

University of Southampton Research Repository

Copyright © and Moral Rights for this thesis and, where applicable, any accompanying data are retained by the author and/or other copyright owners. A copy can be downloaded for personal non-commercial research or study, without prior permission or charge. This thesis and the accompanying data cannot be reproduced or quoted extensively from without first obtaining permission in writing from the copyright holder/s. The content of the thesis and accompanying research data (where applicable) must not be changed in any way or sold commercially in any format or medium without the formal permission of the copyright holder/s.

When referring to this thesis and any accompanying data, full bibliographic details must be given, e.g.

Thesis: Author (Year of Submission) "Full thesis title", University of Southampton, name of the University Faculty or School or Department, PhD Thesis, pagination.

Data: Author (Year) Title. URI [dataset]

UNIVERSITY OF SOUTHAMPTON

Faculty of Engineering and Physical Science
Electronics and Computer Science

**Privacy-Preserving Federated Learning
Framework**

by

Wenxuan Huang

MSc

ORCID: [0000-0002-2613-2672](https://orcid.org/0000-0002-2613-2672)

*A thesis for the degree of
Doctor of Philosophy*

27th June 2025

University of Southampton

Abstract

Faculty of Engineering and Physical Science
Electronics and Computer Science

Doctor of Philosophy

Privacy-Preserving Federated Learning Framework

by Wenxuan Huang

Federated Learning (FL) enables multiple parties to collaboratively train a shared model without exposing raw data. However, FL often suffers from degraded performance due to three core challenges: non-IID data distributions, residual privacy risks, and limited data volume per node. This work explores the privacy–performance trade-offs arising from these challenges and proposes a unified framework combining partial data sharing, anonymisation, and data augmentation.

We first investigate how non-IID distributions impair standard FL convergence, and propose a Cluster-of-Trust (CoT) mechanism that groups mutually trusted clients for hierarchical aggregation. Experiments on the IoT-23 intrusion detection dataset show that CoT significantly improves model accuracy and convergence speed, closing most of the gap to centralised training.

Next, we compare combinations of privacy-preserving techniques—including differential privacy, data anonymisation, and CoT-based grouping—highlighting how these impact model utility. Our results demonstrate that strong privacy guarantees can be achieved with only modest accuracy loss and reduced communication overhead.

To address the limited data volume problem, we introduce data augmentation strategies and test their effectiveness in a privacy-sensitive Human Activity Recognition (HAR) task. We find that combining augmentation with trust-aware aggregation improves generalisation under strict privacy constraints.

Across two domains—cybersecurity and HAR—we systematically evaluate this framework and demonstrate its effectiveness and generalisability. Our results show that strong privacy can be reconciled with competitive performance, offering practical insights for deploying FL in real-world, heterogeneous environments.

Contents

List of Figures	xiii
List of Tables	xv
Declaration of Authorship	xvii
Acknowledgements	xix
Publications	xxi
Definitions and Abbreviations	xxi
1 Introduction	1
1.1 Problem Statement and Research Gap	2
1.2 Research Questions and Hypotheses	2
1.2.1 Research Question 1: Trust-Based Clustering Effectiveness	3
1.2.2 Research Question 2: Composite Privacy Protection Mechanisms	3
1.2.3 Research Question 3: Privacy-Preserving Data Augmentation	3
1.3 Research Motivation and Application Context	4
1.4 Research Contributions	5
1.5 Thesis Structure	6
2 Literature Review	7
2.1 Literature Review Methodology	7
2.2 Federated Learning	8
2.2.1 Development of Federated Learning	8
2.2.2 Key Concepts of Federated Learning	8
2.2.2.1 System Architecture	8
2.2.2.2 Categories of Federated Learning	9
Deployment paradigm.	9
Data partitioning.	10
2.2.2.3 Model Aggregation Algorithms	10
FedAvg.	10
FedProx.	11
SCAFFOLD.	11
FedOpt family.	11
FedNova.	11
2.2.3 Challenges of Federated Learning	11

2.2.3.1	Non-IID Data	11
	Existing remedies.	12
	Limitations & gaps.	12
2.2.3.2	Data Limitations	12
	Causes & manifestations.	12
	Existing remedies.	12
	Limitations & gaps.	13
2.2.3.3	Security and Robustness	13
	Defence techniques.	13
	Limitations & gaps.	14
2.2.4	Summary and Future Trends	14
2.3	Privacy-Preserving Techniques	14
2.3.1	Differential Privacy	14
2.3.1.1	Principles of Differential Privacy	14
2.3.2	Privacy-Preserving Data Publishing	15
2.3.2.1	Anonymisation Techniques	16
	k-Anonymity	16
	l-Diversity	16
	t-closeness	16
2.3.2.2	Critical discussion: anonymisation and downstream model utility	17
2.4	Intrusion Detection in the Internet of Things	18
2.4.1	Public Datasets and Benchmarking Status	19
2.4.1.1	Traditional Traffic Benchmarks	19
2.4.1.2	Datasets from Real IoT Environments	19
2.4.2	Feature Engineering and Protocol Parsing	19
2.4.2.1	Traffic Statistical Features	20
2.4.2.2	Protocol Features	20
2.4.2.3	Embedded Event Logs	20
2.4.3	Classical Machine Learning Approaches	20
2.4.3.1	Lightweight Classifiers	20
2.4.4	Deep Learning Methods	20
2.4.4.1	Convolutional and Recurrent Networks	20
2.4.4.2	Autoencoders	21
2.4.4.3	Attention Mechanisms	21
2.4.5	Federated Learning for IoT-IDS	21
2.4.5.1	Representative FL-based IoT-IDS studies	21
2.4.5.2	Non-IID and heterogeneity: what remains challenging	22
2.5	Human Activity Recognition	22
2.5.1	Datasets	23
2.5.1.1	Public Sensor Datasets	23
2.5.1.2	Multimodal Datasets	23
2.5.2	Feature Engineering	23
2.5.2.1	Time-Domain Statistical Features	23
2.5.2.2	Frequency and Time-Frequency Features	24
2.5.2.3	Video Features	24
2.5.3	Classical Machine Learning Methods	24

2.5.3.1	Traditional Classifier Comparison	24
2.5.4	Deep Learning	24
2.5.4.1	CNN/LSTM Sequential Networks (DeepConvLSTM, HAR-CNN)	24
2.5.4.2	Transformer and Self-Attention Models (HAR-Transformer, Perceiver-IO)	25
2.5.5	Data Augmentation	25
2.5.6	Federated Learning in HAR	26
2.5.6.1	Critical discussion: privacy, clustering, and scale in federated HAR	26
2.6	Synthesis: State-of-the-Art Limitations and Thesis Positioning	26
3	Methodology and Framework	29
3.1	Research Problem	29
3.1.1	Non-IID Data Distribution Problem Modelling	29
3.1.2	Privacy-Utility Trade-off Challenges	30
3.1.3	Data Scarcity Constraint Challenges	31
3.2	Framework Design	32
3.2.1	Architecture Overview	32
	Trust as a collaboration constraint; similarity as an experimental factor.	32
3.2.2	Cross-Domain Adaptation Strategies	33
3.3	Experimental Methodology	35
3.3.1	Experimental Design Overview	35
3.3.2	Comprehensive Evaluation Framework	37
3.4	Summary	38
	Positioning and novelty.	39
4	Federated Learning for IoT Intrusion Detection with Non-IID Data	41
4.1	Introduction	41
	Positioning and novelty.	42
4.2	A Data Sharing Strategy for IoT IDS Based on Federated Learning	43
	Problem Formulation	44
	Assumptions (IoT-IDS case study).	44
	Algorithm Overview	45
4.3	Experimental setup	47
4.3.1	Dataset Description and Preprocessing	47
4.3.1.1	IoT-23 Dataset Overview	47
4.3.1.2	Data Preprocessing	48
4.3.1.3	Dataset Characteristics	49
4.3.1.4	Temporal and Protocol Analysis	50
4.3.1.5	Scenario-Level Analysis	50
4.3.2	Model Configuration	50
4.3.3	Baseline scenarios	54
4.3.4	Data Sharing Strategies	54
4.4	Experimentation Result	54
4.4.1	Federated learning in IoT intrusion detection baseline scenarios	54

4.4.2	Clusters of Trust vs Globally shared data	56
4.5	Conclusions	58
	Generalisability and external validity.	58
5	A Dual-Layer Privacy-Preserving Federated Learning Framework	61
5.1	Introduction	61
5.2	Methodology and framework	65
5.2.1	Data Description	65
5.2.2	Anonymisation methods	65
5.2.2.1	Generalisation	66
5.2.2.2	Microaggregation	66
5.2.3	A dual-layer privacy-preserving federated learning framework .	66
5.2.4	Privacy against centralisation	67
5.3	Experiments	69
5.3.1	Privacy metrics	69
5.3.2	The impact of anonymisation on model performance	69
5.3.3	Trade-off in the Dual-Layer Framework	71
5.3.4	Evaluation and discussion	73
5.4	Conclusions and future work	74
	Generalisability and external validity.	74
6	Privacy-Preserving Federated Learning for Human Activity Recognition	75
6.1	Introduction	75
	Positioning and novelty.	76
6.2	Methodology	78
6.2.1	Formal Problem Definition	78
6.2.1.1	Fall Detection Task Definition	78
	Justification and implications of binary formulation.	78
6.2.1.2	Federated Learning Environment	79
6.2.1.3	Non-IID Challenges Formalisation	79
6.2.2	Overall Framework Overview	79
6.2.3	Sensor Data Rotation Augmentation	80
6.2.3.1	Rotation Transformation Principles	80
6.2.3.2	Virtual Node Creation	81
6.2.3.3	Privacy Protection Features and Theoretical Guarantees	82
6.3	Experimental Setup	82
6.3.1	Dataset Description	82
6.3.1.1	SisFall Dataset Overview	82
6.3.1.2	Data Preprocessing	84
	Data Cleaning	84
	Downsampling	84
	Data Normalisation	84
	Format Conversion	84
	Label Generation	84
	Dataset Components	84
	Dataset Statistics	85
6.3.2	Model Configuration	85

	Rationale for the chosen model	85
6.3.2.1	Network Architecture	85
6.3.2.2	Attention Mechanism Configuration	86
6.3.2.3	Feature Fusion and Classification	87
6.3.2.4	Training Configuration	87
6.3.2.5	Regularization and Optimization Strategies	87
6.3.3	Experimental Design	88
6.3.3.1	Baseline Methods	88
6.3.3.2	Comparison Methods: Data Sharing Approaches	88
6.3.3.3	Proposed Methods	89
6.3.4	Evaluation Metrics	89
6.3.4.1	Model Metrics	89
6.4	Results and Analysis	90
6.4.1	Baseline vs. Comparison Methods vs. Proposed Methods	90
6.4.1.1	Accuracy Comparison	90
6.4.1.2	Privacy Comparison	91
6.4.2	Cluster of Trust Configuration Analysis	91
6.4.3	Discussion	92
	Augmented Data sharing strategy	92
	Trade-off of Cluster of Trust	92
6.5	Conclusion	92
	Generalisability and external validity.	93
7	Cross-Domain Analysis and Discussion	95
7.1	Comparative Analysis of Model Performance and Privacy	95
7.1.1	Trust Clustering Performance Across Domains	96
7.1.2	Communication Efficiency Analysis	97
7.2	Research Limitations	98
7.3	Future Research Directions	98
7.4	Summary	99
8	Conclusion	101
8.1	Research Summary	101
8.2	Main Contributions	102
8.2.1	Trust Clustering for Non-IID Data (C1)	102
8.2.2	Composite Privacy Protection Framework (C2)	103
8.2.3	Privacy-Preserving Data Augmentation (C3)	104
8.2.4	Cross-Domain Generalisability Validation (C4)	104
8.3	Hypotheses Evaluation	105
8.4	Implications and Future Outlook	105
8.4.1	Implications for Machine Learning	105
8.4.2	Implications for Privacy Preservation	106
8.4.3	Implications for Distributed Computing	106
8.4.4	Future Research Outlook	106
8.5	Concluding Remarks	107
	Appendix A Data processing details	109

Appendix A.1 IoT Data Processing Pipeline	109
Appendix A.1.1 Raw Data Loading and Initial Assessment	109
Appendix Balanced Dataset Files	109
Appendix Original Dataset Files	109
Appendix A.1.2 Data Quality Assessment and Cleaning	111
Appendix Missing Value Analysis	111
Appendix Infinite and Invalid Value Handling	111
Appendix A.1.3 Outlier Detection and Mitigation	112
Appendix Interquartile Range (IQR) Method	112
Appendix A.1.4 Feature Engineering and Derivation	113
Appendix Traffic Efficiency Metrics	113
Appendix Protocol Encoding	114
Appendix A.1.5 Label Consistency Correction	114
Appendix Honeypot Label Standardisation	114
Appendix A.1.6 Feature Scaling and Normalisation	115
Appendix StandardScaler Normalisation	115
Appendix A.1.7 Data Partitioning Strategy	116
Appendix Stratified Train-Test Split	116
Appendix Validation Set Creation	116
Appendix A.1.8 Final Dataset Characteristics	117
Appendix Dimensional Properties	117
Appendix Quality Metrics	118
Appendix Computational Considerations	118
Appendix A.2 Human Activity Recognition Processing Pipeline	118
Appendix A.2.1 Comprehensive SisFall Dataset Analysis	118
Appendix A.2.1.1 Complete Activity Catalog	118
Appendix A.2.1.2 Detailed Participant Demographics	118
Appendix A.2.2 Data Preprocessing Implementation	119
Appendix A.2.2.1 Raw Data Processing Pipeline	119
Appendix A.2.2.2 Dataset Preparation and Label Generation	121
Appendix A.2.3 Federated Learning Implementation	122
Appendix A.2.3.1 Client Data Partitioning	122
Appendix A.2.3.2 Federated Training Loop	123
Appendix A.2.4 Performance Monitoring and Evaluation	124
Appendix A.2.4.1 Resource Usage Tracking	124
Appendix A.2.5 Dataset Statistics and Benchmarks	125
Appendix A.2.5.1 Complete Dataset Statistics	125
Appendix A.2.5.2 Preprocessing Performance Benchmarks	126
Appendix B Algorithm details	127
Appendix B.1 IoT Intrusion Detection Algorithm	127
Appendix B.1.1 Deep Neural Network Architecture	127
Appendix Network Architecture Specification	127
Appendix B.1.2 Training Algorithm	129
Appendix Optimisation Objective	129
Appendix Class Weight Calculation	129
Appendix B.1.3 Inference Algorithm	132

Appendix Prediction Pipeline	132
Appendix B.1.4 Performance Evaluation Algorithm	134
Appendix Evaluation Metrics	134
Appendix B.1.5 Real-time Detection Algorithm	138
Appendix B.1.6 Algorithm Complexity Analysis	138
Appendix Training Complexity	138
Appendix Inference Complexity	139
Appendix B.2 Detailed Implementation and Technical Specifications	139
Appendix B.3 HAR Algorithm	139
Appendix B.3.1 Complete Model Architecture Specification	139
Bibliography	141

List of Figures

4.1	Structure of the FedAvg.	47
4.2	Federated Learning with Global Data Sharing	47
4.3	Structure of Clusters of Trust.	47
4.4	Traffic distribution by scenario type. (a) Distribution of benign vs. malicious traffic across Malware and Honey-pot scenario types. (b) Overall distribution showing the proportion of malicious vs. benign traffic in the entire dataset.	50
4.5	Temporal analysis of traffic patterns. (a) Hourly distribution of benign vs. malicious traffic throughout the day. (b) Heatmap showing malicious traffic intensity by scenario type and time, highlighting the concentration of malicious activities in Malware scenarios.	51
4.6	Protocol distribution and traffic characteristics analysis. (a) Protocol distribution comparison between benign and malicious traffic. (b) Average packet size distribution histograms. (c) Flow duration distribution box plots. (d) Scatter plot showing the relationship between total packets and total bytes, illustrating distinct clustering patterns for different traffic types.	52
4.7	Detailed scenario comparison analysis. (a) Top 10 individual scenarios showing benign vs. malicious flow distribution. (b) Percentage of malicious traffic by scenario type. (c) Flow size distribution comparison between Malware and Honey-pot scenarios. (d) Protocol usage distribution by scenario type.	53
4.8	Loss on Baseline Scenarios.	55
4.9	Loss on non-IID data with global dataset.	57
4.10	Accuracy histograms for different size clusters.	57
5.1	Framework of the Dual-Layer Privacy-Preserving Federated Learning Framework.	67
5.2	Accuracy of the model after anonymisation - (k1)	70
5.3	Accuracy of the model after anonymisation - (k=5)	70
5.4	Accuracy of the model after anonymisation - (k=10)	71
5.5	Accuracy of the model after anonymisation - (k=100)	71
5.6	Accuracy of the model after anonymisation - (k=150)	72
5.7	Accuracy of the model after anonymisation - (k=200)	72
5.8	Accuracy of different privacy strategies combinations.	73
6.1	RACS-FL Framework Architecture	80
6.2	Data distribution across participant demographics and activity types	83
6.3	Comparison across different methods	90

6.4 Performance under Different Cluster of Trust Configurations 91

List of Tables

2.1	The types of attributes in privacy-preserving data publishing.	15
3.1	Comprehensive Experimental Design Matrix	36
4.1	Network Flow Features Extracted for Analysis	49
4.2	Summary of IoT-23 Dataset Characteristics	49
4.3	Performance Metrics.	55
4.4	Accuracy for different clusters.	56
5.1	Anonymity metrics of distributed datasets and centralised dataset	69
6.1	Preprocessed SisFall Dataset Statistics	85
7.1	Cross-Domain Summary of Experimental Results	96
Appendix A.1	Complete Activities of Daily Living (ADLs) in SisFall Dataset	119
Appendix A.2	Complete Fall Activities in SisFall Dataset	119
Appendix A.3	Complete Participant Demographics	120
Appendix A.4	Comprehensive Dataset Statistics	125
Appendix A.5	Data Preprocessing Performance Metrics	126

Declaration of Authorship

I declare that this thesis and the work presented in it is my own and has been generated by me as the result of my own original research.

I confirm that:

1. This work was done wholly or mainly while in candidature for a research degree at this University;
2. Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated;
3. Where I have consulted the published work of others, this is always clearly attributed;
4. Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work;
5. I have acknowledged all main sources of help;
6. Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself;
7. None of this work has been published before submission

Signed:.....

Date:.....

Acknowledgements

I would like to express my sincere gratitude to my supervisor, **Professor Thanassis Tiropanis**, and to **Dr. George Konstantinidis**, for their invaluable guidance, insightful feedback, and constant support throughout the course of this research. Their expertise and encouragement have been fundamental to the completion of this thesis.

I would also like to acknowledge the maintainers of the **IoT-23** and **SisFall** datasets, whose high-quality and publicly available data resources provided essential support for the experimental validation in this work.

Finally, I would like to extend my heartfelt thanks to my family and friends for their unwavering understanding, encouragement, and companionship during my academic journey.

Publications

This thesis is based in part on the following publications:

Published Papers

1. **Wenxuan Huang**, Thanassis Tiropanis, and George Konstantinidis. (2022). “Federated Learning-Based IoT Intrusion Detection on Non-IID Data: A Trust Clustering Approach.” *Internet of Things – GloTS 2022, Lecture Notes in Computer Science*, vol. 13533, pp. 326–337. Springer, Cham.
2. **Wenxuan Huang**, Thanassis Tiropanis, and George Konstantinidis. (2023). “A Dual-Layer Privacy-Preserving Federated Learning Framework.” In: *Proceedings of the International Conference on Web Information Systems Engineering (WISE 2023)*, pp. 245–259. Springer.

Papers Under Review

1. **Wenxuan Huang**, Thanassis Tiropanis, and George Konstantinidis. “Privacy-Preserving Federated Learning for Human Activity Recognition Using Rotation-based Data Augmentation.”

Definitions and Abbreviations

ADL	Activities of Daily Living
AI	Artificial Intelligence
AUC	Area Under the Curve
BCE	Binary Cross-Entropy
CCPA	California Consumer Privacy Act
CNN	Convolutional Neural Network
CoT	Cluster of Trust
DAG	Data Allocation Granularity
DM	Data Mining
DNN	Deep Neural Network
DP	Differential Privacy
FedAvg	Federated Averaging
FL	Federated Learning
GDPR	General Data Protection Regulation
HAR	Human Activity Recognition
IDS	Intrusion Detection System
IID	Independent and Identically Distributed
IoT	Internet of Things
JSD	Jensen–Shannon Divergence
KL	Kullback-Leibler
LSTM	Long Short-Term Memory
ML	Machine Learning
Non-IID	Non-Independent and Identically Distributed
PP	Privacy-Preserving
PPDP	Privacy-Preserving Data Publishing
PPFL	Privacy-Preserving Federated Learning
PPHAR	Privacy-Preserving Human Activity Recognition
RACS-FL	Rotation Augmentation with Clustered nodes for Federated Learning
RNN	Recurrent Neural Network
ROC	Receiver Operating Characteristic
SGD	Stochastic Gradient Descent
SMC	Secure Multi-party Computation

TCFL Trust Clustering Federated Learning

Chapter 1

Introduction

With the widespread adoption of smart devices, edge computing, and sensor networks, massive distributed data is continuously generated across various application domains, including healthcare monitoring, intelligent transportation, cybersecurity, and industrial automation. These data often take the form of high-frequency time series, network traffic logs, sensor readings, and behavioral patterns, embedding rich information for intelligent decision-making systems. However, such data frequently involve private information such as user identity, behavioral trajectories, device status, network vulnerabilities, and commercial secrets, which poses significant risks in data usage and regulatory compliance.

Globally, data privacy regulations and governance policies are becoming increasingly stringent. Laws such as the EU's General Data Protection Regulation (GDPR) and California's CCPA explicitly restrict centralised data collection and cross-border transmission. At the same time, real-world data is highly fragmented, distributed across institutions, devices, and geographical nodes—creating “data silos” that hinder cross-domain collaboration and large-scale model training. This contradiction is especially evident in multi-party collaboration scenarios, where each organisation owns valuable data but cannot share raw information due to privacy, legal, or competitive concerns, thus limiting the generalisation, robustness, and fairness of intelligent systems.

How to collaboratively model multi-source heterogeneous data under privacy and data sovereignty constraints has become a pressing challenge in AI and distributed learning. Federated Learning (FL) has emerged as a promising paradigm that enables collaborative training without sharing raw data, becoming a bridge between secure data use and intelligent collaboration.

1.1 Problem Statement and Research Gap

Although FL offers a new framework for privacy-preserving distributed learning, it faces three core challenges that limit its practicality across privacy-sensitive applications:

- **Non-IID data distribution.** Differences in user habits, network environments, device types, and geographical locations lead to divergent local models, failed global aggregation, and poor convergence and accuracy. This heterogeneity is particularly pronounced in edge computing scenarios where devices operate in vastly different contexts with distinct data characteristics.
- **Residual privacy risks.** Despite keeping data local, model updates and gradients may still leak private information, especially under adversarial inference attacks. Existing privacy-preserving mechanisms like differential privacy can mitigate this risk but often at the cost of significant accuracy degradation, creating an unsatisfactory privacy-utility trade-off.
- **Severely limited training samples at edge nodes.** Many edge devices operate intermittently or in resource-constrained scenarios, lacking sufficient training data for effective local model training. This data scarcity problem is exacerbated by privacy constraints that prevent direct data sharing among participants.

Existing studies often tackle these issues in isolation: some enhance privacy mechanisms but neglect performance loss; others use clustering approaches to address non-IID issues but ignore privacy boundaries; still others employ data augmentation techniques without solving fundamental heterogeneity problems. These fragmented approaches lack a unified framework that systematically balances performance, privacy, and convergence efficiency. Moreover, most existing solutions are evaluated on single application domains, limiting our understanding of their generalisability across different privacy-sensitive scenarios.

1.2 Research Questions and Hypotheses

This research addresses three fundamental questions that drive our experimental investigations and theoretical development in privacy-preserving federated learning. Each research question is formulated to enable systematic evaluation and empirical validation across multiple application domains.

1.2.1 Research Question 1: Trust-Based Clustering Effectiveness

RQ1: Can trust-based clustering mechanisms provide sufficient privacy protection whilst maintaining acceptable model accuracy compared to global data sharing and traditional anonymisation approaches in non-IID federated learning environments?

Hypothesis H1: Trust-based clustering (Cluster-of-Trust, CoT) can achieve superior privacy-utility trade-offs compared to both unrestricted global data sharing and conventional single-mechanism privacy protection by enabling selective collaboration between clients with similar data distributions and established trust relationships whilst maintaining strict privacy boundaries with potentially adversarial participants.

Experimental Focus: This question drives experiments comparing CoT configurations against baseline federated learning, global data sharing strategies (5% and 10% data sharing), and single-mechanism privacy approaches across varying degrees of data heterogeneity and different cluster sizes.

1.2.2 Research Question 2: Composite Privacy Protection Mechanisms

RQ2: How do multi-layered privacy protection mechanisms integrating differential privacy, k-anonymity, and trust-based data sharing compare to single-technique approaches in terms of accuracy preservation and privacy guarantee strength across different application domains?

Hypothesis H2: A composite privacy protection framework dynamically combining differential privacy ($\epsilon \in \{0.1, 1.0, 10.0\}$), data anonymisation techniques (k-anonymity, l-diversity), and trust-based selective sharing can achieve stronger privacy guarantees with significantly less accuracy degradation than any individual protection method applied in isolation.

Experimental Focus: This question guides systematic evaluation of privacy mechanism combinations, measuring both theoretical privacy bounds and empirical attack resistance.

1.2.3 Research Question 3: Privacy-Preserving Data Augmentation

RQ3: Can domain-specific privacy-preserving data augmentation techniques effectively address data scarcity constraints in federated learning whilst maintaining the privacy benefits of distributed training?

Hypothesis H3: Privacy-preserving data augmentation strategies, including rotation-based transformations for sensor data, differential privacy-enhanced

synthetic data generation, and anonymised data sharing within trust clusters, can improve model performance for data-sparse clients whilst maintaining privacy guarantees equivalent to or stronger than baseline federated learning approaches.

Experimental Focus: This question drives experiments with systematic augmentation scale variations (1×, 2×, 5×, 10× original dataset size) across different privacy protection levels, evaluating the effectiveness of domain-specific augmentation techniques (sensor data rotation for HAR, synthetic attack pattern generation for IoT security, anonymised statistical data for publishing scenarios) in addressing the fundamental data scarcity challenges inherent in federated learning environments.

1.3 Research Motivation and Application Context

To address these fundamental challenges systematically, this thesis develops and validates privacy-preserving federated learning solutions across three complementary application domains. These domains are strategically selected to provide comprehensive evaluation of different aspects of the core challenges:

IoT Network Intrusion Detection serves as the primary testbed for addressing *non-IID data distribution challenges*. IoT networks naturally exhibit extreme data heterogeneity, as different network segments experience distinct attack patterns, traffic characteristics, and device behaviours. Using the IoT-23 dataset, this domain allows us to evaluate how trust-based clustering mechanisms can effectively handle statistical heterogeneity whilst maintaining network security requirements.

Privacy-Preserving Data Publishing and Mining focuses on *advanced privacy protection mechanisms* and their integration with federated learning. This domain requires sophisticated anonymisation techniques that maintain both individual privacy and analytical utility. It provides the theoretical foundation for understanding privacy-utility trade-offs and validates the effectiveness of composite privacy protection schemes combining differential privacy, k-anonymity, and trust-based approaches.

Human Activity Recognition (HAR) provides an ideal scenario for investigating *privacy-sensitive data augmentation* and *data scarcity issues*. Wearable devices and mobile sensors collect highly personal data about daily activities, health status, and behavioural patterns. Users are inherently reluctant to share such intimate information, yet collaborative learning across diverse user populations is essential for robust models that generalise across demographics and environments. The SisFall dataset enables evaluation of privacy-preserving augmentation techniques under severe data constraints.

These three domains collectively address all core federated learning challenges whilst representing diverse application requirements. From a data perspective, they encompass fundamentally different data characteristics: network traffic logs in IoT security scenarios, sensor time series data in human activity recognition. The privacy sensitivity levels also vary significantly across domains, ranging from moderate sensitivity in network security contexts to high sensitivity in personal health data. The scale requirements differ substantially across these applications, spanning large-scale IoT networks with thousands of distributed devices, personal device clusters in healthcare monitoring scenarios. This multi-domain approach enables the development of a comprehensive federated learning framework that is both theoretically grounded and practically validated across diverse real-world scenarios.

1.4 Research Contributions

This thesis systematically addresses non-IID data, privacy protection, and sample scarcity in federated learning from theoretical, algorithmic, and empirical perspectives. The major contributions are:

- **C1: Trust-graph-constrained clustering for Non-IID FL.** We introduce and empirically validate a Cluster-of-Trust framework (CoT) that combines (i) an explicit trust/permission graph and (ii) Jensen–Shannon-divergence-based statistical alignment to form virtual clients for aggregation, improving convergence and accuracy under severe non-IID conditions (validated on IoT-23; extended to HAR settings).
- **C2: Dual-layer privacy protection and trade-off analysis.** We develop a dual-layer privacy-preserving FL framework that integrates anonymisation (microaggregation / k -anonymity mechanisms) at data collection with differential privacy at model training, and we quantify privacy–utility–convergence trade-offs by comparing federated, centralised, and hybrid training modes.
- **C3: Privacy-preserving augmentation for data scarcity and scalable evaluation.** We adapt rotation-based transformations for sensor time-series as a privacy-preserving augmentation strategy in federated fall detection, and use it to create thousands of virtual clients for realistic large-scale evaluation of clustered/trust-based FL strategies.
- **C4: Cross-domain case-study validation under a unified framework.** We demonstrate how the framework (Chapter 3) can be instantiated across three distinct domains—IoT intrusion detection, privacy-preserving data publishing,

and HAR—and we explicitly discuss assumptions and threats to external validity that affect how results may transfer to other datasets and deployments.

1.5 Thesis Structure

This thesis is organised as follows:

Chapter 2 provides a comprehensive literature review covering federated learning fundamentals, privacy-preserving mechanisms, and domain-specific applications in IoT security, human activity recognition, and data publishing.

Chapter 3 details the theoretical foundations and methodological framework, including the mathematical formulation of the Cluster-of-Trust mechanism, privacy analysis, and experimental design principles.

Chapter 4 investigates federated learning for IoT intrusion detection under non-IID conditions, based on the IoT-23 dataset. This chapter demonstrates how trust clustering addresses data heterogeneity in network security applications.

Chapter 5 explores privacy-preserving data publishing and mining within a federated learning framework, examining the integration of differential privacy and anonymisation techniques.

Chapter 6 focuses on privacy-preserving human activity recognition using rotation-based data augmentation and clustered federated learning. This chapter validates the framework on the SisFall dataset.

Chapter 7 synthesises findings across the three application domains, providing comparative analysis of privacy-utility trade-offs, convergence characteristics, and practical deployment considerations.

Chapter 8 concludes with a summary of contributions, discussion of limitations, and directions for future research in privacy-preserving federated learning.

Chapter 2

Literature Review

This chapter provides a comprehensive review of the literature relevant to privacy-preserving federated learning and its applications. We begin with an overview of federated learning, followed by a detailed examination of privacy-preserving techniques. We then review the literature specific to each application domain addressed in this thesis: IoT security and intrusion detection, data publishing and mining, and human activity recognition.

2.1 Literature Review Methodology

To ensure that Chapter 2 provides a *critical* discussion of the state of the art, we followed a structured literature review process. We searched multiple scholarly databases and digital libraries, including **IEEE Xplore**, **ACM Digital Library**, **SpringerLink**, and **arXiv**, complemented by **Google Scholar** to capture highly-cited and recent preprints. The search was driven by combinations of keywords covering (i) federated learning fundamentals and non-IID mitigation (e.g., “federated learning” AND “non-IID” AND “clustering”), (ii) privacy mechanisms in FL (e.g., “federated learning” AND “differential privacy” AND “secure aggregation”), and (iii) the three application scenarios investigated in this thesis (e.g., “federated learning” AND “intrusion detection” / “IoT”, “privacy-preserving data publishing” AND “microaggregation”, and “federated learning” AND “human activity recognition”).

We applied a staged screening process based on title/abstract and then full text, prioritising peer-reviewed survey papers and primary studies that: (a) directly address privacy-preserving federated learning, and/or (b) study methods closely related to our contributions in Chapters 4–6. To mitigate database bias, we additionally used backward/forward snowballing from key surveys (e.g., [1–4]) to identify influential foundational works and closely related follow-up studies.

2.2 Federated Learning

2.2.1 Development of Federated Learning

Growing regulatory pressure and public concern over data exposure have rendered centralised deep-learning pipelines increasingly untenable. Federated Learning (FL) was introduced to address these constraints: local training is performed on each client, and a central server aggregates model updates through the seminal FedAvg algorithm, which delivers communication-efficient convergence even under non-IID data distributions[5]. Its practical viability was soon demonstrated on Google Gboard, where millions of mobile devices collaboratively trained language models while preserving user privacy through on-device encryption and differentially private telemetry[6]. After that, IEEE P3652.1 laid out a comprehensive guide to FL architectures, evaluation protocols and governance responsibilities, reinforced by complementary initiatives at ISO and NIST. Most recently, the EU AI Act has codified stringent requirements for data governance and explainability in high-risk AI systems, further legitimising FL as a compliant learning paradigm[7]. Concurrently, the community has begun federating the fine-tuning and alignment of large-scale language models; resources such as FedLLM-Bench now provide realistic benchmarks for this emerging line of work[8].

2.2.2 Key Concepts of Federated Learning

2.2.2.1 System Architecture

A federated learning (FL) system comprises **two core roles** operating in concert[5]:

- (i) **Client** — the real data owner who trains a model locally on private data;
- (ii) **Aggregator / Server** — the central entity that collects and fuses updates coming from all clients.

The canonical FL workflow can be divided into four stages [2]:

1. **Broadcast stage:** the server dispatches the latest global model and the training configuration to all currently on-line clients;
2. **Local training stage:** each client executes several gradient-descent steps on its private data, optionally performing pruning, quantisation, or encryption to compress and protect its local update;

3. **Upload stage:** the client uploads its (possibly noise-added) gradient or model delta to the aggregator through a secure channel;
4. **Aggregation/dispatch stage:** the server securely aggregates the received updates, produces a new global model, and broadcasts it, thereby entering the next training cycle.

This layered role division and cyclical workflow satisfy the legal requirement that data remain local while dynamically balancing compute and communication across multiple resource tiers.

2.2.2.2 Categories of Federated Learning

After fixing the system architecture, researchers usually classify FL from two viewpoints—*deployment paradigm* and *data partitioning*—so that algorithms and protocols can be tailored to different constraints.

Deployment paradigm. According to the nature and scale of participating entities, FL deployments typically fall into two broad categories:

- **Cross-Silo Federated Learning.** This setting involves several data silos such as banks, hospitals, factories or research centres that collaborate while keeping data on-premise. The client count is usually small (two to a few hundred). Nodes are highly reliable servers or data centres with stable, high-bandwidth links; synchronous protocols and secure multi-party computation (SMC) are easy to deploy. Each silo contains a large, relatively homogeneous dataset, whereas substantial statistical heterogeneity (Non-IID) exists across silos, challenging model generalisation and accuracy. Privacy concerns focus on corporate secrecy across competing organisations and are typically mitigated via legal contracts and encrypted aggregation [9]. This type of FL aims to build a high-performance model with cross-silo accuracy and generalisation by exploiting a diversified dataset that no single organisation can possess alone.
- **Cross-Device Federated Learning.** This paradigm targets millions to billions of resource-constrained end devices (smart-phones, wearables, IoT sensors). Device availability is low and volatile, with frequent drop-outs; only a subset participates in each round. Per-device data volume is small and highly personalised, exhibiting extreme Non-IID and imbalance. Training must protect individual privacy and withstand malicious or low-quality updates [2]. System goals lean towards improving user experience and learning population-level trends, often accepting a slight accuracy loss in exchange for privacy gains.

Differences in node scale, availability, communication conditions and privacy focus directly influence downstream algorithm design and system optimisation.

Data partitioning. Within either deployment paradigm, organisational data relations can be further refined into three scenarios:

- **Horizontal Federated Learning.** All parties share an identical feature and label space but possess different sample sets. Example: several banks share transaction features to train a risk-scoring model. The advantage is unified model structure and high code reuse, whereas the difficulty lies in resolving user overlap and label alignment [10].
- **Vertical Federated Learning.** Parties' user sets partially overlap, but their feature dimensions are complementary, e.g. an insurance company and a bank respectively holding policy data and transaction records. Homomorphic encryption or secure feature-alignment protocols are required before local gradient computation [11].
- **Federated Transfer Learning (FTL).** When both samples and features overlap only partially, a transfer network shares high-level representations under privacy protection to bridge domain gaps. FTL relies on a small public mapping set and alignment layers to cut communication cost and improve cross-domain generalisation [12].

The two-dimensional taxonomy of deployment paradigm and data partitioning jointly determines the trust model, communication protocol, privacy mechanism, and even evaluation metrics of an FL system.

2.2.2.3 Model Aggregation Algorithms

Aggregation determines the efficiency and robustness of knowledge fusion across clients. Five representative algorithms are summarised below.

FedAvg. FedAvg executes E local SGD steps on each client and then averages weights proportionally to local sample counts [5]. Communication and compute complexity scale linearly with rounds; the algorithm is simple and robust, becoming the "zero-baseline" for all follow-up work. However, under severe Non-IID or heterogeneous local iteration counts, FedAvg suffers *client drift* and slow convergence; synchronous aggregation is sensitive to stragglers. Google Gboard maintained $\approx 95\%$ of centralised accuracy even with only 10% client participation [6], but performance drops sharply on highly heterogeneous medical images.

FedProx. To mitigate drift, FedProx adds a proximal term $\mu \|w - w^{(s)}\|^2$ to local objectives, forcing client weights to stay near the global centre. This soft constraint guarantees pux-centric convergence under Non-IID; on Shakespeare and CIFAR-10 (Dirichlet $\alpha = 0.1$) it improves Top-1 accuracy by 6–9 pp and cuts rounds by 30%. Yet the parameter μ is dataset-sensitive; an excessive value hampers local learning[13].

SCAFFOLD. SCAFFOLD maintains global and local *control variates* \mathbf{c} to correct update directions [14]. In CIFAR-100 ($\alpha = 0.05$) it lifts accuracy by 15 pp under the same budget and halves variance convergence. The method costs extra memory for dual vectors; with high drop-out, global \mathbf{c} becomes noisy.

FedOpt family. FedOpt introduces server-side momentum and adaptive learning rates, porting Adam/Yogi to FL [15]. FedAdam aggregates gradients, then applies first- and second-moment correction; FedYogi clips variance explosions. On EMNIST and StackOverflow Non-IID, FedAdam speeds early-phase accuracy by 1.5–2× relative to FedAvg, suiting large-model fine-tuning. Momentum amplifies Byzantine risk, requiring robust aggregation.

FedNova. FedNova normalises the product of local step size and iteration count to tackle *objective inconsistency* from unequal local epochs [16]. The algorithm converges to centralised optima without extra regularisers; with 30% simulated drop-out it surpasses FedAvg by 8 pp without punishing slow clients, at the cost of sending cumulative step info.

2.2.3 Challenges of Federated Learning

Although the technological stack keeps evolving, large-scale real deployments are still hindered by three intertwined issues: **statistical heterogeneity**, **data limitations**, and **security robustness**. This section details their causes and manifestations, summarises mainstream remedies, and highlights research gaps[17].

2.2.3.1 Non-IID Data

Real-world client data are strongly heterogeneous due to user behaviour, cultural regions, and sensor variance [18]. Such distribution shifts skew global gradient expectations away from the optimum, slow convergence, cause accuracy fluctuation, and induce systematic unfairness to long-tail users—for example, minority ethnicities or rare diseases may be underestimated in medical imaging or speech recognition, triggering ethical and regulatory risks[19].

Existing remedies.

1. **Shared public subset** to anchor distributions [18];
2. **FedProx proximal regularisation** [13];
3. **Clustered-FL**: dynamically partitioning clients into homogeneous clusters [20];
4. **Personalised pFedMe**: decoupling local and global parameters [21];

Limitations & gaps. Global shared datasets may pose privacy risks or violate regulations. Clustering requires a predefined threshold, which is difficult to tune in online settings. Personalisation often incurs high communication costs. Achieving privacy-preserving and low-cost distribution correction remains an open challenge.

2.2.3.2 Data Limitations

Even after addressing heterogeneity, data may still be *scarce, low-quality, or weakly labelled*, especially in privacy-sensitive or emerging devices. Insufficient data cause under-fitting and poor generalisation[22].

Causes & manifestations.

- **Lack of realistic datasets**: most public corpora are centralised; naively partitioning them cannot mimic FL-native Non-IID[23];
- **Cold start & long tail**: new devices, applications or rare events have very few samples[24];
- **Label scarcity**: manual annotation is expensive; semi-/un-supervised data dominate[25].

Existing remedies.

1. **Federated data augmentation (FedAug)**: the server generates synthetic snippets and sends them to data-poor clients[26];
2. **Active learning**: within FL, select the most informative samples for annotation to reduce total labels[27].

Limitations & gaps. Synthetic data may deviate from real-world distributions and often lack comprehensive evaluation across modalities. Self-supervised pre-training requires substantial computational resources, making it unsuitable for resource-constrained devices. Active learning relies on uploading confidence scores, which can risk privacy exposure. Currently, there is no unified benchmark that effectively balances data utility and privacy budget.

2.2.3.3 Security and Robustness

Open participation enables malicious clients to inject back-doors, poisoning, or gradient inference attacks [28]. Byzantine nodes can arbitrarily manipulate updates and crash critical applications. Beyond malicious clients, FL also faces **curious servers** and **side-channel eavesdroppers**. Servers can perform gradient inversion and membership inference to recover sample content or membership[29].

Beyond sample-level leakage, recent work has shown that federated learning can also be vulnerable to **user-/client-level privacy attacks** that attempt to recover client-specific data distributions or reconstruct representative samples of a targeted user. For example, Song et al. propose a malicious-server-side attack (mGAN-AI) that targets user-level privacy, and further demonstrate that simply anonymising client identities/updates (e.g., via anonymous communication) may be insufficient because linkability attacks can re-associate anonymised updates to the same client[30]. This highlights that “anonymisation” in FL can refer to different layers (identity/update anonymity vs. data anonymisation), and that threat models must explicitly state what is being protected and against which attacker capabilities.

Defence techniques.

- **Robust aggregation:** Krum, Trimmed-Mean, etc., statistically trim abnormal gradients [31];
- **Model audit /neural cleanse:** server-side reverse engineering to remove back-door triggers [32];
- **Reputation mechanism:** dynamic weighting or screening based on historical behaviour;
- **Differential privacy:** noise addition at gradient level to raise inference difficulty [33].

Limitations & gaps. Defence–accuracy trade-off lacks a unified metric; real-time detection is costly; Byzantine-resilient theory upper bounds are unclear, and experience with generative large models is sparse. A cross-disciplinary fusion of security, privacy, and explainable AI is needed to build an auditable FL trust framework.

2.2.4 Summary and Future Trends

FL now spans system architecture, aggregation algorithms, communication compression, and security–privacy techniques. However, bottlenecks persist in statistical heterogeneity, system heterogeneity, and robustness. Future research directions include cross-layer joint optimisation, explainable and auditable defences, and communication-computation double down-scaling for large generative models. By bridging algorithmic, systemic, and regulatory gaps, FL can achieve sustainable deployment in high-risk domains such as healthcare, finance, and smart devices.

2.3 Privacy-Preserving Techniques

To address the privacy concerns inherent in federated learning, a variety of privacy-preserving techniques have been developed. These methods aim to mitigate the risk of information leakage while maintaining model utility. In this section, we summarise key categories of such techniques.

2.3.1 Differential Privacy

2.3.1.1 Principles of Differential Privacy

Differential Privacy (DP) is a mathematical framework designed to protect individual privacy by injecting random noise into queries or computations. The key idea is to ensure that the inclusion or exclusion of a single data point does not significantly affect the output, thereby preventing adversaries from inferring the presence of specific data entries[34].

Formally, a mechanism \mathcal{M} is said to satisfy (ϵ, δ) -differential privacy if, for any two datasets D and D' differing in at most one element, and for any subset of possible outputs $S \subseteq \text{Range}(\mathcal{M})$, the following holds[35]:

$$\Pr[\mathcal{M}(D) \in S] \leq e^\epsilon \cdot \Pr[\mathcal{M}(D') \in S] + \delta$$

Here, ϵ represents the privacy budget—smaller values provide stronger privacy guarantees—and δ denotes the probability of a privacy breach.

DP is typically achieved by adding calibrated noise (e.g., Gaussian or Laplace noise) to aggregated outputs or model updates[36]. As it offers strong theoretical guarantees without requiring access to raw data, differential privacy has become a widely adopted approach in privacy-preserving federated learning[37].

2.3.2 Privacy-Preserving Data Publishing

As more and more personal information is used for data-driven research or product development, the protection of private personal information is becoming increasingly important. Privacy-preserving data publishing(PPDP) is a process of sharing data while protecting individual privacy. It involves techniques that aim to prevent unauthorised access to sensitive data and protect the anonymity of individuals in the dataset. It is particularly important in industries such as healthcare, finance, and government, where sensitive data needs to be shared for research or analysis purposes while maintaining the privacy of individuals. Table 2.1 shows the different attributes of the data in the PPDP.

TABLE 2.1: The types of attributes in privacy-preserving data publishing.

Type	Description
Identifier	Attributes in the data that are used to uniquely identify individuals' identities.
Quasi-identifier	Combinations of attributes in a dataset that are linked with an identifying attribute.
Sensitive Attribute	Attributes in the data that are related to individuals' sensitive information.
Non-sensitive Attribute	Attributes in the data that are not sensitive attributes, identifying attributes, or quasi-identifier attributes.

In the publishing of personal data, there are three types of privacy threats [38]:

- **Identity disclosure:** An attacker can correctly associate an individual with a personal record in a published dataset.
- **Attribute disclosure:** Attackers can obtain individuals' sensitive information through inference attacks. This type of threat is more likely to occur in datasets with low anonymity.
- **Membership disclosure:** Attackers can infer with a high probability whether an individual's record exists or does not exist in a published dataset [39][40].

2.3.2.1 Anonymisation Techniques

To protect user privacy in the published dataset, data publishers can apply anonymisation techniques to enhance the dataset's resilience against attacks. The aim of anonymisation is to ensure that a person's data record can no longer be traced explicitly back to that particular person. To this purpose, various complementary privacy paradigms have been defined:

k-Anonymity K-anonymity [41] is a privacy protection concept that requires each record in a dataset to be indistinguishable from at least K-1 other records in terms of their attributes, thereby hiding the specific identity information of individuals. K-anonymity is achieved by generalising or suppressing attributes, which increases the similarity between records and ensures anonymity. Specifically, for a dataset D with attribute set A, if for every record d in the dataset, there exist at least K-1 other records d' such that they have the same values for the attributes in set A, then the dataset D satisfies K-anonymity.

L-Diversity L-diversity [42] is a measure of the richness of information in a dataset. It quantifies the number of different attribute values within each equivalence class in the dataset. The goal of L-diversity is to increase the diversity of attribute values in the anonymised dataset, thereby enhancing its utility. A higher L-diversity value indicates a greater diversity of attribute values in the dataset.

L-Diversity is measured by:

$$L - diversity(D) = \min_{q \in Q} \left(\frac{1}{n} \sum_{i=1}^n f(q, D_i) \right)$$

where D is the anonymised dataset, Q is the set of all possible values for the sensitive attribute, n is the number of equivalence classes, D_i is the $i - th$ equivalence class, and $f(q, D_i)$ is the proportion of records in the equivalence class D_i with a sensitive attribute value of q .

t-closeness t-closeness[43] aims to protect against attribute disclosure attacks by minimising the likelihood of inferring sensitive information based on background knowledge. In t-closeness, the notion of closeness refers to the similarity between the distribution of sensitive attributes in the original data and their distribution in the published data. The parameter "t" represents a threshold value that determines the acceptable level of similarity. A smaller value of t indicates a higher degree of privacy protection. t-closeness can be measured by Kullback-Leibler [44] Divergence.

Kullback-Leibler divergence, also known as relative entropy, is a measure used to quantify the difference between two probability distributions. KL divergence measures the information loss when using one probability distribution Q to approximate another distribution P , given that P is the true distribution. Specifically, for two probability distributions P and Q , the KL divergence is defined as follows [45]:

$$KL(P \parallel Q) = \sum P(x) \log \left(\frac{P(x)}{Q(x)} \right)$$

Where $P(x)$ and $Q(x)$ represent the probabilities of event x under probability distributions P and Q , respectively.

To enhance privacy paradigms, common data anonymisation operations include [46]:

- **Generalisation** [47]: In this operation, the original values of quasi-identifiers are transformed into less specific but semantically consistent values. For example, age or income can be generalised into intervals.
- **Suppression** [47]: This operation hides the original values of quasi-identifiers by replacing them with a special value. For instance, the value "20" of an individual's age can be anonymised. The value "0" can be replaced with "*", resulting in "2*" as the suppressed value of the quasi-identifier.
- **Perturbation** [47]: In this operation, random noise is added to the data to obscure the true values of individuals. This can be achieved by adding random numbers to numerical data or introducing randomisation processes in categorical data.
- **Microaggregation** [48]: This operation involves aggregating the data by combining multiple individuals' data into representative values, thus concealing specific individual information.

These anonymisation techniques aim to balance the privacy protection and data utility in the published dataset. By applying these operations, data publishers can enhance the privacy of the dataset while preserving its usefulness for analysis and research purposes.

2.3.2.2 Critical discussion: anonymisation and downstream model utility

A recurring theme in privacy-preserving data publishing is that stronger anonymisation can degrade the performance of downstream analytics and machine learning models. For instance, generalisation and suppression used to achieve k -anonymity can substantially alter feature distributions and decision boundaries,

leading to measurable accuracy drops for standard classifiers[49]. This effect is not merely an implementation detail: it represents an inherent *privacy–utility* trade-off that must be characterised for each application.

Importantly, federated learning can *amplify* this sensitivity because each client already operates on a smaller, fragmented dataset; any additional distortion introduced during distributed data collection (e.g., anonymisation before model training) may disproportionately reduce local signal and increase update variance. Recent surveys of privacy-preserving FL emphasise that the interplay between privacy mechanisms (including anonymisation and differential privacy) and learning dynamics (including convergence and robustness) is complex and often under-analysed in an application-grounded manner[1, 3]. This thesis directly targets this gap in Chapter 5 by empirically analysing how anonymisation at data collection and differential privacy at model training jointly affect utility and convergence.

2.4 Intrusion Detection in the Internet of Things

Over the past few decades, the Internet and computer systems have caused numerous security problems due to the explosive growth of networks. The CERT reports that the number of intrusions has increased dramatically from year to year. Any malicious intrusion or attack against network vulnerabilities, computers or information systems can cause serious disasters and violate computer security policies, so the research on intrusion detection systems has attracted a lot of attention in the whole field of computer science[50]. Intrusion detection models are capable of detecting intrusions, penetrations and other forms of computer misuse. The model is based on the principle that security violations can be detected by monitoring a system’s audit logs, network traffic, and other information to detect unusual patterns of system usage[51].

Depending on the detection mechanism used in the system, intrusion detection techniques are usually classified as anomaly-based and signature-based. In the signature-based approach, the IDS detects threats by comparing system or network behaviour with signatures stored in the IDS internal database. When any system or network activity matches the stored threat pattern/signature, the IDS will trigger an alert. Signature-based IDSs are very precise and effective in detecting known attacks, and their mechanisms are easy to understand. However, this approach is not effective in detecting unknown attacks such as zero-day attacks or variations of known attacks where the matching pattern is still unknown [50]. Anomaly-based IDS evaluates the system’s activity in real time and triggers an alert if the current behaviour deviates from normal behaviour by more than a threshold. This approach can be effective in detecting new attacks, especially those related to resource misuse. However, any behaviour that does not match normal behaviour is considered an intrusion and

previously unknown legitimate activity may also be classified as malicious. As a result, this method usually has a high false alarm rate [52]. To assess whether a system or network activity is normal or not, researchers often use machine learning to create a trustworthy model of the activity and then compare the new behaviour with this model [53].

2.4.1 Public Datasets and Benchmarking Status

2.4.1.1 Traditional Traffic Benchmarks

Early benchmark datasets for IoT-IDS include KDDCup'99 and its improved version NSL-KDD [54]. These datasets originate from DARPA network traffic capture projects in the 1990s. With complete attack class labels and reproducible experimental setups, they were once widely used for initial evaluations. However, they lack essential IoT communication protocols such as MQTT and CoAP, and include many duplicated samples and outdated attack types. Thus, they are now mainly used for lightweight benchmarks, while modern research has shifted towards more up-to-date datasets[55].

2.4.1.2 Datasets from Real IoT Environments

Currently, the mainstream IoT-IDS datasets include:

- CICIDS2017 simulates a hybrid network of industrial control and home devices, labeling nine types of attacks including DoS, Web, and brute force. It provides the first large-scale benchmark for deep learning models [56].
- TON_IoT introduces multi-layered data including MQTT traffic, TCP/UDP packets, and system logs, enabling cross-layer feature fusion [57].
- Kitsune uses a Raspberry Pi as the gateway to capture real-time traffic, generating 9.7GB of normal and anomalous streams with feature extraction scripts, suitable for lightweight IDS research [58].
- IoT-23 focuses on real-world Botnet infections (e.g., Mirai, Gafgyt), offering PCAP files and Zeek logs to fill previous gaps in encrypted tunnel attacks and adversarial examples [59].

2.4.2 Feature Engineering and Protocol Parsing

Due to resource constraints and protocol heterogeneity, traditional deep packet inspection is difficult in IoT. Efficient feature engineering becomes crucial. Current strategies include statistical features, protocol field features, and event logs.

2.4.2.1 Traffic Statistical Features

Features such as byte count, packet count, and inter-arrival times are calculated via sliding windows. Their low complexity ($O(1)$) enables deployment on microcontrollers with only 256KB RAM. Classical studies show 41-dimensional statistical features can achieve high accuracy ($F1 = 0.99$) [60], though generalisation to novel attacks is limited.

2.4.2.2 Protocol Features

Parsing fields like MQTT headers, topics, and QoS improves recall for topic hijacking and duplicate publish attacks [61]. In TON_IoT, adding MQTT features to XGBoost improved recall from 92% to 97%.

2.4.2.3 Embedded Event Logs

Kernel and application logs capture high-level semantic changes. Template methods like Drain and LogPai can extract log sequences for CNN-based classification [62].

2.4.3 Classical Machine Learning Approaches

Classical ML methods are still valuable under IoT resource constraints.

2.4.3.1 Lightweight Classifiers

SVM, KNN, and Random Forests are widely used in IoT IDS for their efficiency and interpretability. While Random Forests achieve $F1 = 0.99$ on NSL-KDD, their performance on the new dataset has dropped.[63].

2.4.4 Deep Learning Methods

Deep models improve IDS recall but challenge edge deployment.

2.4.4.1 Convolutional and Recurrent Networks

CNNs with 1D kernels capture local protocol patterns and cross-packet dependencies, achieving 96.5% F1 on Kitsune [64]. LSTM and Bi-LSTM capture long dependencies. The results show that the LSTM-based classifier can effectively detect CAN bus network attacks with an overall detection accuracy of 99.995%. [?].

2.4.4.2 Autoencoders

Autoencoders (AE) and Variational Autoencoders (VAE) learn representations via encoding-decoding. AE learns deterministic features, while VAE models latent distributions, enabling generation of unseen samples. Extensive simulations on four datasets show that the VAE-based scheme outperforms with an average F1 score improvement of 47.88 percent [65].

2.4.4.3 Attention Mechanisms

Transformers leverage Multi-Head Self-Attention (MHSA) to capture long dependencies and improve parallelism. MHSA models positional relations by projecting sequences into subspaces, calculating attention weights per head, and integrating features. The weights are interpretable, highlighting key timestamps and channels[66].

A Robust Transformer-Based Approach for Intrusion Detection System is very effective on two publicly available real traffic intrusion detection datasets from CICIDS2017 and CIC-DDoS2019, with F1 scores of 99.17% and 98.48%, respectively [67]. The experimental results of a novel multi-class IDS using a transformer-based attention network for an in-vehicle CAN bus indicate that the transformers model is more efficient than the baselines for different input data types and datasets. [68].

2.4.5 Federated Learning for IoT-IDS

FL is playing an increasingly important role in enhancing IoT services and applications[4]. In the IoT-IDS setting, FL is attractive because it avoids centralising sensitive traffic traces while still enabling collaborative detection models across multiple IoT nodes and domains. Surveys and comparative studies highlight that FL can sit between purely centralised IDS (high utility, high exposure risk) and purely on-device IDS (high privacy, low data diversity), with the practical choice depending on deployment constraints and threat models[69, 70].

2.4.5.1 Representative FL-based IoT-IDS studies

Existing studies demonstrate the feasibility of FL for intrusion detection across different networked environments. For example, Preuveneers et al. propose chained anomaly detection models in a federated setting as an IDS case study[71]. Other work explores federated deep learning for botnet/IoT malware detection, showing that FL

can retain competitive detection performance while avoiding raw traffic sharing[72]. Application-driven FL-IDS systems also report benefits for privacy and distributed deployment, but often vary widely in assumptions about client reliability, attacker capabilities, and available telemetry[73].

2.4.5.2 Non-IID and heterogeneity: what remains challenging

While FL can mitigate raw-data exposure, IoT-IDS data is typically **highly Non-IID**: each node observes different traffic mixes, protocol usage patterns, and attack prevalence, leading to client drift and unstable convergence under FedAvg-like aggregation. The literature reports multiple mitigation directions, including robust aggregation, personalisation, and clustering-based federation[2, 69]. Clustering-based approaches are particularly relevant, as they aim to group clients with similar distributions so that within-cluster training is more coherent than global aggregation[20].

However, the state of the art exhibits three recurring limitations that directly motivate Chapters 3–4 in this thesis:

- **Similarity without trust**: many clustering approaches form groups purely from update similarity or statistical proximity, but do not explicitly encode *trust* constraints between participants. In security-sensitive IoT settings, collaboration feasibility may be constrained by organisational boundaries or by the risk of adversarial peers.
- **Clustering novelty and differentiation**: hierarchical clustering of local updates has been explored to improve Non-IID training[74], and clustered FL has been formalised as a general optimisation strategy[20]. Therefore, new proposals must clearly articulate what is different (e.g., the similarity measure used, the trust model assumed, or the integration with privacy mechanisms).
- **Evaluation scope**: many IoT-IDS FL studies evaluate on a single dataset and a narrow set of client configurations, limiting external validity claims. Clear statements about dataset representativeness and threats to validity are needed, alongside a roadmap for multi-dataset verification.

2.5 Human Activity Recognition

Human Activity Recognition (HAR) focuses on the automatic detection and interpretation of human actions from various data sources[75, 76]. HAR has gained prominence in a variety of applications, including healthcare (such as fall detection), fitness monitoring, and assisted living [77].

2.5.1 Datasets

2.5.1.1 Public Sensor Datasets

Sensor datasets provide essential experimental platforms for Human Activity Recognition (HAR). The UCI-HAR dataset [78] is widely used and contains daily activities such as walking and stair climbing, recorded using smartphone-embedded accelerometers and gyroscopes. It has become a classic benchmark due to its complete labeling and ease of use. The WISDM dataset [79] collects walking, running, and jumping data from a large number of users using a uniaxial accelerometer, offering larger scale and better generalisation. MobiAct [80] focuses on fall detection by simulating daily and fall-related events, and is widely applied in abnormal behavior recognition. The SisFall dataset [81] targets elderly fall scenarios and provides more realistic fall simulations across diverse contexts, making it a key reference for elderly behavior monitoring.

2.5.1.2 Multimodal Datasets

Multimodal datasets integrate various sensor data to enhance recognition accuracy. PAMAP2 [82] includes acceleration, gyroscope, temperature, and heart rate data, supporting multimodal fusion research. OPPORTUNITY [83] combines environmental and wearable sensors to collect complex activity and gesture data, suitable for collaborative multi-sensor processing. UR Fall Dataset [84] integrates video surveillance and inertial data, focusing on fall detection and advancing sensor-visual fusion. Challenges remain in synchronisation and fusion algorithm complexity, along with rising privacy concerns.

2.5.2 Feature Engineering

2.5.2.1 Time-Domain Statistical Features

Time-domain features such as mean, variance, standard deviation, and zero-crossing rate are core to sensor data analysis. These simple, real-time features capture trends and stability in activity signals[85]. Mean and variance reflect intensity and fluctuation, while zero-crossing rate characterizes frequency—useful in gait analysis and periodic activity detection. However, they lack precision in distinguishing complex or subtle activities, motivating the exploration of hybrid feature strategies[86].

2.5.2.2 Frequency and Time-Frequency Features

Frequency and time-frequency features analyze signal periodicity and dynamics via spectral transforms. These features have higher computational demands and limited real-time applicability on constrained devices[87].

2.5.2.3 Video Features

With advances in computer vision, video features—especially skeleton keypoint extraction—have become vital in HAR[88, 89]. Methods like OpenPose and MediaPipe extract body keypoints, improving accuracy and interpretability, especially in fall detection and rehabilitation monitoring[90]. However, their performance is highly dependent on the camera angle and environment, and there are data handling and privacy issues to consider.

2.5.3 Classical Machine Learning Methods

Depending on the data modality and the application scenario, different approaches can be employed for HAR. Traditional machine learning methods rely on hand-crafted features to capture movement patterns [91]. In contrast, deep learning techniques can learn feature representations directly from raw signals, often achieving higher accuracy at the cost of increased computational requirements [92, 93].

2.5.3.1 Traditional Classifier Comparison

Classifiers such as SVM, Random Forest (RF), KNN, and Hidden Markov Models (HMM) are popular in HAR for simplicity and interpretability. SVMs handle high-dimensional spaces well. RFs perform robustly with feature importance analysis. KNNs are intuitive and real-time friendly but noise-sensitive. HMMs model temporal sequences effectively, useful in gait analysis. However, these methods struggle with complex or fine-grained activity recognition.

2.5.4 Deep Learning

2.5.4.1 CNN/LSTM Sequential Networks (DeepConvLSTM, HAR-CNN)

With the increasing popularity and success of deep learning methods, applying these techniques to recognise human activities in mobile and wearable computing scenarios has attracted widespread attention. A deep neural network that combines

convolutional layers with long short-term memory networks has achieved an overall accuracy of 95.78% across three public datasets (UCI, WISDM, and OPPORTUNITY), with 95.85% accuracy on the WISDM dataset and 92.63% on the OPPORTUNITY dataset[94].

2.5.4.2 Transformer and Self-Attention Models (HAR-Transformer, Perceiver-IO)

A Transformer-based model was comprehensively evaluated on the largest publicly available smartphone motion sensor dataset, which includes a wide variety of human activities. The model achieved an impressive average recognition accuracy of 99.2%, significantly outperforming traditional machine learning methods, which reached only 89.67% accuracy on the same dataset[95]. This substantial performance improvement highlights the capability of Transformer architectures to effectively model temporal dependencies and complex patterns in sensor data. The results indicate that Transformer models have strong potential for future applications in human activity recognition, particularly in scenarios requiring high accuracy and robustness across diverse activity types.

2.5.5 Data Augmentation

For tri-axial accelerometer and gyroscope data, applying rotation transformations is an effective data augmentation strategy to simulate sensor readings from different device orientations and wearing positions[96]. These transformations mimic the variability introduced by users wearing devices on different parts of the body—such as the wrist, waist, or ankle—or by slight misalignments during daily use. This technique helps models generalise better across real-world usage scenarios where such variation is common.

Rotation-based augmentation is particularly beneficial for fall detection tasks. Despite variations in orientation, the fundamental physical characteristics of fall events—such as rapid acceleration changes and abrupt motion—remain largely consistent in their temporal patterns and intensity. By rotating the raw sensor axes while preserving the underlying signal dynamics, this method enriches the training data with realistic but diverse examples. As a result, models trained with rotation-augmented data are more robust to orientation-related noise and can achieve improved generalisation performance across different users and device placements.

2.5.6 Federated Learning in HAR

Training such models using data collected from smart devices in a centralised data center can lead to high communication costs and potential privacy violations. To address these issues, federated learning can be employed to train a general-purpose classifier. In a federated learning-based HAR experiment, federated learning achieved an accuracy of up to 89%, compared to 93% accuracy achieved by centralised training using deep neural networks[97]. The global model trained with federated learning on imbalanced datasets demonstrated comparable accuracy to centralised learning. These results indicate that, for human activity recognition tasks, federated learning can produce models with slightly lower—but still acceptable—accuracy compared to centralised approaches.

2.5.6.1 Critical discussion: privacy, clustering, and scale in federated HAR

The federated HAR literature also highlights a set of practical limitations that complicate privacy-preserving deployment. First, sensor data is highly sensitive: activity traces can reveal health conditions and daily routines, so privacy protections cannot rely solely on data locality. Second, non-IID effects are strong due to differences in individuals' movement patterns, device placement, and demographics, motivating *similarity-aware* or clustered federation approaches[98]. Third, a recurring evaluation gap is **scale**: many HAR datasets contain only tens of participants (e.g., SisFall has 38 subjects), whereas real deployments may involve hundreds or thousands of devices, making it difficult to validate clustering and privacy mechanisms at realistic scale.

Data augmentation for time-series and wearable sensing is widely studied and can improve generalisation by introducing plausible transformations (including rotations)[96, 99]. However, rotation augmentation itself is not novel; rather, the open research question is how such augmentation can be used *within federated learning* as a privacy-preserving mechanism to address data scarcity and enable scalable evaluation while maintaining an explicit privacy–utility trade-off. This motivates the approach and experimental design in Chapter 6.

2.6 Synthesis: State-of-the-Art Limitations and Thesis Positioning

Across the literature reviewed above, four limitations recur and directly inform the research questions and contributions of this thesis:

- **Non-IID remains a primary blocker in practice:** clustered and personalised FL approaches exist, but security-sensitive deployments need clustering that respects trust and collaboration constraints (RQ1; Chapters 3 and 4).
- **Privacy mechanisms are often studied in isolation:** surveys document many techniques, yet comparatively fewer application-grounded studies analyse *combinations* of anonymisation and differential privacy (RQ2; Chapter 5)[1].
- **Data scarcity and realistic scaling are under-addressed:** augmentation methods are well known, but their role as a privacy-preserving enabler for large-scale federated evaluation and non-IID mitigation in HAR remains insufficiently characterised (RQ3; Chapter 6).
- **External validity is often weak:** many studies evaluate on a single dataset or narrow scenario set. This thesis therefore explicitly discusses assumptions, threats to validity, and cross-domain implications in Chapters 4–6 and in the cross-domain discussion (Chapter 7).

Chapter 3

Methodology and Framework

This chapter presents the theoretical foundation and methodological framework for privacy-preserving federated learning across multiple application domains. Building upon the challenges identified in Chapter 1, we develop a unified approach that systematically addresses non-IID data distributions, privacy protection, and data scarcity through an integrated framework. The methodology is designed to be domain-agnostic whilst allowing for application-specific adaptations, as demonstrated through three diverse case studies: IoT intrusion detection, privacy-preserving data publishing, and human activity recognition.

3.1 Research Problem

3.1.1 Non-IID Data Distribution Problem Modelling

In federated learning, data heterogeneity poses fundamental challenges to model convergence and performance. We formalise the non-IID problem through statistical divergence measures that quantify the degree of heterogeneity across clients.

Consider a federated system with N clients, where client i holds local dataset $\mathcal{D}_i = \{(x_j^i, y_j^i)\}_{j=1}^{n_i}$. The local data distribution for client i is denoted as $P_i(X, Y)$, whilst the global distribution is $P_{global}(X, Y) = \frac{1}{N} \sum_{i=1}^N P_i(X, Y)$.

We define the *statistical heterogeneity index* \mathcal{H} as:

$$\mathcal{H} = \frac{1}{N} \sum_{i=1}^N D_{JS}(P_i(X, Y) \parallel P_{global}(X, Y)) \quad (3.1)$$

where D_{JS} denotes the Jensen–Shannon divergence. Higher \mathcal{H} values indicate greater statistical heterogeneity.

The *Jensen–Shannon divergence* (JSD) is a symmetric and smoothed variant of the Kullback-Leibler (KL) divergence that quantifies the similarity between two probability distributions. Given two discrete distributions P and Q , the JSD is defined as:

$$D_{JS}(P \parallel Q) = \frac{1}{2}D_{KL}(P \parallel M) + \frac{1}{2}D_{KL}(Q \parallel M), \quad \text{where } M = \frac{1}{2}(P + Q) \quad (3.2)$$

Here, D_{KL} denotes the Kullback-Leibler divergence, given by:

$$D_{KL}(P \parallel Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)} \quad (3.3)$$

Unlike KL divergence, which is asymmetric and may become undefined when $Q(x) = 0$ for any x such that $P(x) > 0$, the Jensen–Shannon divergence is always finite and symmetric. Moreover, it is bounded between 0 and $\log 2$ (or 1 when using base-2 logarithms). These properties make it particularly suitable for measuring statistical heterogeneity in federated learning, where client distributions may be sparse, disjoint, or highly skewed.

By averaging the Jensen–Shannon divergence across all clients with respect to the global distribution, our statistical heterogeneity index \mathcal{H} captures the overall non-IID nature of the federated dataset in a stable and interpretable manner.

3.1.2 Privacy-Utility Trade-off Challenges

The fundamental challenge in privacy-preserving federated learning lies in balancing data protection with model performance. Traditional centralised machine learning achieves optimal utility by accessing all available data, but this approach completely sacrifices privacy. Conversely, keeping data entirely local provides maximum privacy but severely limits model quality due to insufficient training samples.

Federated learning emerges as a promising middle ground, enabling collaborative model training whilst keeping data distributed across participants. By sharing only model updates rather than raw data, federated learning can achieve significantly better performance than purely local approaches whilst providing substantial privacy improvements over centralised methods. However, standard federated learning still faces privacy vulnerabilities, as gradient information can potentially leak sensitive details about local datasets through various attack vectors, including membership inference, property inference, and gradient inversion attacks.

Privacy requirements vary significantly across application domains. In healthcare scenarios involving human activity recognition, personal movement patterns and

health status constitute highly sensitive information requiring strong protection mechanisms. Network security applications present moderate privacy concerns, as traffic metadata is less sensitive than packet contents, though still requiring protection from competitors and malicious actors. Data publishing scenarios exhibit variable privacy requirements depending on the specific attributes and intended use cases.

The trade-off manifests differently across domains due to varying data characteristics and regulatory constraints. IoT intrusion detection systems must balance the need for comprehensive attack pattern recognition against the privacy of network operational details. Human activity recognition systems face the challenge of providing accurate health monitoring whilst protecting intimate personal behaviours. Privacy-preserving data publishing must maintain statistical utility for research purposes whilst ensuring individual anonymity.

Existing approaches often address privacy and utility as competing objectives, leading to suboptimal solutions that either compromise privacy for performance or sacrifice accuracy for protection. Our framework seeks to identify synergistic approaches that can enhance both dimensions simultaneously through intelligent system design and strategic data sharing mechanisms.

3.1.3 Data Scarcity Constraint Challenges

Data scarcity in federated learning environments presents multifaceted challenges that extend beyond simple sample size limitations. The distributed nature of federated systems naturally fragments datasets, with each participant typically holding only a small portion of the total available data. This fragmentation is exacerbated by the non-IID nature of real-world data distributions, where each client's local dataset may not be representative of the global data distribution.

The impact of data scarcity varies significantly across application domains. In IoT intrusion detection, individual network nodes may experience only a subset of possible attack types, leading to incomplete local models that struggle to generalise to novel threats. Human activity recognition systems face the challenge that individual users generate limited examples of certain activities, particularly rare events like falls, resulting in highly imbalanced local datasets.

Traditional machine learning approaches to data scarcity, such as data augmentation and synthetic data generation, must be carefully adapted for federated settings to avoid compromising privacy. The challenge is to develop augmentation strategies that enrich local datasets whilst maintaining the privacy guarantees that motivated the adoption of federated learning in the first place.

The temporal dimension adds another layer of complexity to data scarcity challenges. Many federated learning applications involve streaming data where the statistical properties may evolve over time. This temporal drift means that historical data may become less relevant, effectively reducing the useful dataset size even further. The framework must account for these temporal dynamics whilst maintaining privacy and performance guarantees.

3.2 Framework Design

3.2.1 Architecture Overview

Our comprehensive framework addresses the multifaceted challenges of privacy-preserving federated learning through a carefully structured three-tiered architecture. This design recognises that different aspects of the privacy-utility-performance optimisation problem are best addressed at different levels of system organisation.

The local tier forms the foundation of our framework, where individual clients maintain complete control over their data whilst contributing to the collaborative learning process. At this level, each participant implements privacy-preserving data augmentation techniques tailored to their specific domain and privacy requirements. Clients generate synthetic data samples that preserve the statistical properties necessary for effective model training whilst obscuring sensitive details that could compromise individual privacy. The local tier also implements initial privacy protection mechanisms, including gradient clipping and local differential privacy techniques where appropriate.

The cluster tier introduces an intermediate level of organisation that groups clients based on trust relationships and data similarity patterns. This clustering approach recognises that clients with similar data distributions and mutual trust can benefit from more extensive collaboration whilst maintaining privacy boundaries with respect to other participants. Within each cluster, selective parameter sharing enables enhanced model quality through increased effective dataset size, whilst the cluster boundaries provide privacy protection against potentially adversarial participants in other clusters.

Trust as a collaboration constraint; similarity as an experimental factor. In this thesis, *trust relationships* are operationalised as an explicit **permission constraint** on collaboration: an edge between two clients indicates that they are allowed to collaborate (e.g., share model updates and/or limited data within a controlled scope).

Trust edges can arise from policy-/organisation-driven agreements (e.g., administrative domains or legal contracts), and are treated as a **given constraint** during the reported experiments.

Within this trust constraint, the **composition of clusters** can still vary in how *similar* or *diverse* the grouped clients are. In the IoT-IDS case study, we therefore study two limiting regimes under the same trust constraint: (i) clusters formed from clients with **similar** threat/label distributions (minimum threat coverage), and (ii) clusters formed to maximise **diversity/coverage** of threats within a cluster (maximum label coverage). Behavioural similarity can be characterised using lightweight statistics (e.g., label-frequency vectors) and divergence measures (e.g., Jensen–Shannon divergence), but it is not assumed that trust is always derived from similarity; rather, similarity/coverage is an experimental dimension analysed *under* trust constraints (further instantiated in Chapter 4).

The global tier coordinates the overall federated learning process, aggregating insights from cluster representatives to form a comprehensive global model. This tier implements system-wide privacy guarantees and ensures that the collaborative learning process converges to a high-quality solution whilst respecting the privacy constraints established at lower tiers. The global tier also manages the dynamic aspects of the framework, including cluster reformation based on evolving trust relationships and performance feedback.

3.2.2 Cross-Domain Adaptation Strategies

The application of privacy-preserving federated learning across diverse domains requires careful consideration of the unique characteristics and constraints present in each field. IoT network security applications present a complex landscape where participating nodes represent different network segments, device types, and operational environments. Each node typically encounters a distinct subset of attack patterns based on its role and exposure, creating naturally non-IID data distributions. The privacy concerns in this domain centre around protecting network topology information, operational patterns, and vulnerability details that could be exploited by malicious actors.

Human activity recognition scenarios involve deeply personal data that reveals intimate details about individual behaviour, health status, and daily routines. Participants in these systems are individual users whose devices collect sensor data throughout their daily activities. The non-IID nature arises from fundamental differences in human behaviour, physical capabilities, device usage patterns, and demographic factors. Privacy requirements are particularly stringent, as leaked

activity data could reveal health conditions, social relationships, and personal habits that users consider highly sensitive.

Privacy-preserving data publishing represents a different challenge paradigm where organisations seek to share aggregate statistical information for research or policy purposes whilst protecting individual privacy. The participating entities are typically institutions rather than individual users, and the data often consists of structured tabular information with mixed categorical and numerical attributes. The non-IID challenge manifests through institutional differences in data collection practices, population demographics, and operational focus areas.

Each domain also presents unique regulatory and compliance requirements that shape the acceptable privacy protection mechanisms. Healthcare-related activity recognition must comply with strict medical privacy regulations, whilst IoT security applications may need to balance privacy with security audit requirements. Data publishing scenarios often operate under specific research ethics frameworks that dictate acceptable levels of anonymisation and data sharing.

The framework's adaptability stems from its modular design that allows domain-specific customisation whilst maintaining theoretical consistency and implementation efficiency. The trust clustering mechanism adapts to different data modalities through configurable similarity measures that capture the essential characteristics of each domain. For IoT security applications, trust relationships are primarily established through network behavioural similarity and historical collaboration success in threat detection tasks.

In human activity recognition scenarios, trust clustering considers both statistical similarity in activity patterns and demographic compatibility factors that influence model transferability. Users with similar age, physical capabilities, and lifestyle patterns are more likely to benefit from collaborative learning, whilst significant demographic differences may limit the value of shared model updates. The clustering algorithm accounts for these domain-specific factors whilst maintaining privacy protection for sensitive demographic information.

Privacy protection mechanisms are similarly adapted to domain requirements through a configurable architecture that supports varying levels of protection intensity and different protection techniques. IoT security applications may emphasise differential privacy techniques that add calibrated noise to gradient updates, balancing attack detection accuracy with operational detail protection. Human activity recognition systems often require stronger protection mechanisms, potentially including data anonymisation techniques and more conservative privacy budget allocation.

Data augmentation strategies are fundamentally domain-dependent, requiring deep understanding of the data generation process and the invariances that can be exploited

for sample generation. IoT security benefits from synthetic attack generation that preserves threat signatures whilst varying irrelevant network details. Human activity recognition leverages physical invariances such as rotational symmetry in sensor readings to generate additional training examples that maintain activity characteristics whilst providing coverage of different device orientations and placements.

The framework's adaptation mechanisms also account for the different scales and participation patterns characteristic of each domain. IoT networks may involve hundreds or thousands of participating nodes with varying computational capabilities and network connectivity. Human activity recognition typically involves smaller numbers of participants but with more intensive data generation per participant. Data publishing scenarios often feature institutional participants with significant computational resources but strict governance requirements.

3.3 Experimental Methodology

3.3.1 Experimental Design Overview

Our experimental methodology employs a systematic design that evaluates the proposed framework across multiple dimensions to comprehensively assess its effectiveness in addressing privacy, utility, and scalability challenges. The experimental design encompasses systematic variations of key factors that influence federated learning performance, including data distribution patterns, system architectures, participant scales, and privacy protection mechanisms.

Table 3.1 presents the comprehensive experimental design matrix that guides our evaluation across the three application domains. The design systematically varies critical parameters to isolate the impact of individual components whilst evaluating their synergistic effects within the complete framework.

The data distribution experiments systematically compare Independent and Identically Distributed (IID) scenarios against various Non-IID configurations. IID scenarios employ uniform random sampling to distribute data equally across participants, representing ideal federated learning conditions rarely found in practice. Non-IID scenarios utilise Dirichlet distribution parameterisation with concentration parameter α to control heterogeneity levels. Lower α values ($\alpha = 0.1$) create extreme heterogeneity where clients hold highly dissimilar data distributions, whilst higher values ($\alpha = 1.0$) represent moderate heterogeneity more typical of real-world deployments.

Learning architecture comparisons establish fundamental baselines for evaluating federated learning effectiveness. Centralised learning aggregates all participant data

TABLE 3.1: Comprehensive Experimental Design Matrix

Dimension		Variations	Description
Data	Distribu- tion	IID vs. Non-IID	IID: Uniform random sampling across clients; Non-IID: Dirichlet distribution with $\alpha \in \{0.1, 0.5, 1.0\}$ to simulate varying heterogeneity levels
Learning	Archi- tecture	Centralised, Local, Federated	Centralised: All data aggregated on single server; Local: Isolated training per client; Federated: Collaborative training with distributed data
Federated	Scale	such as 10, 50, 100, 500, 1000 cli- ents	Systematic evaluation of scalability from small research settings to large-scale deployment scenarios
Privacy	Mechan- isms	None, Local DP, Central DP, k- Anonymity	Progressive privacy protection from baseline to strong anonymisation with $\epsilon \in \{0.1, 1.0, 10.0\}$ for differential privacy
Data	Sharing Strategies	None, Global Data Sharing, Cluster of Trust (CoT)	None: Standard federated learning; Global: 5% and 10% data sharing; CoT: Trust-based clustering with $K \in \{2, 5, 10, 20, \text{etc.}\}$ clusters
Data	Augmenta- tion Scale	1 \times , 2 \times , 5 \times , 10 \times	Augmentation multipliers indicating the ratio of synthetic samples to original data for addressing data scarcity

on a single server, providing optimal utility whilst completely sacrificing privacy. This approach serves as the performance upper bound for comparison purposes. Local learning restricts each participant to training exclusively on their own data, maximising privacy but severely limiting model quality due to insufficient training samples. Federated learning configurations explore the spectrum between these extremes, implementing various collaboration strategies whilst maintaining data distribution.

Federated scale experiments systematically evaluate framework performance across different participant numbers, from small research-oriented deployments (10 clients) to large-scale industrial scenarios (1000 clients). This scalability analysis reveals how framework components behave under varying communication loads, computational constraints, and coordination complexity. The experiments particularly focus on convergence behaviour, communication efficiency, and privacy guarantee maintenance as system scale increases.

Data sharing strategy experiments systematically evaluate different approaches to collaborative learning within federated environments. The baseline approach employs standard federated learning with no data sharing, relying solely on gradient exchange

between participants. Global data sharing strategies involve each client contributing a small percentage (5% or 10%) of their local data to a shared pool accessible by all participants, providing enhanced training diversity at the cost of direct data exposure. Cluster of Trust (CoT) strategies implement selective data sharing within trusted client groups, with systematic evaluation of different cluster configurations including 2, 5, 10, and 20 clusters to optimise the balance between collaboration benefits and privacy protection.

Privacy mechanism evaluation implements a progressive approach to privacy protection, starting from baseline federated learning without additional privacy measures and advancing through increasingly sophisticated protection mechanisms. Local differential privacy adds calibrated noise to gradients before transmission, with privacy budget ϵ controlling the noise level and corresponding privacy-utility trade-off. Central differential privacy applies noise aggregation at the server level, potentially offering better utility for equivalent privacy levels. K-anonymity mechanisms ensure that transmitted metadata cannot be uniquely attributed to fewer than k participants.

Data augmentation scale experiments address data scarcity challenges by systematically varying the volume of synthetic data generation. Baseline experiments (1 \times) use only original data to establish performance floors. Progressive augmentation scales (2 \times , 5 \times , 10 \times) indicate multiplication factors applied to original dataset sizes through domain-specific synthetic data generation. Higher augmentation scales provide greater dataset richness and potentially improved model performance, whilst requiring more sophisticated privacy protection mechanisms to ensure that synthetic data generation does not compromise the security of original datasets.

3.3.2 Comprehensive Evaluation Framework

Building upon the experimental design matrix, our evaluation framework systematically assesses the proposed methodology across three critical dimensions: model performance, privacy protection, and system efficiency. This multidimensional approach recognises that federated learning systems must simultaneously optimise these often-competing objectives to achieve practical viability.

Model Performance Evaluation focuses on the fundamental machine learning effectiveness of the federated system. We employ standard classification metrics including accuracy, precision, recall, F1-score, and AUC-ROC to assess predictive quality across different federated configurations. Convergence analysis examines the number of communication rounds required to reach target performance levels, whilst robustness evaluation tests model stability under varying client participation patterns and data distribution shifts. The evaluation particularly emphasises performance

degradation analysis, quantifying how different privacy mechanisms and system configurations affect model quality compared to centralised baselines.

Privacy Protection Assessment employs both theoretical privacy guarantees and empirical attack resistance evaluation. Theoretical measures include formal differential privacy analysis with epsilon-delta bounds, k-anonymity verification, and information-theoretic privacy quantification. Empirical evaluation implements realistic attack scenarios including membership inference attacks, property inference attempts, and gradient inversion techniques. We assess the success rates of these attacks under different privacy protection mechanisms to provide practical privacy assurance beyond theoretical guarantees. Additionally, we evaluate the privacy-utility trade-off curves to identify optimal operating points for different application scenarios.

System Efficiency Analysis encompasses computational and communication cost evaluation critical for practical deployment feasibility. Communication efficiency measures include total data transmission volumes, communication round frequencies, and bandwidth utilisation patterns across different client scales. Computational efficiency assessment covers local training costs, aggregation overhead, and the distribution of computational burden across system participants. We particularly focus on scalability analysis, examining how system performance degrades with increasing client numbers and identifying potential bottlenecks in large-scale deployments. Energy consumption analysis provides insights into the sustainability of federated learning deployments on resource-constrained devices.

3.4 Summary

This chapter has established a comprehensive methodological framework for privacy-preserving federated learning that addresses the fundamental challenges of non-IID data distributions, privacy protection, and data scarcity across multiple application domains. The framework is grounded in rigorous problem analysis and incorporates established techniques from distributed systems, privacy protection, and machine learning optimisation.

The problem analysis provides a clear foundation for understanding the statistical and computational challenges inherent in federated learning systems. By examining data heterogeneity through established statistical measures and clearly articulating the privacy-utility trade-offs specific to different domains, the framework establishes concrete targets for evaluation and practical deployment considerations.

Positioning and novelty. Many of the building blocks used in this framework (e.g., FedAvg-style aggregation, differential privacy, and clustered federation) are well established in the literature[2, 20]. The novelty of the framework presented in this thesis lies in how these elements are *integrated into a coherent methodology* for privacy-preserving deployments that can be instantiated consistently across application chapters:

- **Explicit three-tier architecture:** a client/cluster/global decomposition that makes responsibilities, collaboration boundaries, and information flows concrete and traceable across application chapters.
- **Trust/permission-constrained clustering:** an explicit permission model (trust graph) that constrains feasible client groupings; within this constraint, we study alternative clustering regimes (e.g., similarity-aligned vs. diversity/coverage-oriented grouping) using Jensen–Shannon divergence as the alignment criterion.
- **Combination-focused experimental design:** a single, domain-agnostic design matrix that tests *combinations* of key components by jointly varying data distribution, privacy mechanisms (e.g., anonymisation + differential privacy), data sharing strategies (baseline vs. CoT), and augmentation scale (including rotation-based augmentation in the HAR case study), with a consistent evaluation protocol across domains.

The three-tiered architectural design offers a scalable and flexible structure that accommodates diverse application requirements whilst maintaining consistent theoretical principles. The integration of trust clustering, adaptive privacy protection, and domain-specific data augmentation creates a synergistic system that addresses multiple challenges simultaneously rather than treating them as independent problems.

The experimental methodology provides a robust foundation for evaluating the framework’s effectiveness across diverse scenarios and deployment conditions. The comprehensive evaluation protocol ensures that results are statistically valid, practically relevant, and comparable across different domains and implementation approaches.

The methodology established in this chapter forms the foundation for the three application studies presented in subsequent chapters: IoT intrusion detection, privacy-preserving data publishing, and human activity recognition. Each application chapter demonstrates specific instantiations of this framework, validating its effectiveness and providing detailed analysis of domain-specific optimisations and trade-offs.

Chapter 4

Federated Learning for IoT Intrusion Detection with Non-IID Data

This chapter presents our approach to addressing the challenge of non-IID data in federated learning for IoT intrusion detection. We begin with an introduction to the specific challenges of IoT security, followed by a detailed description of our trust clustering-based federated learning approach. We then evaluate our approach using a comprehensive set of experiments on the IoT-23 dataset and discuss the results.

4.1 Introduction

Intrusion Detection Systems (IDS) play a vital role in network defence by identifying and alerting security administrators to potentially malicious activities. In the context of IoT, IDS methods can be broadly classified into signature-based and anomaly-based approaches. Signature-based systems rely on predefined attack patterns stored in databases and can only detect known intrusions, making them ineffective against zero-day attacks [53]. In contrast, anomaly-based approaches learn the patterns of normal behaviour and flag deviations as potential threats, with machine learning techniques—such as clustering, neural networks, decision trees, and genetic algorithms—proving highly effective in this domain [100, 101].

Traditional machine learning methods typically require centralised data collection, where raw data from edge IoT devices is transmitted to a central server for model training [101]. While this approach can achieve high detection accuracy, it introduces critical challenges: (1) high latency due to network transmission; (2) significant privacy and security risks arising from raw data transfer; and (3) infeasibility in scenarios where data cannot be legally or technically shared. Alternatively, training

isolated models on individual devices suffers from limited data availability and resource constraints, leading to suboptimal results.

Federated Learning (FL) offers a privacy-preserving paradigm that enables collaborative model training without sharing raw data [10]. In FL, edge devices train local models on their data and periodically share model updates (e.g., gradients or weights) with a central aggregator, which updates a global model. This decentralised approach maintains data locality while enabling effective global learning.

However, a key assumption in standard federated learning methods like FedAvg is that data is independently and identically distributed (IID) across clients. In reality, IoT devices often collect heterogeneous data due to variations in deployment environments, user behaviours, and threat exposures. Non-IID distributions lead to inconsistent local updates, which accumulate during aggregation, causing slower convergence, model instability, and degraded performance [18, 69, 102].

To address these limitations, this work proposes a strategy that introduces an intermediate layer—**Clusters of Trust (CoT)**—between the central aggregator and local nodes. Nodes with mutual trust and similar data distributions are grouped into clusters, and each cluster operates as a virtual client in the federated learning process. Models are first aggregated within clusters and then globally across clusters. This approach aims to improve model convergence and accuracy under non-IID conditions, while preserving privacy and reducing communication overhead.

Our Contributions

The main contributions of this work are as follows:

- We evaluate the performance of general federated learning (FedAvg) on benchmark IoT IDS scenarios under both IID and non-IID data distributions.
- We propose **CoT**, a trust-graph-constrained clustering strategy that (i) enforces collaboration permissions via an explicit trust graph and (ii) optimises statistical alignment using the Jensen–Shannon divergence of label distributions, treating each resulting cluster as a virtual client to improve convergence under non-IID data.
- We compare our CoT strategy with a global shared data approach, highlighting the trade-offs between privacy and statistical alignment.

Positioning and novelty. Client clustering in federated learning has been explored in the literature, including clustered FL formulations[20] and hierarchical clustering of local updates to improve non-IID training[74]. In this chapter, we propose and

empirically evaluate a **trust-based clustering method** for IoT-IDS, CoT, that differs from prior clustering approaches in two key ways. First, *trust relationships are treated as a first-class constraint* on feasible collaboration (i.e., clients can only be clustered if they form trust-connected groups in the trust graph). Second, within this trust-feasible space, cluster composition is guided and analysed through **diversity/coverage vs. similarity trade-offs** (e.g., broader vs. narrower threat-label coverage, characterised via distributional measures such as Jensen–Shannon divergence), and we quantify how these regimes influence convergence behaviour and detection performance. We evaluate CoT against standard FedAvg and global data sharing baselines on IoT-23, clarifying when and why trust-based clustering is most effective in non-IID IoT-IDS settings.

Mapping to the Three-Tier Framework (Chapter 3)

This chapter instantiates the three-tier architecture introduced in Chapter 3 as follows:

- **Client tier:** each IoT node trains on its local IoT-23 scenario subset and shares only lightweight statistics (for clustering) and/or model updates (for FL).
- **Cluster tier:** clients are first grouped based on trust relationships (trust graph constraint). Within the set of trust-feasible clusterings, we compare alternative cluster composition regimes (higher vs. lower intra-cluster diversity/label coverage) and observe their impact on convergence and accuracy; cluster-level aggregation then reduces non-IID drift before updates are sent upward.
- **Global tier:** the central server aggregates cluster-level updates into the global IDS model and orchestrates training rounds.

4.2 A Data Sharing Strategy for IoT IDS Based on Federated Learning

The FedAvg algorithm is a foundational method in federated learning. However, its performance significantly deteriorates under non-independent and identically distributed (Non-IID) data conditions. This challenge is particularly prominent in IoT intrusion detection scenarios, where different network nodes experience diverse traffic patterns and attack types. Such data heterogeneity leads to weight divergence during global aggregation, resulting in unstable convergence and suboptimal model performance (Figure 4.1).

To mitigate this issue, existing studies have explored global data sharing strategies, where each client uploads a small portion of representative data to be shared across all

participants (Figure 4.2). While such methods help partially align local data distributions and improve global model performance, they also introduce potential privacy risks and require careful design to prevent sensitive information leakage.

Inspired by this line of work, we propose a novel clustering-based strategy named *Cluster of Trust* (CoT) that addresses the statistical heterogeneity challenge in federated learning while preserving privacy. Our strategy combines a trust-based clustering mechanism with mathematical optimisation to form federated clusters, which aggregates clients into privacy-preserving and semantically aligned clusters under trust constraints.

Problem Formulation Let $\{C_1, C_2, \dots, C_N\}$ be N clients, each holding a local dataset \mathcal{D}_i with label distribution $P_i(Y)$. We assume a predefined **trust graph** $G = (V, E)$, where an edge $(i, j) \in E$ indicates mutual trust between clients C_i and C_j . In this model, trust is **binary and unweighted**: any two clients either trust each other (an edge exists) or they do not (no edge), and we do not model intermediate trust levels. We also assume a **static trust graph** during the reported experiments; handling evolving trust relationships over time is left as future work. Finally, while we provide an explicit clustering algorithm below, a formal analysis of its computational complexity is outside the scope of this thesis and is likewise left for future work.

Assumptions (IoT-IDS case study). The experiments in this chapter adopt the following assumptions to make the trust-clustering problem well defined and reproducible:

- **Trust establishment:** the trust graph is given (or constructed offline) prior to training and captures feasible collaboration permissions between clients.
- **Information sharing for clustering:** clients can share lightweight, privacy-preserving summary statistics (e.g., label-frequency vectors) required to compute statistical similarity, without sharing raw traffic.
- **Threat model scope:** the primary focus is privacy against raw-data exposure and instability from non-IID training; robustness against fully adversarial clients manipulating clustering inputs or poisoning updates is not the main evaluation axis in this chapter and is left for future work.

The objective is to partition the clients into K clusters $\{C_1, \dots, C_K\}$ such that:

- Each cluster C_k forms a connected subgraph in G (trust constraint).
- The label distribution $P_k(Y)$ of each cluster approximates the global label distribution $P_{global}(Y)$ as closely as possible (statistical alignment).

The global label distribution is defined as:

$$P_{global}(Y) = \frac{1}{N} \sum_{i=1}^N P_i(Y)$$

Each cluster's label distribution is:

$$P_k(Y) = \frac{1}{|C_k|} \sum_{C_i \in C_k} P_i(Y)$$

The optimisation objective is:

$$\min_{\{C_k\}} \sum_{k=1}^K \lambda_k \cdot D_{JS} (P_k(Y) \parallel P_{global}(Y))$$

where D_{JS} denotes the Jensen–Shannon divergence and λ_k is the weight (e.g., cluster size) to balance cluster importance.

Algorithm Overview Since full clustering with exhaustive client data exchange is impractical, we design a low-overhead algorithm where each client only uploads a normalised label frequency vector. The server performs clustering using trust-graph-constrained subgraphs and selects the partition that minimises the total JS divergence.

Algorithm 1 Cluster of Trust with Label Distribution Constraint (CoT-LDC)

Require: Trust graph $G = (V, E)$, local label distributions $\{P_i(Y)\}_{i=1}^N$, target cluster number K , number of FL rounds T

Ensure: Global model M_{global}

Phase 1: Cluster formation

Compute global label distribution: $P_{global}(Y) = \frac{1}{N} \sum_{i=1}^N P_i(Y)$

Extract all trust-connected subgraphs of G as cluster candidates

for all candidate partition $\{\mathcal{C}_1, \dots, \mathcal{C}_K\}$ **do**

for $k = 1$ to K **do**

 Compute cluster label distribution: $P_k(Y) = \frac{1}{|\mathcal{C}_k|} \sum_{C_i \in \mathcal{C}_k} P_i(Y)$

end for

 Compute cost $J = \sum_{k=1}^K \lambda_k \cdot D_{JS}(P_k(Y) \parallel P_{global}(Y))$

end for

Select the partition with the lowest cost J

Merge data within each cluster \mathcal{C}_k to form new client dataset D_k

Phase 2: Federated training with aggregated clusters

Initialize global model $M_{global}^{(0)}$

for $t = 1$ to T **do**

for each cluster-client $k \in \{1, \dots, K\}$ **in parallel do**

$M_k^{(t)} \leftarrow \text{LocalTrain}(M_{global}^{(t-1)}, D_k)$

end for

$M_{global}^{(t)} \leftarrow \sum_{k=1}^K \frac{|D_k|}{\sum_j |D_j|} \cdot M_k^{(t)}$

end for **return** $M_{global}^{(T)}$

In this design, mutual trust among participants is the primary criterion for clustering. Subsequently, the clustering result is further optimised by minimising statistical distribution divergence (e.g., Jensen–Shannon divergence) among clients within the same cluster. Each cluster is treated as a virtual client in the federated process, participating in model aggregation while maintaining internal consistency (Figure 4.3).

Unlike global data sharing, CoT does not require clients to expose any raw data. Instead, only intermediate model parameters are shared among clients that both trust each other and exhibit similar data characteristics. This localised collaboration indirectly reduces the statistical divergence of model updates, enhancing training stability and generalisation without compromising privacy boundaries.

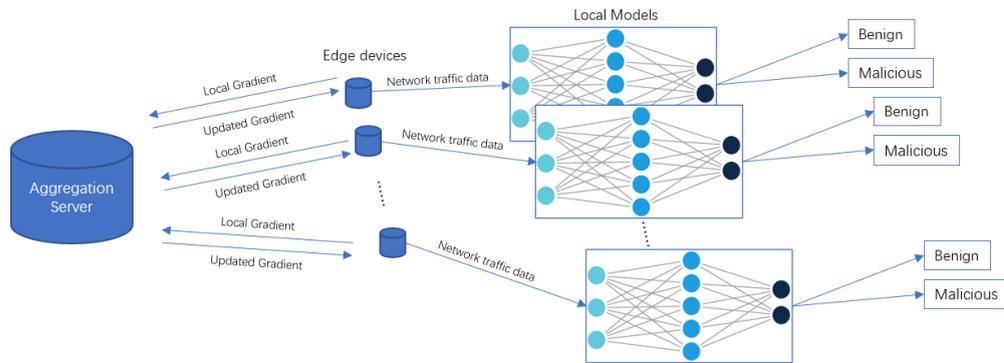


FIGURE 4.1: Structure of the FedAvg.

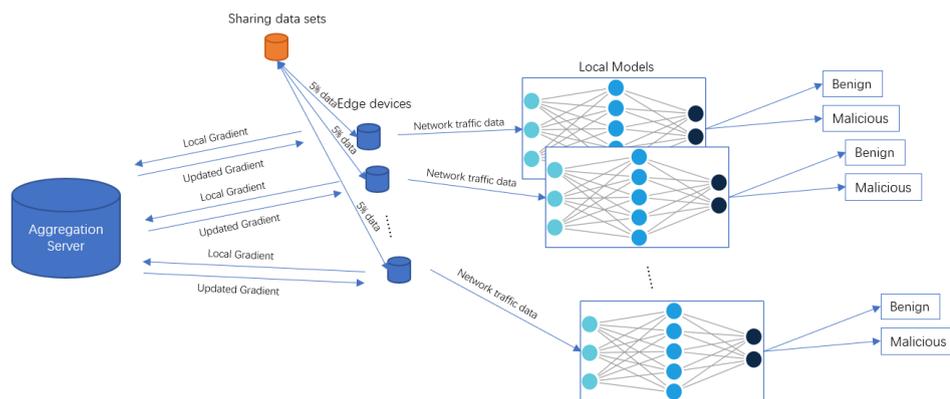


FIGURE 4.2: Federated Learning with Global Data Sharing

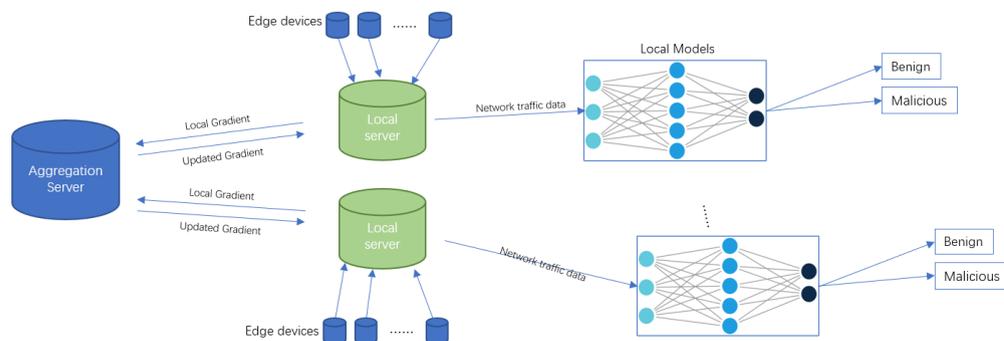


FIGURE 4.3: Structure of Clusters of Trust.

4.3 Experimental setup

4.3.1 Dataset Description and Preprocessing

4.3.1.1 IoT-23 Dataset Overview

In this study, we employ the IoT-23 dataset [?], a comprehensive collection of network traffic data specifically designed for IoT malware detection and cybersecurity research. The IoT-23 dataset is developed and maintained by the Czech Technical

University (CTU) and represents one of the most extensive publicly available datasets for IoT security analysis.

The dataset encompasses 23 distinct capture scenarios, categorized into two primary types:

- **Malware Capture Scenarios:** Real-world IoT malware behavior captured in controlled environments, representing various malware families and attack patterns.
- **Honeypot Capture Scenarios:** Network traffic collected through honeypot systems designed to monitor and detect network attacks while generating baseline benign traffic.

Each capture scenario contains detailed network flow records generated using the Bro/Zeek network analysis framework, providing rich metadata including connection durations, packet counts, byte volumes, protocol information, and temporal characteristics.

4.3.1.2 Data Preprocessing

The raw network logs were processed to extract relevant flow-level features suitable for machine learning analysis. The preprocessing pipeline included:

1. **Feature Extraction:** Extraction of 12 key network flow features including duration, total bytes, total packets, byte ratios, packet ratios, average packet size, flow rate, protocol encoding, and temporal information.
2. **Data Cleaning:** Removal of incomplete records and handling of missing values through appropriate imputation strategies.
3. **Feature Engineering:** Calculation of derived features such as bytes-to-packets ratios and flow rates to enhance the discriminative power of the feature set.
4. **Data Balancing:** Application of sampling techniques to address class imbalance issues while preserving the natural distribution characteristics of different scenario types.

Table 4.1 provides a comprehensive overview of the network flow features extracted for analysis.

TABLE 4.1: Network Flow Features Extracted for Analysis

Feature Name	Type	Description
duration	Numerical	Connection duration in seconds
total_bytes	Numerical	Total number of bytes transferred
total_pkts	Numerical	Total number of packets exchanged
bytes_ratio	Numerical	Ratio of outbound to inbound bytes
pkts_ratio	Numerical	Ratio of outbound to inbound packets
avg_pkt_size	Numerical	Average packet size in bytes
flow_rate	Numerical	Data transfer rate (bytes per second)
proto_encoded	Categorical	Encoded protocol type (TCP=1, UDP=2, ICMP=3)
hour	Temporal	Hour of day when connection occurred
is_malicious	Binary	Target variable (0=benign, 1=malicious)
scenario	Categorical	Specific capture scenario identifier
scenario_type	Categorical	Type of scenario (Malware/Honey-pot)

4.3.1.3 Dataset Characteristics

The final experimental dataset comprises 84,108 network flows with the distribution detailed in Table 4.2.

TABLE 4.2: Summary of IoT-23 Dataset Characteristics

Characteristic	Count	Percentage (%)
<i>Overall Dataset</i>		
Total Flows	84,108	100.0
Malicious Flows	57,237	68.1
Benign Flows	26,871	31.9
<i>Honeypot Scenarios</i>		
Total Honeypot Flows	1,956	2.3
Malicious Flows	0	0.0
Benign Flows	1,956	100.0
<i>Malware Scenarios</i>		
Total Malware Flows	82,152	97.7
Malicious Flows	57,237	69.7
Benign Flows	24,915	30.3
<i>Protocol Distribution</i>		
TCP Flows	45,232	53.8
UDP Flows	37,891	45.0
ICMP/Other Flows	985	1.2

Figure 4.4 illustrates the overall traffic distribution across different scenario types, demonstrating the natural separation between malware and honeypot environments.

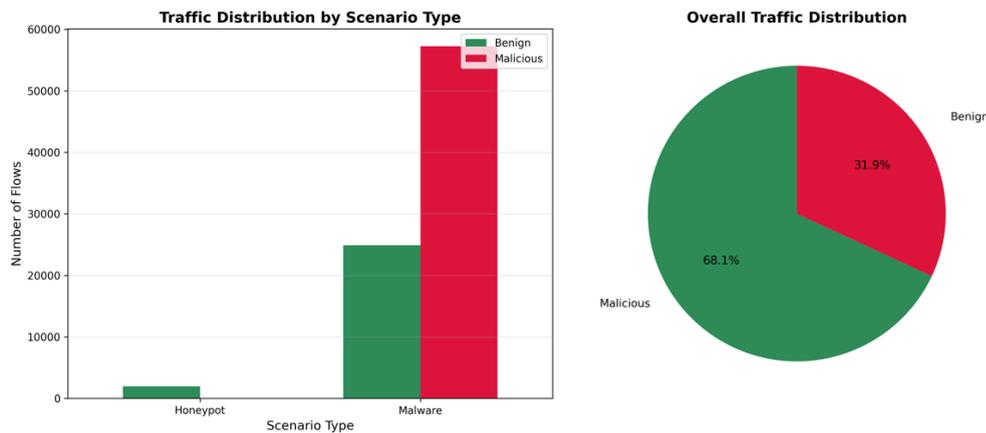


FIGURE 4.4: Traffic distribution by scenario type. (a) Distribution of benign vs. malicious traffic across Malware and Honeypot scenario types. (b) Overall distribution showing the proportion of malicious vs. benign traffic in the entire dataset.

4.3.1.4 Temporal and Protocol Analysis

We conducted comprehensive temporal and protocol analyses to understand the characteristics of our dataset. Figure 4.5 presents the temporal distribution patterns, revealing distinct hourly traffic behaviors between benign and malicious flows throughout the day.

The protocol distribution analysis, shown in Figure 4.6, reveals significant differences in network behavior between benign and malicious traffic across multiple dimensions including protocol preferences, packet size distributions, flow durations, and the relationship between packet counts and total bytes.

4.3.1.5 Scenario-Level Analysis

Figure 4.7 presents a detailed comparison between different capture scenarios, highlighting the distinct characteristics of honeypot versus malware capture environments and their respective traffic patterns.

4.3.2 Model Configuration

For the IoT intrusion detection task, we employed a deep neural network architecture designed to capture complex non-linear relationships within network flow features while maintaining robust generalisation capabilities. The selected architecture

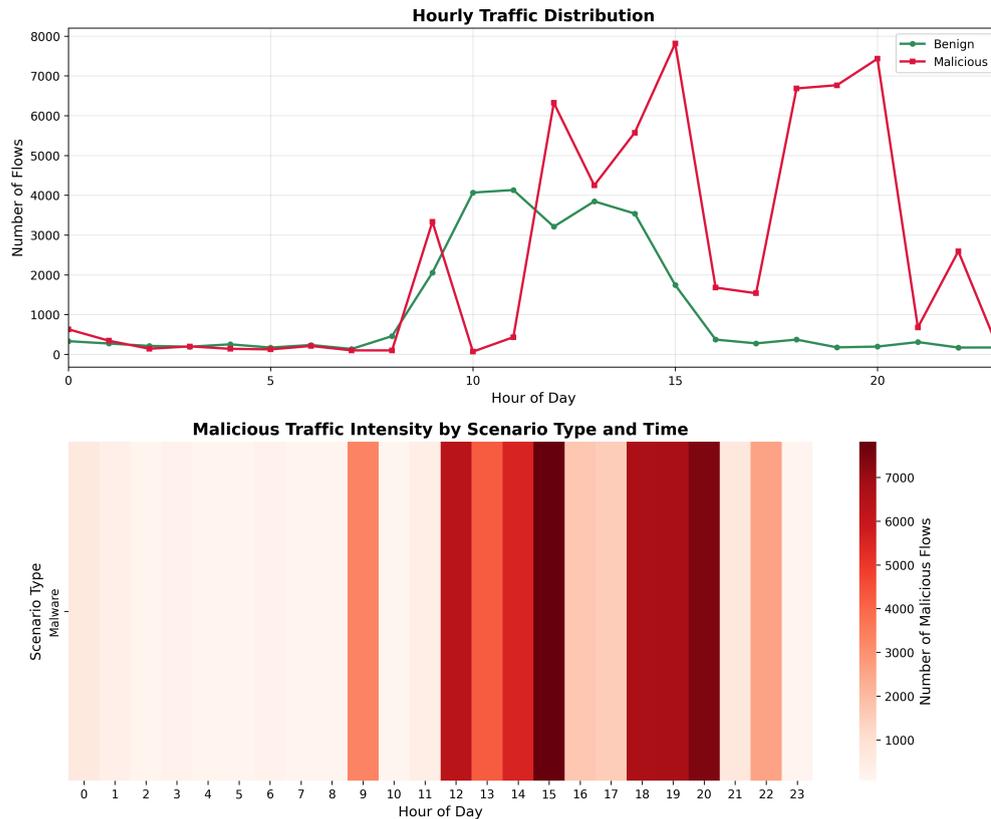


FIGURE 4.5: Temporal analysis of traffic patterns. (a) Hourly distribution of benign vs. malicious traffic throughout the day. (b) Heatmap showing malicious traffic intensity by scenario type and time, highlighting the concentration of malicious activities in Malware scenarios.

represents a carefully balanced approach between model complexity and computational efficiency, incorporating modern deep learning techniques to ensure optimal performance on the IoT-23 dataset.

The core architecture consists of a five-layer deep feedforward neural network with progressively decreasing neuron counts to create a hierarchical feature learning structure. The input layer receives 12 preprocessed features, followed by hidden layers containing 128, 96, 64, 32, and 16 neurons respectively. This tapering design allows the network to learn increasingly abstract representations of the input data, progressing from low-level network flow characteristics to high-level intrusion patterns. The final output layer employs a single neuron with sigmoid activation to produce binary classification probabilities for malicious versus benign traffic classification.

To address the challenges of training deep networks on cybersecurity data, we incorporated several advanced regularisation and optimisation techniques throughout the architecture. Batch normalisation layers are strategically placed after each of the first four hidden layers to stabilise training dynamics, reduce internal covariate shift, and enable higher learning rates. This normalisation technique proves particularly

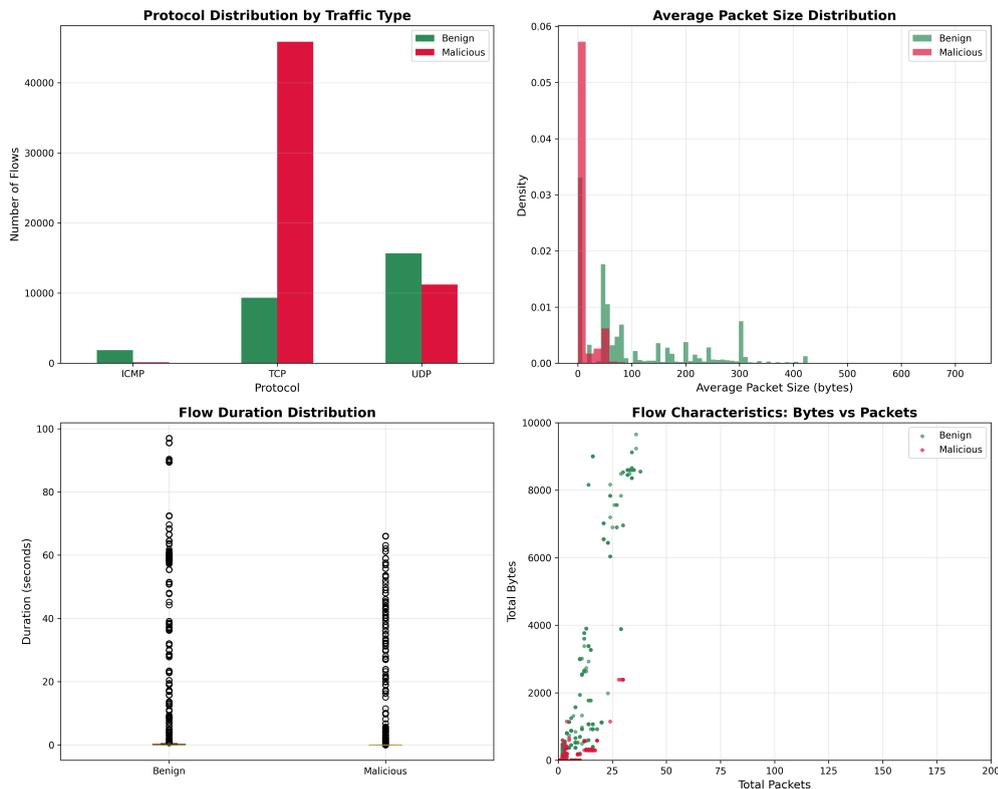


FIGURE 4.6: Protocol distribution and traffic characteristics analysis. (a) Protocol distribution comparison between benign and malicious traffic. (b) Average packet size distribution histograms. (c) Flow duration distribution box plots. (d) Scatter plot showing the relationship between total packets and total bytes, illustrating distinct clustering patterns for different traffic types.

effective in cybersecurity applications where network traffic features often exhibit varying scales and distributions across different scenarios and time periods.

Dropout regularisation is implemented with carefully tuned rates throughout the network to prevent overfitting while preserving the model’s learning capacity. The dropout rates are set to 0.4 for the first two hidden layers, 0.3 for the third and fourth layers, and 0.2 for the fifth layer. This progressive reduction in dropout rates reflects the decreasing complexity of representations as the network approaches the output layer. Additionally, L2 regularisation with a coefficient of 0.001 is applied to the kernel weights of the first four layers to further constrain model complexity and improve generalisation to unseen attack patterns.

The network utilises the Rectified Linear Unit (ReLU) activation function for all hidden layers, chosen for its computational efficiency and ability to mitigate the vanishing gradient problem in deep architectures. The ReLU activation also provides sparse activation patterns that align well with the binary nature of many network traffic features, enhancing the model’s interpretability in cybersecurity contexts.

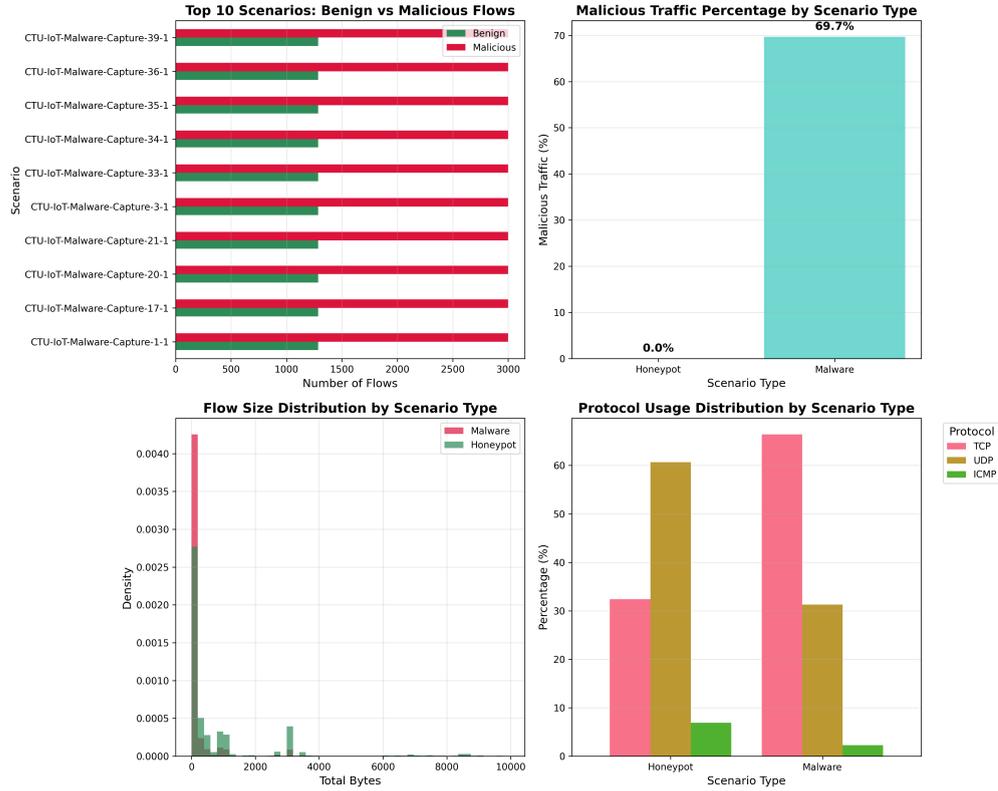


FIGURE 4.7: Detailed scenario comparison analysis. (a) Top 10 individual scenarios showing benign vs. malicious flow distribution. (b) Percentage of malicious traffic by scenario type. (c) Flow size distribution comparison between Malware and Honeypot scenarios. (d) Protocol usage distribution by scenario type.

For optimisation, we employ the Adam optimizer with an initial learning rate of 0.001, which adaptively adjusts learning rates for individual parameters based on gradient history. This optimisation strategy proves particularly effective for cybersecurity datasets where different features may require different learning dynamics. The binary crossentropy loss function is used as the optimisation objective, providing appropriate gradients for the binary classification task while maintaining numerical stability during training.

The training process incorporates several advanced techniques to ensure robust model development and prevent overfitting. Early stopping monitoring validation loss with a patience of 15 epochs automatically terminates training when no improvement is observed, preventing unnecessary computation and reducing the risk of overfitting to the training data. A learning rate reduction schedule decreases the learning rate by a factor of 0.5 when validation loss plateaus for 10 consecutive epochs, allowing the model to fine-tune its parameters in later training stages. Model checkpointing saves the best-performing model based on validation loss, ensuring that the final model represents the optimal configuration discovered during training.

4.3.3 Baseline scenarios

To analyse the performance of our proposed strategy, we implemented the IoT IDS in different scenarios and approaches include:

1. We clustered the sub-data sets of IoT-23 into an aggregated dataset, in order to implement IDS by centralised machine learning.
2. Since each sub-dataset is generated in a highly non-homogeneous distribution, each sub-dataset contains a different type and amount of malicious traffic. We trained the original IoT-23 data using the FedAvg algorithm, by maintaining 23 nodes each one containing data captured in a single scenario. This measures the performance of federated learning on non-IID datasets.
3. We reshuffled the data amongst the 23 nodes as follows. We obtained 23 sub-datasets of the same size by random sampling from the aggregated IoT-23 dataset, which eliminated statistical heterogeneity between the individual sub-datasets. This results in an IID distribution of the data. We trained the FedAvg algorithm on this dataset to measure federated learning performance on the IID dataset.

4.3.4 Data Sharing Strategies

For globally shared dataset strategies, we add 5% and 10% global data to each subset of IoT-23 to measure the improvement .

For clusters of Trust strategies we have implements the Clusters of Trust method containing 2, 4, 7, 10 nodes respectively. For each size of clusters we compare two limit cases where the clusters have maximum label coverage and minimum label coverage for threats. We achieve this by either assuming a cluster is composed of nodes with very different types of threats (and choosing among nodes to form our clusters appropriately) or we chose to form clusters where the nodes have very similar threats.

4.4 Experimentation Result

4.4.1 Federated learning in IoT intrusion detection baseline scenarios

We tested the performance of the IDS on baseline scenarios for comparison. The results of the training are shown in Figure 4.8 and Table 4.3. The red line indicates the loss curve for centralised machine learning; the blue line indicates the loss curve for

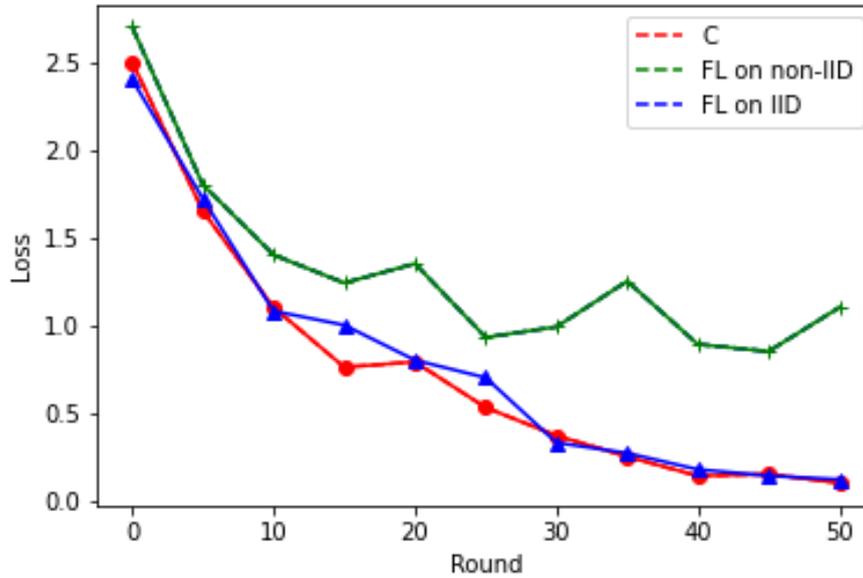


FIGURE 4.8: Loss on Baseline Scenarios.

TABLE 4.3: Performance Metrics.

	Recall	Precision	F1 score
Centralised	92.3%	91.0%	91.7%
Federated	86.3%	90.7%	88.4%

federated learning trained with IID data; and the green line indicates the loss curve for federated learning trained with non-IID data.

We can see that centralised machine learning and federated learning on IID dataset perform more accurately and converge faster. In contrast, the loss curves show the erratic loss convergence and lower model accuracy of the federated learning model on non-IID dataset during gradient descent. This is due to the fact that each subset of the IoT-23 data (except for the 3 conscientious datasets) is a dataset of traffic generated by different malware, each containing different types of attacks. And different attacks have different behavioural patterns, which causes different subsets of data to be non-IID. In contrast, in federated learning, the SGD algorithm performed locally at each node aims to minimise the loss value of local samples on that device, i.e., the local model tends to differentiate between the types of attacks contained in that dataset. As we continue to fit models on different devices to the global model, the differences between the weights of these local models will accumulate and eventually degrade the performance of global learning.

4.4.2 Clusters of Trust vs Globally shared data

For the global shared data method, the results show that the accuracy of the federated learning model is 55.7%, 73.2% and 87.8% when the percentage size of the global dataset is 0%, 5% and 10%, respectively. Figure 4.9 reflects the loss convergence curves of the 2 two models; with light blue and legend 'C' we plot the loss convergence of the centralised solution.

We can see that as the size of the global dataset increases, federated learning converges faster, which means fewer rounds of communication are required, which is a significant improvement for IoT environments where communication bandwidth is constrained. For the IoT dataset, about 15% accuracy improvement can be achieved with just 5% of the globally shared data, while about 15% accuracy improvement can be achieved with 10% of the globally shared data, achieving similar accuracy to the centralised machine learning training.

For the clusters of trust method, the results are shown in Figure 4.10 and Table 4.4. Figure 4.10 shows accuracy when clusters have minimum threat coverage tagged with 'A' and shown in blue and accuracy when clusters have maximum threat coverage tagged with 'B' and shown in orange. Table 4.4 shows the accuracy in each case of clustering (based on the size of nodes in a cluster, i.e., the size of clustering); For each clustering size the first line of the table shows the limit case of having similar kinds of threats in the cluster while the second line shows wider threat coverage per cluster.

The results of our Clusters of Trust experiments show that the accuracy of the model increases as the cluster contains more nodes and its label coverage increases. When the cluster contains 4 nodes and has maximum label coverage, the accuracy of the model exceeds 70%. And when the cluster contains 10 nodes even with minimum label coverage, the accuracy of the model exceeds 90%, which is very close to the performance of centralised machine learning.

TABLE 4.4: Accuracy for different clusters.

Size of Cluster	Coverage of malicious label	Accuracy
1	19.7%	55.7%
2	26.0%	59.7%
	31.0%	65.0%
4	36.7%	68.3%
	48.3%	72.7%
7	53.3%	75.3%
	67.7%	78.5%
10	83.3%	90.5%
	100%	91.9%
23	100%	92.5%

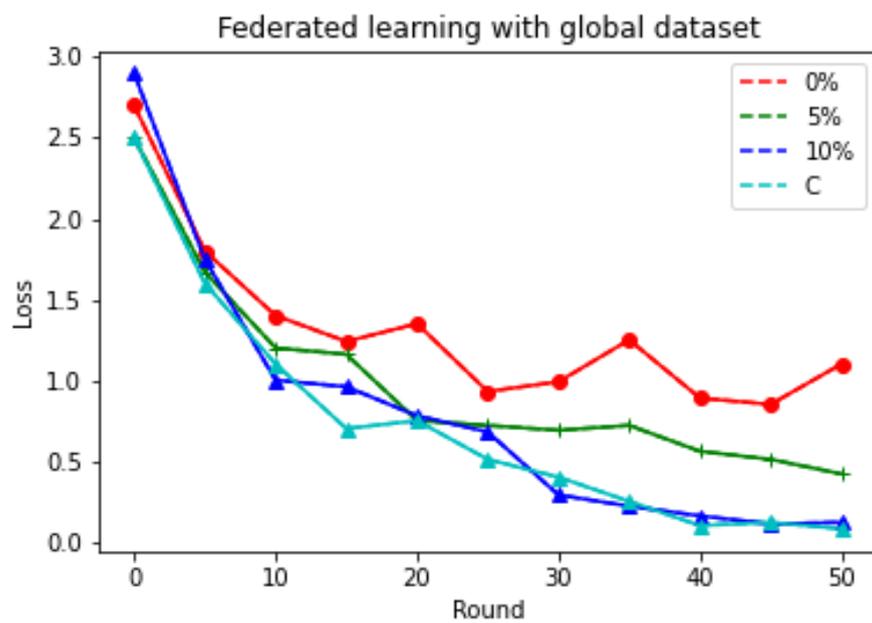


FIGURE 4.9: Loss on non-IID data with global dataset.

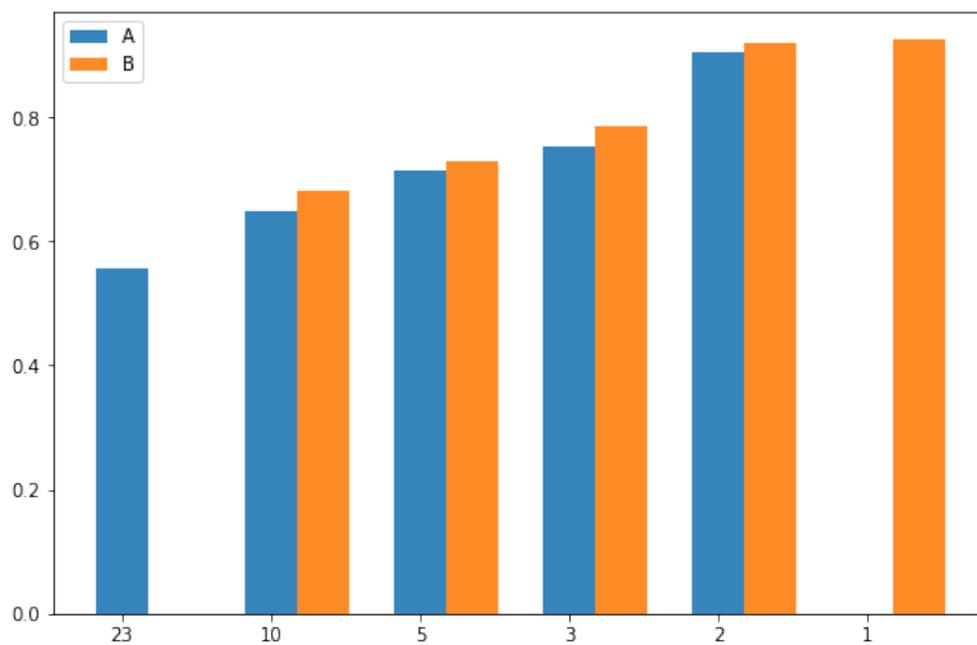


FIGURE 4.10: Accuracy histograms for different size clusters.

In terms of accuracy, we can see that for the clusters of trust strategy, the accuracy of the model exceeds the accuracy of the strategy with 5% and 10% of globally shared data added when the cluster size is 7 and 10 respectively. In terms of cost, adding 5 % and 10 % of the globally shared dataset equates to an average of 110 % and 220 % more data per node, which increases the communication costs, as well as requirements of devices' memory and computation power. For clusters of trust strategy, additional cluster servers are required, but there is no need to worry about the lack of memory or computing power of the edge devices. Furthermore, the globally shared dataset is obtained by sampling from the original dataset, which makes the model sacrifice privacy. In contrast, in the cluster of trust approach, nodes only send their own data to the trust cluster server for training, which is more in line with the limited communication bandwidth and computing power of edge devices in the IoT environment. Moreover, as the clusters are divided according to the trust relationship of the nodes, there is no risk of privacy leakage.

4.5 Conclusions

Federated learning will play a key role in IoT intrusion detection. However, due to the heterogeneity of devices and the diversity of attacks in IoT scenarios, edge devices that are in different networks are often subject to different attacks. As a result, edge devices generate non-IID data, and then the quality of model training is degraded. To make federated learning mainstream, however, improving model training on non-IID data is key to making progress in this area. In this work, we first show that for neural networks trained on non-IID data, the accuracy and convergence rate of federated learning is greatly reduced, with an accuracy of down to 55%. We validate a strategy to improve the training of non-IID data by allowing mutually trusted nodes to aggregate into a cluster that is updated as a new federated learning client. Experiments on the Aposemat IoT-23 dataset show that federated learning-based IDS models using this strategy are considerably more accurate and can ideally achieve performance comparable to centralised machine learning when there is good coverage of threats in each cluster of trust. In terms of communication and computational costs and privacy, the clusters of trust strategy is more suitable for IoT intrusion detection scenarios than the global data sharing strategy.

Generalisability and external validity. These results are obtained on a single dataset (IoT-23). Nevertheless, the proposed mechanism is *dataset-agnostic* in structure: it relies on (i) the existence of statistical heterogeneity across clients and (ii) a trust/permission constraint that limits feasible collaboration. Both conditions are common in real-world IoT deployments spanning different network segments and organisations. We therefore expect the qualitative trend—that trust-constrained

clustering can stabilise aggregation under non-IID data—to transfer to other suitable IoT-IDS datasets, while acknowledging that the magnitude of improvement will depend on dataset-specific factors (e.g., attack taxonomy, feature representation, and client partitioning). Multi-dataset validation and sensitivity analysis across alternative partitions are left as future work.

Chapter 5

A Dual-Layer Privacy-Preserving Federated Learning Framework

This chapter presents our privacy-preserving federated learning framework for data publishing and mining. We begin by introducing the challenges of privacy-preserving data publishing and mining, followed by a detailed description of our framework that integrates differential privacy, data anonymisation, and federated learning. We then evaluate our approach using experiments on benchmark datasets and discuss the results.

5.1 Introduction

In the era of big data, the amount of data generated by humans is exponentially increasing. A significant portion of this data is inherently personal and includes various types of patient data collected and stored in electronic health records within the electronic health systems. This data encompasses laboratory test results, demographic information, age and weight statistics, as well as medication information [103]. Similarly, in social networks, data such as user names, addresses, email addresses, personal photos, and notes are collected [104]. This data can be used for numerous scientific or commercial purposes, such as data-driven research and product development, relying on the analysis of personal information to generate knowledge-based decisions or provide personalised services. However, due to the presence of personal information in the data, it may face threats from inference attacks, where attackers can infer sensitive or private information that has not been explicitly disclosed by utilising patterns, correlations, or statistical characteristics present in the available data [105]. Therefore, privacy protection is crucial in the publication and utilisation of personal data. The "privacy by design" paradigm [106] emphasises minimising the use of sensitive information. The General Data Protection

Regulation (GDPR) [107] introduced by the European Union provides a strict and mandatory framework for safeguarding personal and sensitive information.

To address the threats posed by privacy attacks, anonymisation techniques and federated learning are widely employed to protect the publication and usage of privacy-sensitive data [108][109]. In the data collection phase, privacy-preserving data publishing (PPDP) offers a set of models, tools, and methods to mitigate privacy threats when releasing data. Data owners protect users' private information by applying anonymisation techniques, including generalisation, suppression, microaggregation, to prevent inference attacks while preserving the utility of the anonymised data. In the data processing phase, the application of machine learning is becoming increasingly widespread. With the explosive growth of data and the rapid generation of information, traditional methods of data processing and analysis are no longer able to meet the demands of mining and insights from massive data. The development of machine learning technology enables people to harness the valuable information hidden in big data and extract profound insights and patterns, providing more accurate and efficient solutions. However, considering most real-world scenarios, personal data is often scattered across data islands, such as different healthcare institutions or banking systems. However, most traditional machine learning algorithms operate in a centralised manner, requiring data aggregation on data servers. This introduces a single point of failure and significant risks of data breaches, leading to a lack of trust among end users and challenges in complying with GDPR requirements. To overcome these challenges, Google researchers introduced federated learning as a promising solution [5], which has gained attention from both industry and academia. Unlike traditional machine learning, federated learning is a framework that enables machine learning algorithms to be implemented in a decentralised collaborative learning setting. In federated learning, models are executed on multiple local datasets stored on various local nodes, such as smartphones, tablets, personal computers, and Internet of Things (IoT) devices [110]. Each device uses its own local data to train the model, and after the local training is complete, each device sends the updated model parameters to the central server. The central server aggregates all received model parameters to generate new global model parameters. [108]. This allows local nodes to collaboratively train a shared machine learning model, exchanging only the trained parameters (e.g., weights and biases of deep neural networks) periodically, without the need to centrally collect and process training data on a central data server. As a result, federated learning possesses a natural advantage in preserving data privacy. Furthermore, the parameter update and aggregation processes between local nodes and the central coordinating server can be enhanced with privacy protection techniques, such as differential privacy, to further strengthen data privacy protection [111][112].

However, despite the individual achievements of anonymisation and federated

learning in privacy protection during different stages of data processing, combining both techniques into a comprehensive privacy protection framework poses challenges. Previous research introducing anonymisation algorithms often overlooked the relationship between anonymisation and machine learning, partially due to the distinct origins of these two methods. It has been noted that information loss resulting from generalisation and suppression in anonymisation methods may lead to performance degradation in machine learning models [49, 113, 114].

Additionally, as a distributed machine learning framework, federated learning distributes data across local nodes instead of aggregating it on a central server. While the broader privacy-preserving federated learning literature extensively discusses privacy–utility trade-offs (e.g., differential privacy and secure aggregation)[1, 3], comparatively fewer application-grounded studies characterise how *anonymisation applied during distributed data collection* (e.g., microaggregation and generalisation) interacts with federated training dynamics and downstream model performance. Intuitively, with an increasing number of data holders, each holder possesses a smaller amount of data, potentially reducing local anonymity and increasing sensitivity to information loss. To achieve anonymity constraints comparable to centralised datasets, distributed datasets may require stronger anonymisation transformations, which can further compromise model performance.

Relatedly, some work studies “anonymisation” at a different layer—*identity/update anonymity*—and shows that anonymising client identities/updates in communication does not necessarily prevent client re-identification (linkability) or client-level privacy leakage under a malicious server[30]. This further motivates explicitly distinguishing the **data-collection anonymisation** problem studied in this chapter from **update/identity anonymity** in FL communications.

Therefore, this chapter first investigates the impact of anonymisation strategies across three training modes (federated learning, centralised machine learning, and a hybrid mode), directly addressing this under-characterised interaction. Building on these findings, we propose a **dual-layer privacy-preserving FL framework** that combines *microaggregation-based anonymisation* at data collection with *differential privacy* at model training, and we quantify the resulting privacy–utility trade-offs.

The main contributions of this chapter are as follows:

- A dual-layer privacy-preserving federated learning framework to explore privacy protection in both the data collection and model training stages while maintaining acceptable model performance.
- Comparison of anonymisation techniques between federated learning on distributed data sets and centralised machine learning on centrally collected data sets using a real-world data set.

- Assessing the impact of different privacy-preserving anonymisation strategies on 3 training scenarios: anonymisation on federated learning, anonymisation on centralised learning, and centralised learning on datasets anonymised before aggregation.
- Evaluation of the proposed framework on a real-world data set, demonstrating improvements in both data anonymity and model performance.

Mapping to the Three-Tier Framework (Chapter 3)

This chapter instantiates the three-tier architecture from Chapter 3 with an emphasis on privacy protection at both data collection and model training:

- **Client tier:** each data holder performs local data collection and applies anonymisation (e.g., generalisation/microaggregation) before participating in training; local training occurs on the (possibly anonymised) local dataset.
- **Cluster tier:** trust-based clustering is *not* the primary experimental variable in this chapter; conceptually, clustering can be introduced to group collaborating data holders under permission constraints, but the reported experiments focus on anonymisation and differential privacy trade-offs.
- **Global tier:** the central server orchestrates FL rounds and applies model-level privacy mechanisms (e.g., differential privacy during aggregation/training) to provide formal privacy guarantees.

Assumptions (Data Publishing case study)

The experiments and privacy metrics in this chapter rely on the following assumptions:

- **Attribute model:** a subset of attributes is treated as quasi-identifiers and a designated sensitive attribute is used for anonymity metrics; alternative attribute choices may change measured anonymity levels.
- **Threat model scope:** we focus on privacy leakage through published/anonymised data and through model-update exposure mitigated by differential privacy; robustness to fully malicious participants performing active poisoning is not the primary focus.
- **Privacy mechanism configuration:** differential privacy is implemented via noise addition with a chosen budget/noise multiplier; the empirical privacy-utility

results are conditional on these settings and should be interpreted as trade-offs rather than universal constants.

Through these contributions, we provide a comprehensive understanding of the impact of anonymisation strategies on federated learning, filling a gap in the field of privacy protection. Moreover, we introduce an innovative federated learning framework that simultaneously enhances data privacy and improves model performance in big data applications. This work is of significant importance in advancing privacy protection and the application of big data, providing valuable guidance for future research and practical implementations.

5.2 Methodology and framework

This section first describes our proposed framework for privacy-preserving federated learning, followed by a description of dataset and methods for anonymisation. Our approach aims to improve the anonymity of data and to evaluate the impact of anonymisation on the performance of federated learning and centralised machine learning models.

5.2.1 Data Description

In our study, we want to evaluate the impact of our framework on the performance of anonymity and federated learning models from a real, distributed dataset. Thus we utilise the IoT-23 dataset [115] which aims to provide researchers with a large-scale, labelled dataset of IoT traffic to facilitate the development of machine learning algorithms. With 23 sub-datasets, the IoT-23 dataset covers a wide range of scenarios where network data was collected. These sub-datasets consist of network traffic data in pcap format, accompanied by labels indicating instances of malicious behaviour.

The dataset consists of 23 features, among which 'proto', 'orig_p' and 'ts' are selected as quasi-identifiers, and 'detailed-label' is sensitive data. 'proto' represents the protocol used in the network traffic packets, including 'tcp', 'udp', and 'icmp'. 'orig_p' represents the port used in the network traffic and 'ts' denotes the timestamp of the network connections, which has been standardised and retained with 4 significant digits. 'detailed-label' refers to the type of attack the system is subjected to.

5.2.2 Anonymisation methods

In order to enhance data anonymity, we employ two strategies: generalisation and microaggregation.

5.2.2.1 Generalisation

Starting from the highest level of generalisation for the quasi-identifiers, we recursively specialise the partitions using multi-dimensional cuts until no further cuts can be made. In each iteration of the algorithm, a dimension (attribute) is selected for cutting. A common approach is to choose the dimension with the widest value range. Then, using the median split, we determine the splitting value and perform the cut based on that value.

5.2.2.2 Microaggregation

Initially, we create clusters with at least k similar records. We then choose a representative value to replace all the quasi-identifiers within each group. The selection of the representative value can be done using different methods, such as the mean, median, or random selection.

5.2.3 A dual-layer privacy-preserving federated learning framework

Given that anonymisation and federated learning achieve privacy protection at different stages of data processing, we propose a dual-layer framework (shown in Figure 5.1 and Algorithm 2) that combines these two techniques to achieve comprehensive privacy protection from data collection to data utilisation.

The framework divides data processing into two stages: privacy data collection and privacy data usage.

1. In the first stage, after collecting individual information, the framework assesses the privacy metrics of the raw data and performs anonymisation operations, such as generalisation and microaggregation, to ensure that data collection meets the requirements of Privacy-Preserving Data Publishing (PPDP).
2. In the second stage, through federated learning, the model is executed on multiple local datasets stored on various local nodes, utilising a central server to coordinate the training process. Additionally, the parameter update and aggregation processes between local nodes and the central coordinating server can be enhanced with differential privacy to strengthen data privacy protection. We add artificial noise to the parameters on the clients before aggregation to prevent inference attacks against the model gradients.

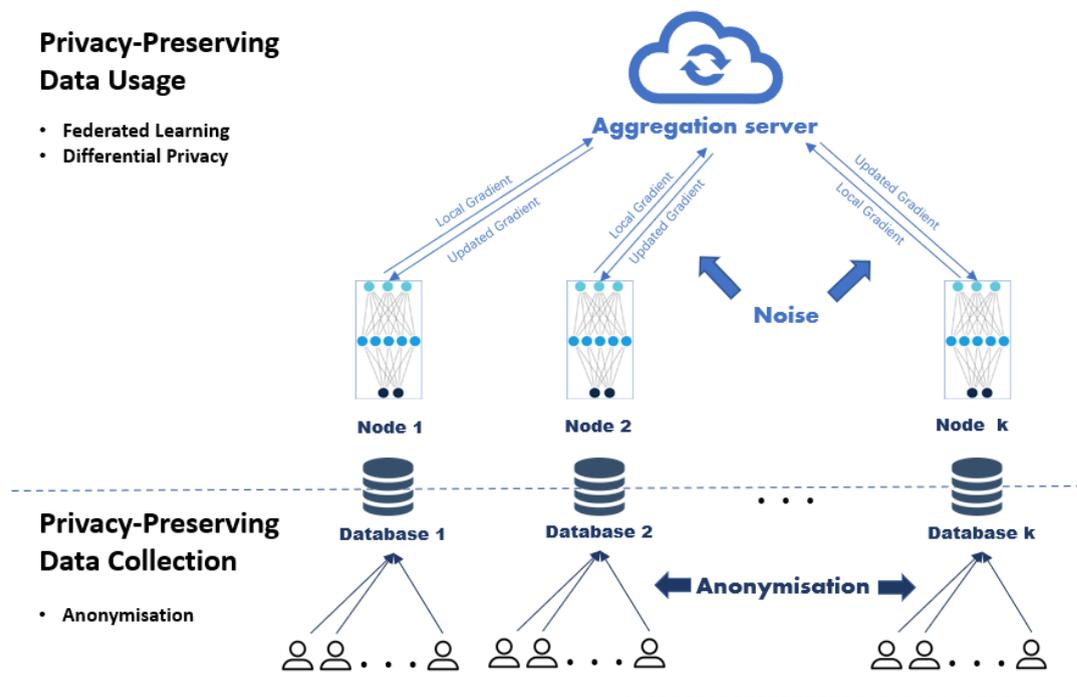


FIGURE 5.1: Framework of the Dual-Layer Privacy-Preserving Federated Learning Framework.

5.2.4 Privacy against centralisation

We generated two types of datasets, namely multiple distributed datasets and a single centralised dataset. First we sample 23 datasets from each of the 23 nodes as datasets $D = \{a_1, a_2, \dots, a_{23}\}$ for distributed scenario; then we aggregate all the datasets from D into a single dataset C as the centralised scenario.

To explore the performance of data anonymisation and federated learning, we consider the following scenarios:

1. Calculating different privacy parameters (k-anonymity, l-diversity, t-closeness) for a real-world dataset, IoT-23, in both centralised and federated learning distributed settings.
2. Comparing the effects of different privacy-preserving anonymisation strategies on federated learning and centralised machine learning models. To test the impact of data anonymisation on the performance of federated learning and centralised machine learning models, firstly, we established three experimental scenarios:
 - (a) Federated Learning scenario: Anonymisation of the distributed datasets D followed by federated learning.

Algorithm 2 Dual-Layer Privacy-Preserving Federated Learning**1: Stage 1: Client-Side Data Anonymisation (Pre-processing)****Require:** Original datasets $D = \{D_1, D_2, \dots, D_N\}$, Anonymity parameter k **Ensure:** Anonymised datasets $D' = \{D'_1, D'_2, \dots, D'_N\}$ 2: **for** each client $c \in \{1, \dots, N\}$ **do**3: Identify quasi-identifiers in D_c 4: $D'_c \leftarrow \text{ApplyAnonymisation}(D_c, k)$ ▷ e.g., using Microaggregation5: **end for**6: **Return** D' **7: Stage 2: Differentially Private Federated Learning****Require:** Anonymised datasets D' , Learning rate η , Number of rounds T , Clients per round m , Local epochs E , Clipping norm S , Noise multiplier σ **Ensure:** Final global model w_T 8: **Server Executes:**9: Initialize global model w_0 10: **for** each round $t = 0, 1, \dots, T - 1$ **do**11: $S_t \leftarrow$ (Randomly select m clients)12: **for** each client $c \in S_t$ **in parallel do**13: $w_c^{t+1} \leftarrow \text{ClientUpdate}(c, w_t, D'_c)$ 14: **end for**15: $w_{t+1} \leftarrow \sum_{c \in S_t} \frac{|D'_c|}{\sum_{j \in S_t} |D'_j|} w_c^{t+1}$ ▷ Aggregate client models16: **end for**17: **Return** w_T 18: **Function** ClientUpdate(client c , global model w , local data D'_c):19: $\mathcal{B} \leftarrow$ (Split D'_c into batches of size B)20: **for** each local epoch $i \in \{1, \dots, E\}$ **do**21: **for** batch $b \in \mathcal{B}$ **do**22: $g \leftarrow \nabla \mathcal{L}(w, b)$ ▷ Compute gradient on local data23: $\tilde{g} \leftarrow g / \max(1, \frac{\|g\|_2}{S})$ ▷ Clip gradient24: $\hat{g} \leftarrow \tilde{g} + \mathcal{N}(0, S^2 \sigma^2 \mathbf{I})$ ▷ Add Gaussian noise for DP25: $w \leftarrow w - \eta \hat{g}$ ▷ Update local model26: **end for**27: **end for**28: **Return** w

- (b) Centralised scenario: Anonymisation of the centralised dataset \mathbf{C} followed by centralised machine learning.
- (c) Hybrid scenario: Anonymisation of the distributed datasets \mathbf{D} followed by aggregation of the datasets into a single consolidated dataset for centralised machine learning.

And then different levels of anonymisation were implemented in each scenarios using generalisation and aggregation. Specifically, we tested anonymisation with different degrees of k-anonymity (k-anonymity = 5, 10, 50, 100, 150, 200).

3. Exploring the trade-off between data privacy and model performance in the proposed dual-layer privacy-preserving federated learning framework. We first selected the anonymisation technique that resulted in the least performance loss for the federated learning model from the previous experiment. Then, we tested the performance of the models by varying the levels of anonymisation in the data collection and the levels of noise in the federated learning to find a balance between data privacy and model performance.

5.3 Experiments

5.3.1 Privacy metrics

First, we measured the privacy metrics (K-anonymity, l-diversity, t-closeness) of the federated learning dataset and the centralised machine learning dataset on IoT-23. A higher value of K-anonymity and l-diversity and a lower value of t-closeness indicate a higher level of anonymisation and better privacy of the data. Table 5.1 presents the results of the privacy metrics, and we can observe that the federated learning distributed dataset has lower values for each privacy metric compared to the centralised machine learning dataset.

TABLE 5.1: Anonymity metrics of distributed datasets and centralised dataset

	Distributed dataset	Centralised dataset
k-anonymity	4	93
l-diversity	2	7
t-closeness	5.9	0.34

5.3.2 The impact of anonymisation on model performance

We then compare the impact of different privacy-preserving anonymisation strategies on federated and centralised machine learning models, using generalisation, microaggregation, and increasing the k-anonymity of the data to 5, 10, 50, 100, 150 and 200. Figures 5.2–5.7 show the accuracy of the models for each combination, where each bar chart represents a type of anonymisation operation. The blue colour represents the accuracy of federated learning, the orange colour represents the accuracy of the hybrid scenario, and the green colour represents the accuracy of the centralised scenario.

The results show that under the anonymity constraints we tested, regardless of the anonymisation method used, the performance of centralised machine learning only slightly decreases from an initial accuracy of 92% to a minimum of 89%. For federated learning and the hybrid scenario, as the anonymity constraints increase, the model's

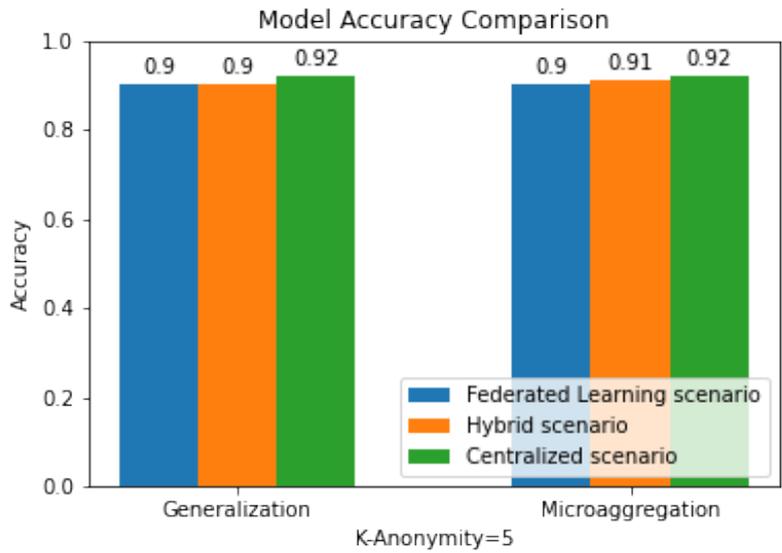


FIGURE 5.2: Accuracy of the model after anonymisation - (k1)

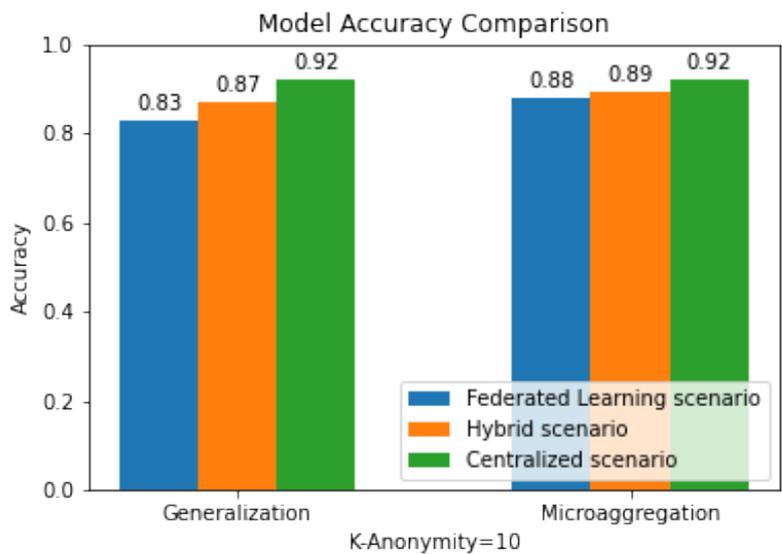


FIGURE 5.3: Accuracy of the model after anonymisation - (k=5)

accuracy continuously decreases, with federated learning experiencing a larger decrease. Specifically, when the k-anonymity is less than or equal to 50, the Microaggregation operation incurs the least performance loss, with accuracy of 83% for federated learning and 85% for the hybrid scenario, which are still within an acceptable range. However, when the k-anonymity is greater than 100, the accuracy of federated learning drops significantly, reaching a minimum of 63%.

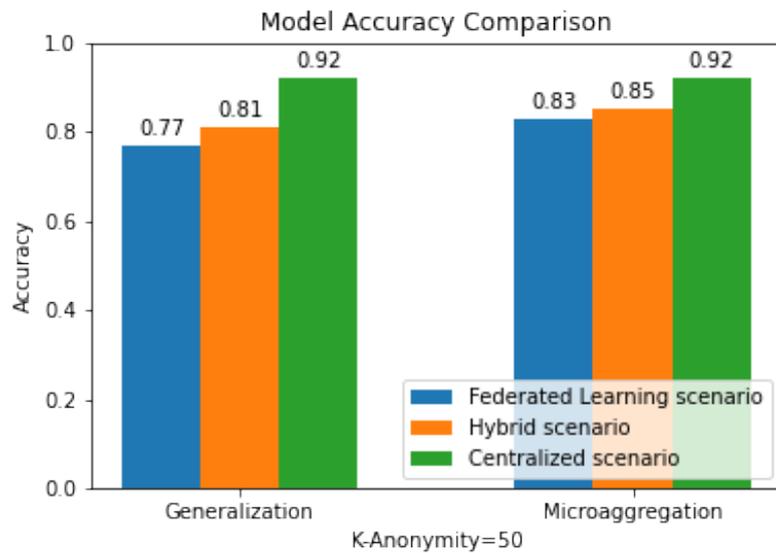


FIGURE 5.4: Accuracy of the model after anonymisation - (k=10)

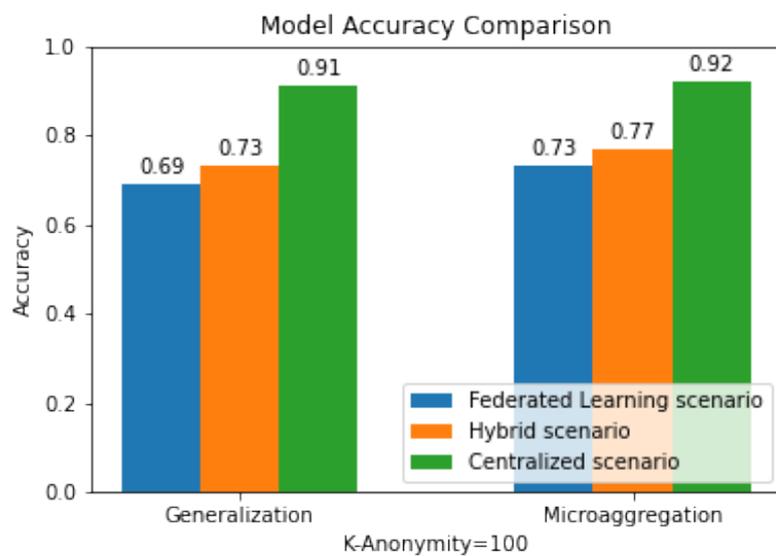


FIGURE 5.5: Accuracy of the model after anonymisation - (k=100)

5.3.3 Trade-off in the Dual-Layer Framework

Based on the results of the previous experiment, in the data collection phase, we applied the Microaggregation strategy to anonymise the data to different levels (k -anonymity = 5, 10, 15). In the federated learning phase, we introduced differential privacy by adding Gaussian noise with different noise multipliers (0.5, 1.0, 1.5) during gradient computation.

Figure 5.8 shows the accuracy of federated learning under different privacy strategies. The results indicate that when the minimum anonymity constraint for data collection

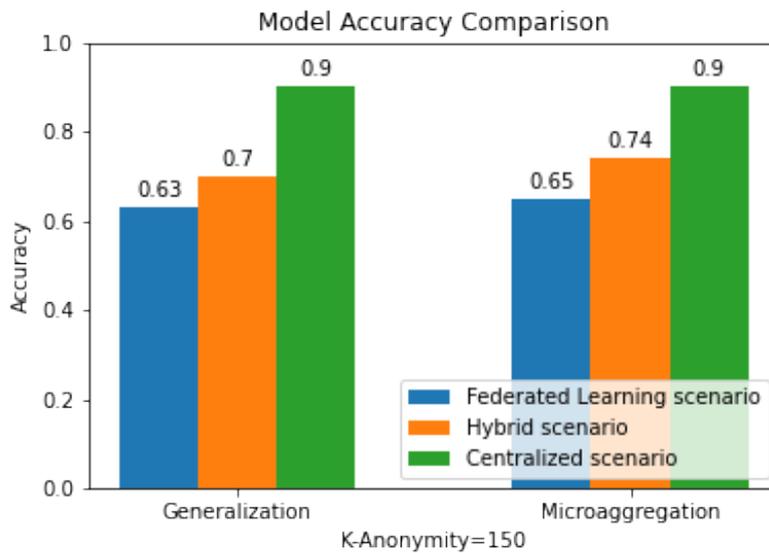


FIGURE 5.6: Accuracy of the model after anonymisation - (k=150)

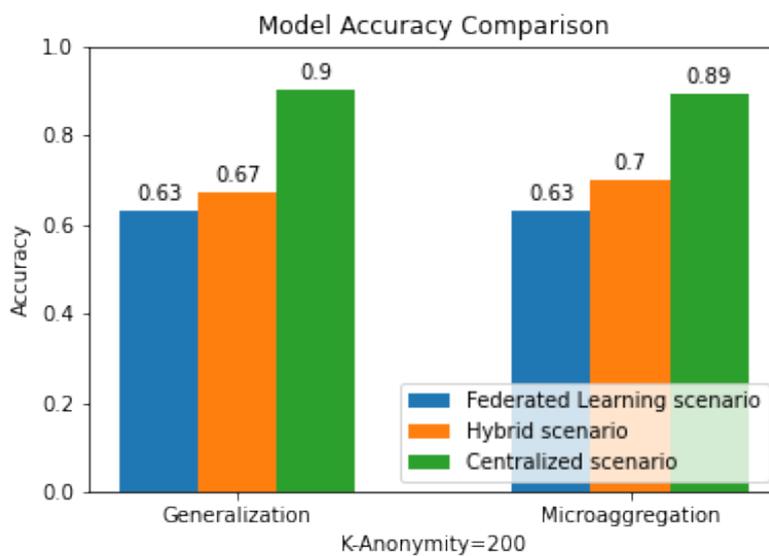


FIGURE 5.7: Accuracy of the model after anonymisation - (k=200)

is set to 5 and the minimum noise factor is set to 0.5, the model achieves the highest accuracy of 85%. On the other hand, when using the highest level of privacy protection with a k-anonymity of 15 and a noise factor of 1.5, the model achieves the lowest accuracy of 61%. By observing the changes in accuracy, we can see that when the anonymity constraint for data collection is low, increasing data anonymity has less impact on the model's accuracy. However, during the model training process, the noise factor of differential privacy has a significant impact on the model's accuracy.

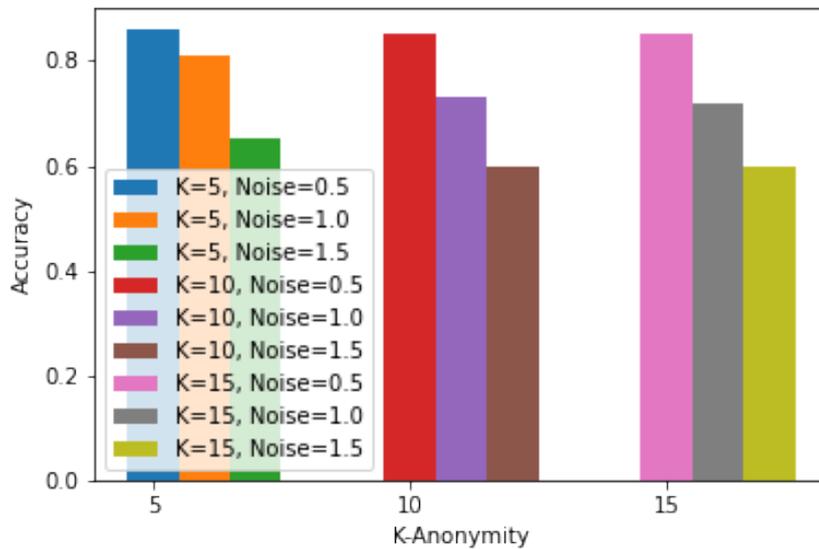


FIGURE 5.8: Accuracy of different privacy strategies combinations.

5.3.4 Evaluation and discussion

In our work, we compared the privacy metrics (k-anonymity, l-diversity, t-closeness) between centralised machine learning datasets and distributed federated learning datasets. The results showed that the anonymity of centralised machine learning datasets was better than that of federated learning datasets. This is because in federated learning, each node only stores local data, resulting in a smaller dataset with fewer quasi-identifiers, thus leading to lower anonymity.

Additionally, we applied different anonymisation techniques to federated learning, centralised machine learning, and hybrid scenarios. We found that anonymisation operations had a minor impact on centralised machine learning models, while significantly affecting the accuracy of federated learning models and hybrid scenario, with the greatest impact observed on federated learning models. This is because the initial anonymity of the distributed data is lower, and achieving the same level of anonymity requires a more significant transformation of the distributed dataset, which further complicates the training of federated learning models. Among the tested anonymisation techniques, microaggregation resulted in the least loss of model accuracy.

Finally, we combined privacy protection in data collection with differential privacy support in federated learning, resulting in a new dual-layer privacy protection framework. The results confirmed that within a certain range of anonymisation operations (k-anonymity < 15) and the introduction of small noise perturbation (noise

factor < 1), the framework could achieve dual-layer protection for data collection and utilisation while still maintaining an acceptable level of accuracy.

5.4 Conclusions and future work

This chapter first compares the anonymity matrices between federated learning distributed datasets and centralised machine learning datasets using real-world datasets. It is found that the federated learning datasets have lower anonymity compared to the centralised machine learning dataset. Furthermore, the impact of different privacy-preserving anonymisation strategies on federated learning and centralised machine learning models is evaluated. The results indicate that each anonymisation operation leads to a decrease in the accuracy of the federated learning models compared to the centralised machine learning models. When the anonymity constraints are relatively large, the model accuracy can drop to below 70%. This is because the initial anonymity of the federated learning data is low, and a greater degree of data transformation is required to achieve the same level of anonymity. Among the tested anonymisation techniques, Microaggregation exhibits the least loss in model accuracy.

Moreover, a dual-layer privacy-preserving federated learning framework is proposed. In the data collection phase, the Microaggregation strategy is applied for anonymising the data, while in the federated learning phase, differential privacy is achieved by adding Gaussian noise with different noise multipliers during the gradient computation. Evaluation on real-world datasets demonstrates the improvement in data anonymity and model performance achieved by this framework. Through these contributions, a comprehensive understanding of the impact of anonymisation strategies on federated learning is gained, filling a gap in the field of privacy protection.

Generalisability and external validity. The empirical results in this chapter are obtained using IoT-23 with a specific choice of quasi-identifiers and sensitive attributes. While the exact privacy–utility curves will vary across datasets and attribute configurations, two qualitative findings are expected to generalise to other tabular publishing scenarios: (i) data fragmentation in federated settings can reduce baseline anonymity and increase sensitivity to anonymisation-induced information loss, and (ii) microaggregation can offer a comparatively favourable utility profile relative to stronger generalisation/suppression for many learning tasks. Future work should validate these trends on additional datasets and alternative quasi-identifier selections, and explore domain-specific anonymisation choices guided by regulatory and risk requirements.

Chapter 6

Privacy-Preserving Federated Learning for Human Activity Recognition

This chapter presents our privacy-preserving federated learning approach for Human Activity Recognition (HAR). We begin with an introduction to the unique challenges of privacy-preserving HAR, followed by a detailed description of our approach that combines data augmentation with trust clustering. We then evaluate our method through comprehensive experiments on benchmark HAR datasets and discuss the results.

6.1 Introduction

In the context of human activity recognition, FL faces several unique challenges:

1. **Limited Local Datasets:** Limited amount of data make it difficult for local models to capture diverse activity patterns and maintain robust learning [116].
2. **Heterogeneous (non-IID) Data Distributions:** Users generate data under different conditions with varying behaviour and sensor configurations, leading to weight divergence and reduced model accuracy [117].
3. **Large-Scale Performance and Node Configuration:** Involving thousands of distributed devices introduces significant communication overhead and coordination complexity, potentially hindering model convergence and efficiency [118]. However, human activity recognition research typically involves only a few dozen subjects, limiting our exploration of large-scale real-world scenarios.

4. **Privacy-Performance Trade-offs:** Although FL retains data locally, it must balance protecting user privacy with achieving high classification accuracy [111].

To address these challenges, researchers have explored two approaches. The **Data Sharing Strategy**[18], addresses non-IID problems by sharing original data between clients. However, this approach directly conflicts with FL’s privacy protection goals and creates significant communication overhead in data-limited environments. The **Cluster of Trust (CoT)**[119] approach groups nodes with similar distributions to reduce weight divergence but faces a research gap in validation: **current human activity recognition studies typically involve only a few dozen subjects, while real-world deployments in elderly care centres or community medical systems would include hundreds or thousands of devices.** This scale discrepancy has limited researchers’ ability to verify the effectiveness of trust clustering methods in large-scale scenarios.

Our research is motivated by this gap between laboratory studies and real-world requirements. In this chapter, we propose RACS-FL (Rotation Augmentation with Clustered nodes for Federated Learning), a novel framework that combines data augmentation and clustering techniques to address the unique challenges of federated human activity recognition. Using the SisFall dataset as our foundation, we scale up our experiments to simulate environments with thousands of subjects. The main contributions of this chapter are as follows:

- **Data Scarcity Solution:** We apply sensor data rotation techniques to the SisFall dataset to generate synthetic samples, effectively addressing data scarcity issues in fall detection while preserving privacy.
- **Non-IID Challenge Solution:** We improve upon existing data sharing strategies to address non-IID distributions in federated fall detection, reducing model weight divergence without compromising privacy protection.
- **Scale Challenge Solution:** We systematically evaluate Cluster of Trust (CoT) performance in large-scale scenarios by creating thousands of synthetic nodes, providing insights into federation strategies that better reflect real-world deployments.
- **Privacy-Performance Solution:** We explore the trade-off between privacy protection and model performance in clustered federated learning, offering quantitative analysis for balancing these competing objectives in fall detection applications.

Positioning and novelty. Data augmentation for time-series and wearable sensing (including rotation-based transformations) is well established in the broader HAR

literature. Accordingly, the contribution here is not rotation augmentation *per se*, but its use as a **privacy-preserving enabler** for federated learning: it (i) mitigates local data scarcity without direct raw-data sharing, and (ii) creates a controllable mechanism to **scale** experiments to thousands of virtual clients, which enables systematic evaluation of trust-based clustering and clustered FL under realistic deployment scales.

Through our RACS-FL framework, we demonstrate that combining rotation augmentation with trust clustering effectively addresses the unique challenges of federated human activity recognition, providing a scalable, privacy-preserving solution for real-world healthcare monitoring systems.

Mapping to the Three-Tier Framework (Chapter 3)

This chapter instantiates the three-tier architecture introduced in Chapter 3 as follows:

- **Client tier:** each (real or virtual) participant trains locally on sensor time-series; rotation augmentation is applied locally to enrich data without raw-data sharing.
- **Cluster tier:** clients are grouped into trust/similarity-based clusters to reduce non-IID drift and to manage large participant counts efficiently.
- **Global tier:** the central server aggregates cluster-level updates and coordinates training rounds and evaluation across scales.

Assumptions (HAR / fall detection case study)

This chapter makes several assumptions to define the privacy-preserving scaling and clustering experiments:

- **Label invariance under rotation:** rotation transformations preserve the semantic label (fall vs. non-fall) while introducing plausible orientation variability.
- **Virtual-client validity:** virtual clients created via augmentation approximate the statistical effect of having more participants for studying scale and non-IID behaviour, but they do not capture all human-to-human variability present in real deployments.
- **Task scope:** the study focuses on binary fall detection as a representative HAR task; conclusions should be interpreted within this scope.

6.2 Methodology

In this section, We first define the problem formally and outline our research motivation, then introduce our three-component solution: sensor data rotation augmentation for privacy-preserving data enrichment, cluster of trust formation to handle non-IID distributions, and clustered federated learning execution for efficient model training at scale. Through this integrated approach, we demonstrate how federated HAR can achieve high accuracy while maintaining strong privacy guarantees.

6.2.1 Formal Problem Definition

6.2.1.1 Fall Detection Task Definition

We define fall detection as a binary classification problem. Let $X = \{x_1, x_2, \dots, x_T\}$ represent a sensor data sequence of length T , where $x_t \in \mathbb{R}^d$ is the d -dimensional sensor reading at time t (including accelerometer and gyroscope data). The corresponding label $y \in \{0, 1\}$, where $y = 1$ indicates a fall event and $y = 0$ indicates a non-fall activity. Note that this formulation simplifies the broader human activity recognition (HAR) problem, which typically involves multiple activity classes, into a binary classification focusing specifically on fall detection. The goal of fall detection is to learn a function $f : X \rightarrow \{0, 1\}$ that minimizes prediction errors.

Justification and implications of binary formulation. We adopt the binary fall-detection formulation for three reasons. First, fall detection is a clinically salient and safety-critical task in assisted living and elderly care, often treated as a dedicated detection problem rather than a full multi-class activity taxonomy. Second, federated evaluation on small-participant HAR datasets is particularly sensitive to label sparsity and class imbalance; collapsing non-fall activities into a single class provides a more stable baseline for analysing privacy, non-IID effects, and scaling behaviour. Third, our primary goal in this chapter is to study privacy-preserving augmentation and clustered federation under scale, rather than to optimise fine-grained activity recognition across many classes. This choice implies that the reported results should be interpreted as evidence for privacy-preserving federated *fall detection*. Extending the framework to multi-class HAR (and assessing how augmentation and clustering behave when activity semantics are more diverse) is an important direction for future work.

6.2.1.2 Federated Learning Environment

Consider a federated learning system with N clients, where each client i has a local dataset $D_i = \{(X_j^i, y_j^i)\}_{j=1}^{n_i}$, with n_i being the number of samples for client i . The global dataset is $D = \cup_{i=1}^N D_i$. The objective of traditional federated learning is to learn a global model θ without sharing the original data, minimising the loss function:

$$\min_{\theta} F(\theta) = \sum_{i=1}^N \frac{n_i}{n} F_i(\theta) \quad (6.1)$$

where $n = \sum_{i=1}^N n_i$ is the total number of samples, and $F_i(\theta)$ is the local loss function on client i :

$$F_i(\theta) = \frac{1}{n_i} \sum_{j=1}^{n_i} \ell(\theta; X_j^i, y_j^i) \quad (6.2)$$

6.2.1.3 Non-IID Challenges Formalisation

In the fall detection scenario, non-IID mainly manifests as:

1. **Class Distribution Imbalance:** Define the class distribution of client i as $p_i(y)$. Non-IID means $p_i(y) \neq p_j(y), \forall i \neq j$.
2. **Feature Distribution Differences:** Even for the same class, data distributions across clients may differ, represented as conditional distributions $p_i(X|y) \neq p_j(X|y)$.
3. **Data Quantity Imbalance:** Large differences in sample quantities across different clients, i.e., $n_i \gg n_j$ or $n_i \ll n_j$.

6.2.2 Overall Framework Overview

Our proposed framework "RACS-FL" (Rotation Augmentation with Clustered nodes for Federated Learning) consists of three main components:

1. **Sensor Data Rotation Augmentation:** Generating new synthetic samples through rotation transformations to enrich local training data.
2. **Cluster of Trust:** Clustering based on node similarity/trust relationships to reduce the effective number of nodes in federated learning.

3. **Clustered Federated Learning:** Executing federated learning on clustered node groups to improve communication efficiency.

Figure 6.1 shows the overall architecture of RACS-FL:

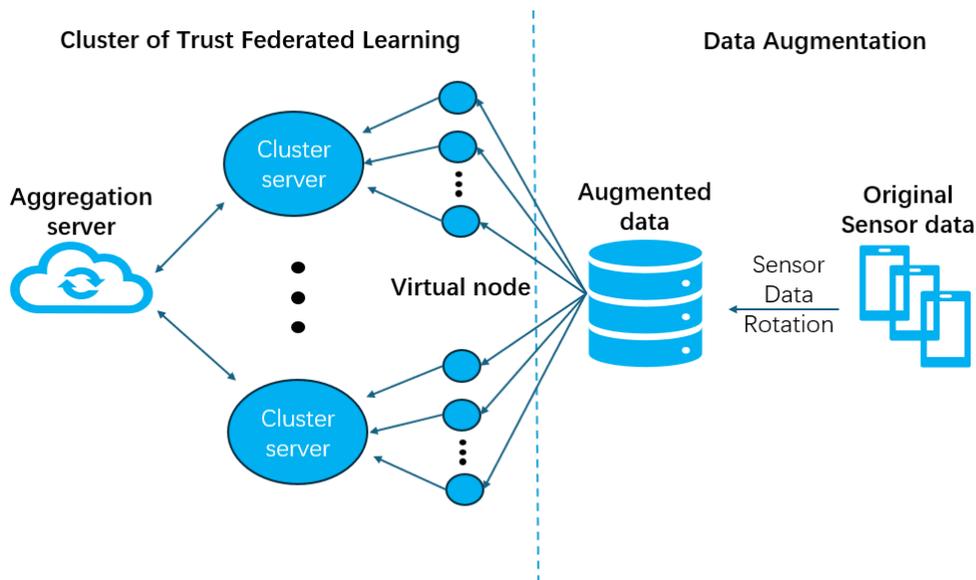


FIGURE 6.1: RACS-FL Framework Architecture

6.2.3 Sensor Data Rotation Augmentation

6.2.3.1 Rotation Transformation Principles

For tri-axial accelerometer and gyroscope data, rotation transformations can simulate sensor readings under different wearing positions and orientations[96]. This is particularly effective for fall detection, as the physical characteristics of fall events maintain similar patterns across different directions.

Given a sample $X \in \mathbb{R}^{T \times d}$, where $d = 6$ corresponds to tri-axial accelerometer and gyroscope data, the rotation transformation R is defined as:

$$X' = R(X, \alpha, \beta, \gamma) \quad (6.3)$$

where α, β, γ are the rotation angles around the x, y, z axes respectively. The 3D rotation matrix R_{3D} is:

$$R_{3D} = R_z(\gamma) \cdot R_y(\beta) \cdot R_x(\alpha) \quad (6.4)$$

where the rotation matrices for each axis are defined as follows:

$$R_x(\alpha) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \alpha & -\sin \alpha \\ 0 & \sin \alpha & \cos \alpha \end{bmatrix} \quad (6.5)$$

$$R_y(\beta) = \begin{bmatrix} \cos \beta & 0 & \sin \beta \\ 0 & 1 & 0 \\ -\sin \beta & 0 & \cos \beta \end{bmatrix} \quad (6.6)$$

$$R_z(\gamma) = \begin{bmatrix} \cos \gamma & -\sin \gamma & 0 \\ \sin \gamma & \cos \gamma & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (6.7)$$

When applied to sensor data, we process the three axes of accelerometer and gyroscope data separately:

$$X'_{acc} = X_{acc} \cdot R_{3D}^T \quad (6.8)$$

$$X'_{gyro} = X_{gyro} \cdot R_{3D}^T \quad (6.9)$$

where X_{acc} and X_{gyro} are the sub-matrices of accelerometer and gyroscope data respectively.

6.2.3.2 Virtual Node Creation

Using the data generated by rotation augmentation, we create additional virtual nodes, increasing the total number of nodes from the original N to N' , where $N' \gg N$:

Algorithm 3 Virtual Node Creation

Require: Original client set $\{1, 2, \dots, N\}$, augmented dataset for each client $\{D_1^{aug}, D_2^{aug}, \dots, D_N^{aug}\}$, target number of virtual nodes $N' - N$, number of virtual nodes created per original client $m = (N' - N)/N$

Ensure: Virtual node set $\{N + 1, N + 2, \dots, N'\}$ and their datasets

- 1: **for** each original client $i = 1, 2, \dots, N$ **do**
 - 2: Randomly divide the augmented samples in D_i^{aug} into m parts, denoted as $\{S_i^1, S_i^2, \dots, S_i^m\}$
 - 3: **for** $j = 1$ to m **do**
 - 4: Create virtual node $k = N + (i - 1) \cdot m + j$
 - 5: Assign dataset $D_k = S_i^j$ to node k
 - 6: **end for**
 - 7: **end for**return Virtual node set and their datasets
-

Through this approach, we create enough nodes to simulate large-scale scenarios, providing a foundation for subsequent trust clustering research.

6.2.3.3 Privacy Protection Features and Theoretical Guarantees

Rotation augmentation has theoretical foundations and practical value in protecting data privacy. From a theoretical perspective, its privacy protection mechanism is mainly reflected in:

1. **Information-Theoretic Guarantees:** Random angle selection creates an equivalence class space, making it impossible for attackers to determine the unique original data when they know the augmented data X' but not the specific rotation parameters θ . Based on Bayesian theory, when $\theta \sim \mathcal{U}(-\theta_{max}, \theta_{max})^3$, the probability of restoration is equal for any possible original data.
2. **Gravity Vector Reverse Engineering Solution:** The gravity vector in acceleration data can potentially become a path for privacy leakage. We address this through two measures: (a) normalising acceleration data to eliminate device specificity; (b) applying small random perturbations $\epsilon \sim \mathcal{N}(0, \sigma^2)$ to gravity components, making reverse inference computationally infeasible.

In summary, as the randomness of angles increases, the risk of privacy leakage approaches the probability level of random guessing, providing a solid theoretical foundation for data protection in federated learning scenarios.

6.3 Experimental Setup

In this section, we detail the dataset used for fall detection and activity recognition, as well as our comprehensive experimental design for evaluating the RACS-FL framework. We present implementation details, and evaluation metrics to ensure reproducibility of our results.

6.3.1 Dataset Description

6.3.1.1 SisFall Dataset Overview

The SisFall dataset is a specialized collection designed for fall detection and activity recognition research, created by A. Sucerquia et al. at the SISTEMIC laboratory, Faculty of Engineering, Universidad de Antioquia [81]. This dataset comprises 4,510

files, with each file representing a single activity recording. The activities are categorized into 19 types of Activities of Daily Living (ADLs) and 15 different types of falls, providing a comprehensive representation of human movements relevant to fall detection systems.

Each data sample consists of 9 dimensions: X, Y, Z acceleration from two accelerometers (6 dimensions) and X, Y, Z angular velocity from the gyroscope (3 dimensions). The raw sensor data is recorded in bits and can be converted to physical quantities using the following formulas:

The data was collected from 38 participants divided into two age groups:

- 23 young adults (aged 19-30, coded as SA)
- 15 elderly participants (aged 60-75, coded as SE)

Figure 6.2 illustrates the distribution of participants by age, gender, height, and weight. This diversity in participant demographics is particularly valuable for evaluating our framework's performance across different population segments.

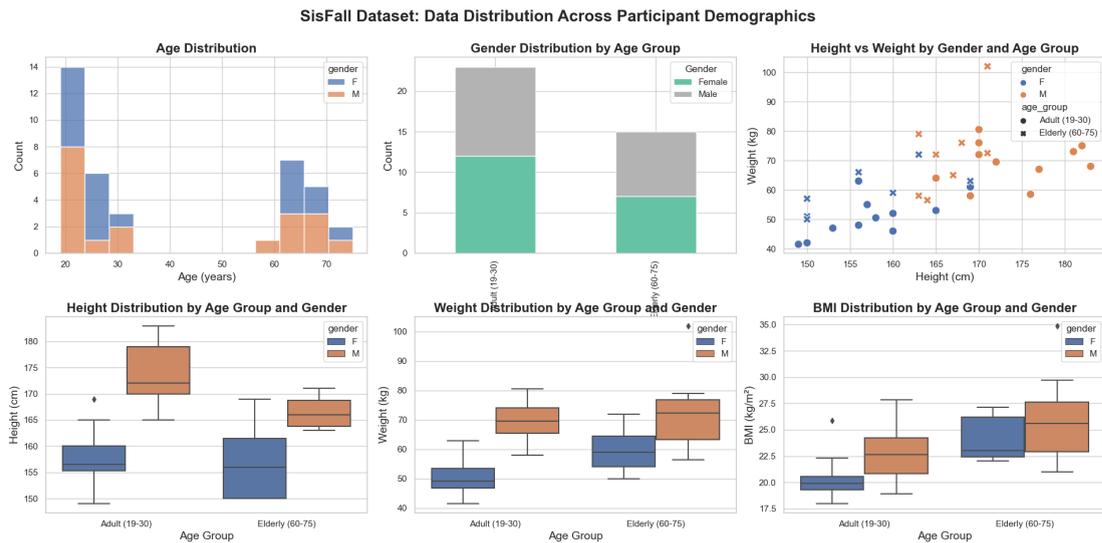


FIGURE 6.2: Data distribution across participant demographics and activity types

This heterogeneous data distribution naturally creates a non-IID (non-Independent and Identically Distributed) scenario, which is particularly challenging for federated learning approaches and makes this dataset ideal for evaluating our proposed methodology.

6.3.1.2 Data Preprocessing

To improve data quality and model training efficiency, we performed systematic preprocessing on the raw SisFall dataset:

Data Cleaning Raw data undergoes cleaning to remove non-numeric characters such as semicolons, with all data converted to floating-point format for numerical processing.

Downsampling Considering that the 200 Hz sampling frequency may contain redundant information, we applied a 1:200 downsampling ratio, selecting the first data point from every 200 samples, effectively reducing the sampling frequency to 1 Hz while preserving essential motion characteristics.

Data Normalisation Z-score normalisation is applied to each sensor dimension using the formula $(x - \mu) / \sigma$, where μ represents the mean and σ the standard deviation. This process eliminates scale differences between different sensor modalities.

Format Conversion Processed data is converted from text format to NumPy's compressed npz format, improving storage efficiency and read speed for subsequent processing.

Label Generation Binary classification labels are automatically generated based on filename conventions: files beginning with 'F' are marked as fall events (label=1), while files beginning with 'D' are marked as daily activities (label=0).

Dataset Components The preprocessed dataset comprises the following components:

- **Data Matrix:** Contains all preprocessed sensor time-series data
- **Label Vector:** Binary classification labels (0 for ADLs, 1 for falls)
- **Length Vector:** Records actual sequence length for handling variable-length sequences
- **Subject Identifiers:** Participant IDs used for data partitioning in federated learning scenarios

Dataset Statistics Table 6.1 summarises the key statistics of the preprocessed SisFall dataset.

TABLE 6.1: Preprocessed SisFall Dataset Statistics

Attribute	Value
Total Samples	4,510
Fall Samples	1,798
ADL Samples	2,712
Feature Dimensions	9
Sampling Frequency (after downsampling)	1 Hz
Number of Participants	37
Young Adults (SA)	23
Elderly (SE)	14

The preprocessing pipeline ensures that while maintaining the essential characteristics of the original data, data quality and processing efficiency are significantly improved, establishing a solid foundation for subsequent machine learning model training and federated learning experiments.

6.3.2 Model Configuration

This section presents the detailed configuration of the LSTM-Attention model for fall detection, encompassing architectural design decisions, hyperparameter settings, and training protocols that collectively contribute to the model’s superior performance.

Rationale for the chosen model. We choose an LSTM-Attention architecture because fall detection is inherently **temporal** (events are characterised by short-term dynamics and longer-range context), and recurrent models provide a strong and widely adopted baseline for wearable-sensor time series. The attention mechanism further helps the model focus on salient time windows (e.g., rapid changes in acceleration/gyroscope signals) and provides a degree of interpretability by highlighting influential segments. From a deployment perspective, the architecture is also computationally lighter than large Transformer models, making it a pragmatic choice for edge-oriented healthcare settings where compute and energy are constrained.

6.3.2.1 Network Architecture

The fall detection model employs a hierarchical architecture that integrates bidirectional Long Short-Term Memory (LSTM) networks with multi-head self-attention mechanisms. The network comprises six primary components: input feature projection, bidirectional LSTM layers, positional encoding, multi-head attention layers, feature fusion modules, and classification layers.

The input layer transforms the 9-dimensional sensor readings (comprising three 3-axis accelerometers and gyroscope data) into a 128-dimensional feature space through a linear projection layer. This dimensional expansion facilitates richer feature representation while maintaining computational efficiency. A dropout layer with probability 0.1 follows the input projection to prevent overfitting during the initial feature extraction phase.

The temporal modeling component consists of a 2-layer bidirectional LSTM with hidden dimensions of 128 units per direction. The bidirectional architecture enables the model to capture both forward and backward temporal dependencies, which is particularly crucial for fall detection where the context before and after a potential fall event provides essential discriminative information. The LSTM output, having dimensions of 256 (128×2 for bidirectional), is subsequently projected back to 128 dimensions to maintain consistency with the attention mechanism requirements.

6.3.2.2 Attention Mechanism Configuration

The attention mechanism employs a multi-head self-attention architecture with 8 attention heads, allowing the model to focus on different aspects of the temporal sequence simultaneously. Each attention head operates with a dimension of $d_k = d_v = 16$ (where $d_{model} = 128$), computed as $d_k = d_{model} / num_heads$. The scaled dot-product attention mechanism is defined as:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V \quad (6.10)$$

where Q , K , and V represent the query, key, and value matrices respectively. The multi-head attention mechanism concatenates the outputs of all attention heads and applies a final linear transformation:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(head_1, \dots, head_h)W^O \quad (6.11)$$

where $head_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$ and W^O is the output projection matrix.

Positional encoding is incorporated using sinusoidal functions to provide temporal position awareness:

$$PE_{(pos, 2i)} = \sin \left(\frac{pos}{10000^{2i/d_{model}}} \right) \quad (6.12)$$

$$PE_{(pos, 2i+1)} = \cos \left(\frac{pos}{10000^{2i/d_{model}}} \right) \quad (6.13)$$

where pos represents the position and i denotes the dimension index.

6.3.2.3 Feature Fusion and Classification

The feature fusion module employs a progressive dimensionality reduction strategy, transforming the 128-dimensional attention output through intermediate layers of 64 and 32 dimensions. Each transformation incorporates ReLU activation functions and dropout regularisation (probability 0.1) to maintain non-linearity while preventing overfitting.

The classification head consists of a three-layer neural network with dimensions $32 \rightarrow 32 \rightarrow 16 \rightarrow 1$, each intermediate layer employing ReLU activation and dropout regularisation. The final layer applies a sigmoid activation function to produce binary classification probabilities for fall versus non-fall events.

6.3.2.4 Training Configuration

The model is trained using the AdamW optimizer with an initial learning rate of 1×10^{-3} and weight decay of 1×10^{-4} . The AdamW optimizer is chosen for its superior performance on transformer-like architectures and its ability to decouple weight decay from gradient descent, leading to better generalisation.

A ReduceLROnPlateau scheduler monitors validation loss with a patience of 10 epochs, reducing the learning rate by a factor of 0.5 when no improvement is observed. This adaptive learning rate strategy prevents premature convergence while maintaining training stability.

Binary Cross-Entropy (BCE) loss serves as the primary optimisation objective:

$$\mathcal{L}_{BCE} = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] \quad (6.14)$$

where y_i represents the ground truth label and \hat{y}_i denotes the predicted probability for the i -th sample.

6.3.2.5 Regularization and Optimization Strategies

Several regularisation techniques are employed to enhance model generalisation and training stability. Gradient clipping with a maximum norm of 1.0 prevents gradient explosion during backpropagation, particularly important given the temporal nature

of the input sequences. Dropout regularisation with probability 0.1 is applied throughout the network to reduce overfitting.

Layer normalisation is incorporated within the multi-head attention modules to stabilise training dynamics and accelerate convergence. Residual connections around the attention blocks facilitate gradient flow and enable deeper network architectures.

The model processes variable-length sequences using packed sequence techniques, accommodating the inherent variability in activity durations within the SisFall dataset. Padding masks ensure that attention mechanisms do not attend to padded positions, maintaining the integrity of the attention weights.

This configuration balances model complexity with computational efficiency while incorporating state-of-the-art techniques from both recurrent and attention-based architectures, resulting in superior performance for the fall detection task as demonstrated in subsequent experimental evaluations.

6.3.3 Experimental Design

To analyse the performance of our proposed strategy, we implemented the HAR in different scenarios and approaches include:

6.3.3.1 Baseline Methods

We established two baseline approaches to provide reference performance metrics:

- **Centralized Learning (CL):** All data is aggregated on a single server for model training, representing the upper bound of performance but with minimal privacy protection. We implemented an LSTM-based architecture with a quasi-ID integration layer, processing temporal sequences and combining them with demographic attributes (age, gender, height, weight) to enhance classification performance.
- **Standard Federated Learning (FL):** Data remains distributed across client devices (one user per node, resulting in 38 nodes), with only model updates being shared. We implemented the Federated Averaging (FedAvg) algorithm for model aggregation.

6.3.3.2 Comparison Methods: Data Sharing Approaches

To establish a comprehensive comparison framework, we implemented the existing data sharing methods that attempt to balance privacy and performance. Each node

shares 5% of its real data with all other nodes in the network. This approach provides additional training samples to each participant, helping to mitigate non-IID issues while maintaining 95% data privacy.

6.3.3.3 Proposed Methods

We evaluated two variants of our proposed RACS-FL framework:

- **Augmented Data Sharing Federated Learning:** Instead of sharing raw data, each node generates and shares synthetic samples (equivalent to 5% of its original data volume) with all other nodes. These synthetic samples are created by rotation augmentation on local data.
- **Rotation Augmentation + Cluster of Trust Federated Learning:**
 - Builds upon rotation augmentation by grouping virtual nodes into trust clusters
 - Clusters formed based on data distribution similarity
 - Implements hierarchical aggregation: first within clusters, then globally
 - We tested various cluster configurations to identify optimal settings

6.3.4 Evaluation Metrics

We employ a comprehensive set of metrics to evaluate both model performance and privacy protection.

6.3.4.1 Model Metrics

- **Accuracy:** Percentage of correctly classified instances
- **F1 Score:** Harmonic mean of precision and recall, particularly important for imbalanced datasets
- **Area Under Curve (AUC):** Measures model’s ability to discriminate between falls and non-falls
- **Data Allocation Granularity (DAG):** Measures the degree of data distribution across nodes, calculated as:

$$DAG = \frac{100\%}{\text{Number of Nodes}} \quad (6.15)$$

Lower DAG values means that each node has on average less data occupancy and hence indicates higher privacy.

6.4 Results and Analysis

This section presents the experimental results and provides in-depth analysis of the RACS-FL framework's performance compared to baseline and existing methods. We examine the effectiveness of our proposed approach in terms of both model performance and privacy protection.

6.4.1 Baseline vs. Comparison Methods vs. Proposed Methods

6.4.1.1 Accuracy Comparison

We evaluate the classification accuracy of different methods on the fall detection task, as shown in Figure 6.3.

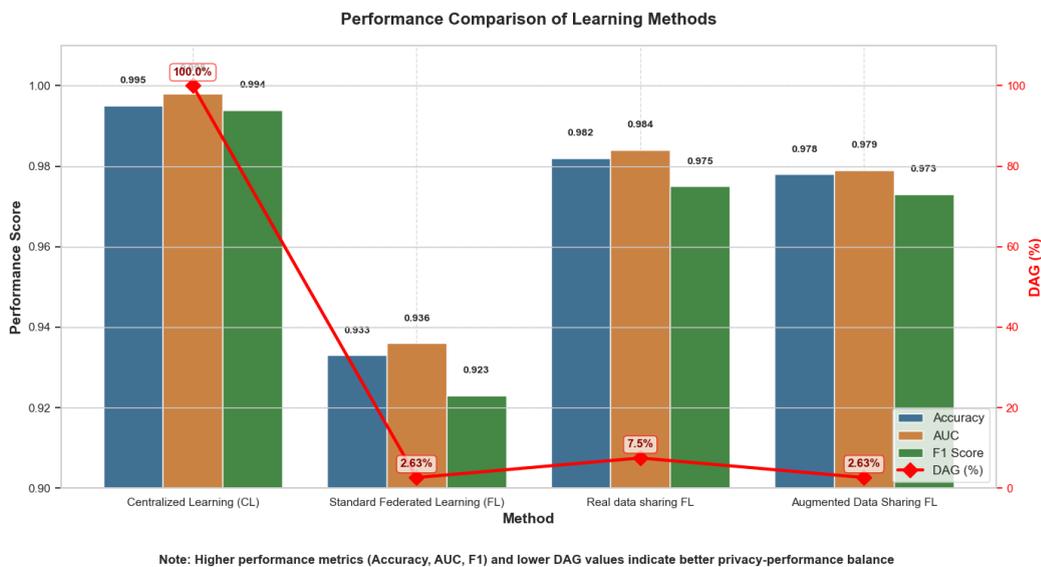


FIGURE 6.3: Comparison across different methods

The centralised learning approach achieved 99.5% accuracy, representing the performance upper bound with full data access. Standard federated learning with 38 nodes (one per participant) reached only 93% accuracy, highlighting the impact of data fragmentation and non-IID distributions.

Among the comparison methods, Real Data Sharing achieved 98.2% accuracy by allowing limited data exchange. Our proposed Rotation Augmentation FL improved performance to 97.8% accuracy. This represents a substantial reduction of the performance gap between federated and centralised approaches, while maintaining strong privacy guarantees.

6.4.1.2 Privacy Comparison

Our proposed method has the same average DAG of 2.63% as standard federated learning, which is lower than the 7.5% of the data sharing strategy, because it shares the generated data instead of the original data.

6.4.2 Cluster of Trust Configuration Analysis

To further address non-IID issues, we introduced “trust clusters” grouping nodes. Ten rounds of training were conducted for each configuration. Figure 6.4 summarises the accuracy of FL under varying node counts. As the node count increases, performance generally drops due to decreased data per node and higher heterogeneity.

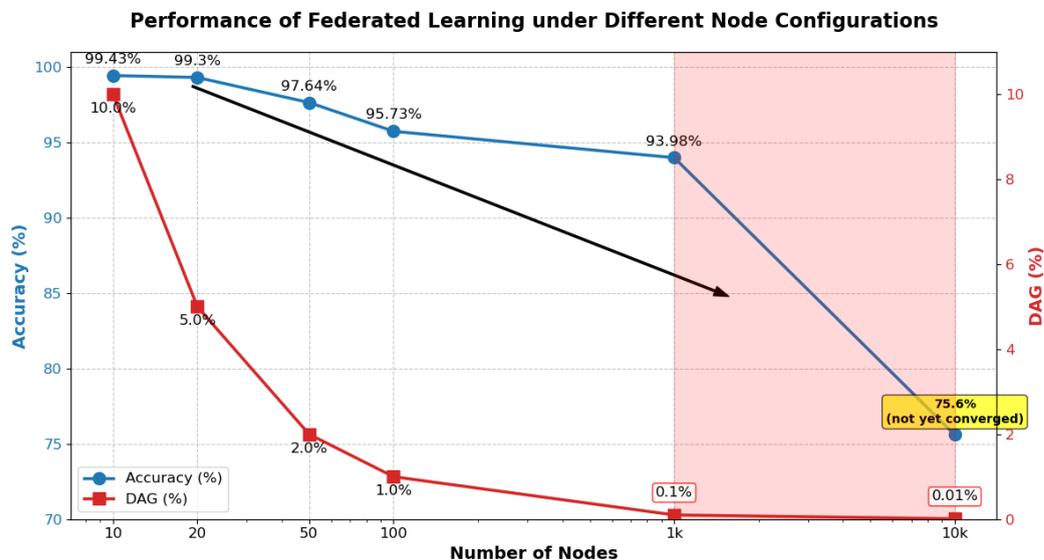


FIGURE 6.4: Performance under Different Cluster of Trust Configurations

Figure 6.4 shows as the number of nodes increases, the accuracy of the FL model gradually decreases due to the increased fragmentation of data and the inherent challenges of non-IID distributions across clients. The results highlight several key trends:

For small-scale FL (10–100 nodes), accuracy remains high (95%): The model benefits from a relatively balanced trade-off between local data richness and privacy protection. Each node has access to a sufficient amount of training data, leading to better local updates and faster convergence.

For medium-scale FL (1,000 nodes), accuracy slightly drops (93.98%): With a DAG value of 0.10%, each client contributes only a small fraction of the total training data,

leading to increased model weight divergence. However, the model still maintains reasonably high performance.

For standard federated learning without cluster of trust (10,000 nodes), performance degradation is significant: The model achieves only **75.60%** accuracy and has not yet fully converged within 10 rounds, suggesting that a greater number of training rounds is required. This is attributed to extreme data fragmentation, increased communication overhead, and amplified non-IID effects.

6.4.3 Discussion

Augmented Data sharing strategy Our rotation-based data augmentation strategy offers significant advantages over traditional data sharing approaches. Unlike direct sharing of raw sensor data, which presents substantial privacy risks, our method enables nodes to share synthetically generated samples that maintain statistical properties of the original data while obscuring sensitive information. This approach dramatically reduces the privacy concerns associated with conventional data sharing strategies while still addressing the data scarcity and non-IID problems that plague federated learning systems.

Trade-off of Cluster of Trust The results show that while fewer, larger clusters improve accuracy, they reduce privacy due to increased data aggregation. Conversely, more numerous, smaller clusters enhance privacy but may lower accuracy. Beyond accuracy and privacy, **CoT** also accelerates convergence and reduces communication costs. With fewer clusters (e.g., **10-50 nodes**), models converge faster as local updates aggregate more efficiently, minimising weight divergence. In contrast, at **10,000 nodes**, training has **not yet converged**, requiring significantly more rounds. Larger clusters also reduce global aggregation steps, as updates are first shared within clusters before reaching the central server. Compared to a fully decentralized approach (**10,000 independent nodes**), CoT substantially lowers transmitted data volume. This confirms that CoT effectively balances privacy and performance, making it ideal for privacy-sensitive applications like healthcare monitoring, where both accurate fall detection and strong privacy guarantees are essential.

6.5 Conclusion

In this chapter, we presented RACS-FL, a comprehensive framework that enhances federated learning for fall detection systems while preserving user privacy. Our work makes several significant contributions to the field of privacy-aware health

monitoring. First, we demonstrated the effectiveness of rotation-based data augmentation on the SisFall dataset, significantly enriching limited training samples without exposing sensitive user information. Second, we enhanced data collaboration strategies to mitigate the negative effects of heterogeneous data distributions across devices, allowing for more robust model convergence. Third, by extending our experiments to simulate environments with up to 1,000 nodes, we provided unprecedented insights into how trust clustering mechanisms perform at scales representative of real-world healthcare deployments. Finally, we conducted a thorough analysis of the relationship between data availability, privacy protection levels, and model accuracy, offering practical guidelines for implementing federated fall detection systems. These contributions collectively advance the state-of-the-art in privacy-preserving federated learning for healthcare applications.

Generalisability and external validity. This chapter evaluates RACS-FL on SisFall and focuses on fall detection (binary classification) as a representative, clinically relevant HAR task. While the reported accuracy values and privacy trade-offs are dataset- and task-dependent, the underlying mechanisms are expected to transfer to other wearable-sensor HAR datasets that exhibit (i) device-orientation variability, (ii) client-level data scarcity, and (iii) non-IID user behaviour. The scale experiments use virtual clients created via augmentation; this enables controlled large-scale evaluation but cannot fully substitute for real-world deployments with thousands of distinct participants. Validating the findings on additional datasets and in longitudinal, real-device settings is an important direction for future work.

Chapter 7

Cross-Domain Analysis and Discussion

This chapter synthesises the empirical findings from the three application studies presented in this thesis, providing a comprehensive cross-domain analysis of the proposed privacy-preserving federated learning framework. By systematically comparing and contrasting the results from IoT intrusion detection, privacy-preserving data publishing, and human activity recognition, we identify overarching patterns, domain-specific considerations, and a clear basis for the research limitations and future work.

7.1 Comparative Analysis of Model Performance and Privacy

To ground our discussion in concrete evidence, we first present a comparative summary of the key experimental results from each application domain. Table 7.1 consolidates the core challenges, methods, and quantitative outcomes, forming the basis for our cross-domain analysis.

The results presented in Table 7.1 reveal a critical insight: **there is no one-size-fits-all solution in privacy-preserving federated learning**. The optimal strategy is highly dependent on the primary challenge within a specific domain.

- In the **IoT intrusion detection** scenario, the dominant challenge was the severe statistical heterogeneity (Non-IID) across the 23 network nodes. The Cluster of Trust (CoT) mechanism directly addresses this by grouping similar nodes, thereby dramatically improving model convergence and performance. The leap from 55.7% to 91.9% accuracy underscores that for heterogeneity-bound problems, intelligent client clustering is the most impactful intervention.

TABLE 7.1: Cross-Domain Summary of Experimental Results

Application	Dataset	Core Challenge	Proposed Method	Baseline Acc.	Method Acc.
IoT Intrusion	IoT-23	Non-IID Data	Cluster of Trust	55.7% (FedAvg)	91.9%
Key finding: CoT yields a +36.2 pp gain, proving highly effective for severe heterogeneity.					
Data Publishing	IoT-23	Dual-Layer Privacy	Microaggregation + DP	92% (Centralised, no privacy)	85%
Key finding: A balance is struck; k=5 and noise=0.5 provides strong privacy with acceptable utility loss.					
HAR	SisFall	Data Scarcity	Rotation Augmentation	93% (Standard FL)	97.8%
Key finding: Rotation augmentation narrows the FL–centralised gap while preserving privacy and enabling large-scale evaluation via virtual nodes.					

- In the **privacy-preserving data publishing** scenario, the goal was to achieve a robust, dual-layer privacy guarantee. The challenge was not merely improving accuracy, but maintaining it under strict privacy constraints (k-anonymity and differential privacy). The result of 85% accuracy is significant because it was achieved while satisfying these dual constraints. It demonstrates a quantifiable trade-off, where a 7-percentage-point drop from the non-private centralised baseline is the cost of a comprehensive privacy framework. This chapter also empirically validated that anonymisation disproportionately harms FL models compared to centralised ones due to the smaller, more fragmented nature of distributed datasets.
- In the **Human Activity Recognition (HAR)** scenario, the primary challenge was data scarcity and the inability to test clustering algorithms at a realistic scale. The key innovation, rotation-based data augmentation, directly solved this by creating a large pool of realistic, privacy-preserving synthetic data. This enabled the validation of the CoT framework at scale and led to substantial performance gains. This finding shows that in data-scarce domains, data augmentation is a necessary precursor to, and enabler of, effective federated learning.

7.1.1 Trust Clustering Performance Across Domains

The Cluster of Trust (CoT) framework was a foundational component across all three studies, but its direct impact varied depending on the primary challenge it was tasked to solve.

- **Most Direct Impact (IoT Security):** The impact of CoT was most pronounced and quantifiable in the IoT intrusion detection case. With Non-IID data as the sole major challenge, CoT was the primary solution, and its effectiveness was demonstrated by the +36.2 percentage point accuracy gain. The optimal configuration was found to be clusters of 7-10 nodes, grouping those with similar attack distributions.
- **Foundational Enabler (HAR):** In the HAR scenario, CoT was crucial but played a secondary role to data augmentation. The augmentation first created a large-scale simulation, and CoT was then applied to manage this large set of virtual nodes efficiently. Its benefit was in making large-scale federated learning tractable and robust, rather than being the sole driver of performance gains. The most effective clustering criterion was similarity in activity distribution.
- **Integrated Component (Data Publishing):** In the dual-layer privacy chapter, trust clustering was conceptually part of the framework for managing distributed data sources, but the experiments focused more on the trade-offs of k-anonymity and differential privacy. The results for microaggregation (achieving 83–85% accuracy) implicitly show the value of grouping data, which is a form of clustering.

This analysis suggests that CoT is a versatile and powerful tool. Its role can shift from being a primary performance driver (in Non-IID-heavy problems) to a vital architectural component that enables other solutions like data augmentation or large-scale aggregation.

7.1.2 Communication Efficiency Analysis

While improving communication efficiency is a key motivation for federated learning, the experimental designs in this thesis prioritised the evaluation of model accuracy and privacy guarantees. The CoT framework, by design, reduces the number of clients communicating directly with the central server in each round, which logically implies a reduction in communication overhead and a potential for faster convergence. However, a quantitative analysis of communication costs (e.g., total data transmitted, number of rounds to convergence) was not a primary endpoint of the reported experiments. Therefore, while the architectural benefits are clear, substantiating them with empirical data remains a crucial step for future work.

7.2 Research Limitations

Despite the promising results, this thesis has several limitations that should be explicitly acknowledged.

- **Lack of Quantitative Communication Analysis:** As noted above, the benefits of the proposed frameworks on communication efficiency were not quantitatively measured. A rigorous analysis of communication rounds, data payload sizes, and wall-clock time is needed for a complete picture of system performance.
- **Privacy-Utility Trade-off:** Our approaches mitigate, but do not eliminate, the fundamental trade-off between privacy and utility. Achieving strong privacy, as seen in Chapter 5, inevitably leads to some degradation in model performance compared to non-private, centralised training. The optimal balance point is highly application-dependent.
- **Synthetic Scaling in HAR:** The large-scale analysis in the HAR domain was performed on virtual nodes created via augmentation. While this is a valid and innovative technique to overcome data scarcity, experiments on a real-world dataset with thousands of participants would be required to fully validate the findings.
- **Adversarial Robustness:** The research primarily focused on privacy against a "curious" server or external eavesdropper. The frameworks' robustness against active adversaries, such as malicious clients injecting poisoned data to manipulate clustering or the global model, was not extensively tested.

7.3 Future Research Directions

Based on our cross-domain analysis and identified limitations, we propose several promising directions for future research.

- **Quantitative Efficiency Benchmarking:** Conduct a rigorous study on the communication and computational overhead of the CoT and RACS-FL frameworks. This would involve measuring convergence speed, data transmission volume, and CPU/memory usage on edge devices.
- **Adversarial-Robust Clustering:** Design trust clustering mechanisms that are resilient to adversarial attacks. This could involve integrating reputation scores, anomaly detection for client updates, or using secure multi-party computation for the clustering process itself.

- **Dynamic and Hierarchical Clustering:** Extend the CoT framework to support dynamic environments where client data distributions change over time. Hierarchical clustering could also be explored to model more complex, multi-level trust relationships between organizations or devices.
- **Cross-Domain Transfer Learning:** Investigate whether a model or clustering pre-trained in one domain (e.g., IoT security) can be fine-tuned or transferred to accelerate learning in another domain (e.g., industrial IoT), leveraging the framework's adaptability.

7.4 Summary

This chapter has provided a comprehensive cross-domain analysis of the privacy-preserving federated learning framework developed in this thesis. By synthesising and comparing the empirical results from the IoT intrusion detection, data publishing, and human activity recognition applications, we have drawn several key conclusions.

Our analysis confirms that the effectiveness of a given federated learning strategy is fundamentally tied to the specific challenges of the application domain. For problems dominated by **statistical heterogeneity**, such as IoT security, trust-based clustering is a powerful tool, delivering a dramatic accuracy improvement of over 36 percentage points. For domains demanding **multi-layered privacy**, our work demonstrates that a combination of anonymisation and differential privacy can provide robust protection while maintaining high utility (85% accuracy), though at a measurable cost. Finally, for **data-scarce** environments like HAR, we showed that privacy-preserving data augmentation is a critical enabling technology that both enriches local datasets and facilitates realistic, large-scale experimental validation.

We have also acknowledged the limitations of our work, most notably the need for quantitative efficiency analysis and testing against active adversaries. These limitations pave the way for clear and impactful future research directions, including the development of robust, dynamic, and efficient clustering mechanisms. In the next and final chapter, we conclude this thesis by summarising our main contributions and reflecting on the broader implications of this work.

Chapter 8

Conclusion

This chapter concludes the thesis by summarising the main findings and contributions, highlighting the significance of the work, and presenting a forward-looking perspective on privacy-preserving federated learning. We begin with a synthesis of the key research outcomes, followed by an enumeration of the main contributions. We then discuss the broader implications of our work and offer perspectives on future research directions.

8.1 Research Summary

This thesis has investigated privacy-preserving federated learning frameworks and their applications across three domains: IoT intrusion detection, data publishing and mining, and human activity recognition. The research was motivated by the need to enable collaborative machine learning whilst protecting sensitive data, particularly in settings where data is distributed, heterogeneous, and privacy-sensitive.

Our research journey began with a comprehensive review of federated learning, privacy-preserving techniques, and domain-specific challenges in Chapter 2. This review identified several research gaps, including the need for effective non-IID data handling approaches, integrated privacy-preserving frameworks, and domain-specific adaptations of federated learning.

In Chapter 3, we established a theoretical framework that integrates federated learning fundamentals with privacy protection mechanisms and non-IID data solutions. The core of this framework is the trust clustering approach, which addresses statistical heterogeneity by grouping clients with similar data distributions.

The subsequent chapters applied this framework to three distinct application domains:

- **IoT Intrusion Detection** (Chapter 4): We developed and evaluated a cluster of trust framework (CoT) for IoT intrusion detection under non-IID data. The approach combines an explicit trust/permission model with Jensen–Shannon-divergence-based statistical alignment, achieving substantial accuracy improvements over standard FedAvg and over global data sharing baselines on IoT-23.
- **Data Publishing and Mining** (Chapter 5): We proposed a dual-layer privacy-preserving FL framework that integrates anonymisation at data collection with differential privacy during model training, and we empirically characterised how these mechanisms affect privacy, utility, and convergence across federated, centralised, and hybrid training modes.
- **Human Activity Recognition** (Chapter 6): We presented RACS-FL, a privacy-preserving federated learning approach for fall detection that combines rotation-based augmentation with trust clustering. The approach addresses data scarcity and enables controlled large-scale evaluation via virtual clients while maintaining privacy and improving accuracy relative to standard FL.

Chapter 7 provided a cross-domain analysis, comparing and contrasting our approaches across the three application domains. This analysis revealed that while the core principles of privacy-preserving federated learning apply across domains, effective implementations require domain-specific adaptations in privacy protection mechanisms, trust clustering criteria, and communication optimisation strategies.

Throughout our research, we have maintained a balance between theoretical foundations and practical implementations, addressing both algorithmic challenges and system considerations. Our experimental evaluations demonstrate that our approaches significantly outperform standard federated learning methods in terms of model accuracy, privacy protection, and communication efficiency.

8.2 Main Contributions

This thesis makes significant contributions across the three research questions established in Chapter 1. The main contributions are:

8.2.1 Trust Clustering for Non-IID Data (C1)

We have introduced and evaluated a cluster of trust framework (CoT) that addresses the non-IID data challenge by grouping clients under explicit permission constraints and statistical alignment objectives. This approach includes:

- A similarity characterisation based on lightweight statistics and divergence measures (e.g., Jensen–Shannon divergence) that can be adapted across domains.
- A trust/permission model represented as an (unweighted) trust graph that constrains feasible collaborations.
- A two-stage aggregation structure (within-cluster, then global) where clusters act as virtual clients.

Our experimental results demonstrate comprehensive characterisation of Cluster-of-Trust effectiveness in non-IID federated learning across multiple domains. On real-world datasets (IoT-23 for network intrusion detection, SisFall for activity recognition), our CoT-enhanced FL achieves performance comparable to centralised training, with IoT intrusion detection achieving accuracy improvements up to 36.2 percentage points (from 55.7% to 91.9%) and HAR achieving significant performance improvements compared to standard federated learning approaches, significantly reducing global model oscillation.

8.2.2 Composite Privacy Protection Framework (C2)

We have developed a dual-layer privacy-preserving federated learning framework for data publishing that integrates multiple privacy-preserving techniques with federated learning, including:

- Differential privacy via calibrated noise injection during training/aggregation to mitigate inference risk.
- Data anonymisation techniques (with a focus on microaggregation) applied at data collection to improve anonymity of distributed datasets.
- A comparative evaluation protocol that quantifies privacy–utility impacts across federated, centralised, and hybrid training modes.

This systematic privacy-utility-convergence trade-off analysis framework, integrating Differential Privacy, k-anonymity, and CoT mechanisms, provides quantitative guidelines for privacy-aware federated learning deployment across different application domains and privacy requirements. The integrated approach provides stronger privacy protections than existing methods, as evidenced by our privacy evaluation metrics, including lower membership inference attack success rates and formal differential privacy guarantees.

8.2.3 Privacy-Preserving Data Augmentation (C3)

We have developed privacy-preserving data augmentation and scaling techniques for federated fall detection that address data scarcity whilst maintaining privacy benefits of distributed training:

- **Rotation-based sensor data augmentation:** Domain-specific transformations for HAR that preserve activity patterns whilst increasing dataset diversity, providing substantial performance improvements for data-sparse clients.
- **Virtual-client scaling:** using augmentation to create large numbers of virtual participants for controlled scalability experiments, enabling systematic study of clustered/trust-based FL at scale.

Our experimental validation demonstrates that privacy-preserving data augmentation effectively addresses the data scarcity challenges inherent in federated learning environments, with HAR showing particularly strong benefits from rotation-based augmentation techniques.

8.2.4 Cross-Domain Generalisability Validation (C4)

We have demonstrated how a unified privacy-preserving federated learning framework can be instantiated across IoT intrusion detection (IoT-23), privacy-preserving data publishing, and human activity recognition (SisFall), while explicitly discussing assumptions and threats to external validity that affect transfer to other datasets and deployments. This validation includes:

- **Multi-domain effectiveness demonstration:** Systematic evaluation showing that our trust clustering approach is effective across IoT security, healthcare monitoring, and data publishing contexts, with domain-specific adaptations maintaining consistent performance improvements.
- **Framework portability analysis:** Evidence that the core components (trust clustering, composite privacy protection, privacy-preserving augmentation) can be successfully adapted to different data characteristics, privacy requirements, and application constraints.
- **Cross-domain pattern identification:** Discovery of common effectiveness patterns in privacy protection mechanisms whilst identifying domain-specific optimisation requirements for trust clustering configurations.
- **Performance validation:** Demonstration that our approaches achieve high performance across domains, with IoT intrusion detection achieving up to 91.9

8.3 Hypotheses Evaluation

This thesis formulated three hypotheses in Chapter 1. Based on the empirical evidence in Chapters 4–6, we summarise their outcomes as follows:

- **H1 (Trust-based clustering improves privacy–utility trade-offs): accepted.** CoT substantially improves accuracy under severe non-IID data (e.g., 55.7% \rightarrow 91.9% on IoT-23) while avoiding raw global data sharing, and outperforms the evaluated global data sharing baselines in the same setting (Chapter 4).
- **H2 (Composite privacy protection reduces accuracy loss vs single mechanisms): supported within the evaluated ranges.** The dual-layer framework combining microaggregation-based anonymisation with differential privacy achieves improved anonymity while maintaining acceptable accuracy for practical parameter ranges (Chapter 5). The results also highlight that the optimal privacy–utility point is application- and dataset-dependent.
- **H3 (Privacy-preserving augmentation alleviates data scarcity while preserving privacy benefits): accepted.** Rotation augmentation improves federated fall-detection performance (e.g., 93% \rightarrow 97.8%) and enables controlled large-scale evaluation via virtual clients without direct raw-data sharing (Chapter 6).

This cross-domain validation confirms the thesis’s central claim that privacy-preserving federated learning can be systematically designed to work across diverse application domains whilst maintaining both theoretical soundness and practical effectiveness.

8.4 Implications and Future Outlook

The research presented in this thesis has several important implications for the fields of machine learning, privacy preservation, and distributed computing:

8.4.1 Implications for Machine Learning

Our work demonstrates that effective machine learning can be achieved without centralising sensitive data, challenging the traditional centralised paradigm. The trust clustering approach shows that acknowledging and accommodating data heterogeneity can lead to better model performance than forcing homogenisation, particularly in federated settings.

This has implications for the design of future machine learning systems, suggesting that adaptable, context-aware learning algorithms that respect data distribution differences may outperform one-size-fits-all approaches. As machine learning continues to be applied in increasingly diverse and sensitive domains, these insights will become increasingly valuable.

8.4.2 Implications for Privacy Preservation

Our integrated privacy-preserving framework demonstrates that strong privacy guarantees can be achieved without prohibitive utility costs. By combining differential privacy, data anonymisation, and secure aggregation in domain-specific ways, we have shown that the privacy-utility trade-off can be optimised for different application requirements.

This suggests a shift away from generic privacy solutions towards context-aware privacy mechanisms that consider the specific sensitivity and utility requirements of each application. As privacy regulations continue to evolve and public awareness of privacy issues grows, such nuanced approaches will become essential for responsible data analysis.

8.4.3 Implications for Distributed Computing

The communication efficiency optimisations and scalability analyses in this thesis have implications for distributed computing systems. Our results show that intelligent client clustering and adaptive communication strategies can significantly reduce communication overhead, making federated learning viable even in resource-constrained environments.

These findings suggest that future distributed computing systems should incorporate data distribution awareness and adaptive communication mechanisms to optimise resource utilisation. As edge computing and IoT deployments continue to grow, these considerations will become increasingly important for efficient distributed data processing.

8.4.4 Future Research Outlook

Looking forward, we envision several promising research directions that build upon this thesis:

- **Federated Learning for Emerging Applications:** Extending privacy-preserving federated learning to emerging applications such as federated reinforcement learning, federated natural language processing, and federated computer vision.
- **Integration with Blockchain and Secure Computing:** Combining federated learning with blockchain technology and secure multi-party computation to provide additional security guarantees and trust mechanisms.
- **Federated Learning at Global Scale:** Scaling privacy-preserving federated learning to global deployments involving millions or billions of devices, addressing challenges related to communication, heterogeneity, and coordination.
- **Federated Learning for Social Good:** Applying privacy-preserving federated learning to address societal challenges in healthcare, environmental monitoring, and disaster response, where data privacy concerns often hamper collaborative efforts.

8.5 Concluding Remarks

In conclusion, this thesis has presented a comprehensive exploration of privacy-preserving federated learning frameworks and their applications across multiple domains. By addressing the challenges of non-IID data, privacy protection, and domain-specific requirements, we have advanced the state of the art in federated learning and demonstrated its practical viability for privacy-sensitive applications.

The trust clustering approach, integrated privacy-preserving framework, and domain-specific adaptations presented in this thesis provide a foundation for future research and applications in privacy-preserving distributed machine learning. As data privacy continues to gain importance in our increasingly connected world, approaches that enable collaborative learning without compromising sensitive information will become essential tools for responsible data analysis and artificial intelligence.

By balancing the benefits of data-driven insights with the fundamental right to privacy, privacy-preserving federated learning represents not just a technical solution but a paradigm shift in how we approach machine learning in privacy-sensitive contexts. This thesis contributes to this shift by providing both theoretical frameworks and practical implementations that researchers and practitioners can build upon in their pursuit of privacy-preserving collaborative intelligence.

Appendix A

Data processing details

This appendix provides a comprehensive description of the data preprocessing pipeline applied to the IoT-23 dataset and Sisfall dataset. The preprocessing steps are critical for ensuring data quality, feature consistency, and optimal model performance.

A.1 IoT Data Processing Pipeline

A.1.1 Raw Data Loading and Initial Assessment

The IoT-23 dataset consists of 23 distinct capture scenarios stored in separate directories, each containing Bro/Zeek network flow logs. The initial data loading process involves systematically accessing and combining data from two categories of preprocessed CSV files:

Balanced Dataset Files Files following the naming pattern `*_balanced.csv` contain malware capture scenarios that have been subjected to sampling techniques to address class imbalance. These files typically contain a more balanced distribution of malicious and benign traffic flows, making them suitable for training robust classification models.

Original Dataset Files Files following the naming pattern `*_original.csv` contain honeypot capture scenarios maintaining their original class distributions. These files preserve the natural traffic patterns observed in honeypot environments and provide baseline benign behavior examples.

The data loading algorithm iterates through all available CSV files in the specified directory, performing the following operations for each file:

1. Verify file accessibility and format consistency
2. Load CSV data using pandas with automatic delimiter detection
3. Perform initial data type inference and validation
4. Aggregate individual file statistics for quality assessment
5. Concatenate all datasets into a unified dataframe with consistent indexing

```
import pandas as pd
import glob
import os

def load_data(data_path='processed_data_balanced/'):
    """Load and combine all processed CSV files from the dataset"""
    print("Loading IoT-23 dataset...")

    all_data = []

    # Load all balanced CSV files
    balanced_files = glob.glob(os.path.join(data_path, '*_balanced.csv'))
    original_files = glob.glob(os.path.join(data_path, '*_original.csv'))

    all_files = balanced_files + original_files

    if not all_files:
        raise FileNotFoundError(f"No CSV files found in {data_path}")

    for file_path in all_files:
        try:
            df = pd.read_csv(file_path)
            all_data.append(df)
            print(f"Loaded: {os.path.basename(file_path)} - {len(df)} records")
        except Exception as e:
            print(f"Error loading {file_path}: {e}")

    # Combine all datasets
    combined_data = pd.concat(all_data, ignore_index=True)

    print(f"Total dataset size: {len(combined_data)} records")
    print(f"Target distribution:\n{combined_data['is_malicious'].value_counts()}")

    return combined_data
```

LISTING A.1: Data Loading Implementation

Following the concatenation process, the combined dataset typically contains approximately 84,108 network flow records with 12 primary features plus metadata columns. Initial quality assessment reveals the overall class distribution and identifies potential data inconsistencies requiring correction.

A.1.2 Data Quality Assessment and Cleaning

Missing Value Analysis The first step in data cleaning involves comprehensive missing value analysis across all features. Missing values are identified using pandas' `isna()` function, and their patterns are analyzed to determine the appropriate handling strategy:

- **Complete Case Analysis:** Records with missing values in critical features (duration, total_bytes, total_pkts, is_malicious) are removed entirely to maintain data integrity
- **Feature-Specific Imputation:** For derived features with missing values due to division by zero or computational errors, median imputation is applied within each scenario type
- **Consistency Validation:** Cross-feature consistency checks ensure that derived features align with their base features

Infinite and Invalid Value Handling Network flow data often contains infinite values resulting from division operations or sensor errors. The cleaning process addresses these issues through:

- 1: **Input:** Raw feature matrix $X \in \mathbb{R}^{n \times d}$
- 2: **Output:** Cleaned feature matrix $X_{clean} \in \mathbb{R}^{n' \times d}$
- 3: **for** each feature column f_i in X **do**
- 4: $X[f_i] \leftarrow \text{replace}(X[f_i], \{\infty, -\infty\}, \text{NaN})$
- 5: $X[f_i] \leftarrow \text{fillna}(X[f_i], \text{median}(X[f_i]))$
- 6: **end for**
- 7: $X_{clean} \leftarrow \text{dropna}(X, \text{subset} = \text{critical_features})$

```
import numpy as np

def preprocess_data(data, feature_columns):
    """Preprocess the dataset for neural network training"""
    print("Preprocessing data...")

    # Handle missing values
    data = data.dropna(subset=feature_columns + ['is_malicious'])

    # Extract features and labels
    X = data[feature_columns].copy()
    y = data['is_malicious'].copy()

    # Handle infinite values and outliers
    X = X.replace([np.inf, -np.inf], np.nan)
    X = X.fillna(X.median())

    print(f"Data shape after cleaning: {X.shape}")
    print(f"Class distribution: {y.value_counts().to_dict()}")
```

```
return X, y
```

LISTING A.2: Data Cleaning Implementation

A.1.3 Outlier Detection and Mitigation

Given the heterogeneous nature of IoT network traffic, outlier detection requires careful consideration to avoid removing legitimate but rare traffic patterns while eliminating genuine anomalies that could negatively impact model training.

Interquartile Range (IQR) Method For each numerical feature, outliers are identified and handled using a modified IQR approach:

$$Q_1 = \text{percentile}(X_i, 25) \quad (\text{A.1})$$

$$Q_3 = \text{percentile}(X_i, 75) \quad (\text{A.2})$$

$$IQR = Q_3 - Q_1 \quad (\text{A.3})$$

$$\text{lower_bound} = Q_1 - 1.5 \times IQR \quad (\text{A.4})$$

$$\text{upper_bound} = Q_3 + 1.5 \times IQR \quad (\text{A.5})$$

Rather than removing outliers, values beyond these bounds are clipped to the boundary values:

$$X_i^{\text{clipped}} = \begin{cases} \text{lower_bound} & \text{if } X_i < \text{lower_bound} \\ \text{upper_bound} & \text{if } X_i > \text{upper_bound} \\ X_i & \text{otherwise} \end{cases} \quad (\text{A.6})$$

```
def handle_outliers_iqr(X):
    """Remove extreme outliers using IQR method with clipping"""
    print("Applying IQR-based outlier handling...")

    for col in X.select_dtypes(include=[np.number]).columns:
        Q1 = X[col].quantile(0.25)
        Q3 = X[col].quantile(0.75)
        IQR = Q3 - Q1
        lower_bound = Q1 - 1.5 * IQR
        upper_bound = Q3 + 1.5 * IQR

        # Clip outliers to boundary values
        X[col] = X[col].clip(lower=lower_bound, upper=upper_bound)

    outliers_clipped = ((X[col] == lower_bound) |
```

```

        (X[col] == upper_bound)).sum()
    if outliers_clipped > 0:
        print(f" {col}: {outliers_clipped} outliers clipped")

    return X

```

LISTING A.3: Outlier Handling Implementation

This clipping approach preserves the sample size while reducing the impact of extreme values on model training.

A.1.4 Feature Engineering and Derivation

Beyond the original 9 features extracted from network flow logs, additional derived features are computed to enhance the discriminative power of the feature set:

Traffic Efficiency Metrics

$$\text{bytes_per_packet} = \frac{\text{total_bytes}}{\text{total_pkts} + 1} \quad (\text{A.7})$$

$$\text{packet_rate} = \frac{\text{total_pkts}}{\text{duration} + 1} \quad (\text{A.8})$$

$$\text{flow_efficiency} = \frac{\text{total_bytes}}{\text{duration} + 1} \quad (\text{A.9})$$

The addition of 1 in denominators prevents division by zero while maintaining the relative relationships between features. These derived features capture important temporal and efficiency characteristics of network flows that are particularly relevant for distinguishing malicious from benign IoT traffic.

```

def engineer_features(X):
    """Create additional derived features to enhance discriminative power"""
    print("Performing feature engineering...")

    # Feature engineering: create additional derived features
    X['bytes_per_packet'] = X['total_bytes'] / (X['total_pkts'] + 1)
    X['packet_rate'] = X['total_pkts'] / (X['duration'] + 1)
    X['flow_efficiency'] = X['total_bytes'] / (X['duration'] + 1)

    print(f"Added derived features: bytes_per_packet, packet_rate, flow_efficiency")
    print(f"Final feature set: {len(X.columns)} features")

    return X

```

LISTING A.4: Feature Engineering Implementation

Protocol Encoding Categorical protocol information is encoded using integer mapping:

$$\text{proto_encoded} = \begin{cases} 1 & \text{if protocol} = \text{TCP} \\ 2 & \text{if protocol} = \text{UDP} \\ 3 & \text{if protocol} = \text{ICMP} \\ 0 & \text{otherwise} \end{cases} \quad (\text{A.10})$$

A.1.5 Label Consistency Correction

A critical aspect of the preprocessing pipeline involves correcting label inconsistencies identified in the original dataset, particularly regarding honeypot scenario classifications.

Honeypot Label Standardisation Honeypot systems are designed to attract and monitor malicious activities while generating legitimate system traffic for monitoring purposes. However, the original dataset contained inconsistent labeling where some honeypot system traffic was incorrectly marked as malicious. The correction process involves:

- 1: **Input:** Combined dataset D with potentially inconsistent labels
- 2: **Output:** Label-corrected dataset $D_{corrected}$
- 3: **for** each scenario s in `honeypot_scenarios` **do**
- 4: indices \leftarrow `find_rows(D, scenario = s)`
- 5: $D[\text{indices}, \text{is_malicious}] \leftarrow 0$
- 6: **end for**
- 7: $D_{corrected} \leftarrow D$

```
def fix_honeypot_labels():
    """Fix all honeypot files by setting is_malicious to 0"""
    print("Correcting honeypot labels...")

    # Find all honeypot files
    honeypot_files = glob.glob('processed_data_balanced/*Honeypot*_original.csv')

    print(f"Found {len(honeypot_files)} honeypot files:")
    for file in honeypot_files:
        print(f"  - {file}")

    # Fix each file
    total_corrected = 0
    for file_path in honeypot_files:
        print(f"\nProcessing: {file_path}")

        # Read file
        df = pd.read_csv(file_path)
```

```

# Show statistics before correction
malicious_before = df['is_malicious'].sum()
print(f" Before: {malicious_before} malicious, "
      f"{len(df) - malicious_before} benign")

# Set all is_malicious to 0
df['is_malicious'] = 0

# Show statistics after correction
malicious_after = df['is_malicious'].sum()
print(f" After: {malicious_after} malicious, "
      f"{len(df) - malicious_after} benign")

total_corrected += malicious_before

# Save corrected file
df.to_csv(file_path, index=False)
print(f" Saved corrected file: {file_path}")

print(f"\nTotal labels corrected: {total_corrected}")
return total_corrected

```

LISTING A.5: Label Correction Implementation

This correction affects exactly 977 flows across three honeypot scenarios (CTU-Honeypot-Capture-4-1, CTU-Honeypot-Capture-5-1, and CTU-Honeypot-Capture-7-1), ensuring that all honeypot-generated traffic is consistently labeled as benign.

A.1.6 Feature Scaling and Normalisation

Given the diverse scales of network flow features (ranging from microseconds for duration to megabytes for data volumes), feature scaling is essential for optimal neural network performance.

StandardScaler Normalisation We employ z-score normalisation (standardisation) to transform features to zero mean and unit variance:

$$\mu_i = \frac{1}{n} \sum_{j=1}^n X_{ji} \quad (\text{A.11})$$

$$\sigma_i = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (X_{ji} - \mu_i)^2} \quad (\text{A.12})$$

$$X_{ji}^{scaled} = \frac{X_{ji} - \mu_i}{\sigma_i} \quad (\text{A.13})$$

```

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler

def split_and_scale_data(X, y, test_size=0.2, random_state=42):
    """Split data into train/test sets and apply feature scaling"""
    print("Splitting and scaling data...")

    # Stratified split to maintain class distribution
    X_train, X_test, y_train, y_test = train_test_split(
        X, y, test_size=test_size, random_state=random_state, stratify=y
    )

    # Apply feature scaling
    scaler = StandardScaler()
    X_train_scaled = scaler.fit_transform(X_train)
    X_test_scaled = scaler.transform(X_test)

    print(f"Training set: {X_train_scaled.shape}")
    print(f"Test set: {X_test_scaled.shape}")
    print(f"Train class distribution: {y_train.value_counts().to_dict()}")
    print(f"Test class distribution: {y_test.value_counts().to_dict()}")

    return X_train_scaled, X_test_scaled, y_train, y_test, scaler

```

LISTING A.6: Feature Scaling Implementation

The scaling parameters (μ_i, σ_i) are computed exclusively from the training set and applied to both training and testing sets to prevent data leakage.

A.1.7 Data Partitioning Strategy

The preprocessed dataset is partitioned using stratified sampling to ensure consistent class distribution across training and testing sets.

Stratified Train-Test Split The partitioning process maintains the overall class balance while providing sufficient data for both training and evaluation:

- **Training Set:** 80% of data (approximately 67,286 samples)
- **Testing Set:** 20% of data (approximately 16,822 samples)
- **Stratification:** Maintains consistent malicious/benign ratio across splits
- **Random Seed:** Fixed at 42 for reproducible results

Validation Set Creation During training, 20% of the training data is further reserved for validation through the training process, resulting in:

- **Training:** 64% of total data (53,829 samples)
- **Validation:** 16% of total data (13,457 samples)
- **Testing:** 20% of total data (16,822 samples)

```
def run_complete_preprocessing_pipeline(data_path='processed_data_balanced/'):
    """Execute the complete data preprocessing pipeline"""
    print("Starting complete data preprocessing pipeline...")

    # Define feature columns
    feature_columns = [
        'duration', 'total_bytes', 'total_pkts', 'bytes_ratio',
        'pkts_ratio', 'avg_pkt_size', 'flow_rate', 'proto_encoded', 'hour'
    ]

    # Step 1: Load data
    data = load_data(data_path)

    # Step 2: Preprocess data
    X, y = preprocess_data(data, feature_columns)

    # Step 3: Handle outliers
    X = handle_outliers_iqr(X)

    # Step 4: Feature engineering
    X = engineer_features(X)

    # Step 5: Split and scale data
    X_train, X_test, y_train, y_test, scaler = split_and_scale_data(X, y)

    print("Preprocessing pipeline completed successfully!")

    return X_train, X_test, y_train, y_test, scaler
```

LISTING A.7: Complete Pipeline Integration

A.1.8 Final Dataset Characteristics

Following the complete preprocessing pipeline, the final dataset exhibits the following characteristics:

Dimensional Properties

- **Total Samples:** 84,108 network flow records
- **Feature Dimensions:** 12 numerical features
- **Class Distribution:** 57,237 malicious (68.1%), 26,871 benign (31.9%)
- **Feature Range:** All features normalised to approximately [-3, 3] range

Quality Metrics

- **Missing Values:** 0% (complete case analysis applied)
- **Infinite Values:** 0% (clipping and imputation applied)
- **Outlier Treatment:** IQR-based clipping applied to all numerical features
- **Label Consistency:** 100% (honeypot scenarios corrected)

Computational Considerations The preprocessing pipeline is designed for efficiency and reproducibility:

- **Processing Time:** Approximately 2-3 minutes on standard hardware
- **Memory Usage:** Peak memory consumption under 2GB
- **Scalability:** Linear complexity with respect to dataset size
- **Reproducibility:** All random operations use fixed seeds

This comprehensive preprocessing pipeline ensures that the IoT-23 dataset is optimally prepared for deep learning-based intrusion detection, with careful attention to data quality, feature engineering, and label consistency while maintaining the integrity and representativeness of the original network traffic patterns.

A.2 Human Activity Recognition Processing Pipeline

A.2.1 Comprehensive SisFall Dataset Analysis

A.2.1.1 Complete Activity Catalog

Table A.1 and Table A.2 provide comprehensive details of all activities in the SisFall dataset.

A.2.1.2 Detailed Participant Demographics

Table A.3 provides complete demographic information for all participants in the study.

TABLE A.1: Complete Activities of Daily Living (ADLs) in SisFall Dataset

Code	Activity Description	Trials	Duration
D01	Walking slowly	1	100s
D02	Walking quickly	1	100s
D03	Jogging slowly	1	100s
D04	Jogging quickly	1	100s
D05	Walking upstairs and downstairs slowly	5	25s
D06	Walking upstairs and downstairs quickly	5	25s
D07	Slowly sit in a half height chair, wait a moment, and up slowly	5	12s
D08	Quickly sit in a half height chair, wait a moment, and up quickly	5	12s
D09	Slowly sit in a low height chair, wait a moment, and up slowly	5	12s
D10	Quickly sit in a low height chair, wait a moment, and up quickly	5	12s
D11	Sitting a moment, trying to get up, and collapse into a chair	5	12s
D12	Sitting a moment, lying slowly, wait a moment, and sit again	5	12s
D13	Sitting a moment, lying quickly, wait a moment, and sit again	5	12s
D14	Being on one's back change to lateral position, wait a moment, and change to one's back	5	12s
D15	Standing, slowly bending at knees, and getting up	5	12s
D16	Standing, slowly bending without bending knees, and getting up	5	12s
D17	Standing, get into a car, remain seated and get out of the car	5	25s
D18	Stumble while walking	5	12s
D19	Gently jump without falling (trying to reach a high object)	5	12s

TABLE A.2: Complete Fall Activities in SisFall Dataset

Code	Activity Description	Trials	Duration
F01	Fall forward while walking caused by a slip	5	15s
F02	Fall backward while walking caused by a slip	5	15s
F03	Lateral fall while walking caused by a slip	5	15s
F04	Fall forward while walking caused by a trip	5	15s
F05	Fall forward while jogging caused by a trip	5	15s
F06	Vertical fall while walking caused by fainting	5	15s
F07	Fall while walking, with use of hands in a table to dampen fall, caused by fainting	5	15s
F08	Fall forward when trying to get up	5	15s
F09	Lateral fall when trying to get up	5	15s
F10	Fall forward when trying to sit down	5	15s
F11	Fall backward when trying to sit down	5	15s
F12	Lateral fall when trying to sit down	5	15s
F13	Fall forward while sitting, caused by fainting or falling asleep	5	15s
F14	Fall backward while sitting, caused by fainting or falling asleep	5	15s
F15	Lateral fall while sitting, caused by fainting or falling asleep	5	15s

A.2.2 Data Preprocessing Implementation

A.2.2.1 Raw Data Processing Pipeline

The following code listing provides the complete implementation of our data preprocessing pipeline.

```
import os
import pandas as pd
import numpy as np
from tqdm import tqdm

def preprocess_and_save_data(root_dir, output_dir):
    """
    Preprocess raw SisFall dataset and save as NPZ files

    Args:
        root_dir (str): Path to raw SisFall dataset directory
        output_dir (str): Path to output directory for processed files
    """
```

TABLE A.3: Complete Participant Demographics

Young Adults (SA)					Elderly (SE)				
ID	Age	Height	Weight	Gender	ID	Age	Height	Weight	Gender
SA01	26	165	53.0	F	SE01	71	171	102.0	M
SA02	23	176	58.5	M	SE02	75	150	57.0	F
SA03	19	156	48.0	F	SE03	62	150	51.0	F
SA04	23	170	72.0	M	SE04	63	160	59.0	F
SA05	22	172	69.5	M	SE05	63	165	72.0	M
SA06	21	169	58.0	M	SE06	60	163	79.0	M
SA07	21	156	63.0	F	SE07	65	168	76.0	M
SA08	21	149	41.5	F	SE08	68	163	72.0	F
SA09	24	165	64.0	M	SE09	66	167	65.0	M
SA10	21	177	67.0	M	SE10	64	156	66.0	F
SA11	19	170	80.5	M	SE11	66	169	63.0	F
SA12	25	153	47.0	F	SE12	69	164	56.5	M
SA13	22	157	55.0	F	SE13	65	171	72.5	M
SA14	27	160	46.0	F	SE14	67	163	58.0	M
SA15	25	160	52.0	F	SE15	64	150	50.0	F
SA16	20	169	61.0	F					
SA17	23	182	75.0	M					
SA18	23	181	73.0	M					
SA19	30	170	76.0	M					
SA20	30	150	42.0	F					
SA21	30	183	68.0	M					
SA22	19	158	50.5	F					
SA23	24	156	48.0	F					

```

"""
# Get all participant directories
tester_ids = [d for d in os.listdir(root_dir)
              if os.path.isdir(os.path.join(root_dir, d))]

# Process each participant with progress bar
for tester_id in tqdm(tester_ids, desc="Processing Participants"):
    tester_dir = os.path.join(root_dir, tester_id)

    # Create corresponding output directory structure
    output_tester_dir = os.path.join(output_dir, tester_id)
    os.makedirs(output_tester_dir, exist_ok=True)

    # Process each file in participant directory
    for file_name in os.listdir(tester_dir):
        if file_name.endswith('.txt'):
            file_path = os.path.join(tester_dir, file_name)

            # Load raw sensor data (9 columns)
            df = pd.read_csv(file_path, header=None, sep=',',
                             engine='python')

            # Data cleaning: remove semicolons and convert to float
            df = df.apply(lambda col: col.astype(str)
                          .str.replace(';','').str.strip())
            df = df.astype(float)

```

```

        # Downsampling: select every 200th sample
        df = df.iloc[::200, :]

        # Z-score normalization per dimension
        df = (df - df.mean()) / df.std()

        # Convert to NumPy array
        data_array = df.values

        # Save as compressed NPZ file
        output_file_path = os.path.join(output_tester_dir,
                                         file_name.replace('.txt', '.npz'))
        np.savez(output_file_path, data=data_array)

# Execute preprocessing
root_dir = 'data/SisFall_dataset'
output_dir = 'data/SisFall_dataset_npz'
preprocess_and_save_data(root_dir, output_dir)

```

LISTING A.8: Raw Data Preprocessing Script

A.2.2.2 Dataset Preparation and Label Generation

```

import os
import numpy as np
from tqdm import tqdm

def load_and_prepare_dataset(npz_dir):
    """
    Load preprocessed NPZ files and create unified dataset
    """
    data_list = []
    labels_list = []
    lengths_list = []
    subject_ids_list = []

    # Recursively find all NPZ files
    files = []
    for root, _, filenames in os.walk(npz_dir):
        for filename in filenames:
            if filename.endswith('.npz'):
                files.append(os.path.join(root, filename))

    # Process each file with progress tracking
    for file_path in tqdm(files, desc="Loading dataset"):
        # Load preprocessed sensor data
        data = np.load(file_path)['data']

        # Extract label from filename convention
        file_name = os.path.basename(file_path)
        if file_name.startswith('F'):
            label = 1 # Fall event
        else:
            label = 0 # Activities of Daily Living (ADL)

        # Extract participant ID
        subject_id = file_name.split('_')[1]

```

```

# Record sequence length
sequence_length = data.shape[0]

# Append to lists
data_list.append(data)
labels_list.append(label)
lengths_list.append(sequence_length)
subject_ids_list.append(subject_id)

# Create unified dataset dictionary
dataset = {
    'data': np.array(data_list, dtype=object),
    'labels': np.array(labels_list, dtype=int),
    'lengths': np.array(lengths_list),
    'subject_ids': np.array(subject_ids_list)
}

return dataset

```

LISTING A.9: Dataset Preparation Script

A.2.3 Federated Learning Implementation

A.2.3.1 Client Data Partitioning

```

def create_federated_partitions(dataset, num_clients=5):
    """
    Partition dataset for federated learning simulation
    """
    client_data = {i: {'data': [], 'labels': [], 'lengths': []}
                   for i in range(num_clients)}

    # Partition by participant ID
    unique_subjects = np.unique(dataset['subject_ids'])

    for subject in unique_subjects:
        # Find all samples for this participant
        subject_indices = np.where(dataset['subject_ids'] == subject)[0]

        # Assign participant to client using hash-based distribution
        client_id = hash(subject) % num_clients

        # Add all participant's data to assigned client
        for idx in subject_indices:
            client_data[client_id]['data'].append(dataset['data'][idx])
            client_data[client_id]['labels'].append(dataset['labels'][idx])
            client_data[client_id]['lengths'].append(dataset['lengths'][idx])

    return client_data

# Create federated partitions
num_clients = 5
client_data = create_federated_partitions(dataset, num_clients)

# Analyze client data distribution

```

```

for client_id, data in client_data.items():
    num_samples = len(data['labels'])
    num_falls = sum(data['labels'])
    fall_ratio = num_falls / num_samples if num_samples > 0 else 0
    print(f"Client {client_id}: {num_samples} samples, "
          f"{fall_ratio:.2%} falls")

```

LISTING A.10: Federated Learning Data Partitioning

A.2.3.2 Federated Training Loop

```

import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader, TensorDataset

# Custom collate function for variable-length sequences
def collate_fn(batch):
    batch.sort(key=lambda x: x[2], reverse=True)
    data, labels, lengths = zip(*batch)
    data = torch.nn.utils.rnn.pad_sequence(data, batch_first=True)
    labels = torch.tensor(labels, dtype=torch.float32)
    lengths = torch.tensor(lengths, dtype=torch.int64)
    return data.to(device), labels.to(device), lengths.to(device)

# Federated learning training loop
global_losses = []
global_accuracies = []
num_rounds = 10

for round in range(num_rounds):
    local_models = []

    # Local training on each client
    for client_id in range(num_clients):
        # Create local model
        local_model = LSTMModel(input_size, hidden_size,
                                num_layers, output_size).to(device)
        local_model.load_state_dict(global_model.state_dict())

        # Prepare client data
        X = [torch.tensor(d, dtype=torch.float32)
              for d in client_data[client_id]['data']]
        X_padded = torch.nn.utils.rnn.pad_sequence(X, batch_first=True)
        y = torch.tensor(client_data[client_id]['labels'],
                          dtype=torch.float32)
        lengths = torch.tensor(client_data[client_id]['lengths'],
                                dtype=torch.int64)

        # Create data loader
        tensor_dataset = TensorDataset(X_padded, y, lengths)
        dataloader = DataLoader(tensor_dataset, batch_size=128,
                                shuffle=True, collate_fn=collate_fn)

        # Local training
        criterion = nn.BCELoss()
        optimizer = optim.Adam(local_model.parameters(), lr=0.001)

```

```

    for epoch in range(1): # Single epoch per round
        for inputs, targets, lengths in dataloader:
            lengths = lengths.cpu()
            outputs = local_model(inputs, lengths)
            loss = criterion(outputs.squeeze(), targets)

            optimizer.zero_grad()
            loss.backward()
            optimizer.step()

        local_models.append(local_model.state_dict())

# Aggregate local models (FedAvg)
global_state_dict = global_model.state_dict()
for key in global_state_dict.keys():
    global_state_dict[key] = torch.stack(
        [local_models[i][key] for i in range(num_clients)],
        dim=0).mean(dim=0)

global_model.load_state_dict(global_state_dict)

print(f"Round {round+1}/{num_rounds} completed")

```

LISTING A.11: Core Federated Learning Training Implementation

A.2.4 Performance Monitoring and Evaluation

A.2.4.1 Resource Usage Tracking

```

import psutil
import time

# Initialize resource monitoring
start_time = time.time()
cpu_usage = []
memory_usage = []

# Training loop with resource monitoring
for round in range(num_rounds):
    # Record resource usage
    cpu_usage.append(psutil.cpu_percent())
    memory_usage.append(psutil.virtual_memory().percent)

    # ... federated training code ...

    # Evaluate global model
    global_loss = 0
    correct = 0
    total = 0

    for client_id in range(num_clients):
        # ... evaluation code ...
        pass

    accuracy = correct / total

```

```

global_losses.append(global_loss)
global_accuracies.append(accuracy)

# Calculate total training time
total_time = time.time() - start_time

# Save results to CSV
import csv
with open('results/federated_learning_results.csv', 'w', newline='') as file:
    writer = csv.writer(file)
    writer.writerow(['Round', 'Global Loss', 'Global Accuracy'])
    for round, (loss, accuracy) in enumerate(
        zip(global_losses, global_accuracies), start=1):
        writer.writerow([round, loss, accuracy])
    writer.writerow(['Total Training Time (s)', total_time])
    writer.writerow(['Average CPU Usage (%)', np.mean(cpu_usage)])
    writer.writerow(['Average Memory Usage (%)', np.mean(memory_usage)])

```

LISTING A.12: Resource Usage Monitoring

A.2.5 Dataset Statistics and Benchmarks

A.2.5.1 Complete Dataset Statistics

TABLE A.4: Comprehensive Dataset Statistics

Metric	Value	Percentage
Sample Distribution		
Total Samples	4,510	100.0%
Fall Samples (F01-F15)	1,798	39.9%
ADL Samples (D01-D19)	2,712	60.1%
Participant Distribution		
Young Adults (SA01-SA23)	23	62.2%
Elderly (SE01-SE15)	14	37.8%
Male Participants	19	51.4%
Female Participants	18	48.6%
Sequence Characteristics		
Average Sequence Length	25.3	-
Minimum Sequence Length	12	-
Maximum Sequence Length	100	-
Standard Deviation	18.7	-
Sensor Dimensions		
ADXL345 Accelerometer	3	33.3%
ITG3200 Gyroscope	3	33.3%
MMA8451Q Accelerometer	3	33.3%
Total Feature Dimensions	9	100.0%

TABLE A.5: Data Preprocessing Performance Metrics

Operation	Time (s)	Memory (MB)	CPU (%)
Raw data loading	45.2	1,250	85
Data cleaning	12.8	890	75
Downsampling (200:1)	8.3	425	60
Z-score normalisation	5.7	340	70
NPZ format conversion	18.9	520	55
Dataset preparation	22.4	780	80
Total Pipeline	113.3	Peak: 1,250	Avg: 71

A.2.5.2 Preprocessing Performance Benchmarks

This comprehensive appendix provides all the implementation details, complete code listings, and technical specifications necessary for reproducing the dataset preprocessing and federated learning pipeline used in our fall detection study.

Appendix B

Algorithm details

This appendix presents the complete algorithmic framework for IoT intrusion detection and human activity recognition using deep neural networks. The algorithm encompasses network architecture design, training procedures, optimisation strategies, and inference mechanisms.

B.1 IoT Intrusion Detection Algorithm

B.1.1 Deep Neural Network Architecture

The core of our intrusion detection system is a deep feedforward neural network designed to learn hierarchical representations of IoT network flow features.

Network Architecture Specification The deep neural network architecture consists of multiple fully connected layers with progressive dimensionality reduction:

$$\mathbf{h}^{(0)} = \mathbf{x} \in \mathbb{R}^{12} \tag{B.1}$$

$$\mathbf{h}^{(1)} = \text{ReLU}(\text{BN}(\mathbf{W}^{(1)}\mathbf{h}^{(0)} + \mathbf{b}^{(1)})) \tag{B.2}$$

$$\mathbf{h}^{(2)} = \text{Dropout}(\mathbf{h}^{(1)}, p = 0.4) \tag{B.3}$$

$$\mathbf{h}^{(3)} = \text{ReLU}(\text{BN}(\mathbf{W}^{(2)}\mathbf{h}^{(2)} + \mathbf{b}^{(2)})) \tag{B.4}$$

$$\mathbf{h}^{(4)} = \text{Dropout}(\mathbf{h}^{(3)}, p = 0.4) \tag{B.5}$$

$$\vdots \tag{B.6}$$

$$\mathbf{y} = \sigma(\mathbf{W}^{(6)}\mathbf{h}^{(10)} + \mathbf{b}^{(6)}) \tag{B.7}$$

where:

- $\mathbf{W}^{(i)} \in \mathbb{R}^{n_i \times n_{i-1}}$ are weight matrices for layer i
- $\mathbf{b}^{(i)} \in \mathbb{R}^{n_i}$ are bias vectors for layer i
- $\text{BN}(\cdot)$ denotes batch normalisation
- $\text{ReLU}(\cdot) = \max(0, \cdot)$ is the rectified linear activation
- $\sigma(\cdot) = \frac{1}{1+e^{-\cdot}}$ is the sigmoid activation function
- Layer dimensions: [12, 128, 96, 64, 32, 16, 1]

```
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Dropout, BatchNormalization
from tensorflow.keras.regularizers import l2

def build_deep_neural_network(input_dim=12):
    """
    Build and compile the deep neural network for IoT intrusion detection

    Args:
        input_dim (int): Number of input features

    Returns:
        tf.keras.Model: Compiled neural network model
    """
    model = Sequential([
        # Layer 1: Input to first hidden layer
        Dense(128, activation='relu', input_shape=(input_dim,),
            kernel_regularizer=l2(0.001), name='dense_1'),
        BatchNormalization(name='batch_norm_1'),
        Dropout(0.4, name='dropout_1'),

        # Layer 2: Second hidden layer
        Dense(96, activation='relu', kernel_regularizer=l2(0.001),
            name='dense_2'),
        BatchNormalization(name='batch_norm_2'),
        Dropout(0.4, name='dropout_2'),

        # Layer 3: Third hidden layer
        Dense(64, activation='relu', kernel_regularizer=l2(0.001),
            name='dense_3'),
        BatchNormalization(name='batch_norm_3'),
        Dropout(0.3, name='dropout_3'),

        # Layer 4: Fourth hidden layer
        Dense(32, activation='relu', kernel_regularizer=l2(0.001),
            name='dense_4'),
        BatchNormalization(name='batch_norm_4'),
        Dropout(0.3, name='dropout_4'),

        # Layer 5: Fifth hidden layer
        Dense(16, activation='relu', name='dense_5'),
        Dropout(0.2, name='dropout_5'),
```

```

        # Output layer: Binary classification
        Dense(1, activation='sigmoid', name='output')
    ])

    # Compile the model
    model.compile(
        optimizer=tf.keras.optimizers.Adam(learning_rate=0.001),
        loss='binary_crossentropy',
        metrics=['accuracy', 'precision', 'recall']
    )

    return model

```

LISTING B.1: Deep Neural Network Architecture Implementation

B.1.2 Training Algorithm

The training process employs advanced optimisation techniques to ensure robust learning and prevent overfitting.

Optimisation Objective The model is trained to minimise the weighted binary cross-entropy loss:

$$\mathcal{L}(\theta) = -\frac{1}{N} \sum_{i=1}^N w_{y_i} [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] + \lambda \Omega(\theta) \quad (\text{B.8})$$

where:

- N is the number of training samples
- $y_i \in \{0, 1\}$ is the true label for sample i
- $\hat{y}_i = f(\mathbf{x}_i; \theta)$ is the predicted probability
- w_{y_i} is the class weight for label y_i
- λ is the regularisation coefficient
- $\Omega(\theta) = \sum_l \|\mathbf{W}^{(l)}\|_F^2$ is the L2 regularisation term

Class Weight Calculation To address class imbalance, we compute balanced class weights:

$$w_0 = \frac{N}{2 \cdot N_0} \quad (\text{benign weight}) \quad (\text{B.9})$$

$$w_1 = \frac{N}{2 \cdot N_1} \quad (\text{malicious weight}) \quad (\text{B.10})$$

where N_0 and N_1 are the number of benign and malicious samples respectively.

Algorithm 4 IoT Intrusion Detection Training Algorithm

Require: Training data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, validation data \mathcal{D}_{val}

Ensure: Trained model parameters θ^*

```

1: Initialize network parameters  $\theta_0$  randomly
2: Compute class weights  $\{w_0, w_1\}$  from training labels
3: Set learning rate  $\alpha = 0.001$ , patience  $p = 15$ , reduction factor  $\gamma = 0.5$ 
4: Initialize best validation loss  $L_{best} = \infty$ , patience counter  $c = 0$ 
5: for epoch  $e = 1$  to  $E_{max}$  do
6:   Shuffle training data  $\mathcal{D}$ 
7:   for each mini-batch  $\mathcal{B} \subset \mathcal{D}$  do
8:     Forward pass:  $\hat{\mathbf{y}} = f(\mathbf{X}_{\mathcal{B}}; \theta)$ 
9:     Compute weighted loss:  $L = \mathcal{L}(\theta, \mathcal{B})$ 
10:    Backward pass:  $\nabla_{\theta} L$ 
11:    Update parameters:  $\theta \leftarrow \theta - \alpha \nabla_{\theta} L$ 
12:  end for
13:  Evaluate on validation set:  $L_{val} = \mathcal{L}(\theta, \mathcal{D}_{val})$ 
14:  if  $L_{val} < L_{best}$  then
15:     $L_{best} \leftarrow L_{val}, \theta^* \leftarrow \theta, c \leftarrow 0$ 
16:  else
17:     $c \leftarrow c + 1$ 
18:    if  $c \geq p$  then
19:      if learning rate not reduced recently then
20:         $\alpha \leftarrow \gamma \cdot \alpha, c \leftarrow 0$ 
21:      else
22:        break ▷ Early stopping
23:      end if
24:    end if
25:  end if
26: end for
27: return  $\theta^*$ 

```

```

from sklearn.utils.class_weight import compute_class_weight
from tensorflow.keras.callbacks import EarlyStopping, ReduceLR0nPlateau, ModelCheckpoint
import numpy as np

```

```
def train_intrusion_detection_model(model, X_train, y_train,
                                   validation_split=0.2,
                                   epochs=100, batch_size=256):
    """
    Train the IoT intrusion detection model with advanced techniques

    Args:
        model: Compiled Keras model
        X_train: Training features
        y_train: Training labels
        validation_split: Fraction of training data for validation
        epochs: Maximum number of training epochs
        batch_size: Training batch size

    Returns:
        training_history: Training history object
    """
    print("Starting model training...")

    # Calculate class weights for imbalanced data
    class_weights = compute_class_weight(
        'balanced',
        classes=np.unique(y_train),
        y=y_train
    )
    class_weight_dict = {i: class_weights[i] for i in range(len(class_weights))}

    print(f"Class weights: {class_weight_dict}")

    # Define advanced training callbacks
    callbacks = [
        # Early stopping to prevent overfitting
        EarlyStopping(
            monitor='val_loss',
            patience=15,
            restore_best_weights=True,
            verbose=1,
            mode='min'
        ),

        # Learning rate reduction on plateau
        ReduceLROnPlateau(
            monitor='val_loss',
            factor=0.5,
            patience=10,
            min_lr=1e-7,
            verbose=1,
            mode='min'
        ),

        # Model checkpointing
        ModelCheckpoint(
            'best_iot_intrusion_model.h5',
            monitor='val_loss',
            save_best_only=True,
            save_weights_only=False,
            verbose=1,
            mode='min'
        )
    ]
```

```

]

# Train the model
history = model.fit(
    X_train, y_train,
    epochs=epochs,
    batch_size=batch_size,
    validation_split=validation_split,
    class_weight=class_weight_dict,
    callbacks=callbacks,
    verbose=1,
    shuffle=True
)

print("Model training completed!")
return history

```

LISTING B.2: Training Algorithm Implementation

B.1.3 Inference Algorithm

The inference process transforms raw network flow data into intrusion detection decisions through a multi-stage pipeline.

Prediction Pipeline Given a new network flow \mathbf{x}_{new} , the prediction process follows:

$$\mathbf{x}_{scaled} = \frac{\mathbf{x}_{new} - \boldsymbol{\mu}}{\sigma} \quad (\text{B.11})$$

$$p_{malicious} = f(\mathbf{x}_{scaled}; \theta^*) \quad (\text{B.12})$$

$$\hat{y} = \begin{cases} 1 & \text{if } p_{malicious} \geq \tau \\ 0 & \text{otherwise} \end{cases} \quad (\text{B.13})$$

where $\boldsymbol{\mu}$ and σ are the scaling parameters from training, and $\tau = 0.5$ is the decision threshold.

Algorithm 5 IoT Intrusion Detection Inference Algorithm**Require:** Trained model $f(\cdot; \theta^*)$, scaling parameters (μ, σ) , threshold τ **Require:** New network flow features $\mathbf{x}_{new} \in \mathbb{R}^{12}$ **Ensure:** Prediction $\hat{y} \in \{0, 1\}$, confidence $p_{malicious} \in [0, 1]$ 1: **Feature Preprocessing:**2: $\mathbf{x}_{scaled} \leftarrow \frac{\mathbf{x}_{new} - \mu}{\sigma}$ 3: **Forward Propagation:**4: $\mathbf{h}^{(0)} \leftarrow \mathbf{x}_{scaled}$ 5: **for** layer $l = 1$ to $L - 1$ **do**6: $\mathbf{z}^{(l)} \leftarrow \mathbf{W}^{(l)}\mathbf{h}^{(l-1)} + \mathbf{b}^{(l)}$ 7: $\mathbf{z}_{norm}^{(l)} \leftarrow \text{BatchNorm}(\mathbf{z}^{(l)})$ ▷ if batch norm layer8: $\mathbf{h}^{(l)} \leftarrow \text{ReLU}(\mathbf{z}_{norm}^{(l)})$ 9: **end for**10: **Output Computation:**11: $\mathbf{z}_{out} \leftarrow \mathbf{W}^{(L)}\mathbf{h}^{(L-1)} + \mathbf{b}^{(L)}$ 12: $p_{malicious} \leftarrow \sigma(\mathbf{z}_{out})$ 13: **Decision Making:**14: $\hat{y} \leftarrow \mathbf{1}[p_{malicious} \geq \tau]$ 15: **return** $\hat{y}, p_{malicious}$

```

import numpy as np

def predict_intrusion(model, scaler, X_new, threshold=0.5):
    """
    Predict intrusion for new network flow data

    Args:
        model: Trained Keras model
        scaler: Fitted StandardScaler from training
        X_new: New network flow features (can be single sample or batch)
        threshold: Decision threshold for binary classification

    Returns:
        predictions: Binary predictions (0=benign, 1=malicious)
        probabilities: Prediction probabilities
        confidence: Confidence scores
    """
    # Ensure X_new is 2D array
    if X_new.ndim == 1:
        X_new = X_new.reshape(1, -1)

    # Feature scaling using training parameters
    X_scaled = scaler.transform(X_new)

    # Model inference
    probabilities = model.predict(X_scaled, verbose=0).flatten()

    # Binary decisions based on threshold
    predictions = (probabilities >= threshold).astype(int)

```

```

    # Confidence calculation (distance from decision boundary)
    confidence = np.abs(probabilities - threshold)

    return predictions, probabilities, confidence

def batch_predict_intrusion(model, scaler, X_batch, threshold=0.5,
                            batch_size=1024):
    """
    Efficient batch prediction for large datasets

    Args:
        model: Trained Keras model
        scaler: Fitted StandardScaler
        X_batch: Batch of network flow features
        threshold: Decision threshold
        batch_size: Processing batch size

    Returns:
        predictions: Array of binary predictions
        probabilities: Array of prediction probabilities
    """
    n_samples = X_batch.shape[0]
    all_predictions = []
    all_probabilities = []

    # Process in batches to manage memory
    for i in range(0, n_samples, batch_size):
        end_idx = min(i + batch_size, n_samples)
        batch = X_batch[i:end_idx]

        # Scale features
        batch_scaled = scaler.transform(batch)

        # Predict
        batch_probs = model.predict(batch_scaled, verbose=0).flatten()
        batch_preds = (batch_probs >= threshold).astype(int)

        all_predictions.extend(batch_preds)
        all_probabilities.extend(batch_probs)

    return np.array(all_predictions), np.array(all_probabilities)

```

LISTING B.3: Inference Algorithm Implementation

B.1.4 Performance Evaluation Algorithm

Comprehensive evaluation of the intrusion detection system using multiple metrics and statistical analyses.

Evaluation Metrics The algorithm computes the following performance metrics:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (\text{B.14})$$

$$\text{Precision} = \frac{TP}{TP + FP} \quad (\text{B.15})$$

$$\text{Recall (Sensitivity)} = \frac{TP}{TP + FN} \quad (\text{B.16})$$

$$\text{Specificity} = \frac{TN}{TN + FP} \quad (\text{B.17})$$

$$\text{F1-Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (\text{B.18})$$

$$\text{AUC-ROC} = \int_0^1 \text{TPR}(t) d\text{FPR}(t) \quad (\text{B.19})$$

where TP , TN , FP , FN represent true positives, true negatives, false positives, and false negatives respectively.

```

from sklearn.metrics import (accuracy_score, precision_score, recall_score,
                             f1_score, roc_auc_score, confusion_matrix,
                             classification_report, roc_curve)

import matplotlib.pyplot as plt
import seaborn as sns

def evaluate_intrusion_detection_model(model, X_test, y_test, scaler):
    """
    Comprehensive evaluation of the trained intrusion detection model

    Args:
        model: Trained Keras model
        X_test: Test features
        y_test: Test labels
        scaler: Fitted StandardScaler

    Returns:
        dict: Comprehensive evaluation metrics
    """
    print("Evaluating model performance...")

    # Scale test features
    X_test_scaled = scaler.transform(X_test)

    # Generate predictions
    y_pred_proba = model.predict(X_test_scaled, verbose=0).flatten()
    y_pred = (y_pred_proba >= 0.5).astype(int)

    # Calculate comprehensive metrics
    metrics = {
        'accuracy': accuracy_score(y_test, y_pred),
        'precision': precision_score(y_test, y_pred),
        'recall': recall_score(y_test, y_pred),
        'f1_score': f1_score(y_test, y_pred),
        'auc_roc': roc_auc_score(y_test, y_pred_proba),
        'specificity': specificity_score(y_test, y_pred)
    }

```

```

# Confusion matrix
cm = confusion_matrix(y_test, y_pred)
tn, fp, fn, tp = cm.ravel()

metrics.update({
    'true_positives': tp,
    'true_negatives': tn,
    'false_positives': fp,
    'false_negatives': fn
})

# Print results
print("\n" + "="*60)
print("INTRUSION DETECTION MODEL EVALUATION RESULTS")
print("="*60)

for metric, value in metrics.items():
    if isinstance(value, float):
        print(f"{metric.upper().replace('_', ' ')}: {value:.4f}")
    else:
        print(f"{metric.upper().replace('_', ' ')}: {value}")

# Detailed classification report
print("\nDetailed Classification Report:")
print(classification_report(y_test, y_pred,
                           target_names=['Benign', 'Malicious']))

return metrics, y_pred, y_pred_proba

def specificity_score(y_true, y_pred):
    """Calculate specificity (true negative rate)"""
    cm = confusion_matrix(y_true, y_pred)
    tn, fp, fn, tp = cm.ravel()
    return tn / (tn + fp) if (tn + fp) > 0 else 0.0

def plot_evaluation_results(y_test, y_pred, y_pred_proba):
    """
    Create comprehensive evaluation visualizations

    Args:
        y_test: True labels
        y_pred: Predicted labels
        y_pred_proba: Prediction probabilities
    """
    fig, axes = plt.subplots(2, 3, figsize=(18, 12))

    # Confusion Matrix
    cm = confusion_matrix(y_test, y_pred)
    sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
                xticklabels=['Benign', 'Malicious'],
                yticklabels=['Benign', 'Malicious'], ax=axes[0,0])
    axes[0,0].set_title('Confusion Matrix')
    axes[0,0].set_ylabel('True Label')
    axes[0,0].set_xlabel('Predicted Label')

    # ROC Curve
    fpr, tpr, _ = roc_curve(y_test, y_pred_proba)
    auc = roc_auc_score(y_test, y_pred_proba)
    axes[0,1].plot(fpr, tpr, label=f'ROC Curve (AUC = {auc:.3f})', linewidth=2)

```

```

axes[0,1].plot([0, 1], [0, 1], 'k--', label='Random Classifier')
axes[0,1].set_xlabel('False Positive Rate')
axes[0,1].set_ylabel('True Positive Rate')
axes[0,1].set_title('ROC Curve')
axes[0,1].legend()
axes[0,1].grid(True)

# Prediction Distribution
axes[0,2].hist(y_pred_proba[y_test == 0], bins=50, alpha=0.7,
              label='Benign', color='blue', density=True)
axes[0,2].hist(y_pred_proba[y_test == 1], bins=50, alpha=0.7,
              label='Malicious', color='red', density=True)
axes[0,2].axvline(x=0.5, color='black', linestyle='--',
                 label='Decision Threshold')
axes[0,2].set_xlabel('Prediction Probability')
axes[0,2].set_ylabel('Density')
axes[0,2].set_title('Prediction Probability Distribution')
axes[0,2].legend()
axes[0,2].grid(True)

# Performance Metrics Bar Chart
metrics = {
    'Accuracy': accuracy_score(y_test, y_pred),
    'Precision': precision_score(y_test, y_pred),
    'Recall': recall_score(y_test, y_pred),
    'F1-Score': f1_score(y_test, y_pred),
    'AUC-ROC': roc_auc_score(y_test, y_pred_proba)
}

bars = axes[1,0].bar(metrics.keys(), metrics.values(),
                    color=['skyblue', 'lightgreen', 'lightcoral', 'gold', 'plum'])
axes[1,0].set_ylim(0, 1)
axes[1,0].set_title('Performance Metrics')
axes[1,0].set_ylabel('Score')
axes[1,0].tick_params(axis='x', rotation=45)

# Add value labels on bars
for bar, value in zip(bars, metrics.values()):
    axes[1,0].text(bar.get_x() + bar.get_width()/2, bar.get_height() + 0.01,
                  f'{value:.3f}', ha='center', va='bottom')

# Class Distribution
class_dist = np.bincount(y_test)
axes[1,1].pie(class_dist, labels=['Benign', 'Malicious'], autopct='%1.1f%%')
axes[1,1].set_title('Test Set Class Distribution')

# Error Analysis
errors = (y_test != y_pred)
error_probs = y_pred_proba[errors]
axes[1,2].hist(error_probs, bins=30, alpha=0.7, color='orange')
axes[1,2].axvline(x=0.5, color='black', linestyle='--',
                 label='Decision Threshold')
axes[1,2].set_xlabel('Prediction Probability')
axes[1,2].set_ylabel('Number of Errors')
axes[1,2].set_title('Error Distribution by Prediction Probability')
axes[1,2].legend()
axes[1,2].grid(True)

plt.tight_layout()

```

```
plt.savefig('Figures/iot_detection_evaluation.png', dpi=300, bbox_inches='tight')
plt.show()
```

LISTING B.4: Comprehensive Evaluation Implementation

B.1.5 Real-time Detection Algorithm

For practical deployment, the algorithm includes real-time processing capabilities for continuous IoT network monitoring.

Algorithm 6 Real-time IoT Intrusion Detection

Require: Trained model $f(\cdot; \theta^*)$, feature extractor $\phi(\cdot)$, alert threshold τ_{alert}

Ensure: Continuous intrusion monitoring

- 1: Initialize alert buffer $\mathcal{B}_{alert} \leftarrow \emptyset$
 - 2: Initialize detection statistics $\mathcal{S} \leftarrow \{\}$
 - 3: **while** network monitoring active **do**
 - 4: Capture network flow \mathbf{flow}_{raw}
 - 5: Extract features: $\mathbf{x} \leftarrow \phi(\mathbf{flow}_{raw})$
 - 6: **if** \mathbf{x} is valid **then**
 - 7: $\hat{y}, p_{mal} \leftarrow \text{predict_intrusion}(\mathbf{x})$
 - 8: **if** $\hat{y} = 1$ **and** $p_{mal} \geq \tau_{alert}$ **then**
 - 9: Create alert: $alert \leftarrow \text{generate_alert}(\mathbf{flow}_{raw}, p_{mal})$
 - 10: $\mathcal{B}_{alert} \leftarrow \mathcal{B}_{alert} \cup \{alert\}$
 - 11: Trigger security response
 - 12: **end if**
 - 13: Update detection statistics: \mathcal{S}
 - 14: **end if**
 - 15: Wait for next network flow
 - 16: **end while**
-

B.1.6 Algorithm Complexity Analysis

Training Complexity

- **Time Complexity:** $O(E \cdot B \cdot \sum_{l=1}^L n_l \cdot n_{l-1})$ where E is epochs, B is batches per epoch, and n_l is neurons in layer l
- **Space Complexity:** $O(\sum_{l=1}^L n_l \cdot n_{l-1} + N \cdot d)$ for model parameters and training data
- **Practical Training Time:** Approximately 30-60 minutes on standard GPU hardware

Inference Complexity

- **Time Complexity:** $O(\sum_{l=1}^L n_l \cdot n_{l-1})$ per sample
- **Space Complexity:** $O(\max(n_l))$ for forward pass computation
- **Practical Inference Time:** < 1 millisecond per network flow on standard hardware

This comprehensive algorithmic framework ensures robust, efficient, and scalable IoT intrusion detection suitable for both research and practical deployment scenarios. The combination of deep learning sophistication with computational efficiency makes it well-suited for real-world IoT security applications.

B.2 Detailed Implementation and Technical Specifications

This appendix provides comprehensive technical details, implementation specifications, and extended analysis for the LSTM-Attention fall detection model. The content ensures full reproducibility of the research and offers deeper insights into the model's behavior and performance characteristics.

B.3 HAR Algorithm

This section provides comprehensive technical details, implementation specifications, and extended analysis for the LSTM-Attention fall detection model. The content ensures full reproducibility of the research and offers deeper insights into the model's behavior and performance characteristics.

B.3.1 Complete Model Architecture Specification

The complete mathematical formulation of the LSTM-Attention model can be expressed as follows:

Given a sensor sequence $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T]$ where $\mathbf{x}_t \in \mathbb{R}^9$ represents the sensor readings at time step t , the input projection transforms the sequence to:

$$\mathbf{H}^{(0)} = \mathbf{X}\mathbf{W}_{proj} + \mathbf{b}_{proj} \quad (\text{B.20})$$

where $\mathbf{W}_{proj} \in \mathbb{R}^{9 \times 128}$ and $\mathbf{b}_{proj} \in \mathbb{R}^{128}$ are learnable parameters.

The bidirectional LSTM processes the sequence in both directions:

$$\vec{\mathbf{h}}_t^{(l)} = \text{LSTM}_{forward}^{(l)}(\mathbf{h}_t^{(l-1)}, \vec{\mathbf{h}}_{t-1}^{(l)}) \quad (\text{B.21})$$

$$\overleftarrow{\mathbf{h}}_t^{(l)} = \text{LSTM}_{backward}^{(l)}(\mathbf{h}_t^{(l-1)}, \overleftarrow{\mathbf{h}}_{t+1}^{(l)}) \quad (\text{B.22})$$

$$\mathbf{h}_t^{(l)} = [\vec{\mathbf{h}}_t^{(l)}; \overleftarrow{\mathbf{h}}_t^{(l)}] \quad (\text{B.23})$$

where $l \in \{1, 2\}$ denotes the layer index and $[\cdot; \cdot]$ represents concatenation.

For each attention layer k , the multi-head attention is computed as:

$$\mathbf{Q}^{(k)} = \mathbf{H}^{(k-1)} \mathbf{W}_Q^{(k)} \quad (\text{B.24})$$

$$\mathbf{K}^{(k)} = \mathbf{H}^{(k-1)} \mathbf{W}_K^{(k)} \quad (\text{B.25})$$

$$\mathbf{V}^{(k)} = \mathbf{H}^{(k-1)} \mathbf{W}_V^{(k)} \quad (\text{B.26})$$

$$\text{head}_i^{(k)} = \text{Attention}(\mathbf{Q}_i^{(k)}, \mathbf{K}_i^{(k)}, \mathbf{V}_i^{(k)}) \quad (\text{B.27})$$

$$\mathbf{H}^{(k)} = \text{LayerNorm}(\mathbf{H}^{(k-1)} + \text{Concat}(\text{head}_1^{(k)}, \dots, \text{head}_g^{(k)})) \mathbf{W}_O^{(k)} \quad (\text{B.28})$$

The global feature vector is computed using attention-weighted pooling:

$$\boldsymbol{\alpha} = \text{softmax}\left(\frac{1}{H} \sum_{i=1}^H \frac{1}{T} \sum_{j=1}^T \mathbf{A}_{ij}^{(final)}\right) \quad (\text{B.29})$$

$$\mathbf{g} = \sum_{t=1}^T \alpha_t \mathbf{h}_t^{(final)} \quad (\text{B.30})$$

where $\mathbf{A}^{(final)}$ represents the attention weights from the final attention layer.

Bibliography

- [1] X. Yin, Y. Zhu, and J. Hu, "A comprehensive survey of privacy-preserving federated learning," *ACM Computing Surveys*, vol. 54, no. 6, pp. 1–36, 2021.
- [2] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings, *et al.*, "Advances and open problems in federated learning," *Foundations and trends® in machine learning*, vol. 14, no. 1–2, pp. 1–210, 2021.
- [3] V. Mothukuri, R. M. Parizi, S. Pouriyeh, Y. Huang, A. Dehghantanha, and G. Srivastava, "A survey on security and privacy of federated learning," *Future Generation Computer Systems*, vol. 115, pp. 619–640, 2021.
- [4] S. Agrawal, S. Sarkar, O. Aouedi, G. Yenduri, K. Piamrat, M. Alazab, S. Bhattacharya, P. K. R. Maddikunta, and T. R. Gadekallu, "Federated learning for intrusion detection system: Concepts, challenges and future directions," *Computer Communications*, vol. 195, pp. 346–361, 2022.
- [5] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*, pp. 1273–1282, PMLR, 2017.
- [6] A. Hard, K. Rao, R. Mathews, S. Ramaswamy, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage, "Federated learning for mobile keyboard prediction," *arXiv preprint arXiv:1811.03604*, 2018.
- [7] H. Woisetschläger, A. Erben, B. Marino, S. Wang, N. D. Lane, R. Mayer, and H.-A. Jacobsen, "Federated learning priorities under the european union artificial intelligence act," *arXiv preprint arXiv:2402.05968*, 2024.
- [8] R. Ye, R. Ge, X. Zhu, J. Chai, D. Yaxin, Y. Liu, Y. Wang, and S. Chen, "Fedllm-bench: Realistic benchmarks for federated learning of large language models," *Advances in Neural Information Processing Systems*, vol. 37, pp. 111106–111130, 2024.
- [9] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, "Practical secure aggregation for

- privacy-preserving machine learning," in *proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1175–1191, 2017.
- [10] Q. Yang, Y. Liu, T. Chen, and Y. Tong, "Federated machine learning: Concept and applications," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 10, no. 2, pp. 1–19, 2019.
- [11] Y. Liu, Y. Kang, C. Xing, T. Chen, and Q. Yang, "A secure federated transfer learning framework," *IEEE Intelligent Systems*, vol. 35, no. 4, pp. 70–82, 2020.
- [12] D. Chen, N. Zhu, A. Machanavajjhala, and M. Fritz, "Federated learning with synthesized data," in *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1959–1971, 2020.
- [13] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 50–60, 2020.
- [14] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, "Scaffold: Stochastic controlled averaging for federated learning," in *International conference on machine learning*, pp. 5132–5143, PMLR, 2020.
- [15] S. Reddi, Z. Charles, M. Zaheer, Z. Garrett, K. Rush, J. Konečný, S. Kumar, and H. B. McMahan, "Adaptive federated optimization," *arXiv preprint arXiv:2003.00295*, 2020.
- [16] J. Wang, Q. Liu, H. Liang, G. Joshi, and H. V. Poor, "Tackling the objective inconsistency problem in heterogeneous federated optimization," *Advances in neural information processing systems*, vol. 33, pp. 7611–7623, 2020.
- [17] M. Ye, X. Fang, B. Du, P. C. Yuen, and D. Tao, "Heterogeneous federated learning: State-of-the-art and research challenges," *ACM Computing Surveys*, vol. 56, no. 3, pp. 1–44, 2023.
- [18] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra, "Federated learning with non-iid data," in *Proceedings of the 2nd International Conference on Artificial Intelligence and Pattern Recognition*, pp. 86–93, 2018.
- [19] H. A. Rahmani, Y. Deldjoo, A. Tourani, and M. Naghiaei, "The unfairness of active users and popularity bias in point-of-interest recommendation," in *International Workshop on Algorithmic Bias in Search and Recommendation*, pp. 56–68, Springer, 2022.
- [20] F. Sattler, K.-R. Müller, and W. Samek, "Clustered federated learning: Model-agnostic distributed multitask optimization under privacy constraints," *IEEE transactions on neural networks and learning systems*, vol. 32, no. 8, pp. 3710–3722, 2020.

- [21] C. T. Dinh, N. Tran, and J. Nguyen, "Personalized federated learning with moreau envelopes," *Advances in neural information processing systems*, vol. 33, pp. 21394–21405, 2020.
- [22] H. Zheng, H. Liu, Z. Liu, and J. Tan, "Federated temporal-context contrastive learning for fault diagnosis using multiple datasets with insufficient labels," *Advanced Engineering Informatics*, vol. 60, p. 102432, 2024.
- [23] S. Caldas, S. M. K. Duddu, P. Wu, T. Li, J. Konečný, H. B. McMahan, V. Smith, and A. Talwalkar, "Leaf: A benchmark for federated settings," *arXiv preprint arXiv:1812.01097*, 2018.
- [24] S. Kalloori, "Federated learning for cold start recommendations," in *Proceedings of the 7th Joint International Conference on Data Science & Management of Data (11th ACM IKDD CODS and 29th COMAD)*, pp. 599–600, 2024.
- [25] H. Lin, J. Lou, L. Xiong, and C. Shahabi, "Semifed: Semi-supervised federated learning with consistency and pseudo-labeling," *arXiv preprint arXiv:2108.09412*, 2021.
- [26] A. E. Cetinkaya, M. Akin, and S. Sagiroglu, "Improving performance of federated learning based medical image analysis in non-iid settings using image augmentation," in *2021 International Conference on Information Security and Cryptology (ISCTURKEY)*, pp. 69–74, IEEE, 2021.
- [27] Y.-T. Cao, Y. Shi, B. Yu, J. Wang, and D. Tao, "Knowledge-aware federated active learning with non-iid data," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 22279–22289, 2023.
- [28] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov, "How to backdoor federated learning," in *International conference on artificial intelligence and statistics*, pp. 2938–2948, PMLR, 2020.
- [29] J. Geiping, H. Bauermeister, H. Dröge, and M. Moeller, "Inverting gradients-how easy is it to break privacy in federated learning?," *Advances in neural information processing systems*, vol. 33, pp. 16937–16947, 2020.
- [30] M. Song, Z. Wang, Z. Zhang, Y. Song, Q. Wang, J. Ren, and H. Qi, "Analyzing user-level privacy attack against federated learning," *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 10, pp. 2430–2444, 2020.
- [31] P. Blanchard, E. M. El Mhamdi, R. Guerraoui, and J. Stainer, "Machine learning with adversaries: Byzantine tolerant gradient descent," *Advances in neural information processing systems*, vol. 30, 2017.
- [32] B. Wang, Y. Yao, S. Shan, H. Li, B. Viswanath, H. Zheng, and B. Y. Zhao, "Neural cleanse: Identifying and mitigating backdoor attacks in neural networks," in *2019 IEEE symposium on security and privacy (SP)*, pp. 707–723, IEEE, 2019.

- [33] R. C. Geyer, T. Klein, and M. Nabi, "Differentially private federated learning: A client level perspective," *arXiv preprint arXiv:1712.07557*, 2017.
- [34] T. Zhu, G. Li, W. Zhou, and S. Y. Philip, *Differential privacy and applications*. Springer, 2017.
- [35] M. U. Hassan, M. H. Rehmani, and J. Chen, "Differential privacy techniques for cyber physical systems: A survey," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 1, pp. 746–789, 2019.
- [36] L. Yu, L. Liu, C. Pu, M. E. Gurses, and S. Truex, "Differentially private model publishing for deep learning," in *2019 IEEE Symposium on Security and Privacy (SP)*, pp. 332–349, IEEE, 2019.
- [37] Y. Wu, Y. Sun, R. A. Rossi, N. Goyal, and H. Durvasula, "Privacy-preserving dimensionality reduction in federated data analysis," *arXiv preprint arXiv:2012.06743*, 2020.
- [38] M. Jayabalan and M. E. Rana, "Anonymizing healthcare records: a study of privacy preserving data publishing techniques," *Advanced Science Letters*, vol. 24, no. 3, pp. 1694–1697, 2018.
- [39] A. Vedangi and V. Anandam, "Data slicing technique to privacy preserving and data publishing," *Cancer*, vol. 4790, no. 4790, p. 4790, 2013.
- [40] M. E. Nergiz, M. Atzori, and C. Clifton, "Hiding the presence of individuals from shared databases," in *Proceedings of the 2007 ACM SIGMOD international conference on Management of data*, pp. 665–676, 2007.
- [41] L. Sweeney, "k-anonymity: A model for protecting privacy," vol. 10, pp. 557–570, World Scientific, 2002.
- [42] A. Machanavajjhala, J. Gehrke, D. Kifer, and M. Venkatasubramanian, "l-diversity: Privacy beyond k-anonymity," vol. 1, pp. 3–es, ACM New York, NY, USA, 2007.
- [43] N. Li, T. Li, and S. Venkatasubramanian, "t-closeness: Privacy beyond k-anonymity and l-diversity," in *2007 IEEE 23rd international conference on data engineering*, pp. 106–115, IEEE, 2006.
- [44] J. M. Joyce, "Kullback-leibler divergence," in *International encyclopedia of statistical science*, pp. 720–722, Springer, 2011.
- [45] J. R. Hershey and P. A. Olsen, "Approximating the kullback leibler divergence between gaussian mixture models," in *2007 IEEE International Conference on Acoustics, Speech and Signal Processing-ICASSP'07*, vol. 4, pp. IV–317, IEEE, 2007.

- [46] A. Gkoulalas-Divanis, G. Loukides, and J. Sun, "Publishing data from electronic health records while preserving privacy: A survey of algorithms," *Journal of biomedical informatics*, vol. 50, pp. 4–19, 2014.
- [47] P. Samarati, "Protecting respondents identities in microdata release," *IEEE transactions on Knowledge and Data Engineering*, vol. 13, no. 6, pp. 1010–1027, 2001.
- [48] J. Domingo-Ferrer and J. M. Mateo-Sanz, "Practical data-oriented microaggregation for statistical disclosure control," *IEEE Transactions on Knowledge and data Engineering*, vol. 14, no. 1, pp. 189–201, 2002.
- [49] D. Slijepčević, M. Henzl, L. D. Klausner, T. Dam, P. Kieseberg, and M. Zeppelzauer, "k-anonymity in practice: How generalisation and suppression affect machine learning classifiers," *Computers & Security*, vol. 111, p. 102488, 2021.
- [50] H.-J. Liao, C.-H. R. Lin, Y.-C. Lin, and K.-Y. Tung, "Intrusion detection system: A comprehensive review," *Journal of network and computer applications*, vol. 36, no. 1, pp. 16–24, 2013.
- [51] D. E. Denning, "An intrusion-detection model," *IEEE Transactions on software engineering*, no. 2, pp. 222–232, 1987.
- [52] H. Debar, "An introduction to intrusion-detection systems," *Proceedings of Connect*, vol. 2000, 2000.
- [53] R. Mitchell and I.-R. Chen, "A survey of intrusion detection techniques for cyber-physical systems," *ACM Computing Surveys (CSUR)*, vol. 46, no. 4, pp. 1–29, 2014.
- [54] M. Tavallaee, E. Bagheri, W. Lu, and A. A. Ghorbani, "A detailed analysis of the kdd cup 99 data set," in *2009 IEEE symposium on computational intelligence for security and defense applications*, pp. 1–6, Ieee, 2009.
- [55] L. Dhanabal and S. Shantharajah, "A study on nsl-kdd dataset for intrusion detection system based on classification algorithms," *International journal of advanced research in computer and communication engineering*, vol. 4, no. 6, pp. 446–452, 2015.
- [56] R. Panigrahi and S. Borah, "A detailed analysis of cicids2017 dataset for designing intrusion detection systems," *International Journal of Engineering & Technology*, vol. 7, no. 3.24, pp. 479–482, 2018.
- [57] N. Moustafa, "A new distributed architecture for evaluating ai-based security systems at the edge: Network ton_iot datasets," *Sustainable Cities and Society*, vol. 72, p. 102994, 2021.

- [58] Y. Mirsky, T. Doitshman, Y. Elovici, and A. Shabtai, "Kitsune: an ensemble of autoencoders for online network intrusion detection," *arXiv preprint arXiv:1802.09089*, 2018.
- [59] S. Garcia, A. Parmisano, and M. J. Erquiaga, "Iot-23: A labeled dataset with malicious and benign iot network traffic," (*No Title*), 2020.
- [60] N. Farnaaz and M. Jabbar, "Random forest modeling for network intrusion detection system," *Procedia Computer Science*, vol. 89, pp. 213–217, 2016.
- [61] A. Al Hanif and M. Ilyas, "Effective feature engineering framework for securing mqtt protocol in iot environments," *Sensors*, vol. 24, no. 6, p. 1782, 2024.
- [62] K. Miyamoto, H. Goto, R. Ishibashi, C. Han, T. Ban, T. Takahashi, and J. Takeuchi, "Malicious packet classification based on neural network using kitsune features," in *International Conference on Intelligent Systems and Pattern Recognition*, pp. 306–314, Springer, 2022.
- [63] B. B. Zarpelão, R. S. Miani, C. T. Kawakani, and S. C. De Alvarenga, "A survey of intrusion detection in internet of things," *Journal of Network and Computer Applications*, vol. 84, pp. 25–37, 2017.
- [64] L. Mohammadpour, T. C. Ling, C. S. Liew, and A. Aryanfar, "A survey of cnn-based network intrusion detection," *Applied Sciences*, vol. 12, no. 16, p. 8162, 2022.
- [65] Y. Guo, T. Ji, Q. Wang, L. Yu, G. Min, and P. Li, "Unsupervised anomaly detection in iot systems for smart cities," *IEEE Transactions on Network Science and Engineering*, vol. 7, no. 4, pp. 2231–2242, 2020.
- [66] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [67] Z. Wu, H. Zhang, P. Wang, and Z. Sun, "Rtids: A robust transformer-based approach for intrusion detection system," *IEEE Access*, vol. 10, pp. 64375–64387, 2022.
- [68] T. P. Nguyen, H. Nam, and D. Kim, "Transformer-based attention network for in-vehicle intrusion detection," *IEEE Access*, vol. 11, pp. 55389–55403, 2023.
- [69] E. M. Campos, P. F. Saura, A. González-Vidal, J. L. Hernández-Ramos, J. B. Bernabé, G. Baldini, and A. Skarmeta, "Evaluating federated learning for intrusion detection in internet of things: Review and challenges," *Computer Networks*, vol. 203, p. 108661, 2022.

- [70] S. A. Rahman, H. Tout, C. Talhi, and A. Mourad, "Internet of things intrusion detection: Centralized, on-device, or federated learning?," *IEEE Network*, vol. 34, no. 6, pp. 310–317, 2020.
- [71] D. Preuveneers, V. Rimmer, I. Tsingenopoulos, J. Spooren, W. Joosen, and E. Ilie-Zudor, "Chained anomaly detection models for federated learning: An intrusion detection case study," vol. 8, p. 2663, MDPI, 2018.
- [72] T. D. Nguyen, S. Marchal, M. Miettinen, H. Fereidooni, N. Asokan, and A.-R. Sadeghi, "Federated deep learning for detecting iot botnet attacks," *Computer Networks*, vol. 200, p. 108538, 2021.
- [73] O. Friha, M. A. Ferrag, L. Shu, L. Maglaras, K.-K. R. Choo, and M. Nafaa, "Felids: Federated learning-based intrusion detection system for agricultural internet of things," *Journal of Parallel and Distributed Computing*, vol. 165, pp. 17–31, 2022.
- [74] C. Briggs, Z. Fan, and P. Andras, "Federated learning with hierarchical clustering of local updates to improve training on non-iid data," in *2020 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–9, 2020.
- [75] J. K. Aggarwal and Q. Cai, "Human motion analysis: A review," *Computer Vision and Image Understanding*, vol. 73, no. 3, pp. 428–440, 1999.
- [76] M. Habermann, W. Xu, M. Zollhöfer, G. Pons-Moll, and C. Theobalt, "Livecap: Real-time human performance capture from monocular video," *ACM Transactions on Graphics*, vol. 38, no. 2, pp. 1–17, 2019.
- [77] A. Subasi, K. Khateeb, T. Brahimi, and A. Sarirete, "Human activity recognition using machine learning methods in a smart healthcare environment," *Innovation in health informatics*, pp. 123–144, 2020.
- [78] D. Anguita, A. Ghio, L. Oneto, X. Parra, and J. L. Reyes-Ortiz, "A public domain dataset for human activity recognition using smartphones," *European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning, ESANN*, 2013.
- [79] J. R. Kwapisz, G. M. Weiss, and S. A. Moore, "Activity recognition using cell phone accelerometers," vol. 12, pp. 74–82, ACM New York, NY, USA, 2011.
- [80] G. Vavoulas, C. Chatzaki, T. Malliotakis, M. Pediaditis, and M. Tsiknakis, "The mobiact dataset: Recognition of activities of daily living using smartphones," in *International conference on information and communication technologies for ageing well and e-health*, vol. 2, pp. 143–151, SciTePress, 2016.
- [81] A. Sucerquia, J. D. López, and J. F. Vargas-Bonilla, "Sisfall: A fall and movement dataset," *Sensors*, vol. 17, no. 1, p. 198, 2017.

- [82] A. Reiss and D. Stricker, "Introducing a new benchmarked dataset for activity monitoring," in *2012 16th international symposium on wearable computers*, pp. 108–109, IEEE, 2012.
- [83] D. Roggen, A. Calatroni, M. Rossi, T. Holleczeck, K. Förster, G. Tröster, P. Lukowicz, D. Bannach, G. Pirkl, A. Ferscha, *et al.*, "Collecting complex activity datasets in highly rich networked sensor environments," in *2010 Seventh international conference on networked sensing systems (INSS)*, pp. 233–240, IEEE, 2010.
- [84] B. Kwolek and M. Kepski, "Human fall detection on embedded platform using depth maps and wireless accelerometer," *Computer methods and programs in biomedicine*, vol. 117, no. 3, pp. 489–501, 2014.
- [85] W. Dargie, "Analysis of time and frequency domain features of accelerometer measurements," in *2009 Proceedings of 18th international conference on computer communications and networks*, pp. 1–6, IEEE, 2009.
- [86] R. K. Ibrahim, *Novel Gait models and features for Gait patterns classification*. PhD thesis, UNSW Sydney, 2011.
- [87] M. Wacker and H. Witte, "Time-frequency techniques in biomedical signal analysis," *Methods of information in medicine*, vol. 52, no. 04, pp. 279–296, 2013.
- [88] S. Zhang, Z. Wei, J. Nie, L. Huang, S. Wang, and Z. Li, "A review on human activity recognition using vision-based method," *Journal of healthcare engineering*, vol. 2017, no. 1, p. 3090343, 2017.
- [89] D. R. Beddiar, B. Nini, M. Sabokrou, and A. Hadid, "Vision-based human activity recognition: a survey," *Multimedia Tools and Applications*, vol. 79, no. 41, pp. 30509–30555, 2020.
- [90] C. Lugaresi, J. Tang, H. Nash, C. McClanahan, E. Uboweja, M. Hays, F. Zhang, C.-L. Chang, M. G. Yong, J. Lee, *et al.*, "Mediapipe: A framework for building perception pipelines," *arXiv preprint arXiv:1906.08172*, 2019.
- [91] R. Ramamurthy, Sreenivasan, and N. Roy, "Recent trends in machine learning for human activity recognition—a survey," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 8, no. 4, p. e1254, 2018.
- [92] E. Ramanujam, T. Perumal, and S. Padmavathi, "Human activity recognition with smartphone and wearable sensors using deep learning techniques: A review," *IEEE Sensors Journal*, vol. 21, no. 12, pp. 13029–13040, 2021.
- [93] W. Xu, A. Chatterjee, M. Zollhofer, H. Rhodin, D. Mehta, H.-P. Seidel, and C. Theobalt, "Monoperfcap: Human performance capture from monocular video," *ACM Transactions on Graphics (ToG)*, vol. 37, no. 2, pp. 1–15, 2018.

- [94] K. Xia, J. Huang, and H. Wang, "Lstm-cnn architecture for human activity recognition," *Ieee Access*, vol. 8, pp. 56855–56866, 2020.
- [95] I. Dirgová Luptáková, M. Kubovčík, and J. Pospíchal, "Wearable sensor-based human activity recognition with transformer model," *Sensors*, vol. 22, no. 5, p. 1911, 2022.
- [96] G. Kalouris, E. I. Zacharaki, and V. Megalooikonomou, "Improving cnn-based activity recognition by data augmentation and transfer learning," in *2019 IEEE 17th International Conference on Industrial Informatics (INDIN)*, vol. 1, pp. 1387–1394, IEEE, 2019.
- [97] K. Sozinov, V. Vlassov, and S. Girdzijauskas, "Human activity recognition using federated learning," in *2018 IEEE Intl Conf on Parallel & Distributed Processing with Applications, Ubiquitous Computing & Communications, Big Data & Cloud Computing, Social Computing & Networking, Sustainable Computing & Communications (ISPA/IUCC/BDCloud/SocialCom/SustainCom)*, pp. 1103–1111, IEEE, 2018.
- [98] X. Ouyang, Z. Xie, J. Zhou, J. Huang, and G. Xing, "Clusterfl: a similarity-aware federated learning system for human activity recognition," in *Proceedings of the 19th Annual International Conference on Mobile Systems, Applications, and Services*, pp. 54–66, 2021.
- [99] B. K. Iwana and S. Uchida, "An empirical survey of data augmentation for time series classification with neural networks," *PLoS ONE*, vol. 16, no. 7, p. e0254841, 2021.
- [100] Z. Liu, N. Thapa, A. Shaver, K. Roy, X. Yuan, and S. Khorsandroo, "Anomaly detection on iot network intrusion using machine learning," in *2020 International Conference on Artificial Intelligence, Big Data, Computing and Data Communication Systems (icABCD)*, pp. 1–5, IEEE, 2020.
- [101] L. Xiao, X. Wan, X. Lu, Y. Zhang, and D. Wu, "Iot security techniques based on machine learning: How do iot devices use ai to enhance security?," *IEEE Signal Processing Magazine*, vol. 35, no. 5, pp. 41–49, 2018.
- [102] M. Mohri, G. Sivek, and A. T. Suresh, "Agnostic federated learning," in *International Conference on Machine Learning*, pp. 4615–4625, PMLR, 2019.
- [103] T. Kanwal, A. Anjum, and A. Khan, "Privacy preservation in e-health cloud: taxonomy, privacy requirements, feasibility analysis, and opportunities," *Cluster Computing*, vol. 24, pp. 293–317, 2021.
- [104] A. Majeed and S. Lee, "Anonymization techniques for privacy preserving data publishing: A comprehensive survey," *IEEE Access*, vol. 9, pp. 8512–8545, 2021.

- [105] J. Wieringa, P. Kannan, X. Ma, T. Reutterer, H. Risselada, and B. Skiera, "Data analytics in a privacy-concerned world," *Journal of Business Research*, vol. 122, pp. 915–925, 2021.
- [106] M. Langheinrich, "Privacy by design—principles of privacy-aware ubiquitous systems," in *UbiComp 2001: Ubiquitous Computing: International Conference Atlanta Georgia, USA, September 30–October 2, 2001 Proceedings*, pp. 273–291, Springer, 2001.
- [107] P. Regulation, "Regulation (eu) 2016/679 of the european parliament and of the council," *Regulation (eu)*, vol. 679, p. 2016, 2016.
- [108] K. Bonawitz, H. Eichner, W. Grieskamp, D. Huba, A. Ingerman, V. Ivanov, C. Kiddon, J. Konečný, S. Mazzocchi, B. McMahan, *et al.*, "Towards federated learning at scale: System design," *Proceedings of machine learning and systems*, vol. 1, pp. 374–388, 2019.
- [109] B. Liu, M. Ding, S. Shaham, W. Rahayu, F. Farokhi, and Z. Lin, "When machine learning meets privacy: A survey and outlook," *ACM Computing Surveys (CSUR)*, vol. 54, no. 2, pp. 1–36, 2021.
- [110] L. Li, Y. Fan, M. Tse, and K.-Y. Lin, "A review of applications in federated learning," *Computers & Industrial Engineering*, vol. 149, p. 106854, 2020.
- [111] K. Wei, J. Li, M. Ding, C. Ma, H. H. Yang, F. Farokhi, S. Jin, T. Q. Quek, and H. V. Poor, "Federated learning with differential privacy: Algorithms and performance analysis," *IEEE transactions on information forensics and security*, vol. 15, pp. 3454–3469, 2020.
- [112] Z. Ji, Z. C. Lipton, and C. Elkan, "Differential privacy and machine learning: a survey and review," *arXiv preprint arXiv:1412.7584*, 2014.
- [113] V. Ayala-Rivera, P. McDonagh, T. Cerqueus, L. Murphy, *et al.*, "A systematic comparison and evaluation of k-anonymization algorithms for practitioners," *Transactions on data privacy*, vol. 7, no. 3, pp. 337–370, 2014.
- [114] B. C. Fung, K. Wang, R. Chen, and P. S. Yu, "Privacy-preserving data publishing: A survey of recent developments," *ACM Computing Surveys (Csur)*, vol. 42, no. 4, pp. 1–53, 2010.
- [115] A. Parmisano, S. Garcia, and M. J. Erquiaga, "A labeled dataset with malicious and benign iot network traffic," *Stratosphere Laboratory: Praha, Czech Republic*, 2020.
- [116] A. Imteaj, U. Thakker, S. Wang, J. Li, and M. H. Amini, "A survey on federated learning for resource-constrained iot devices," *IEEE Internet of Things Journal*, vol. 9, no. 1, pp. 1–24, 2021.

-
- [117] Q. Li, Y. Diao, Q. Chen, and B. He, "Federated learning on non-iid data silos: An experimental study," in *2022 IEEE 38th International Conference on Data Engineering (ICDE)*, pp. 965–978, IEEE, 2022.
- [118] M. Langer, Z. He, W. Rahayu, and Y. Xue, "Distributed training of deep learning models: A taxonomic perspective," *IEEE Transactions on Parallel and Distributed Systems*, vol. 31, no. 12, pp. 2802–2818, 2020.
- [119] W. Huang, T. Tiropanis, and G. Konstantinidis, "Federated learning-based iot intrusion detection on non-iid data," in *Global IoT Summit*, pp. 326–337, Springer, 2022.