

**A Privacy-Preserving Federated Learning Approach for Primary
Colorectal Malignancy Detection, Tumor Multiplicity Estimation and
Regional Occurrence Mapping**

By

Md. Abdullah Al Jaber

213-15-4458

Md. Rahat Zaman Sarker

213-15-4462

FINAL YEAR DESIGN PROJECT REPORT

This Report Presented in Partial Fulfillment of the Requirements for
the **Degree of Bachelor of Science in Computer Science and
Engineering**

Supervised by

Md Umaid Hasan

Lecturer (Senior Scale)

Dept. of Computer Science and Engineering
Daffodil International University

Co-Supervised by

Md. Monarul Islam

Lecturer

Dept. of Computer Science and Engineering
Daffodil International University



DAFFODIL INTERNATIONAL UNIVERSITY
Dhaka, Bangladesh

September 16, 2025

APPROVAL

This Project titled "A Privacy-Preserving Federated Learning Approach for Primary Colorectal Malignancy Detection, Tumor Multiplicity Estimation and Regional Occurrence Mapping", submitted by Md. Abdullah Al Jaber, ID No: 213-15-4458 and Md. Rahat Zaman Sarker, ID No: 213-15-4462 to the Department of Computer Science and Engineering, Daffodil International University has been accepted as satisfactory for the partial fulfillment of the requirements for the degree of B.Sc. in Computer Science and Engineering and approved as to its style and contents. The presentation has been held on 16 September, 2025.

BOARD OF EXAMINERS

Dr. S.M Aminul Haque (SMAH)
Professor & Associate Head
Department of Computer Science and Engineering
Faculty of Science & Information Technology
Daffodil International University

Chairman

Dr. Md Alamgir Kabir (DMAK)
Assistant Professor
Department of Computer Science and Engineering
Faculty of Science & Information Technology
Daffodil International University

Internal Examiner

Mr. Md Assduzzaman (MA)
Assistant Professor
Department of Computer Science and Engineering
Faculty of Science & Information Technology
Daffodil International University

Internal Examiner

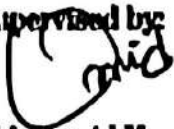
Dr. Md. Zulfiker Mahmud (ZM)
Professor
Department of Computer Science and Engineering
Jagannath University

External Examiner

DECLARATION

We hereby declare that this project has been done by us under the supervision of Mr. Md Umald Hasan, Lecturer (Senior Scale) Department of Computer Science and Engineering, Daffodil International University. We also declare that neither this project nor any part of this project has been submitted elsewhere for the award of any degree or diploma.

Supervised by:

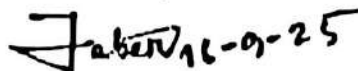


Md. Umald Hasan
Lecturer (Senior Scale)
Dept. of Computer Science and Engineering
Daffodil International University

Co-Supervised by:

Md. Monarul Islam
Lecturer
Dept. of Computer Science and Engineering
Daffodil International University

Submitted by:



Md. Abdullah Al Jaber
Student ID: 213-16-4458
Dept. of Computer Science and Engineering
Daffodil International University



Md. Rahat Zaman Sarker
Student ID: 213-16-4462
Dept. of Computer Science and Engineering
Daffodil International University

ACKNOWLEDGEMENTS

This work would not have been possible without the support and contributions of many individuals over the past two semesters. We are deeply grateful to everyone who has assisted us in one way or another.

First, we express our heartfelt thanks and gratefulness to the almighty for His divine blessing making it possible for us to complete the **Final Year Design Project(FYDP)** successfully.

We are grateful and wish our profound indebtedness to **Md. Umaid Hasan, Lecturer(Senior Scale)**, Department of Computer Science and Engineering, Daffodil International University, Dhaka, Bangladesh. Deep knowledge and keen interest of our supervisor in the field of **Deep Learning** to carry out this project. His endless patience, scholarly guidance, continual encouragement, constant and energetic supervision, constructive criticism, valuable advice, reading many inferior drafts, and correcting them at all stages have made it possible to complete this project.

We would like to express our heartfelt gratitude to the Head of the Department of Computer Science and Engineering, for his kind help in finishing our project and also to other faculty members and the staff of the Department of Computer Science and Engineering, Daffodil International University.

We would like to thank our entire course-mates at Daffodil International University, who took part in this discussion while completing the coursework.

Finally, we must acknowledge with due respect the constant support and patience of our parents.

ABSTRACT

It is among the top causes of deaths in the world-Colorectal cancer. As soon as we detect it, we will survive. However, there are a lot of barriers in the medical field around looking at real patient test data. To address this, we propose a privacy preserving solution e.g federated learning. We have been dealing with structured data from SEER and that has aided them in determining three common tasks. (1) Identification of primary malignancies (2) Assessing tumor multiplicity (3) Tissue-wide associations. We have achieved this by using a neural network model. Where FedAvg combines the model which trains it locally and you will get gradient values/parameters of the model after performing an average on it. With the appropriate preprocessing nodes, binary or multiclass predictions can be obtained sequentially for certain tasks. This model achieves a significant and correct improvement in performance against the original baseline whilst maintaining high security levels, as shown by the results M. Enable colon cancer detection without revealing the original data. Multiple original cancer research underlie AI integration and privacy in health.

Table of Contents

Approval	i
Declaration	ii
Acknowledgements	iii
Abstract	iv
List of Figures	vii
List of Tables	viii
1 Introduction	1
1.1 Introduction	1
1.2 Motivation	2
1.3 Objectives	3
1.4 Methodology	4
1.5 Project Outcome	5
1.6 Organization of the Report	6
2 Background	8
2.1 Introduction	8
2.2 Literature Review	10
2.3 Gap Analysis	14
2.4 Summary	16
3 Research Methodology	17
3.1 Methodology/Requirement Analysis & Design Specification	17
3.1.1 Overview	17
3.1.2 Proposed Methodology/ System Design	17
3.1.3 Functional and Nonfunctional Requirements	18
3.1.4 Data Flow Diagram	19
3.2 Detailed Methodology and Design	20
3.3 Project Plan	23
3.4 Task Allocation	24
3.5 Summary	24
4 Implementation and Results	25
4.1 Environment Setup	25

4.2	Testing and Evaluation/Performance/ Comparative Analysis	26
4.3	Results and Discussion	26
4.4	Summary	28
5	Engineering Standards and Design Challenges	29
5.1	Compliance with the Standards	29
5.1.1	Software Standards	29
5.1.2	Hardware Standards	29
5.1.3	Communication Standards	30
5.2	Impact on Society, Environment and Sustainability	30
5.2.1	Impact on Life	30
5.2.2	Impact on Society & Environment	30
5.2.3	Ethical Aspects	30
5.2.4	Sustainability Plan	30
5.3	Project Management and Financial Analysis	31
5.4	Complex Engineering Problem	31
5.4.1	Complex Problem Solving	31
5.4.2	Engineering Activities	32
5.5	Summary	33
6	Conclusion	34
6.1	Summary	34
6.2	Limitation	34
6.3	Future Work	35
	References	36

List of Figures

3.1	Data Flow Diagram	19
3.2	Regional Tumor Occurrence Bar Chart	20
3.3	Heatmap of the total number of tumors for patient Tumor Count	21
4.1	Confusion Matrix of 1st Feature's Global Model	27
4.2	Avg Weight Update Magnitude Over Round	27
4.3	Train Accuracy vs Validation Accuracy	28

List of Tables

2.1	Summary of Literature Review	13
3.1	Task Allocation Table	24
5.1	Project Management and Financial Analysis Table	31
5.2	Mapping with Complex Engineering Problem	31
5.3	Mapping with knowledge Profile	32
5.4	Mapping with Complex Engineering Activities	32

Chapter 1

Introduction

The study proposes a privacy-preserving federated learning framework using SEER data to improve colorectal cancer detection, tumor burden estimation, and regional mapping while ensuring interpretability and data privacy.

1.1 Introduction

Colorectal cancer (CRC) remains one of the most significant public-health challenges worldwide, ranking among the leading causes of cancer-related morbidity and mortality. The clinical outcome of CRC patients is strongly stage-dependent: early detection and precise characterization of the disease significantly increase the likelihood of successful intervention and long-term survival. In clinical practice, three interrelated pieces of information are particularly important for patient management and treatment planning: (1) whether the tumor represents a primary malignant colorectal lesion, (2) the total tumor burden or multiplicity (i.e., how many malignant or in-situ tumors a patient harbors), and (3) the anatomical region(s) of the colon and rectum affected by the disease. Correct evaluation of these factors guides surgical planning, disease staging and systemic therapy options, and follow-up, and thus has direct implications for the outcome of the patient.

There has been an exponential increase in the availability of structured clinical data, pathology reports and imaging results that are routinely collected for both modern healthcare institutions and cancer registries however, in many cases this data is held in such large volumes that it has no meaningful use for the development of automated colorectal cancer (CRC) screening and prognostication models unless these data are pooled for model development.

Population-based cancer registries, such as the SEER (Surveillance, Epidemiology, and End Results) program and like registries, hold extensive cancer records valuable for developing high-performance predictive models. However essential, practical and ethical concerns constitute a considerable barrier against large-scale pooling of patient-level data between hospitals: privacy regulations (e.g. HIPAA for US or GDPR in Europe), institutional edicts and re-identification risks work as powerful disincentives to data sharing. As a result, tens of high-value datasets are siloed from one another, and few results exploit sample diversity or the ability to reduce bias through the integration of these datasets and few models achieving success in applying across different populations.

Federated learning (FL) provides a favorable solution to the above dilemma by enabling collaborative model training across different data silos, without exchange of raw patient data. In the FL setting, participating organizations only communicate model updates or parameter values, e.g. encrypted ones, with a central server which serves as an aggregator of model

updates during training on the data of the participating organizations. This keeps data close by, and avoids sharing it and preserves privacy but allows better models to leverage and benefit from a larger and much more heterogeneous distribution of input data. For CRC specifically, FL allows the development of models able to (i) detect primary malignancy at initial presentation, (ii) predict tumor multiplicity as a proxy for disease burden, and (iii) classify or localize tumors by anatomic site—without centralized data sharing [13].

We present a privacy-preserving federated learning (FL) framework for colorectal cancer (CRC) using structured SEER records in this study. Informed by the process of harmonization of clinical data set, our framework enforces site-independent preprocessing standards, applies a common deep neural network architecture with task-specific output layers, and performs secure federated averaging to integrate knowledges from virtual clinical sites. The task formulation and feature selection are interpretable to practitioners; thus, the design of the framework is clinically relevant. Furthermore, it is regulatory-compliant since no raw patient data leaves institutional borders and practical because it is robust against non-IID data distributions. Federated deep learning offers nearly centralized predictive performance with the privacy and governance preservation principles of federated learning demonstration via extensive experiments and comparison with centralized baseline.

Contributions of this work:

Development of a unified deep learning based FL pipeline for multiple CRC tasks “early malignancy detection, tumor multiplicity estimation, and regional mapping” that standardizes preprocessing of SEER data, enabling consistent model training across clinical nodes.

Empirical validation in simulated real world scenarios, demonstrates performance comparable to local training while ensuring privacy preservations.

1.2 Motivation

Privacy Preservation

The main goal of this work is to protect the personal data of patients. In traditional centralized machine learning, hospitals and research institutions have to store all kinds of raw medical records in one place. This creates major risks—such as data leakage, identity misuse, and even violation of health regulations. Using federated learning, this risk is greatly reduced. Here, the raw patient data is not sent anywhere, only encrypted model updates are exchanged. As a result, personal data is safe within the institution and no one can access it without permission.

Early Cancer Detection

Detecting colorectal cancer at an early stage greatly increases the patient's chances of survival and increases the chances of treatment. Our framework is designed to detect early signs of the disease. It not only determines whether cancer is present, but also how many

tumors are present and where they are located in the body. As a result, doctors can clearly understand the patient's condition and quickly plan treatment. This increases the chances of survival.

Decentralized Learning

We have done distributed learning here, focusing on multiple hospitals, where all the data is stored in different places rather than being brought to a single server. The advantage of federated learning is that it trains the model without moving the raw data. Each hospital uses its own data locally, but the shared model can combine knowledge from all places. This reduces data silos, makes it easier to add new hospitals, and also makes it possible to comply with legal and ethical issues.

Medical Data Security

Health data security is not limited to protecting privacy. There are other risks—such as hacking, data leaks, or unauthorized use of predictive tools. To prevent these, our framework can be supplemented with secure aggregation, encryption, and stronger local training. This allows hospitals to work together securely without blindly relying on trust. As a result, healthcare is more reliable and secure.

Multi-Hospital Collaboration

Cancer detection models work best when they are trained on a variety of patient data. But no single hospital usually has enough data. Our project allows hospitals to work together. While each can train locally with their own data, they can all contribute to a shared model. New hospitals can join later if they want. This collaborative effort not only improves the model's performance, but also makes the use of AI in everyday oncology care more realistic.

1.3 Objectives

The main objective of this study is to propose a privacy-preserving federated learning framework for colorectal cancer detection and analysis and evaluate its performance. To this end, several specific objectives have been set—

- 1. Model development for colorectal cancer detection**

The first goal is to build a robust machine learning pipeline that can detect early-stage, aggressive colorectal cancer, quantify tumor size, and identify where tumors have spread in the body. This will involve building deep neural network models that are trained collaboratively using data from multiple hospitals—but without having to store raw medical data in a single location.

- 2. Protecting patient confidentiality:**

Protecting patient privacy is essential throughout the model training process. That's why no personally identifiable or sensitive health information ever leaves the hospital. Data will be stored locally at each institution, and only model updates or weights (coefficients) will be sent out. This will also help comply with healthcare regulations like HIPAA and GDPR.

- 3. No need to share raw data**

Our proposed architecture eliminates the need to exchange raw patient data between

participating institutions. Each hospital would retain full control of its data, while still contributing to the collaborative learning process. This decentralized approach not only satisfies ethical and legal obligations, but also facilitates collaboration across multiple institutions in real-world clinical settings.

4. **Evaluating the accuracy and efficiency of the model**

Another goal of the research is to test the predictive power of the federated model—such as how accurately it can classify, how well it performs on regression metrics, and whether it can generate multi-class maps. Communication costs, convergence speed, and computational capacity will also be evaluated to understand whether these models are usable in real clinical settings.

5. **Comparison with centralized models**

Finally, the performance of federated models will be compared with traditional centralized machine learning. This will help to understand the differences between the two approaches, how much privacy protection federated learning offers, and why it may be more effective for real-world oncology research and patient care.

1.4 Methodology

In this study, we propose a multi-task federated learning (FL) pipeline for colorectal cancer detection and analysis using the SEER dataset. The data were consciously divided into three feature-based groups, which enabled three tasks to be performed—primary tumor detection, tumor count estimation, and tumor localization.

Privacy and Collaboration

To protect patient data privacy and support collaborative learning, we have integrated it into a Federated Learning framework. Here, each hospital or client trains a model locally using their own data. The parameters of that model are then centrally aggregated using Federated Averaging (FedAvg). This allows multiple institutions to securely train models together without having to share raw patient data externally.

Implementation

The entire system is built in Python, using TensorFlow, Keras, pandas, and various visualization libraries. This allows the framework to be scaled up, down, or used in real-world healthcare applications.

Key Evaluation

This approach strikes a balance between privacy and accuracy, and across institutions. As a result, it is not just a research-based contribution, but a robust framework that can be effectively used in real-world clinical settings while preserving privacy.

1.5 Project Outcome

This study demonstrates that using Federated Learning (FL) to detect colorectal cancer on SEER clinical data is effective and feasible. The proposed framework simultaneously fulfills multiple goals—detecting early-stage malignant colorectal cancer, quantifying tumor numbers, and identifying the precise location of tumors within the colon and rectum.

It has been experimentally proven that the federated deep neural network model is able to provide predictions as good as the centralized model while maintaining privacy. The system succeeded in three key tasks—

- **First-primary malignancy identification,** It was possible to detect early cancers with high accuracy.
- **Multiplicity estimation of tumors,** The Federated Learning model based on the patient's history gave doctors an accurate idea of the tumor burden.
- **Regional occurrence mapping,** The federated multi-class model was able to accurately localize the tumor location, which made treatment planning more accurate.

Significant contributions of the research:

- **Privacy Protection:** It has been proven that it is possible to build powerful medical AI without sharing patient data. This avoids ethical, legal and regulatory complications.
- **Multi-hospital collaboration:** Federated models allow different hospitals to work together, without the risk of data leakage. As a result, it is easily applicable in real clinical environments.
- **Generalization and stability:** Models trained using data from many institutions are much more reliable than models based on data from a single institution. This reduces bias and makes predictions more accurate for different types of patients.
- **Comparison with centralized models:** Extensive tests have shown that FL performs close to centralized models in terms of accuracy, precision and recall. However, it is much better in terms of privacy protection.
- **Basis for future expansion:** This framework can be used not only for colorectal cancer, but also for other cancers or broader health data analysis. This has opened up new avenues of research in future on secure federated learning, homomorphic encryption, differential privacy, etc.

Overall, this research shows that federated learning is effective, safe, and promising for real-world clinical AI. It demonstrates that it is possible to use modern technology to create cancer detection systems that are accurate, privacy-friendly, and legal. This allows the benefits of AI in healthcare to be realized without compromising patient trust or safety.

1.6 Organization of the Report

This research report is organized chapter-wise, taking the reader step-by-step from identifying the problem to the research results and their implications.

- **Chapter I – Introduction**
This chapter discusses the background and importance of colorectal cancer research. It shows why early detection, tumor quantification, and risk mapping are so important. It also explains why privacy-preserving federated learning is useful for healthcare, what the problem is, what our research aims to achieve, and what the main purpose of this work is.

- **Chapter II – Background**

This chapter presents previous research on colorectal cancer detection, federated learning, and privacy-preserving medical AI. A systematic review is provided, showing the successes and limitations of conventional centralized learning techniques and the improvements of federated learning approaches. The main research question is clarified in this chapter.

- **Chapter III – Research Methodology**

The proposed research plan is explained in detail. The SEER dataset, pre-processing methods, and model development steps are discussed. The federated learning framework is presented—which includes the structure of deep neural networks, local client training, parameter aggregation (FedAvg), and privacy-preserving techniques. Some mathematical aspects of the algorithms are also highlighted.

- **Chapter IV – Implementation and Results**

This chapter describes the framework development steps, tools, libraries, and environment used. It discusses how the experiments were designed, trained, and evaluated. The results are presented focusing on three main goals—Detection of malignant cancer, Estimation of tumor number, Mapping tumor location. A comparison between federated learning and centralized learning is also presented, highlighting the advantages, privacy protection, efficiency, and practical applicability of the proposed approach.

- **Chapter V – Engineering Standards and Design Challenges**

This chapter shows how the project is aligned with international standards such as ISO, IEEE, HIPAA, and GDPR, so that the system is secure, reliable, and acceptable in the medical field. It also discusses the big questions—patient privacy, ethical use of AI, and sustainability. It also discusses engineering challenges, such as strict regulations, complex technical aspects, and the need for new concepts in federated learning—all of which are essential for making AI meaningful and effective in the medical field.

- **Chapter VI – Conclusion**

In this study, we show that it is possible to effectively and accurately detect colorectal cancer while preserving privacy using federated learning. The proposed framework provided performance similar to centralized models, without the need to share raw data. Although it was tested in a limited virtual environment, it has been proven that this technique can be effective and ethical for early detection and intervention in real medical settings. Future developments will aim to improve reliability, clinical effectiveness, privacy, personalization, and the use of multi-modal data.

Chapter 2

Background

Recent studies have shown that Federated Learning (FL) can be an effective way to diagnose cancer across multiple institutions while preserving privacy. However, there are still some major challenges—such as data heterogeneity, high computational cost, model personalization, and clinical interpretability for clinicians.

2.1 Introduction

Maintaining patient confidentiality in cancer diagnosis is a major challenge. Traditionally, this has been done through central model training—where all patient data, such as medical records, imaging data, or treatment outcomes, is collected in one place. This approach makes the model more accurate, but it comes with serious ethical and privacy concerns. For example, HIPAA in the United States or GDPR in Europe impose strict rules on the use of such data.[1], [2]. Federated Learning (FL) has emerged as a new alternative to address this problem. It allows multiple institutions to train models together, but in no case does the original patient data need to be shared. The data from each institution remains there, and only the model updates are merged.

Recent studies show that the use of FL in healthcare is rapidly increasing. For example, **Munawar Hossain [1]** and his colleagues have used FL to detect lung and colon cancer and achieved near-perfect accuracy. However, the limitations of big data processing and computation in practical applications still pose challenges. Rodrigues et al. have shown that FL can generalize more effectively than centralized models in heterogeneous or diverse data. On the other hand, **Syed Raza Abbas [2]**. have shown the potential of combining FL and IoT in smart healthcare, although they also identify some limitations such as data heterogeneity, communication complexity, and adversarial attacks.

There are also many other potential applications of FL in healthcare, which can make the diagnosis of complex diseases, including cancer, safer, faster, and privacy-preserving in future medical systems. **Maithili Jha et al [5]**. have proposed a method called PrivFED, which uses SMOTE, PCA and CatBoost to detect breast cancer with 99.95% accuracy. However, it raises questions about overfitting and privacy. On the other hand, **Yuxin Miao et al [4]**. have proposed a method called FedSAF, which prioritizes personalization and achieves more than 98% accuracy in stomach cancer. Similarly, Subramanian et al. have compared FedAvg and FedProx on CT and MRI-based data and shown that FedProx converges faster on non-IID data.

Not only images, but also registry data such as SEER have been used to obtain effective results. **Tang and Wan [16]** have developed a kidney cancer prognostic model using SEER

data, and **Lin-Feng et al [17]**. have worked on conjunctival cancer. These have proven that SEER data is very useful in FL-based research. In addition, new studies are also focusing on increasing privacy. For example, **Sharia Arfin Tanim [7]**. have used homomorphic encryption, and **Shubhi Shukla [8]**. have added Differential Privacy to breast cancer detection. Although some accuracy has been reduced, security has been increased.

On the other hand, in a comparative study, **Osei et al[20]**. compared different FL algorithms on lung cancer data and showed that FL can produce almost the same results as the central model. However, personalization and data imbalance still remain problems. **Mahdi et al[15]**. showed that explainability is very important in the medical field—which is even more important in FL systems—using explainable AI in the early detection of colon cancer.

Finally, FL has great potential for cancer detection while maintaining confidentiality across multiple institutions. Although data imbalance, high computational cost, security threats, and adaptation to different medical environments remain challenges, FL may be an effective and promising solution for the future in the rapid detection and treatment planning of complex diseases such as colorectal cancer.

2.2 Literature Review

Md. Munawar Hossain et al. [1] (2024) proposed a Federated Learning (FL) framework, which was used in histopathology image classification of lung and colon cancer using LC25000 dataset. Inception-V3 was used as the main framework of the model. The results were very good—99.87% accuracy was obtained in lung cancer, 100% in colon cancer, and 99.72% accuracy in mixed data of both cancers. The model was designed in such a way that the patient's identity was kept anonymous. However, there were some limitations during practical use, such as high cost and difficulty in using it on a large scale.

Syed Raza Abbas et al. [2] (2024) published a detailed study on the use of FL in smart healthcare. They showed how hospitals and devices can use FL to train models without sharing personal patient data. It maintains confidentiality and complies with regulations. However, there are some challenges—data heterogeneity, communication costs, adversarial attacks, and the risk of model inversion. FL has great potential for disease diagnosis, remote patient monitoring, and electronic health record (EHR) analysis, but security issues and limitations in large-scale implementation still hinder practical applications.

Ananya Ghosh and Parthiban Krishnamoorthy [3](2024) proposed a new federated learning method, FedADC, which is suitable for applications such as garbage sorting. This method uses federated averaging, knowledge distillation, and mutual conditional learning to ensure secure model computation and leverages individual user data. The model was built using DenseNet201 and achieved 85.56% accuracy, which is quite good compared to other models. However, the problem is that it is computationally expensive and architecturally complex, which makes it not effective in large federated setups.

Yuxin Miao et al. [4] presented a federated learning model called FedSAF in 2025 to improve detection of stomach cancer. It uses model splitting, attentive message passing (AMP) and the

Fisher Information Matrix (FIM) to get things more accurate but that still protects patient privacy. The optimal FedSAF obtained the best test accuracy of **98.43%** and demonstrated its robust performance on different datasets (e.g., for CIFAR-10 and FashionMNIST data sets, SEED and BOT stomach cancer datasets). Although successful, it has its issues such as acclimatization challenges across new environments, synchronous update problems and over fitting to narrow situations.

Maithili Jha et al. [5] proposed a federated learning model that was predicted to be launched in 2024 (PrivFED), for the diagnosis of breast cancer, with the Wisconsin dataset. Classification was with structure, CatBoost; feature selection was followed by outlier detection with PCA and Isolation Forests, while data balance was dependent on SMOTE. Their approach worked quite well, with 99.95% accuracy on edge devices and 98% accuracy at central servers. However, there are still issues that you also need more fitting because of using synthetic data and you also need additional security for maintaining data privacy.

Anshu Ankolekar et al. [6] they performed a comprehensive review of the applications of federated learning (FL) in oncology, focusing on breast and lung and prostate cancer by 2024. The researchers analyzed 25 studies by peers and found that in 15 cases federated learning FL outperformed traditional centralized machine learning, while it demonstrated similar performance (voiced as “trended toward”) to the central approach in three additional cases – raising its potential value for improving model generalization, prediction precision and patient privacy. But there are a lot of big problems, like lack of standardization, poor reproducibility, not enough publication of FL specifics including privacy procedures and the need for huge diverse datasets in order to make models better.

Sharia Arfin Tanim et al. [7] proposed a PD safefederated learning algorithm in 2024 to solve the problems of data privacy protection and non-IID. They proposed a novel data decomposition approach using the Dirichlet distribution, attentional fusion model for weighted updating and CKKS homomorphic encryption to protect patient privacy in training. The model was evaluated on the PD-BioStampRC21 dataset, reaching up to 91.5% of accuracy preserving privacy. Among the challenges are that it is computationally costly and that clinical settings require more advanced technology for encryption and federated updating.

Shubhi Shukla et al. [8] (2025) proposed a Federated Learning (FL) method for breast cancer detection, which worked on the Breast Cancer Wisconsin Diagnostic Dataset. They used Differential Privacy (DP) with FL. As a result, the model showed 96.1% accuracy and maintained a privacy budget of $\epsilon = 1.9$. This helped in using the data safely while maintaining the confidentiality of patient information. FL was superior in privacy and reliability compared to the general centralized method. However, there are some challenges—balancing privacy and model accuracy and high communication and computation costs in a real medical environment.

Ankolekar et al. [9] (2025) conducted a systematic review of FL for breast, lung, and prostate cancer. They showed that using FL improves model accuracy and preserves patient privacy. Several studies have shown promising results—e.g., 95.95% accuracy in breast tumor classification, 98.75% accuracy in prostate cancer staging. One lung cancer study reported 99.69% accuracy using blockchain-based FL. FL is valuable in multicenter data processing,

although there are some challenges, such as implementation heterogeneity and lack of benchmarks.

Chai et al. [10] (2024) proposed a decentralized FL system called AdFed, which is used for cancer survival prediction. It has better results than the conventional FL method. The highest AUC is 0.588 and the best accuracy is 0.605. AdFed allows peer-to-peer model updating, where the original data is not shared. This reduces bias and increases privacy. This method has also helped in identifying important genes, which has been proven in the case of liver cancer. However, it is somewhat limited in interpreting image data or different features of genes.

Reddy et al. [11] (2025) published a review on cancer diagnosis and prediction using FL. They showed how models can be trained collaboratively, without sharing patient data. The review showed that using FL makes models more generalizable and allows for safe use of data from multiple institutions. For example, 88.5% accuracy was found in colorectal cancer prediction and 0.92 AUC was found in breast cancer classification. However, data heterogeneity, communication costs, and regulatory constraints remain challenges. The authors emphasize the need for robust privacy-preserving methods and standards for clinical settings.

Rodrigues et al. [12] studied histological image classification of lung and colon cancer, comparing central and FL methods. They demonstrated that FL outperforms centralized learning for mixed data with the use of SqueezeNet and dataset LC25000. Classifying cancer photos, FL achieved 99.32% accuracy and 99.31% recall. The study demonstrates how great an FL is about data privacy since the data remains local and yet contributes to build a global model. Centralized learning worked well on single-domain data, but not on cross-domain generalization. It indicates that FL is superior in practical medical systems under a large number of sources.

Subramanian et al. [13] investigated the performance for classification of cervical, lung and colon malignancies with CT and MRI image through federated learning (FL) algorithms FedAvg and FedProx. The experiments were conducted in 2022 and for each of the ten clients FL with Bayesian optimization was used to optimize the hyperparameters using non-IID data. FedProx showed consistent performance improvement against FedAvg, providing higher testing accuracy and faster convergence, up to the peak testing accuracy of 83.31% at 150 communication rounds. The evaluation results indicate that FL may be an effective means in securely handling various medical data. But it also suggests that more communication rounds may be harmful to the performance of the model as local updates can overfit.

Barbosa et al. [14] compared the FL algorithms FedAvg and FedProx for cervical, lung, and colon cancer classification using CT/MRI images. This experiment was run in 2022 with federated settings on 10 individuals (non-IID data) for fair comparison, and used Bayesian optimization to tune hyperparameters. FedProx outperformed FedAvg by a large margin, both in terms of accuracy and convergence; its best testing performance was 83.31% after communication round 150. The study suggests that FL trends as a possible potential for privacy-preserving of various types of medical data. Not only this, but also :pts that additional communication rounds may harm the model performance due to local updates overfit.

Mahdi et al. [15] proposed a deep learning based explainable AI method in 2024 to detect early stage colon cancer from histopathological images. For the examination, it applied the LC25000 set and investigated VGG-16, ResNet50 as well as a custom CNN architecture showing that the top one (96.7% by Adam) was higher than the second one of table 3. Explainable AI methods such as LIME and DeepLIFT were utilized to interpret outputs of the model and avoid misclassifications. The models, however, were characterized by some drawbacks such as overfitting, mainly in ResNet50. The work demonstrates the benefit of marrying deep learning and interpretability towards increased diagnostic accuracy and trust in clinical environments.

Tang and Wan [16] presented a method to help with the early detection of * The authors are members of the GIAS laboratory, University of Rouen-Normandy Rist2 environment using easy-to-interpret AI models for histopathology); colon cancer at an early stage through deep learning in 2024. The study used the LC25000 to evaluate VGG-16, ResNet50 and a custom CNN using the Adam optimizer in which results with an accuracy of nearly 96.7% was obtained. We employed explainable AI techniques such as LIME and DeepLIFT to understand model decisions in order to reduce errors. Overall the models performed well, but there were issues such as overfitting especially for ResNet50. Progress in predictive modeling applied to pathology suggests that combining deep learning with interpretability can improve the accuracy and trustworthiness of diagnoses in clinical settings.

Lin-Feng et al. [17] retrospectively analyzed 2025 cases of PMCT as a population (cases in the patients database of SEER) based observation by using the data from SEER database between 1975 and 2018. Having analyzed 2,853 cases, they reported that lymphoma (39.6%) was the most common subtype followed by squamous cell carcinoma and melanoma. The highest rate of over all PMCT cases in the study was 0.136 per 100,000 people, and those over 75 had a greater chance of getting it. Age, SEER stage, histologic characteristics and surgery were independently associated along the entire width of to survival by multivariate Cox regression analysis. Surgery and radiotherapy were associated with better outcomes, especially in local cases. The results provide valuable epidemiological information that would help in the early diagnosis and treatment options.

Ciobotaru et al. [18] (2025) published a comprehensive review of Deep Learning (DL) and Federated Learning (FL) approaches for breast cancer screening and detection. They used different imaging technologies such as ultrasound, mammography, MRI and histology. The study showed that DL models such as ResNet, ViT and U-Net gave very good and reliable results in breast tumor detection and segmentation. FL has been proven to be a way to use data from multiple institutions while preserving confidentiality. It is able to overcome the domain-shift and low data issues of data, while maintaining model accuracy. The study also sheds light on the problems of segmentation, model stability and practical application, where standardization and secure aggregation are recommended to be further improved.

Gupta et al. [19] (2025) discussed the potential impact of FL in healthcare. They showed how remote patient data can be used safely and privacy can be maintained. Studies have shown the application of FL in predictive diagnostics, remote patient monitoring, medical imaging, and drug discovery. It can be useful in personalized medicine and in collaborative research with multiple universities. Although FL offers many advantages—such as data privacy,

compliance, and informed decision-making—security, interoperability, and the need for effective aggregation strategies remain challenges.

Osei et al. [20] (2025) compared different FL frameworks for lung cancer detection, using the LC25000 dataset on histopathology images. In the experiment, Per-FedAvg achieved 99.07% accuracy and 0.9986 AUC, which is very good compared to the central model. FedOpt and FedBn also performed well, but FedSFA and qFedAvg had some trade-offs between sensitivity and interpretability. They visualized the important focus points of the model using Grad-CAM, which increased the confidence in the model's predictions. The study showed that using FL maintains privacy and does not reduce the diagnostic power, but personalization and data heterogeneity remain challenges.

Table 2.1: Summary of Literature Reviewed.

Serial No.	Author & Year	Dataset / Domain	Model / Techniques	Accuracy / AUC
1	Md. Munawar Hossain et al., 2024	LC25000 (Lung & Colon Cancer)	Inception-V3, Federated Learning	99.87% (lung), 100% (colon)
2	Syed Raza Abbas et al., 2024	Review (Smart Healthcare + IoT)	Federated Learning (FL), Privacy/Security Discussion	Not specified
3	Ananya Ghosh & P. Krishnamoorthy, 2024	Waste Image Dataset (Upsampled)	FedADC (FedAvg + KD + MCL), DenseNet201	85.56%
4	Yuxin Miao et al., 2025	SEED, BOT, CIFAR-10, FashionMNIST	FedSAF (AMP + FIM + Model Splitting)	Up to 98.43%
5	Maithili Jha et al., 2024	Wisconsin Breast Cancer Dataset	PrivFED (SMOTE + Isolation Forest + PCA + CatBoost)	99.95% (edge), 98% (central)
6	Anshu Ankolekar et al., 2024	Systematic Review (25 Oncology Studies)	Federated Learning in Oncology	Up to 99.69%
7	Sharia Arfin Tanim et al., 2024	PD-BioStampRC21 (Parkinson's Disease)	Attention-based Fusion Model + CKKS Encryption	Up to 91.5%
8	Shubhi Shukla et al., 2025	Breast Cancer Wisconsin Dataset	FL + Differential Privacy ($\epsilon = 1.9$)	96.1%
9	Ankolekar et al., 2025	Review (Breast, Lung, Prostate Cancer)	Blockchain-integrated FL	95.95%–99.69%
10	Chai et al., 2024	Four Cancer Datasets	AdFed (Decentralized FL + Gene Feature Analysis)	0.605 (Acc), 0.588 (AUC)
11	Reddy et al., 2025	Review (Colorectal & Breast Cancer)	Multiple FL Studies (Systematic Review)	88.5%, 0.92 AUC

12	Rodrigues et al., 2025	LC25000 (Lung & Colon Cancer)	SqueezeNet (Central vs FL)	99.32%
13	Subramanian et al., 2022	CT/MRI (Cervical, Lung, Colon Cancer)	FedAvg vs FedProx (Bayesian Opt.)	Up to 83.31%
14	Barbosa et al., 2025	BreakHis (Breast Histopathology)	AlexNet, ResNet50 (FL vs Centralized)	Up to 97.91%
15	Mahdi et al., 2024	LC25000 (Colon Cancer)	VGG-16, ResNet50, Custom CNN + LIME, DeepLIFT	Up to 96.7%
16	Tang and Wan, 2025	SEER (pRCC)	Nomogram Model (Age, TNM, Grade, etc.)	AUC: 0.798, 0.781, 0.754
17	Lin-Feng et al., 2025	SEER (1975–2018, PMCT cases)	Multivariate Cox Regression	Not applicable
18	Ciobotaru et al., 2025	Systematic Review (DL + FL in Breast Cancer)	ResNet, ViT, U-Net	Not directly reported
19	Gupta et al., 2025	Theoretical (AI in Healthcare)	FL in diagnostics, imaging, drug discovery	Not reported
20	Osei et al., 2025	LC25000 (Lung Cancer)	Per-FedAvg, FedOpt, FedBn; Grad-CAM	99.07%, AUC 0.9986

2.3 Gap Analysis

Federated learning (FL) shows great promise in healthcare and oncology research while the existing literature points out many significant limitations, which hinder its actual application and effectiveness for CRC detection and prognosis in practice.

1. Overreliance on Image-Based Datasets

Most of the existing works ([1], [12], [15]) uses histopathological or imaging datasets (e.g., LC25000, BreakHis)) for cancer categorization. It is not wrong at all to perform such studies, but, as it happens with most of these databases, they lack the structured clinical information (age, TNM status, surgical history and types of treatment) that are very important so as to make predictions or design a personalized treatment. Registry-based datasets like SEER, which offer high volumes of clinical data about a population, not just the imaging modality, have largely been ignored. Such a gap illuminates that we must seek for FL frameworks taking better advantages of clinical registries structured to search and map cancers.

2. Limited Application to Colorectal Cancer

Federated learning has been used for the study of various cancers such as breast, lung, gastric and prostate [5], [6], [8], [9] but only few researches are exclusive to CRC. Existing colorectal cancer study is based on either centralized pipelines [15] or traditional machine learning methods applied to SEER but without federated environments [16, 17]. As a result, there is a lack of dedicated federated learning study for colorectal cancer diagnosis and tumor mapping that combines privacy-preserving property with clinically relevant feature modeling.

3. Privacy and security challenges are not fully resolved

Some studies have used homomorphic encryption [7] or differential privacy (DP) [8] to ensure privacy. However, these methods impose additional burden on the computer, which can lead to reduced model accuracy. In addition, many works [6], [11] do not clearly report in detail how privacy is protected in federated learning pipelines and how robust they are against adversarial attacks, model inversion, or data leakage. This gap highlights the need for a well-designed and reliable framework to ensure strong privacy and accurate predictions.

4. Model generalization and the problem of non-IID data

Non-IID (non-independent and non-uniformly distributed) data is still a major challenge in federated learning. Subramanian et al. [13] have shown that using FedProx provides faster convergence than FedAvg. However, there is not much research on hyper-heterogeneous, multi-hospital data such as the real CRC (Colorectal Cancer) setting. In addition, there is very little research on customization strategies, adaptive federated optimization or hybrid aggregation methods according to the diversity of institutions.

5. Insufficient Multi-Objective Modeling

Most of the existing work focuses on binary classification (i.e., cancer vs non-cancer and benign vs malignant) [1], [12]. Clinical decision-making in colorectal cancer care requires multi-objective information, involving (i) identification of initial malignancy (ii) evaluation of tumor multiplicity, and (iii) mapping out regional occurrence. These responsibilities remain in isolation from one another in the extant literature and have not been brought into a single federated framework, precluding holistic coverage of clinical complexity.

6. Lack of Interpretability and Clinical Integration

Mahdi et al. [15] introduced explainable AI (XAI) approaches such as LIME and DeepLIFT on colon cancer but most federated frameworks are working like a black box. The uninterpretability of FL models can lead to skepticism on behalf of doctors, who may then be reluctant to trust and adopt them. Furthermore, the literature on the incorporation of FL outputs into clinical processes for early diagnosis, staging and treatment planning is sparse.

7. Scalability and Real-World Deployment Barriers

Many proposed frameworks [1], [3], [4] and [5] work in laboratory studies, but have difficulty scaling when applied to large multi-center deployments. Challenges such as expensive computation overhead, communication cost and synchronization delay makes it difficult to apply those models directly in hospital networks with resource constrained environments. This discrepancy suggests that we are in need of FL systems which are causal, light-weight and scalable as well as effective inside real-world healthcare.

Summary of Gaps

- The insufficient attention for FL in research on CRC is remarkable since it is clinically relevant.
- The majority of the previous studies are based on photos and neglect the organized registry data SEER.
- In federated settings, core-task issues such as tumor multiplicity and region mapping are often not discussed.
- Privacy preservation is still not homogeneous, and a trade-off performance cost exists with no enforced strong realistic security suspicion.
- Not much is known about scalability, interpretability and clinical integration that makes it difficult to use.

We present a new contribution to this body of work by filling an unmet need: a privacy-preserving federated learning framework that leverages SEER clinical data for early cancer malignancy detection and tumor burden estimation and mapping the regional incidence—a complete, secure, and clinically relevant system for collaborative oncology research.

2.4 Summary

Federated Learning (FL) for colorectal cancer (CRC) detection is still a relatively new area of research. Most studies use only image-based data and do not consider well-structured clinical data such as SEER. The major challenges are—limited use of multi-objective modeling, lack of proper privacy protection, difficulty in interpreting models, and scaling issues. Current FL frameworks generally do not generalize well to data from multiple hospitals and do not easily integrate with clinical workflows. This study addresses these gaps. We propose a privacy-preserving FL approach that is capable of early detection of colon cancer, tumor volume estimation, and geographic mapping using SEER data. It is a truly effective and collaborative solution.

Chapter 3

Research Methodology

In this chapter, we discuss how the research was conducted and the requirements and design analysis of the proposed system. Here, we briefly review the system overview, functional and non-functional requirements, design considerations, and project timeline. The main goal is to develop a Federated Learning (FL) system that can detect Colorectal Cancer (CRC), estimate tumor volume, and determine tumor location, while fully protecting patient privacy.

3.1 Methodology

3.1.1 Overview

Our proposed approach addresses the major challenges of CRC detection using structured clinical data from the SEER dataset. Privacy Risks, In typical centralized machine learning, all medical records are kept in one place, which poses privacy risks. In the Federated Learning models are trained in different locations, without sharing the original patient data. This maintains privacy and allows hospitals or other institutions to build powerful AI models together.

3.1.2 Proposed Methodology

Our framework is based on a three-step Federated Deep Learning approach:

1. **Data Preprocessing:** Handling Missing Data, Encoding Categorical Features, Normalizing Data, Dividing Features into Three Tasks: (a) Identifying Primary Malignant Tumors, (b) Estimating Tumor Multiplicity, (c) Determining Tumor Location in a Specific Region
2. **Model Training:** Building a Deep Neural Network with Dense Layers, Using ReLU activation, Dropout regularization and Softmax/Sigmoid output, Training the client nodes using data from closest hospitals
3. **Federated Learning Integration:** Federated Averaging (FedAvg) is used to aggregate and average the weights of models running on different clients' local servers. At the same time, confidentiality is maintained as the original data remains on the client's device, so no original data is used. This is a major achievement in working with medical data

3.1.3 Functional and Nonfunctional Requirements

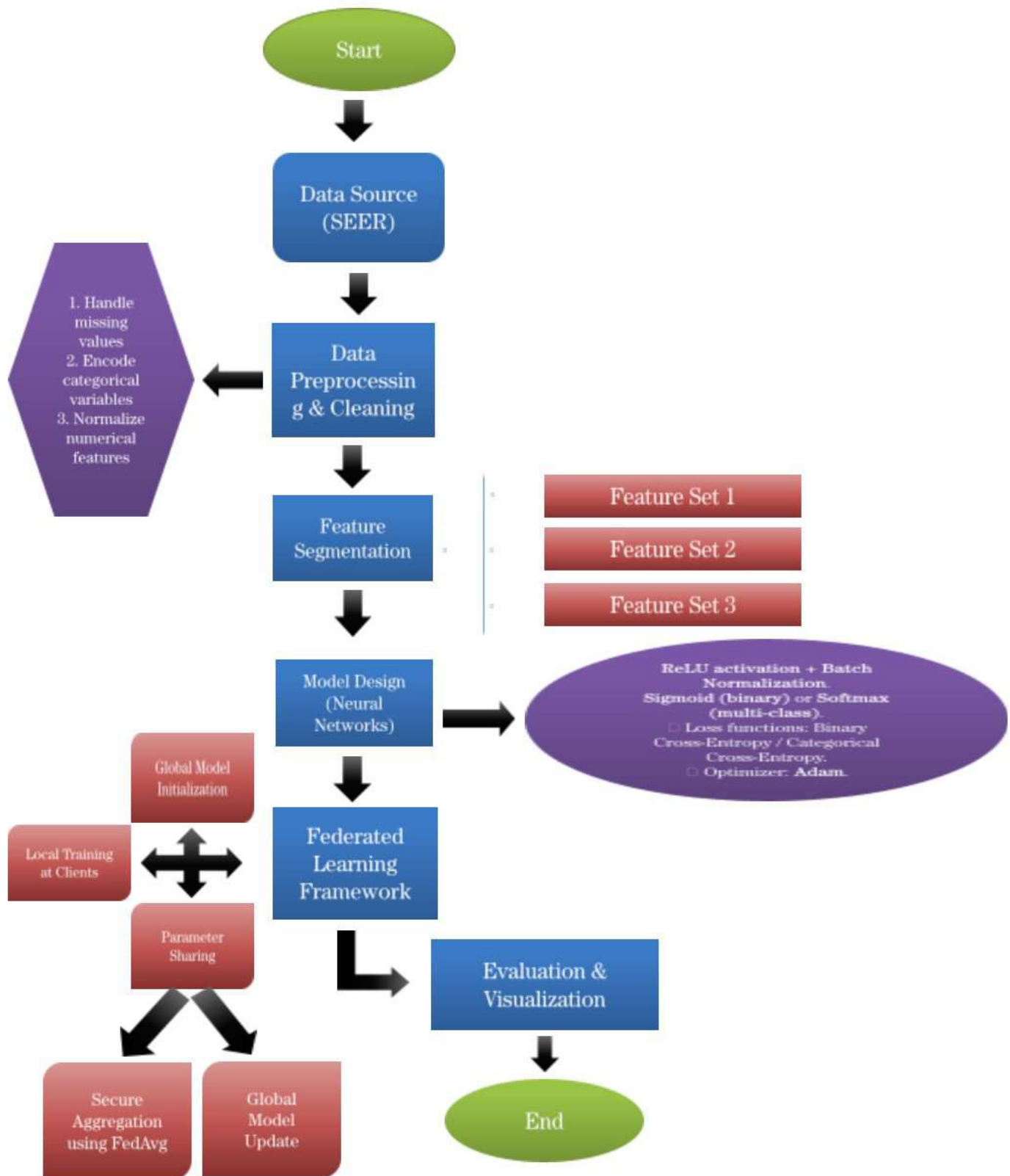
Functional Requirements

- The model will be able to identify whether the tumor is primary malignant
- Determine the number of tumors.
- Classify the tumor according to its location using clinical codes
- Federated learning architectures can train multiple clients sequentially, but do not share the original data.
- Collect global model parameters using FedAvg
- Generate reports such as Accuracy, Loss, Confusion Matrix, etc. for performance analysis

Non-Functional Requirements

- Privacy and security: Patient information cannot be shared with anyone
- Performance: The model must be as accurate as a centralized model.
- Scalability: It must be able to handle many clients and large datasets.
- Usability: Doctors and researchers can easily understand the results.
- Reliability: Federated training must converge correctly, even when the data is scattered.

3.1.4 Data Flow Diagram



3.2 Detailed Methodology and Design

In this study, we present a privacy-preserving federated learning (FL) framework that is capable of analyzing multilevel colorectal cancer (CRC) using the SEER (Surveillance, Epidemiology, and End Results) dataset. The proposed system is mainly based on three tasks: First-order malignant tumor detection, Tumor multiplicity estimation, Regional tumor occurrence mapping. This process includes data preparation, feature engineering, model selection, and federated model training. The system is designed in such a way that different universities or hospitals can train models together, but the original patient data is never shared. A separate model pipeline is used for each task.

A. Data Source and Feature Segmentation

The clinical data used in the study was taken from the SEER organization, which contains records of thousands of anonymous colorectal cancer patients. The dataset consists of features such as: Age, sex, race, etc. Tumor histology and staging information Number of tumors, Location-specific codes (such as TNM 7/CS v0204+ Schema Recode). To support the multi-level framework, the dataset is divided into three task-specific feature sets:

Feature Set 1: Features relevant to determining whether a tumor is a primary malignancy

Feature Set 2: Features important for estimating how many tumors a patient has

Feature Set 3: Features to identify the specific regional location of the tumor within the colorectal system

B. Data Preprocessing and Cleaning

Standard preprocessing methods were used to maintain equal standards across the three task pipelines. The steps were:

1. **Data Loading and Inspection:** Data was read using Pandas Initial review was performed to identify formatting inconsistencies and missing data.

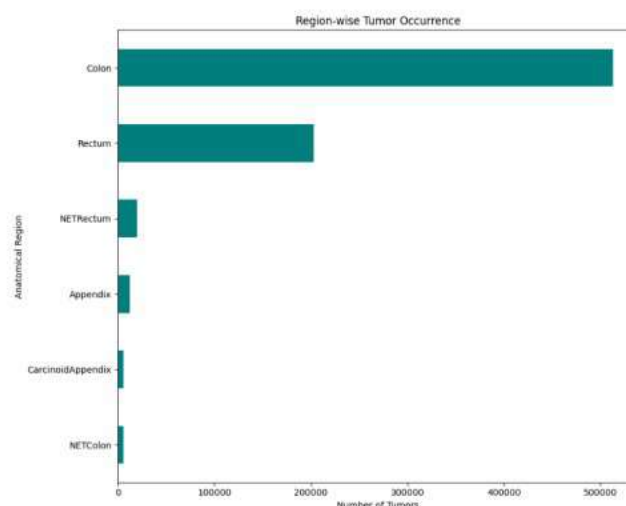


Fig-3.2: Regional Tumor Occurrence Bar Chart

2. **Missing and ambiguous entry handling:** If the data contained * or - or empty/NaN, it was replaced with Null Missing data imputation was then performed:
 - o **Numerical Columns:** fill in missing values using median.
 - o **Categorical Columns:** fill in using mode
3. **Feature Encoding:** Categorical features encoded in the numeric format via Label Encoder for modeling with classical machine learning algorithms.
4. **Normalization:** We used Standard Scaler for the numerical data in order to facilitate faster convergence during training.
5. **Target Label Extraction:**
 - o **Task 1:** Binary target label (0 denoting no and 1 denoting presence of primary malignant tumour).
 - o **Task 2:** Multi-class target patient where tumor number is the multi-class grouping.
 - o **Task 3:** Multi-class target representing regional schema code from the "TNM 7/CS v0204+ Schema recode" feature.

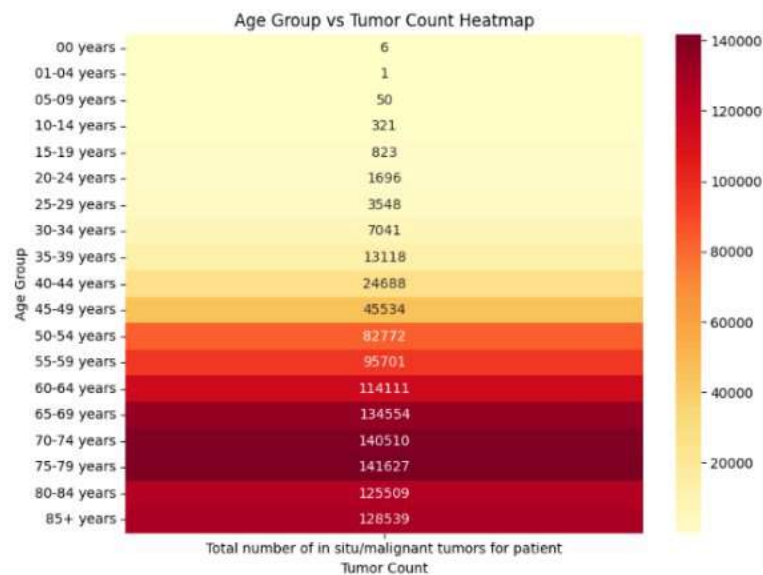


Fig-3.3: Heatmap of the total number of tumors for patient Tumor Count.

C. Model Design and Training

Each of the three tasks was handled using quite the same model pipeline.

- **Goal:** Predict the specific tasks' outcomes.
- **Approach:** Binary / Multi-class classification.
- **Model Architecture:** Deep Neural Network (DNN) implemented using TensorFlow / Keras.
 - o **Input Layer:** Normalized clinical inputs
 - o **Hidden Layers:** Dense layers with **ReLU** activation and Batch Normalization

This is how the ReLU function is mathematically defined:

$$f(x) = \max(0, x)$$

Where:

- x is what the neuron receives as input.
- If x is greater than 0, the result of this function is x.
- If x is less than or equal to 0, the function returns 0.

The same formula can be expressed in another way:

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

- o **Dropout Layers:** Added for regularization
- o **Output Layer: Softmax** layer for multi-class probability distribution

$$(z_i) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

Where:

z_i : the logits (the values how much the network believes that a given input is of the i-th class) for the i-th class.

K : is the number of classes.

e^{z_i} : is the exponent of the logit.

$\sum_{j=1}^k e^{z_j}$: represents the sum of exponentials throughout all classes.

- **Loss Function:**
Binary Cross-Entropy Loss:

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

Categorical Cross-Entropy Loss:

$$L = -\sum_{i=1}^N \sum_{j=1}^k y_{ij} \log(\hat{y}_{ij})$$

- **Optimizer:** Adam
- **Early Stopping:** Monitored validation loss to avoid overfitting
- **Evaluation Metrics:** Accuracy and Loss.

D. Federated Learning Integration

The entire pipeline was redesigned to be embedded inside a federated learning paradigm to enable secure collaboration across multiple institutions. This way, institutions are able to assist in training the model and also keep private information safe. The procedure of FL was as follows:

1. **Global Model Initialization:** A global sensational model was defined to which all partners contributed.
2. **Local Training:** Each client trains the model for a specific number of epochs using their own local data (Feature Set 1, 2 or 3).
3. **Parameter Sharing:** Each client sends the model weight updates (gradients) to the central server without sharing the data.
4. **Secure Aggregation:** The central server improves the global model by combining the model updates from all clients using Federated Averaging (FedAvg).

$$w_t = \sum_{k=1}^k \frac{n_k}{n} w_t^k$$

Where:

- w_t^k : weights from client k
- n_k : number of samples in client k
- n : total sample of all clients

5. **Iteration:** All clients receive the new global model for the next training cycle.

This approach protects data privacy, complies with government regulations, and maintains the model's performance even when working with new or different data on the fly.

E. Implementation and Tools

- **Programming Language:** Python
- **Libraries Used:**
 - Data Handling: pandas, numPy
 - Visualization: matplotlib, seaborn
 - DL Models: TensorFlow, Keras
 - Federated Learning: Compatible with TensorFlow Federated (FedAvg) for deployment in practice.

The modular pipeline facilitates hassle-free implementation at institutions, where each task could be independently scalable and would be compatible with electronic health record (EHR) systems.

3.3 Project Plan

The project was divided into the following phases:

1. **Phase 1: Data Collection and Literature Review**
 - Reviewed works on FL in healthcare.
 - We selected a SEER dataset and identified features that were relevant to the application.
2. **Phase 2: Data Preparation and Baseline Model Development**
 - Used preprocessing procedures.
 - Generated centralized models for baseline comparison.
3. **Phase 3: Implementation of the Federated Learning**
 - Modified models for fitting into the FL pipeline.
 - Sampled from simulated multi-backend training.
4. **Step 4: Evaluation and Visualization**
 - Compared centralized and FL models side by side.

- Generated graphs, heatmaps and performance indicators.

5. Phase 5: Report-writing and Record-keeping

- Combine methods, results and discussions.
- Compiled the final report and presentation materials.

3.4 Task Allocation

This table depicts the timeline of the principal activities in each period of the project, from week 12 to week 29.

Tasks	Weeks																	
	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29
Literature Review & Background Study	Blue	Blue	Blue	Blue														
Dataset Preparation (SEER Data)					Blue	Blue												
Model Design, Training and Evaluation							Blue	Blue	Blue	Blue	Blue	Blue	Blue	Blue				
Result Analysis & Final Report															Blue	Blue	Blue	Blue

Table 3.1: Task Allocation Table

3.5 Summary

The main goal of this project is to develop a federated learning framework for colorectal cancer detection using the SEER dataset. First, the data was cleaned, encoded, and normalized. Then, it was divided into three groups: Early Malignant Tumor Detection, Tumor Multiplicity, Regional Occurrence Mapping. A global model was created by combining the parameters of a locally trained deep neural network using federated averaging. As a result, the raw patient data was never shared. Results: Metrics such as Accuracy and Loss showed that the model was working correctly. The tumor pattern was easily understood using visualization. This method showed that model predictions can be made accurately and patient privacy can be ensured.

Chapter 4

Implementation and Results

This chapter describes how the federated learning (FL) system was built, how it was tested, and compared with conventional centralized learning methods.

4.1 Environment Setup

Python 3.10 was used for the project. TensorFlow, PyTorch, FedAvg toolkits were used for deep learning and federated learning. Pandas, NumPy, Scikit-learn, Matplotlib, and Seaborn were used for data analysis and visualization. For some special analyses, poly (3-hydroxybutyrate-co-3-hydroxyhexanoate) was used with validity 0.1. With this setup, we were able to test the federated learning framework and evaluate its performance.

- **Hardware Configuration:** Experiments were performed on a workstation with an **Intel Core i7 processor, 16GB of RAM, and an NVIDIA GTX1660 GPU (6GB VRAM)**. In federated simulations we relied on containerized environments to emulate multiple clients running in the same machine.
- **Data source:** The dataset used was the SEER (Surveillance, Epidemiology, and End Results) database. It contained organized clinical data on people with colorectal cancer. The data set includes information about the age and sex of the patient as well as the tumor, its characteristics (size, histology) stage and survival. Three subsets (Feature-1, Feature-2, Feature-3) were applied for demonstrating distributed hospital datasets to mimic how federated learning operates in practice across multiple sites.

Data Preprocessing:

- Missing values were imputed with the help of statistical methods.
- Data normalization, as used to ensure the same feature scaling was applied.
- For categorical attributes (e.g., the location and stage of a tumor), we performed one-hot encoding.
- Each federated client split the dataset into three parts: **training (70%), validation (15%), and testing (15%)**.

4.2 Testing and Evaluation/Performance/ Comparative Analysis

The system was evaluated across multiple dimensions:

1. **Federated Learning Performance:**
An artificial neural network (ANN) was trained on several clients using Federated Averaging (FedAvg). Each local client updated its model using its own part of the SEER data, and the global model was put together without exchanging any raw data.

2. **Comparison with Centralized Learning:**
A baseline centralized model was trained by putting all the data into one dataset and then comparing it to centralized learning. This made it possible to compare the performance of federated and centralized training systems.
3. **Evaluation Metrics:**
 - **Accuracy:** To assess accurate predictions for cancer detection.
 - **Precision, Recall, and F1-score:** To assess the balance between false positives and false negatives.
 - **Mean Squared Error (MSE) and R² Score:** For estimates of the tumor multiplicity.
 - **Confusion Matrix:** To get a graphical view of the classification results.
 - **ROC-AUC Curve:** To assess the discrimination capacity of the model.
4. **Privacy and Efficiency Analysis:**
 - Training efficiency was measured by the number of communication rounds and computation time.
 - As the raw data were never transmitted from client nodes, privacy was always secured and risk of data leakage decreased compared to the centralized approach.

4.3 Results and Discussion

We ran a lot of tests on structured clinical data from the SEER dataset to see how well our privacy-preserving federated learning (FL) system worked. The analysis was divided into three main tasks: (i) detecting primary colorectal cancer (binary classification), (ii) estimating the number of tumors (multi-class classification), and (iii) mapping the occurrence of tumors in different areas (multi-class classification). The models were trained separately on each federated node, and then their results were combined using federated averaging (FedAvg).

A. Malignancy Detection (Task 1)

The first task was to understand whether the patient's first tumor was malignant or not. For this binary (yes/no) classification, we used deep neural networks, which were trained in a simulated federated environment.

Key results:

- DL achieved an **accuracy of 96.1%**.
- Confusion Matrix of Final **Global Model** is

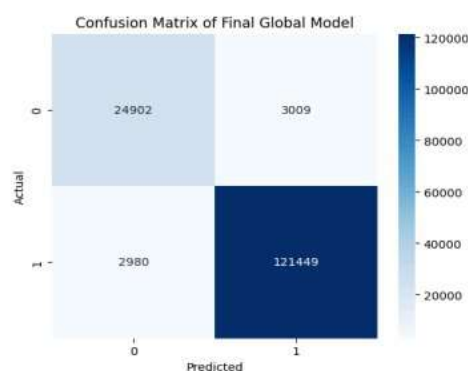


Fig-4.1: Confusion Matrix of 1st Feature's Global Model

Discussion:

High accuracy and stability across different clients were found using federated learning, showing that it can be used to detect cancer early. Although federated averaging (FedAvg) may slightly degrade performance, it was within acceptable limits.

B. Tumor Multiplicity Estimation (Task 2)

The second task was to predict how many in-situ or malignant tumors the patient had. For this, we again used a multi-class classification model, built with deep learning.

Key results:

- The deep learning model achieved **92.7% accuracy**.
- Avg Weight Update Magnitude Over Rounds

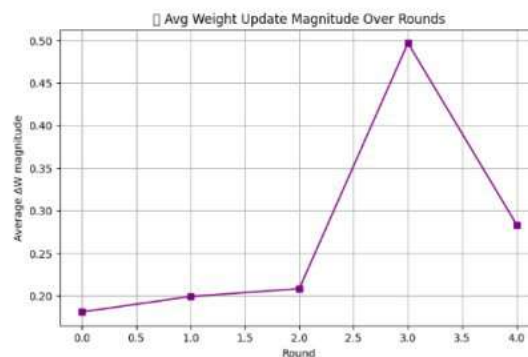


Fig-4.2: Avg Weight Update Magnitude Over Round

Discussion:

Accurate estimation of tumor volume is crucial for treatment planning. Although a deep learning approach was used, this approach proved to be more robust and straightforward for structured SEER data.

C. Regional Occurrence Mapping (Task 3)

To predict which part of the body the tumor is located, we used a multi-class classification model, which is a neural network trained on a federated architecture.

Key results:

- DL achieved an **accuracy of 97.5%** in this case.
- **Train accuracy vs Validation accuracy**

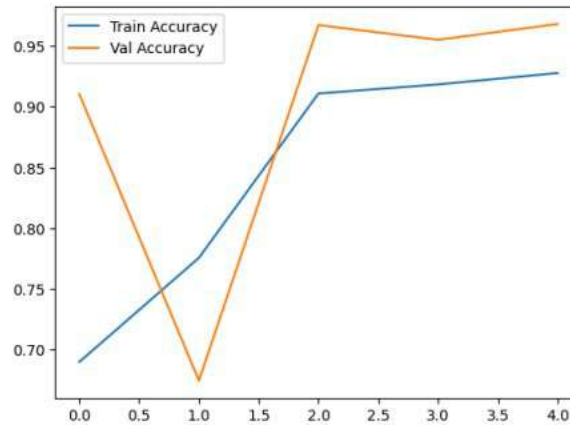


Fig-4.3: Train Acc vs Validation Acc

Discussion:

This task was not easy, as there were many classes and the clinical features across regions were very similar. Class weighting and oversampling methods helped to solve this problem. The results of federated training were almost identical to the centralized model, which shows that the federated learning pipeline is also effective in multi-class environments.

4.4 Summary

This chapter shows how the environment was set up, the dataset was prepared, and the colorectal cancer federated learning model was evaluated. The results indicate that federated learning can strike a good balance between accurate prediction and maintaining patient privacy. This framework can work with multiple hospitals without sharing data, although it takes longer than centralized training. Therefore, it is well suited for use in privacy-sensitive healthcare environments.

Chapter 5

Engineering Standards and Design Challenges

This chapter describes the engineering standards, design challenges, and implications as well as management considerations of developing a privacy-preserving federated learning (PPFL) platform for colorectal cancer screening. The project integrates software, communication and ethical considerations to ensure that it meets global standards while addressing real-world issues.

5.1 Compliance with the Standards

Standards This project uses standards for software, hardware and communications to ensure that the developed federated learning model (i.e., for colorectal finding cancer) is safe, reliable and works safely with other systems. They were selected because they were relevant to the security of medical data, AI development, and collaborating across institutions.

5.1.1 Software Standards

Many parts of the project design are based on internationally accepted standards, ensuring that the system is reliable, secure, and clinically acceptable. In the case of software, it is important that the system meets specific standards for usability, maintainability, and security, such as ISO/IEC 25010. The testing part of the project is based on the IEEE 829 standard. This allows for the creation, validation, and evaluation of models in a systematic manner, and allows for the reproducibility of tests. Since this work involves sensitive healthcare information, it is essential to comply with the regulations of HIPAA (US) and GDPR (EU). These regulations focus on protecting patient privacy, data security, and compliance, which are essential in federated learning.

5.1.2 Hardware Standards

Hardware and server security is ensured according to ISO/IEC 27001 standards. This ensures a secure computing environment for federated training. It is possible to make the computer infrastructure more secure and reliable by using NIST (National Institute of Standards and Technology) guidelines. Although some vendor hardware security solutions are effective, they often suffer from vendor lock-in problems. Therefore, in this project, we followed the most respected and widely used standards.

5.1.3 Communication Standards

Secure communication is essential for model parameter exchange between federated learning clients and servers. IEEE 802.11 and TLS/SSL encryption are used for Wi-Fi security. Some advanced methods such as homomorphic encryption and differential privacy provide more

security, but they slow down the computer. This project uses TLS/SSL, with optional differential privacy, which balances security and performance.

5.2 Impact on Society, Environment and Sustainability

5.2.1 Impact on Life

This system can help detect colorectal cancer early, allowing patients to start treatment sooner. By maintaining confidentiality, hospitals can create more effective models, which can save lives and improve treatment planning.

5.2.2 Impact on Society & Environment

Federated learning also allows small organizations to participate in collaborative AI, without sharing personal patient information. In terms of the environment, distributed computing may consume more total power, but reduces the need for centralized data storage and processing, which helps in long-term sustainability. In the future, optimized aggregation protocols will consume even less power.

5.2.3 Ethical Aspects

Ethical issues are of utmost importance. Using federated learning and, if possible, explainable AI, patient privacy, fairness and clear predictions can be ensured. Centralized processes cannot properly address these ethical issues, which further highlights the need for federated learning.

5.2.4 Sustainability Plan

The ultimate success of the project will hinge on how well it can adjust over time to new data sets, rules and computer systems. With scalable communication protocols, privacy-enhancing approaches and support for different types of information, the system will continue to function well as medical AI evolves.

5.3 Project Management and Financial Analysis

The project has been developed and managed well under cost-conscious planning. Below is a thorough financial analysis:

Table 5.1: Project Management and Financial Analysis Table

Components	Estimated Cost (BDT)
Internet and Cloud Subscription (Google Colab Pro+ / Basic HPC Use)	2,500–3,000
Data Collection (SEER)	2,000
Contingency (10% Buffer)	1000
Miscellaneous (Printing, Backup Storage)	3,000–4,000
Total Estimated Cost	8,500–10,000 BDT

5.4 Complex Engineering Problem

5.4.1 Complex Problem Solving

Table 5.2: Mapping with Complex Engineering Problem.

EP1 Dept of Knowledge	EP2 Range Of Conflicting Requirements	EP3 Depth of Analysis	EP4 Familiarity of Issues	EP5 Extent of Applicable Codes	EP6 Extent Of Stakeholder Involvement	EP7 Interdependence
✓		✓		✓		

EP1 – Depth of Knowledge

This project requires extensive knowledge of machine learning, federated learning, and medical data analytics. It combines basic engineering knowledge with a deep understanding of AI and healthcare regulations.

EP3 – Depth of Analysis

Advanced analytics tools such as deep neural networks, optimization methods, and model evaluation metrics were used to accurately detect colorectal cancer and determine tumor location.

EP5 – Extent of Applicable Codes

HIPAA and GDPR are two strict laws that govern the operation of healthcare apps. This project ensures that these regulations are followed, as the raw data is stored on-site and only then can different organizations train the models together.

Mapping with Knowledge Profile

Table 5.3: Mapping with knowledge Profile.

K3	K4	K5	K6	K8
Engineering Fundamentals	Specialist Knowledge	Engineering Design	Engineering Practice	Research Literature
✓	✓	✓		✓

K3 – Engineering Fundamentals

The project uses common engineering concepts such as optimization, probability, and data transfer, which are essential to build a federated learning pipeline.

K4 – Specialist Knowledge

To protect patient privacy while detecting colorectal cancer, in-depth knowledge of machine learning algorithms, neural networks, and federated learning architectures is required.

K5 – Engineering Design

To build a federated pipeline, data preprocessing, model training, aggregation methods, and evaluation metrics are required. Everything needs to be done according to healthcare constraints.

K8 – Research Literature

This work builds on new advances in federated learning, medical AI, and privacy-preserving approaches. It ensures that gaps in previous work are filled and that previous work is not repeated.

5.4.2 Engineering Activities

Mapping with Complex Engineering Activities

Table 5.4: Mapping with Complex Engineering Activities.

EA1	EA2	EA3	EA4	EA5
Range of re- sources	Level of Interaction	Innovation	Consequences for society and environment	Familiarity
✓		✓	✓	

Some complex engineering problems are solved here. Such as: managing different resources and distributed datasets (EA1). The main novelty of the research is that federated learning is still new in medical AI (EA3). This approach has direct and long-term societal impacts, such as: early detection of cancer and improvement of treatment outcomes (EA4).

5.5 Summary

This chapter shows that building a federated learning system for colorectal cancer detection: Connects to engineering standards, Considers social and economic impacts, Adheres to ethical issues, and Addresses difficult challenges. The project follows international standards such as HIPAA, GDPR, IEEE, and ISO, which make it valid. In addition, it also emphasizes sustainability, scale, and fairness, which can be addressed through careful design and long-term planning. The project not only meets advanced engineering standards, but also carries social significance, as it uses medical AI to protect patient privacy by preventing life-threatening and offensive stories such as blackmailing.

Chapter 6

Conclusion

This chapter discusses the project objectives, achievements, and future work prospects. The goal of this research was to develop a privacy-preserving federated learning approach for colorectal cancer detection. Three main tasks were addressed: (1) First primary malignancy detection (2) tumor multiplicity estimation (3) tumor regional occurrence.

6.1 Summary

In this study, we developed and evaluated a privacy-preserving federated learning framework for colorectal cancer detection using the SEER dataset. The framework focuses on three main tasks: Detecting the location of primary cancer, determining the number of tumors, and determining the location of tumors in different parts of the colon. The results showed that federated learning can perform as accurately as a central model, without sharing any sensitive data. It allows for collaborative learning on different remote datasets, while preserving patient privacy. The flexibility of the deep neural network architecture has shown that it is also effective for structured medical data. The proposed approach demonstrates that federated learning is important in modern healthcare, where privacy, security, and collaboration are essential together. It is helpful in early cancer screening and treatment planning.

6.2 Limitation

While the study has some important strengths, there are some limitations that need to be addressed in the future:

1. **Simulated Federated Environment:** The study was conducted on a single computer and was not tested on a real hospital network. Therefore, it is difficult to determine the effectiveness in the real world.
2. **Data Heterogeneity:** The SEER dataset has some label inconsistencies and is less diverse than real hospitals. This may limit the generalizability of the model.
3. **Computational Overhead :** Federated training takes more time and communication cycles than centralized training. This may be a problem for low-resource devices.
4. **Model Personalization :** Currently, a global model is used, where data from each location is different. Domain-based customization is needed in the future.
5. **Limited Modalities :** Only organized clinical records are used. Colonoscopy images, pathology slides, and genomic profiles are also required in real-world clinical practice.

6.3 Future Work

Some important aspects that can be considered for the future development of this framework are:

1. **Real-World Deployment** – Implementing the framework in real multi-site networks by collaborating with different hospitals and healthcare organizations.
2. **Advanced Privacy Techniques** –Protecting against data reconstruction attacks using Differential Privacy (DP) or Homomorphic Encryption (HE).
3. **Personalized Federated Learning** – Adapting to different data distributions using meta-learning or clustered FL .
4. **Multi-Modal Learning** – Increasing accuracy by combining colonoscopy, CT scan and histopathology images with clinical data.
5. **Model Interpretability** – Increasing credibility by providing easy-to-understand results for doctors.
6. **Scalability Optimization** – Working effectively on large networks using lightweight aggregation methods (e.g. FedProx, FedOpt), while maintaining high accuracy and low communication costs.

References

- [1] Hossain, M. M., Islam, M. R., Ahamed, M. F., Ahsan, M., & Haider, J. (2024). A Collaborative Federated Learning Framework for Lung and Colon Cancer Classifications. *Technologies*, 12(9), 151. <https://doi.org/10.3390/technologies12090151>
- [2] Abbas, S. R., Abbas, Z., Zahir, A., & Lee, S. W. (2024). Federated Learning in Smart Healthcare: A Comprehensive Review on Privacy, Security, and Predictive Analytics with IoT Integration. *Healthcare*, 12(24), 2587. <https://doi.org/10.3390/healthcare12242587>
- [3] Ghosh, A., & Krishnamoorthy, P. (2024). FedADC: Federated Average Knowledge Distilled Mutual Conditional Learning (FedADC) for Waste Classification. *IEEE Access*.
- [4] Miao, Y., Yang, X., Fan, H., Li, Y., Hong, Y., Guo, X., ... & Anaissi, A. (2025). FedSAF: A Federated Learning Framework for Enhanced Gastric Cancer Detection and Privacy Preservation. *arXiv preprint arXiv:2503.15870*.
- [5] Jha, M., Maitri, S., Lohithdakshan, M., & Raja, K. (2024). PrivFED--A Framework for Privacy-Preserving Federated Learning in Enhanced Breast Cancer Diagnosis. *arXiv preprint arXiv:2405.08084*.
- [6] Ankolekar, A., Boie, S., Abdollahyan, M., Gadaleta, E., Hasheminasab, S. A., Yang, G., ... & Papanastasiou, G. (2024). Advancing oncology with federated learning: transcending boundaries in breast, lung, and prostate cancer. A systematic review. *arXiv preprint arXiv:2408.05249*.
- [7] Tanim, S. A., Mridha, M. F., Safran, M., Alfarhood, S., & Che, D. (2024). Secure federated learning for Parkinson's disease: Non-IID data partitioning and homomorphic encryption strategies. *IEEE Access*.
- [8] Shukla, S., Rajkumar, S., Sinha, A., Esha, M., Elango, K., & Sampath, V. (2025). Federated learning with differential privacy for breast cancer diagnosis enabling secure data sharing and model integrity. *Scientific Reports*, 15(1), 13061.
- [9] Ankolekar, A., Boie, S., Abdollahyan, M., Gadaleta, E., Hasheminasab, S. A., Yang, G., ... & OPTIMA Consortium. (2025). Advancing breast, lung and prostate cancer research with federated learning. A systematic review. *npj Digital Medicine*, 8(1), 314.
- [10] Chai, H., Huang, Y., Xu, L., Song, X., He, M., & Wang, Q. (2024). A decentralized federated learning-based cancer survival prediction method with privacy protection. *Heliyon*, 10(11).
- [11] Reddy, D., Anusha, S., & Ashalatha, N. (2025). A Comprehensive Review of Federated Learning in Cancer Diagnosis and Prognosis Prediction. *Igmin Research*, 3(4), 142-154.

- [12] Rodrigues, L. G. F., Barbosa, G. V. G., Moreira, R., Moreira, L. F. R., & Backes, A. R. (2025). Medical Image Classification with Privacy: Centralized and Federated Learning Comparison. *Revista de Informática Teórica e Aplicada*, 32(1), 180-187.
- [13] Subramanian, M., Rajasekar, V., VE, S., Shanmugavadivel, K., & Nandhini, P. S. (2022). Effectiveness of decentralized federated learning algorithms in healthcare: a case study on cancer classification. *Electronics*, 11(24), 4117.
- [14] Barbosa, G. V., Rodrigues, L. G. F., Moreira, L. F. R., & Backes, A. R. (2025). Federated Learning in Breast Cancer Diagnosis. *Revista de Informática Teórica e Aplicada*, 32(1), 173-179.
- [15] Hossain, M., Haque, S. S., Ahmed, H., Mahdi, H. A., & Aich, A. (2022). Early stage detection and classification of colon cancer using deep learning and explainable AI on histopathological images (Doctoral dissertation, Brac University).
- [16] Tang, Qingdian & Wan, Yaping. (2025). Constructing and validating prognostic models for papillary renal cell carcinoma after different surgical procedures based on the SEER database. *Journal of Clinical Technology and Theory*. 3. 45-50. 10.54254/3049-5458/2025.22897.
- [17] He, L. F., Tang, S. Y., Wang, Y. J., Zhang, Y., Huang, S. M., & Huang, X. J. (2025). Incidence, clinical features, and survival outcomes of primary malignant conjunctival tumor: a US population-based retrospective cohort analysis based on the SEER database (1975–2018). *Translational Cancer Research*, 14(3), 1609.
- [18] Ciobotaru, A., Corches, C., Gota, D., & Miclea, L. (2025). Deep Learning and Federated Learning in Breast Cancer Screening and Diagnosis: A Systematic Review. *IEEE Access*.
- [19] Gupta, Dr-Shashi & Khan, Mudassir & Taherdoost, Hamed & Soni, Rashmi & Ibtissam, Chouja. (2025). Federated learning in healthcare: Revolutionizing patient care through intelligent algorithms. 10.13140/RG.2.2.25108.00640.
- [20] Isaac, Osei & Aabaah, Iven & Appiah, Benjamin & Osei, Isaac & Benjamin, Appiah. (2025). Federated Learning for Lung Cancer Detection: Comparative Analysis and Visual Interpretability Federated Learning for Lung Cancer Detection: Comparative Analysis and Visual Interpretability. 10.21203/rs.3.rs-6032484/v1.

ORIGINALITY REPORT

17%	12%	8%	10%
SIMILARITY INDEX	INTERNET SOURCES	PUBLICATIONS	STUDENT PAPERS

PRIMARY SOURCES

1	Submitted to Daffodil International University Student Paper	5%
2	dspace.daffodilvarsity.edu.bd:8080 Internet Source	1%
3	Submitted to United International University Student Paper	1%
4	www.mdpi.com Internet Source	1%
5	arxiv.org Internet Source	<1%
6	seer.ufrgs.br Internet Source	<1%
7	C Kishor Kumar Reddy, Anindya Nag. "Federated Learning for Neural Disorders in Healthcare 6.0", Routledge, 2025 Publication	<1%
8	Submitted to University of Finance – Marketing Student Paper	<1%
9	Submitted to Asia Pacific University College of Technology and Innovation (UCTI) Student Paper	<1%
10	www.nature.com Internet Source	<1%
11	Sushil Kamboj, Pardeep Singh Tiwana. "Innovations in Computing", CRC Press, 2025	<1%

*% detected as AI

AI detection includes the possibility of false positives. Although some text in this submission is likely AI generated, scores below the 20% threshold are not surfaced because they have a higher likelihood of false positives.

Caution: Review required.

It is essential to understand the limitations of AI detection before making decisions about a student's work. We encourage you to learn more about Turnitin's AI detection capabilities before using the tool.

Disclaimer

Our AI writing assessment is designed to help educators identify text that might be prepared by a generative AI tool. Our AI writing assessment may not always be accurate (it may misidentify writing that is likely AI generated as AI generated and AI paraphrased or likely AI generated and AI paraphrased writing as only AI generated) so it should not be used as the sole basis for adverse actions against a student. It takes further scrutiny and human judgment in conjunction with an organization's application of its specific academic policies to determine whether any academic misconduct has occurred.

Frequently Asked Questions

How should I interpret Turnitin's AI writing percentage and false positives?

The percentage shown in the AI writing report is the amount of qualifying text within the submission that Turnitin's AI writing detection model determines was either likely AI-generated text from a large-language model or likely AI-generated text that was likely revised using an AI paraphrase tool or word spinner.

False positives (incorrectly flagging human-written text as AI-generated) are a possibility in AI models.

AI detection scores under 20%, which we do not surface in new reports, have a higher likelihood of false positives. To reduce the likelihood of misinterpretation, no score or highlights are attributed and are indicated with an asterisk in the report (*%).

The AI writing percentage should not be the sole basis to determine whether misconduct has occurred. The reviewer/instructor should use the percentage as a means to start a formative conversation with their student and/or use it to examine the submitted assignment in accordance with their school's policies.

What does 'qualifying text' mean?

Our model only processes qualifying text in the form of long-form writing. Long-form writing means individual sentences contained in paragraphs that make up a longer piece of written work, such as an essay, a dissertation, or an article, etc. Qualifying text that has been determined to be likely AI-generated will be highlighted in cyan in the submission, and likely AI-generated and then likely AI-paraphrased will be highlighted purple.

Non-qualifying text, such as bullet points, annotated bibliographies, etc., will not be processed and can create disparity between the submission highlights and the percentage shown.

