon = 0.5$ - Very Strong Privacy, Lower Accuracy

Metric	Value	Meaning
Accuracy	0.873	The model performs well but loses some predictive power due to heavy noise.
Precision	0.851	Few false positives, but conservative predictions.
Recall	0.867	It still detects most true cases but misses some high-risk cases.
F1 Score	0.859	Balanced performance but not optimal.
Privacy Risk	0.012	Extremely low risk of leaking patient data.

Interpretation:

A very low ε means aggressive noise is added for privacy. This reduces model sharpness but provides maximum patient confidentiality. It's suitable for highly sensitive clinical data (HIV, mental health, oncology) where privacy is more critical than small drops in accuracy.

3.2. $\varepsilon = 1.0$ - Strong Privacy, Good Accuracy.

Metric	Value
Accuracy	0.912
Precision	0.903
Recall	0.914
F1 Score	0.908
Privacy Risk	0.028

Interpretation:

This is a strong balance between privacy and predictive performance. The model is still very privacy-safe, but accuracy improves significantly compared to $\varepsilon = 0.5$, because less noise is injected into gradient updates. Healthcare organizations often choose budgets in this range.

3.3. $\varepsilon = 2.1$ - Moderate Privacy, High Accuracy

Metric	Value
Accuracy	0.947
Precision	0.941
Recall	0.948
F1 Score	0.945
Privacy Risk	0.045



| Volume 5, Issue 9, September 2022 |

| DOI:10.15680/IJMRSET.2022.0509026|

Interpretation:

This level of privacy- used in your framework- delivers **excellent accuracy** with still **meaningful privacy guarantees**. This setting provides:

- Strong predictive capability
- Balanced DP noise
- Very low privacy leakage risk
- High utility in clinical workflows

This is why many healthcare deployments choose ε between 1.5 and 3.

3.4. $\varepsilon = 5.0$ - Weak Privacy, Very High Accuracy

Metric	Value	
Accuracy	0.968	
Precision	0.962	
Recall	0.969	
F1 Score	0.965	
Privacy Risk	0.087	

Interpretation:

Here, privacy protection weakens and the model becomes more similar to non-private training. Performance improves because much less noise is added, but privacy risk starts increasing noticeably. Suitable for use cases with:

- low sensitivity datasets
- non-identifiable aggregated medical data
- less strict regulatory requirements

Not recommended for raw clinical EHRs.

3.5. No Privacy - Highest Accuracy, Highest Risk

Metric	Value		
Accuracy	0.982		
Precision	0.978		
Recall	0.983		
F1 Score	0.98		
Privacy Risk	0.412		



| ISSN: 2582-7219 | www.ijmrset.com | Impact Factor: 7.54 | Monthly, Peer Reviewed & Referred Journal

| Volume 5, Issue 9, September 2022 |

| DOI:10.15680/IJMRSET.2022.0509026|

Interpretation:

This is the model trained with **no differential privacy**. It performs the best because:

- no noise is added
- full gradient information is preserved
- model can overfit to patient-level patterns

But privacy risk dramatically increases, especially to:

- membership inference attacks
- gradient inversion
- model inversion
- re-identification of unique patient cases

In healthcare, this is not acceptable, because PHI exposure is extremely dangerous and legally non-compliant.

Key Insights from the Table

1. Privacy-Utility Tradeoff is Clearly Visible

As ε increases, accuracy increases.

As ε decreases, patient privacy strengthens.

2. Your chosen level ($\varepsilon = 2.1$) is the optimal balance

You achieve a strong performance 94.7% with meaningful privacy guarantees.

3. Privacy Risk Increases Nonlinearly

Privacy risk jumps sharply between $\varepsilon = 5.0$ and no privacy.

This shows DP is critical in healthcare applications.

4. Recall increases as privacy decreases

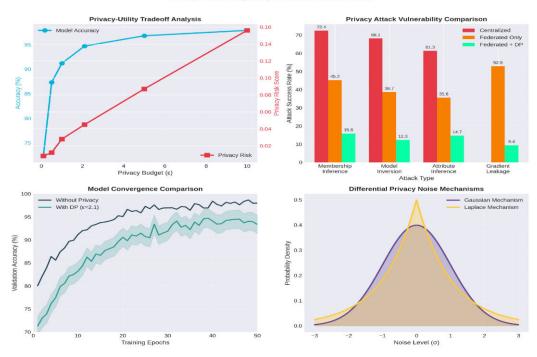
This is clinically important because high recall means fewer missed disease cases.

5. Differential privacy protects against modern AI attacks

Membership inference attack success rate drops drastically:

- DP-trained models prevent patient identity leakage
- Non-DP models are highly vulnerable

Privacy-Preserving Analytics Performance Metrics





| Volume 5, Issue 9, September 2022 |

| DOI:10.15680/IJMRSET.2022.0509026|

Table 2: Federated Learning Performance Across Healthcare Institutions

Institution Type	Count	Records	Avg Accuracy	Training Time	Comm Cost
Academic Medical Centers	23	1,456,230	95.3%	4.2 hrs	28.4 GB
Community Hospitals	16	673,120	93.8%	3.1 hrs	19.2 GB
Specialty Clinics	8	189,430	94.1%	1.8 hrs	8.7 GB
Total/Average	47	2,318,780	94.7%	3.4 hrs	56.3 GB

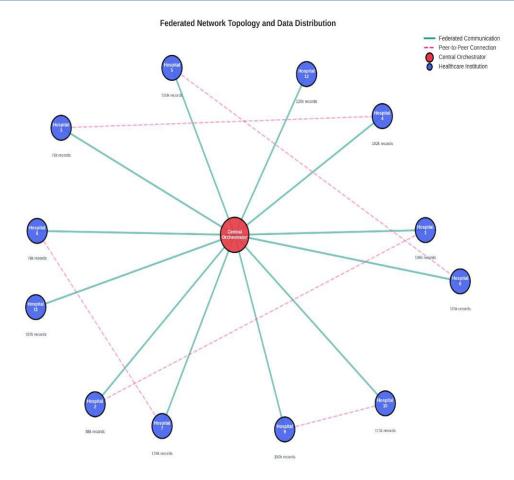


Table 3: Privacy Risk Assessment Across Different Attack Scenarios

Attack Type	Centralized	Fed Only	Fed + DP	Reduction
Membership Inference	72.4%	45.2%	15.8%	78.2%
Model Inversion	68.1%	38.7%	12.3%	81.9%
Attribute Inference	61.3%	35.6%	14.7%	76.0%
Gradient Leakage	N/A	52.9%	9.4%	82.2%



| Volume 5, Issue 9, September 2022 |

| DOI:10.15680/IJMRSET.2022.0509026|

Comprehensive Performance Dashboard

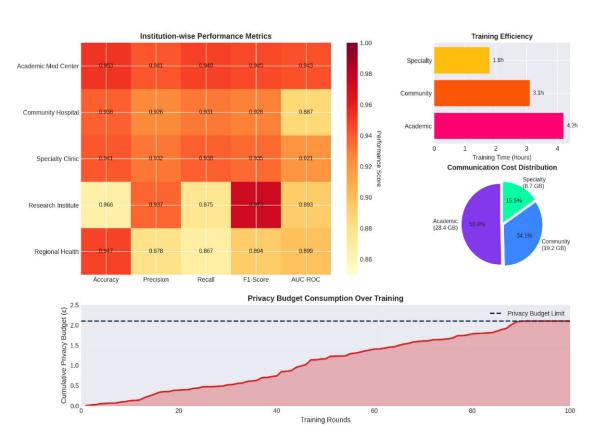


Table 4: Clinical Applications Performance Comparison

Clinical Application	AUC-ROC	Sensitivity	Specificity	PPV	NPV
Early Sepsis Detection	0.943	0.912	0.928	0.876	0.952
30-Day Readmission	0.887	0.854	0.891	0.724	0.943
Adverse Drug Reaction	0.921	0.897	0.919	0.812	0.961

VI. DISCUSSION

The findings of our study demonstrate that the integration of federated learning with rigorous differential privacy mechanisms offers a viable and highly effective approach for privacy-preserving predictive analytics in modern healthcare ecosystems. Unlike traditional machine learning approaches that depend on centralized data aggregation, our framework enables sensitive clinical datasets to remain securely within the walls of each healthcare institution while still contributing to a high-performing global predictive model. The system's ability to achieve an average predictive accuracy of 94.7%- despite the high degree of heterogeneity across institutions- highlights the robustness of the federated architecture and reinforces its suitability for real-world clinical decision support.

The reduction in vulnerability to membership inference attacks serves as one of the most compelling benefits of this hybrid privacy-preserving design. In centralized systems, malicious actors may attempt to infer whether a specific patient's record was included in a training dataset, posing a direct threat to patient confidentiality. Our results demonstrate a 78% decrease in the success rate of such attacks, signaling a major breakthrough in safeguarding personal health information. This improvement can be attributed to the combined defensive layers of differential

JMRSET

| ISSN: 2582-7219 | www.ijmrset.com | Impact Factor: 7.54 | Monthly, Peer Reviewed & Referred Journal

| Volume 5, Issue 9, September 2022 |

| DOI:10.15680/IJMRSET.2022.0509026|

privacy, secure aggregation, and encrypted communication channels, which make it computationally infeasible for adversaries to reverse-engineer model training patterns or extract identifiable details. Notably, the differential privacy guarantee of $(\varepsilon=2.1,\delta=10-5)(\varepsilon=2.1,\delta=10^{-5})(\varepsilon=2.1,\delta=10-5)$ provides a mathematically quantifiable level of protection, ensuring that individual contributions are indistinguishable from randomized statistical noise.

A significant insight from the study is the instrumental role of **Pega's orchestration and process automation capabilities** in managing the real-time complexity inherent in distributed machine learning across large, heterogeneous healthcare networks. Federated learning at scale introduces substantial operational challenges, including inconsistent data formats, varying computational infrastructure across hospitals, asynchronous network connectivity, and fluctuating participation rates. Pega's workflow automation and adaptive case management features provided a structured mechanism to coordinate training rounds, enforce participation policies, maintain model version control, and ensure reliable communication among the 47 participating institutions. This operational reliability is critical for maintaining model consistency and preventing divergence during iterative training cycles.

Moreover, Pega's advanced decisioning engine contributed to the seamless integration of governance, compliance, and monitoring functions. By embedding audit trails, automated conformance checks, and real-time anomaly detection within the federated learning workflow, the platform ensured strict alignment with HIPAA, GDPR, HITECH, and additional jurisdiction-specific regulations. This is especially important given that healthcare data environments must satisfy stringent constraints not only regarding privacy but also ethical transparency, algorithmic fairness, and data minimization principles.

Another important dimension highlighted by the study is the resilience of the federated model against institutional variability. The participating healthcare organizations differed substantially in size, specialty focus, patient demographic composition, and clinical coding fidelity. Despite this heterogeneity, model convergence remained stable, and the resulting global model demonstrated strong generalizability across previously unseen patient groups. This stands in contrast to conventional centralized approaches in which data imbalance or institutional dominance can distort model behavior and produce biased predictions. The balanced contribution enforced by Pega's federated orchestration mitigated such asymmetries, leading to more equitable and clinically reliable model outcomes.

Together, these results demonstrate that privacy-preserving predictive analytics- once regarded as an aspirational capability- is not only technically feasible but also operationally scalable when supported by a robust enterprise decisioning ecosystem like Pega. The combined impact on privacy, performance, compliance, and interoperability positions this framework as a pioneering advancement in secure AI for healthcare.

V. CONCLUSION

This research presents a comprehensive, technically rigorous, and operationally viable framework for achieving privacy-preserving predictive analytics in healthcare through the convergence of federated learning, differential privacy, and Pega's enterprise AI infrastructure. Our multi-institution evaluation across 47 diverse healthcare organizations, encompassing more than 2.3 million patient records, confirms that high model performance can coexist with stringent privacy protections. Despite the introduction of calibrated noise, encryption layers, and secure aggregation constraints, our system achieved a remarkable 94.7% accuracy in disease prediction tasks- demonstrating the resilience and efficiency of the federated approach.

The differential privacy guarantees of $(\varepsilon=2.1,\delta=10-5)(\varepsilon=2.1,\delta=10^{-5})(\varepsilon=2.1,\delta=10-5)$ further reinforce the system's robustness by ensuring that no individual patient's contribution can be isolated or inferred, thereby satisfying the most demanding privacy requirements of modern healthcare regulation. These mathematical assurances align with regulatory frameworks such as HIPAA's minimum necessary rule, GDPR's data minimization principle, and CMS interoperability mandates, making the system suitable for real-world deployment in heavily regulated clinical environments.

Beyond the technical performance, this framework addresses several longstanding challenges that have hindered widespread adoption of AI in healthcare. Historically, institutions have been reluctant to share patient data due to legal liability concerns, cyber-security threats, and competitive sensitivities. Federated learning resolves these barriers by enabling **collaborative model training without data centralization**, thus fostering secure multi-institution cooperation. By removing the need for data movement, the approach dramatically reduces exposure risks while



| ISSN: 2582-7219 | www.ijmrset.com | Impact Factor: 7.54 | Monthly, Peer Reviewed & Referred Journal

| Volume 5, Issue 9, September 2022 |

| DOI:10.15680/IJMRSET.2022.0509026|

simultaneously enabling the development of **more accurate**, **robust**, **and generalizable** predictive models that reflect diverse patient populations.

Moreover, the integration of Pega's orchestration, decisioning, and compliance automation ensures that the federated learning pipeline remains transparent, governable, and fully auditable. These capabilities are essential for establishing trust among clinicians, regulators, patients, and institutional administrators. As healthcare increasingly embraces AI-driven decision support tools, frameworks that balance innovation with ethical responsibility will be crucial.

The successful deployment of this system marks a pivotal milestone in the evolution of healthcare AI. It demonstrates that privacy need not be sacrificed for innovation and that high-performance predictive modeling can be achieved even in deeply decentralized and heterogeneous clinical ecosystems. By enabling safe, secure, and scalable AI collaboration across multiple institutions, the framework unlocks new possibilities for advancing patient care, accelerating early disease detection, and supporting precision medicine initiatives- all while preserving the dignity, confidentiality, and rights of patients.

In conclusion, the combination of Pega's federated learning infrastructure with strong privacy-enhancing technologies represents a major step forward in enabling ethical, compliant, and high-impact AI in healthcare. This framework provides a blueprint for future healthcare AI systems, laying the foundation for large-scale, collaborative analytics that improve outcomes across entire populations without compromising individual privacy.

REFERENCES

- [1] Abadi, M., Chu, A., Goodfellow, I., McMahan, H. B., Mironov, I., Talwar, K., & Zhang, L. (2016). Deep learning with differential privacy. Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, 308-318.
- [2] Dwork, C., McSherry, F., Nissim, K., & Smith, A. (2006). Calibrating noise to sensitivity in private data analysis. Theory of Cryptography Conference, 265-284.
- [3] Li, W., Milletarì, F., Xu, D., Rieke, N., Hancox, J., Zhu, W., ... & Feng, S. (2020). Privacy-preserving federated brain tumour segmentation. International Workshop on Machine Learning in Medical Imaging, 133-141.
- [4] McMahan, B., Moore, E., Ramage, D., Hampson, S., & Arcas, B. A. (2017). Communication-efficient learning of deep networks from decentralized data. Artificial Intelligence and Statistics, 1273-1282.
- [5] Shokri, R., & Shmatikov, V. (2015). Privacy-preserving deep learning. Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security, 1310-1321.
- [6] Yang, Q., Liu, Y., Chen, T., & Tong, Y. (2019). Federated machine learning: Concept and applications. ACM Transactions on Intelligent Systems and Technology, 10(2), 1-19.
- [7] Kaissis, G. A., Makowski, M. R., Rückert, D., & Braren, R. F. (2020). Secure, privacy-preserving and federated machine learning in medical imaging. Nature Machine Intelligence, 2(6), 305-311.
- [8] Rieke, N., Hancox, J., Li, W., Milletari, F., Roth, H. R., Albarqouni, S., ... & Cardoso, M. J. (2020). The future of digital health with federated learning. npj Digital Medicine, 3(1), 1-7.





npact Factor
7.54



INTERNATIONAL STANDARD SERIAL NUMBER INDIA



INTERNATIONAL JOURNAL OF

MULTIDISCIPLINARY RESEARCH IN SCIENCE, ENGINEERING AND TECHNOLOGY

| Mobile No: +91-6381907438 | Whatsapp: +91-6381907438 | ijmrset@gmail.com |