



# An advanced data fabric architecture leveraging homomorphic encryption and federated learning

Sakib Anwar Rieyan <sup>a</sup>, Md. Raisul Kabir News <sup>a</sup>, A.B.M. Muntasir Rahman <sup>a</sup>, Sadia Afrin Khan <sup>a</sup>, Sultan Tasneem Jawad Zaarif <sup>a</sup>, Md. Golam Rabiu Alam <sup>a</sup>, Mohammad Mehedi Hassan <sup>b,\*</sup>, Michele Ianni <sup>c</sup>, Giancarlo Fortino <sup>c</sup>

<sup>a</sup> Department of Computer Science and Engineering, School of Data and Sciences, BRAC University, 66 Mohakhali, Dhaka, 1212, Bangladesh

<sup>b</sup> Department of Information Systems, College of Computer and Information Sciences, King Saud University, Riyadh, 11543, Saudi Arabia

<sup>c</sup> Department of Informatics, Modeling, Electronics, and Systems, University of Calabria, Rende, CS, 87036, Italy

## ARTICLE INFO

### Keywords:

Data fabric  
Federated learning  
Partially homomorphic encryption  
Data fusion  
Data lake

## ABSTRACT

Data fabric is an automated and AI-driven data fusion approach to accomplish data management unification without moving data to a centralized location for solving complex data problems. In a Federated learning architecture, the global model is trained based on the learned parameters of several local models that eliminate the necessity of moving data to a centralized repository for machine learning. This paper introduces a secure approach for medical image analysis using federated learning and partially homomorphic encryption within a distributed data fabric architecture. With this method, multiple users or clients (hospitals/medical data centers) can collaborate in training a machine-learning model without exchanging raw data. The approach complies with laws and regulations such as HIPAA and GDPR, ensuring the privacy and security of the data. The study demonstrates the method's effectiveness through a case study on pituitary tumor classification, achieving a significant accuracy of 83.31%. However, the primary focus of the study is using the data fabric architecture to securely store and analyze medical images while complying with HIPAA and GDPR regulations. The results highlight the potential of these techniques to be applied to other privacy-sensitive domains and contribute to the growing body of research on secure and privacy-preserving machine learning.

## 1. Introduction

Artificial intelligence (AI) has become an integral part of our daily lives, including in the field of healthcare. In order to make the most of AI in healthcare, it is important to have access to large amounts of high-quality data. However, the confidentiality and sensitivity of healthcare data pose significant challenges to its storage and analysis. For example, according to a study [1], almost 5,150 data breaches were reported to OCR (Optical Character Recognition) between October 21, 2009, and December 31, 2022. The volume of healthcare data is often very large, particularly due to the prevalence of image-based data. In addition, the confidentiality of healthcare data is of the utmost importance, as it is often personal and sensitive in nature.

To address these challenges, we have built an advanced data fabric architecture that brings together healthcare centers in a region and stores patient data and diagnoses in a secure and privacy-preserving

manner. Data fabric is a data fusion [2,3] and integration approach to accomplish data management unification through analytics and AI. The proposed approach utilizes federated learning and partially homomorphic encryption [4] to allow for collaborative machine learning on encrypted data, while still maintaining compliance with laws and regulations such as the Health Insurance Portability and Accountability Act (HIPAA) [5] and the General Data Protection Regulation (GDPR) Act 2018 [6].

In the context of advancing healthcare through technological innovation, the work of [7] emerges as a pivotal exploration into the application of deep learning methods within the domain of multi-source heterogeneous data fusion. They compare two fusion approaches: stage-based fusion, which aligns different data sources toward a common goal but misses interaction, and feature-based fusion, which overlooks redundancy among features, affecting correlation. By combining

\* Corresponding author.

E-mail addresses: [sakib.anwar.rieyan@g.bracu.ac.bd](mailto:sakib.anwar.rieyan@g.bracu.ac.bd) (S.A. Rieyan), [md.raisul.kabir.news@g.bracu.ac.bd](mailto:md.raisul.kabir.news@g.bracu.ac.bd) (Md.R.K. News), [a.b.m.muntasir.rahman@g.bracu.ac.bd](mailto:a.b.m.muntasir.rahman@g.bracu.ac.bd) (A.B.M.M. Rahman), [sadia.afrin.khan@g.bracu.ac.bd](mailto:sadia.afrin.khan@g.bracu.ac.bd) (S.A. Khan), [sultan.tasneem.jawad.zaarif@g.bracu.ac.bd](mailto:sultan.tasneem.jawad.zaarif@g.bracu.ac.bd) (S.T.J. Zaarif), [rabiul.alam@bracu.ac.bd](mailto:rabiul.alam@bracu.ac.bd) (Md.G.R. Alam), [mmhassan@ksu.edu.sa](mailto:mmhassan@ksu.edu.sa) (M.M. Hassan), [michele.ianni@unical.it](mailto:michele.ianni@unical.it) (M. Ianni), [giancarlo.fortino@unical.it](mailto:giancarlo.fortino@unical.it) (G. Fortino).

deep learning and data fusion, [7] underscores the potential of harnessing hierarchical features through unsupervised training. This resonates with our work proposing a sophisticated data fabric architecture, emphasizing the fusion of technology and contemporary challenges.

In this study, we have used pituitary tumor classification as a case study, employing various deep-learning models such as VGG16, VGG19, ResNet50, and ResNet152. Our results show promising potential for the use of federated learning [2] and partially homomorphic encryption in secure medical image analysis. Specifically, we achieved good performance with VGG16 and VGG19 models, while ResNet50 and ResNet152 achieved lower accuracy and precision for both classes. However, our custom CNN architecture outperformed all of these pre-trained models in almost every metric that we used. Our findings contribute to the growing body of research on secure and privacy-preserving [8,9] machine learning [10] and demonstrate the potential for these techniques to be applied in other privacy-sensitive domains.

This paper's structure follows a clear sequence: Section 1 introduces the context, addresses healthcare data security challenges, presents an advanced data fabric architecture integrating encryption and federated learning, outlines research objectives and contributions. Section 2 provides background on key concepts: Data Fabric, Federated Learning, and Homomorphic Encryption for medical image analysis. Section 3 explores existing research on Data Fabric architecture, secure healthcare data management, encryption methods, and their applications in healthcare data privacy and analysis. Section 4 outlines the experimental approach, including data description, preprocessing, model architecture, and neural network models for encrypted medical image classification. Section 5 presents study outcomes, including the use of homomorphic encryption, model evaluation on encrypted and unencrypted data, performance metrics for different models in classifying pituitary tumors, and discussing the results' implications. Section 6 summarizes research achievements in developing an advanced data fabric architecture using Partial Homomorphic Encryption and Federated Learning for secure and decentralized machine learning on medical data. It discusses study assumptions, limitations, and potential future directions.

### 1.1. Motivation

In the field of medical image analysis, ensuring the security of sensitive patient data is of utmost importance. However, with the increasing use of machine learning and deep learning techniques for medical image analysis, there is a pressing need for an effective and secure architecture to handle such data. Previous studies have shown that the use of conventional security measures, such as encryption and access control, is not enough to ensure the privacy of patient data in the context of machine learning and deep learning operations. In 2015, Anthem Inc., one of the largest health insurance companies in the United States, suffered a massive data breach that exposed the personal information of nearly 79 million individuals. The attackers gained access to a vast amount of sensitive medical data. The breach was a success because of the Centralized Database Vulnerability, Lack of Encryption, and Slow Detection [1].

In addition, the use of traditional centralized architectures [11,12] for processing medical image data can be slow and resource-intensive, which can further compromise the security of the data [13,14]. Therefore, there is a clear need for a new, advanced data fabric architecture that is specifically designed to handle the unique challenges of securing medical image data while also supporting efficient machine learning and deep learning operations. This research aims to address this gap in the current state of the art by proposing and evaluating a novel architecture that is capable of effectively securing medical image data while also enabling fast and accurate machine learning and deep learning operations.

### 1.2. Research problems

The integration of data into healthcare has the potential to improve the prediction of diseases and epidemics, enhance treatment outcomes, and prevent premature deaths. However, the confidentiality of healthcare data and the complexity of managing large and diverse datasets pose significant challenges to the integration of data into healthcare. Ensuring data security and privacy is of utmost importance, as security breaches in healthcare are on the rise. According to a study [15], there were 3,033 data breaches reported between 2010 and 2019, resulting in the exposure of 255.18 million records of data. Furthermore, as previously mentioned, during last 13 years, there have 5,150 cases of data breaches in this sector according to [1].

Moreover, the substantial volume of healthcare data presents challenges in terms of efficient processing, storage, and communication. Conventional methods may prove inadequate when confronted with the magnitude of the data at hand. In one proposed solution [16], a big data healthcare cloud would host clinical, financial, social, physical, and psychological data from patients in a centralized location. However, proper governance of the data cloud is necessary to effectively work with and analyze complex data.

In this study, we aim to address these challenges by proposing an advanced data fabric architecture that brings together healthcare centers in a region and stores patient data and diagnoses in a secure and privacy-preserving manner using federated learning and partially homomorphic encryption. We demonstrate the effectiveness of our approach using pituitary tumor classification as a case study. However, the primary focus of our work is on the development and evaluation of federated learning and partially homomorphic encryption as tools for secure medical image analysis in the healthcare sector.

The primary objective of this study is to address the research question:

*How effective and practical is the implementation of advanced data fabric architecture using federated learning and partially homomorphic encryption for secure medical image analysis in the healthcare sector?*

### 1.3. Contributions

Through our work, we show a fully-fledged data fabric architecture based on healthcare data can be built whilst complying with privacy regulations and maintaining good accuracy scores. Our primary contribution spans four aspects:

- We propose an advanced data fabric architecture for storing and collaborating and fusing healthcare data in an encrypted form by using Partial Homomorphic Encryption (PHE) and sharing it with other parties without revealing its content. In this architecture, medical images of various patients/clients are encrypted on the client side and these encrypted images are then used as inputs for deep learning models, enabling the models to learn and classify tumors. Subsequently, the system collects the classified tumor data for further analysis and processing. Further processing of data is done in its encrypted state. Here, The raw data was encrypted and generated to local weights using the local FL Model before getting global attention. Therefore even if the data was backtracked the end result will produce nothing but an encrypted image. Thus, this architecture provides a secure and efficient mechanism for processing encrypted data, while preserving data privacy and confidentiality regulations such as HIPAA and GDPR. This transformative approach marks a distinct departure from currently available data architectures that lack such encryption mechanisms.
- Our architecture also encompasses a federated learning framework, allowing multiple clients to collaboratively train machine learning models on their respective data. Unlike the existing

general federated learning frameworks which function on real-time local and global updates, our framework offers the flexibility to modify, scale, merge or select the local model updates before using them into the global model. In this way, the framework we proposed facilitates the systematic exchange of model updates between the local and global models. This innovation grants healthcare organizations the unparalleled ability to securely collaborate on model training while maintaining data privacy. No existing architecture seamlessly integrates federated learning within a data fabric, thus highlighting the exceptional nature of our contribution.

- Moreover, we have tailored a convolutional neural network (CNN) architecture, inspired by VGG16 and VGG19, with a smaller input size, resulting in a reduced parameter size compared to the aforementioned models. This customization enables enhanced efficiency by reducing computational complexity, particularly when leveraging Partially Homomorphic Encryption (PHE) techniques. This optimization sets it apart from previously established architectures that often overlook the symbiotic relationship between encryption and model design.
- We further evaluate the proposed approach by implementing a prototype of the homomorphic encryption-based data fabric and the federated learning framework. The assessment indicates that the suggested method offers an effective and reliable data fusion for sharing and analyzing data securely. The experimental results demonstrate that the proposed approach achieves satisfactory accuracy in the collaborative training of machine learning models, even when the data is encrypted. This pivotal advancement distinguishes our work from current architectures that require disjointed solutions for data management and operational practices.

## 2. Background studies

### 2.1. Data fabric

According to Gartner [17], Data Fabric is a unified and integrated platform that enables data discovery, fusion and integration, management, and access across multiple environments. It provides a consistent and scalable approach to managing data assets that are distributed across various locations, such as on-premises, cloud, and edge computing. Data Fabric helps organizations to simplify and optimize their data management processes, reduce data silos, and enable real-time access to data. It also supports the creation of a self-service data marketplace, allowing users to discover, share, and consume data in a secure and governed manner. Data Fabric is increasingly becoming a critical component of modern data architectures, as organizations seek to manage the growing volume, velocity, and variety of data generated by digital business initiatives.

In our research, we are utilizing homomorphic encryption to classify pituitary tumors from MRI images in our dataset. We have used a Data Fabric architecture to store the weights of different machine-learning models as encrypted data. The ML models are run on client PCs, and the resulting encrypted data is saved in our Data Lake. Using homomorphic encryption, a server PC can perform computations on the encrypted data, allowing for the creation of a homogeneous global model. The server can then provide users with the requested results without compromising the privacy of the MRI images. This approach benefits from the Data Fabric's ability to provide a unified and integrated platform that enables data discovery, integration, management, and access across multiple environments. By utilizing homomorphic encryption and a Data Fabric architecture, we can classify pituitary tumors from MRI images in a privacy-preserving manner, contributing to the development of more secure and privacy-preserving medical imaging technologies.

### 2.1.1. Vanilla architecture of data fabric

[Fig. 1](#) provides a visual overview of the key components and processes of the Vanilla Architecture of Data Fabric.

#### (i) Accessing Data:

- (a) **Data Collecting and Encryption:** Data is a volatile resource, and Medical Data is considered highly sensitive as it can contain personally identifiable information such as names, addresses, dates of birth, and medical records which can be exploited if they fall into wrong hands. To comply with this shortcoming, in this architecture, data is neither collected nor stored in a central server which may possess the risk of data leakage.

Here, firstly, to ensure privacy and reduce data volatility, medical data from various users are first selected and then encrypted with Partially Homomorphic Encryption (PHE). Subsequently, the encrypted data is selected to train the model locally for collecting updated model weights. The data of each user is generated and stored locally, without being transferred to the central server. Instead, the generated model updates are stored and merged for the global model formation.

- (b) **Master Data Management:** Following the generation of local model updates and subsequent merging of data, feature selection is employed as a means of optimizing and enhancing efficiency. By selecting the most relevant features or weights, the dimensionality of the data can be reduced, facilitating ease of analysis. Moreover, feature selection mitigates the risk of overfitting, improves model accuracy, and reduces computational costs, thereby achieving heightened efficiency through the utilization of a reduced training dataset.

FedMax, FedAvg, and FedMin are optimization algorithms used in Federated Learning for feature selection. In all three algorithms, updated model weights are sent to the server/stored for future usage. However, in the case of FedMax, the server/user selects the model with the highest accuracy while it chooses the lowest loss model for FedMin and an average of all models for FedAvg.

In our architecture, we selected FedMax as our feature selection algorithm to select important and relevant model weights as it showed more accuracy and efficiency compared to FedAvg and FedMin. The selected data is collected and kept together as "Master Data".

#### (ii) Managing Life Cycle:

- (a) **Governance:** Data governance is an essential component of data fabric architecture. Data fabric architecture is an approach to data management that enables organizations to manage and process data from multiple sources, locations, and formats. It provides a unified view of data across the organization and supports various data processing requirements, such as data integration, analytics, and artificial intelligence.

Data governance in data fabric architecture refers to the policies, processes, and standards that organizations implement to manage their data assets effectively. Data governance helps organizations ensure that their data is accurate, consistent, and compliant with regulatory requirements. It also helps organizations manage data privacy, security, and access.

The following are some key considerations for data governance in data fabric architecture:

**Data quality:** Data governance policies should include measures to ensure data quality, such as data profiling, data cleansing, and data validation.

**Metadata management:** Data governance policies should include metadata management to ensure that data is properly tagged, categorized, and classified. Metadata helps organizations understand the meaning and context of their data and facilitates data discovery and reuse.

**Data privacy and security:** Data governance policies should include measures to ensure data privacy and security, such as access controls, data encryption, and data masking.

**Data lineage:** Data governance policies should include data lineage to track the origin, transformation, and movement of data across the organization. Data lineage helps organizations understand how data is used and facilitates compliance with regulatory requirements.

**Data ownership and stewardship:** Data governance policies should define data ownership and stewardship to ensure that data is managed and maintained by the appropriate individuals and teams.

(b) **Compliance:** Data compliance refers to adhering to relevant laws, regulations, and industry standards related to the handling, processing, and storage of data. In the context of data fabric, data compliance refers to ensuring that data is managed in accordance with these requirements across the entire data fabric. To ensure data compliance in a data fabric, it is necessary to establish policies and procedures that cover the entire data lifecycle, from data ingestion to archival and deletion. Personal data must be collected, processed, and stored in compliance with privacy regulations such as GDPR, CCPA, HIPAA, etc. In this architecture, the feature selected weights were trained on various models such as VGG16, VGG19, ResNet 50, ResNet 152, etc. and updates are stored in a data lake structurally based on models they were trained on complying with HIPAA regulations.

(iii) **Exposing Data:** Data exposure refers to making data available for consumption and analysis by users or applications within an organization. Exposing data in a data fabric involves providing access to the data for authorized users or applications. There are several ways to expose data in a data fabric, including:

**APIs:** Application Programming Interfaces (APIs) enable applications to access and retrieve data from the data fabric.

**Data Catalogs:** A data catalog provides a searchable inventory of data assets in the data fabric. Users can discover and access data assets through the data catalog.

**Self-Service Analytics:** A self-service analytics platform enables users to create their own queries and reports using the data available in the data fabric.

**Data Virtualization:** Data virtualization enables users to access and combine data from multiple sources as if it were in a single location.

## 2.2. Federated learning

Federated Learning is a distributed machine learning technique that enables multiple clients to collaboratively learn a shared model without exchanging their raw data. This technique has gained popularity in recent years due to its ability to preserve data privacy and security while improving model performance. As shown in Fig. 2, each client trains a local model using its own data and then sends the local model weights to a central server. The central server then aggregates the local model weights to update a global model that is shared among all clients. This process continues iteratively until the global model achieves the desired level of accuracy. According to the report [18], Federated Learning has been successfully applied to various domains, such as speech recognition, natural language processing, and healthcare, where data privacy is a major concern.

## 2.3. Homomorphic encryption

A cryptographic method called homomorphic encryption enables mathematical operations to be carried out on ciphertext without exposing the underlying plaintext. Table 1 offers a comprehensive view of different types of homomorphic encryption along with the distinctions that set them apart. In our research, we used partially homomorphic encryption to encrypt sensitive medical images, specifically brain MRI scans.

Partially homomorphic encryption (PHE) is a type of homomorphic encryption that only supports a limited set of mathematical operations, such as addition or multiplication. By encrypting the medical images using this technique, we were able to process and analyze the data without exposing the sensitive information contained within it. Fig. 3 provides an illustrative overview of Partially Homomorphic Encryption (PHE) Technique at the Pixel Level to Dataset Images.

One major benefit of using partially homomorphic encryption in this context is that it ensures the confidentiality of medical data. As medical information is often highly sensitive and personal, it is important to protect it from unauthorized access. By encrypting the data, we were able to securely process and analyze it without compromising its confidentiality.

In addition, partially homomorphic encryption allows for more efficient processing of the encrypted data. Because the mathematical operations can be performed directly on the ciphertext, there is no need to decrypt the data first, which can be a time-consuming process. This was particularly useful when working with large datasets or when processing data in real time [19].

Overall, our use of partially homomorphic encryption proved to be a successful and effective method for protecting the confidentiality of sensitive medical images while still enabling their analysis.

### 2.3.1. The paillier encryption scheme

As previously mentioned, we are using Partial Homomorphic Encryption for our dataset which follows The Paillier Encryption Scheme. The Paillier encryption scheme [20] is an additively homomorphic cryptosystem based on the computational difficulty of the decisional composite residuosity assumption. The scheme's security relies on the difficulty of factoring the product of two large prime numbers. Key components of the Paillier encryption scheme are:

1. **Key Generation:** The key generation process creates a public key  $(n, g)$  and a private key  $(\lambda, \mu)$ .

- $n$  is the product of two large prime numbers, kept secret and known only to the data owner.
- $g$  is a public system parameter, typically set as  $g = n + 1$ .
- $\lambda$  is Carmichael's totient function [21],  $\lambda = lcm(p-1, q-1)$ , where  $p$  and  $q$  are the large prime factors of  $n$ .
- $\mu$  is the modular multiplicative inverse of  $\lambda$  modulo  $n$ .

2. **Encryption:** Given a plaintext image  $x$ , the encryption process is performed as follows:

$$Enc(x) = (g^x \times x^n) \% n^2 \quad (1)$$

where  $x$  is a random value chosen for each encryption, ensuring probabilistic encryption.

### 2.3.2. Homomorphic operations

1. **Addition:** Addition on encrypted values is akin to combining the original plaintext values after decryption. Given two encrypted images  $Enc(x)$  and  $Enc(y)$ , the homomorphic addition can be performed as:

$$Enc(x + y) = Enc(x) \times Enc(y) \% n^2 \quad (2)$$

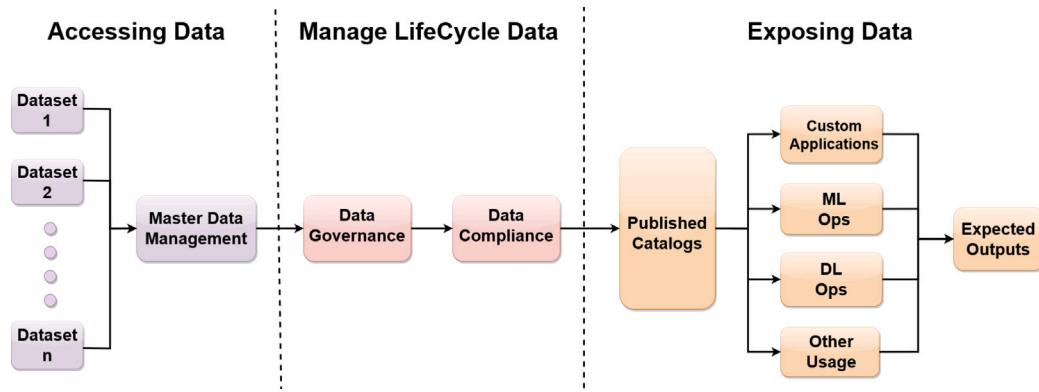


Fig. 1. Key stages of the vanilla data fabric architecture.

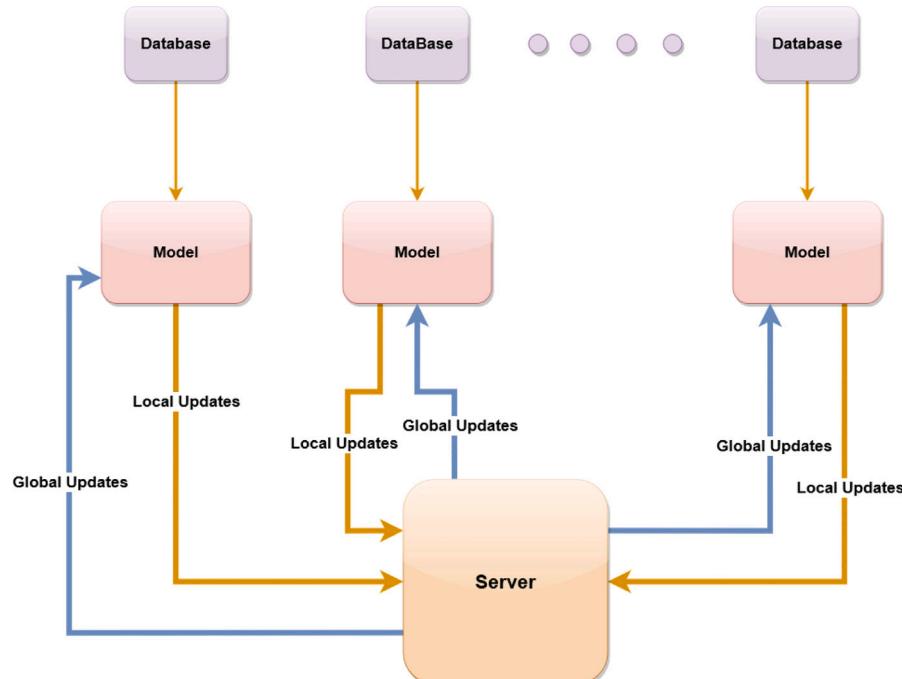


Fig. 2. Workflow of the federated learning model.

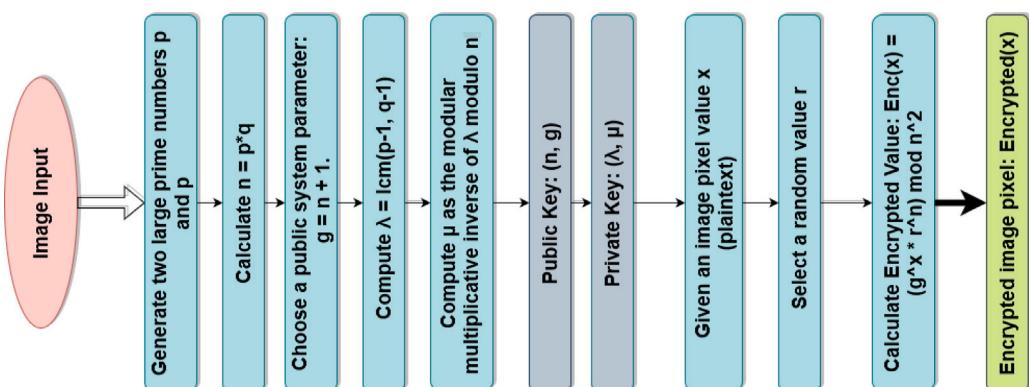


Fig. 3. Application of partially homomorphic encryption (PHE) technique at the pixel level of dataset images.

**Table 1**  
Comparison among different types of Homomorphic Encryption.

Types of Homomorphic Encryption	Partially Homomorphic Encryption (PHE)	Somewhat Homomorphic Encryption (SHE)	Fully Homomorphic Encryption (FHE)
Supported Operations	Addition or Multiplication	Addition & Multiplication	Arbitrary Computations
Security	High	Moderate to High	High
Computational Efficiency	High	Moderate	Low
Computational Intensity	Low	Moderate	High
Encryption	Simple	More Complex than PHE	Very Complex
Encryption Overhead	Low	Moderate	High
Implementation Ability	Easy	Moderate	Difficult

**2. Scalar Multiplication:** Scalar multiplication involves multiplying an encrypted value by a plaintext scalar, equivalent to raising the corresponding plaintext value to that scalar power after decryption. For an encrypted image  $Enc(x)$  and a plaintext scalar  $k$ , the scalar multiplication can be conducted as:

$$Enc(k \times x) = Enc(x)^k \% n^2 \quad (3)$$

### 3. Literature review

Data fabric architecture is a relatively new concept that has already been utilized by notable organizations, including IBM, for data fusion, management, and unification purposes. Despite its potential benefits, there is limited research available on the implementation of this architecture in the healthcare system. Due to the sensitive nature of healthcare, developing a secure data fabric architecture can present challenges. This chapter examines various data architectures, processes, and encryption methods that can be employed to ensure the security of healthcare data.

There are a few works that are related to the architecture we are working on. They are described below:

In [22], the authors describe a proof-of-concept implementation that uses the Hyperledger Fabric framework. They claim that this concept is capable of storing patient records effectively with keeping all the privacy protocols intact. Lastly, they compare the read-write times according to their claim.

In [23], Roehrs et al. divide personal health record (PHR) information into data blocks that may appear to be centrally stored but are actually distributed among participating devices. The authors claim that their proposed openPHR protocol is practical, flexible, and scalable for adoption by multiple organizations. Although the authors provide a detailed architecture, questions have been raised about the feasibility of their approach, especially regarding security and privacy concerns. It is important to note that PHRs are controlled by patients, while electronic health records (EHRs) are managed by healthcare institutions. Nevertheless, EHRs and PHRs are electronically stored and distributed and can be assessed based on metrics such as performance, scalability, privacy protection, and compliance with the GDPR.

In [24], the MeDShare platform shares several similarities with PREHEALTH [22]. However, the authors do not explicitly specify the underlying blockchain framework. Additionally, their emphasis is more on examining the fundamental components of blockchain technology, including data blocks and smart contracts, rather than presenting a practical solution.

Ming and Zhang [25] present an effective privacy-preserving access control (PPAC) strategy for cloud-based EHR systems. Their approach utilizes the cuckoo filter and an innovative attribute-based signcryption (ABSC) mechanism to achieve both anonymity and computational efficiency. The authors offer comprehensive assurances of privacy and conduct thorough performance evaluations for comparison. However, it is uncertain whether their approach complies with the GDPR regulations.

In [26], Fu et al. mainly focus on sharing data among different participants safely. Using Hyperledger Fabric, they propose a more secure decentralized distributed data storage over traditional centralized

data for SKA data due to high management costs and low credible traceability. The SKA Data Management Alliance significantly reduces costs and improves the overall security of its data by adopting this distributed storage system.

In [27], Ram et al. propose a multi-sensor data fusion technology which is effective to use multi-sensors for collecting data from similar target and analyze, the collected data using various computer technologies. Multi-sensor refers to a process, which combines observations from a few different sensors for providing a robust and complete overview environment or any process of interest. Many nodes are included in the IoT system, and timely data fusion is effective for organizations to enhance their decision-making process along with some effective challenges in the future such as imperfect data, inconsistency, conflicting nature of data, and several others.

In [28], the authors present an information fusion model integrated with traffic management. Their approach considers the information filtering processing and diverse membership function fusion, while also covering the traditional methods as well. By considering these aspects, the proposed method aims to enhance the effectiveness of information fusion in traffic management.

Our architecture utilizes federated learning to train deep learning models. In our research, we have come across some related studies that are somewhat aligned with our work. Here are brief descriptions of these studies:

In [29], the authors propose a secured model of federated learning using homomorphic encryption and attempt to classify Covid-19 using X-ray images. They secure the federated process and claim that sensitive information can fall into the hands of attackers if the DL process is not secured. They do not encrypt the dataset but rather encrypt the weight matrix used in federated learning. They claim to achieve an accuracy of 84.00% with a precision of 86.89%.

It is obvious that with the vast growth of artificial intelligence and big data, contradictions between user data policy and data policy also grow proportionally. That is why in [30], the authors come up with a vertical federated learning system for Bayesian machine learning using homomorphic encryption. Their model can be compared with 90% of models trained by a single union server. Additionally, their system can be used in education, finance, medicine, risk controls, and other fields.

Here are some additional studies that are relevant to our research and have been reviewed and summarized below:

In [31], the authors provide insights into the difficulties of medical data analysis and security and propose a solution based on a decentralized architecture. They utilize the Exonum framework. In their proposed architecture, they separate the whole system into two parts - 'Closed Information' and 'Open Information.' In the closed part, encrypted data is stored in a blockchain, while in the open part, non-encrypted service-related data is stored.

Zhang et al. [32] propose a meaningful usage of optimizing Electronic Health Records (EHRs) using big data analytics. Here, they propose an insight into how to improve electronic health records using three methods: Data Collection, Data Storage, and Data Utilization. Firstly, in the data collection method, records are divided into structured and unstructured data. Structured data includes demographics,

**Table 2**  
Comparative analysis of Contemporary Data Architectures in the Healthcare Sector.

Architecture	Technology	Access Level	HIPAA	GDPR	Privacy Preserving
PREHEALTH [22]	Hyperledger Fabric	Private	Not Mentioned	Yes	Yes
OmniPHR [23]	Peer - to - Peer	Private	No	No	No
MeDShare [24]	Agnostic	Open	No	No	No
Access Control Based EHR [25]	AC Scheme	Private	No	No	Yes
Our Proposed Work	Data Fabric	Private	Yes	Yes	Yes

health status, lab results, billing, etc., while surgical videos or diagnosis notes fall under unstructured data. After collecting, they propose a transformation engine where data is moved, cleaned, merged, and validated, and is stored in DBMS, Cloud, or NoSQL. Finally, transformed data is processed using mapping and reduction, and stream computing and in-database analytics are used for generating reporting systems, which help achieve a meaningful usage of EHRs using big data analytics. However, the authors also map out the limitations of this research, and emphasize that it lays the foundation for interesting opportunities in the future.

In [33], the authors introduce a novel big data platform that can redesign modern medical data and bring an effective and quick solution to the healthcare system. They propose a system that is lightning-fast, supports stream processing, and integrates with both NoSQL and RDBMS. This process aims to exploit open-source technologies as much as possible and build the system on top of them. The core of the system is Spark core, and to utilize it, the system uses the Spark framework with real-time graphical image processing. For handling structured data, the authors propose SparkSQL Structured Data and MLlib Machine Learning for classifying the data.

In [34], the authors describe an architecture where an improvised big data model is involved to create a cloud computing environment for healthcare. In this process, huge amounts of data from medical sources and processes are collected in cloud storage, and real-time analysis is done using cloud computing for better accuracy. The Healthcare Data Management Framework is built on Hadoop Clusters and certain key components. Semantic Practitioner, Big data container and processing layer, Query formulator, Batch scheduling, and Data reader are examples of such tools. This architecture is also open to implementing various cryptographic techniques on the Cloud.

### 3.1. Comparative analysis with related works

Our work is different from existing approaches as it enables effective and secure handling of highly sensitive health data through our proposed Data Fabric architecture. Another advantage of this work is that its privacy-preserving features (Homomorphic encryption) are compliant with GDPR [35], as it supports the right to be forgotten, as the actual health data is neither collected nor stored. Rather, after training, the encrypted data which is collected is stored as model updates. Similarly, it is possible to delete a client's data from the data lake upon request. In addition to this, compared to other works, our model also complies with HIPAA guidelines, as it ensures confidentiality, integrity, and availability of personal health information, safeguards data from threats, and protects impermissible access as mentioned here [5].

Table 2 displays a comparison between our data fabric architecture and other existing works that have demonstrated proof-of-concept implementations, which can be implemented and practical for real-world scenarios. In this table, the technology, privacy-preserving features, GDPR and HIPAA compliance assessment results of each existing proposal have been provided.

## 4. Methodology

### 4.1. Data description

In our experimental evaluation, we used the dataset [36] named “Brain MRI Images for Brain Tumor Detection”, available on Kaggle.<sup>1</sup> Fig. 4 presents a subset of samples extracted from the employed dataset, and Fig. 5 provides an encompassing overview of the dataset. The dataset includes a large number of 2D MRI images for the classification of brain tumors. Brain tumors are classified into three types: Benign, Malignant, and Pituitary. However, for our selected dataset, images are classified into two categories; Pituitary Tumor and No Tumor.

From the dataset, 1852 images have been used by us as shown in the figure below. On the images, we performed encryption and machine learning techniques for ensuring privacy and achieving desired results.

For preparing our own data fabric to perform MLOps, an image was selected from the dataset and converted into a Numpy array with its pixel values. After that, homomorphic encryption was performed on the array to make the data fabric encrypted. A brief overview is shown in the following figure.

### 4.2. Data preprocessing

To prepare our data we needed to preprocess the whole dataset. Firstly, we have encrypted our dataset using the Partially Homomorphic Encryption Algorithm. After that, we resized the images to  $128 \times 128$  dimensions. The shape of our dataset became (1852, 128, 128). Then we reshaped our dataset by multiplying the dimensions of individual pixels. Then the shape ultimately (1852, 15376). For machine learning classifiers we have scaled the image dataset to value 0 to 1.

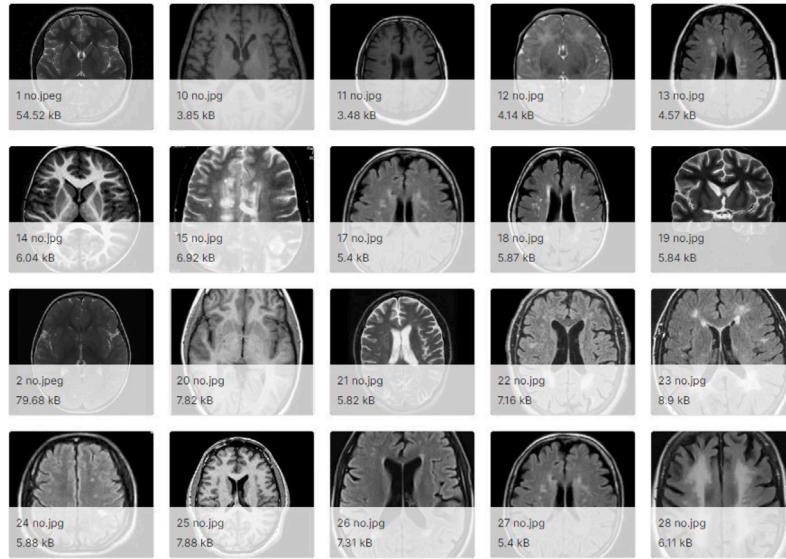
### 4.3. Proposed model

#### 4.3.1. Advanced architecture of data fabric

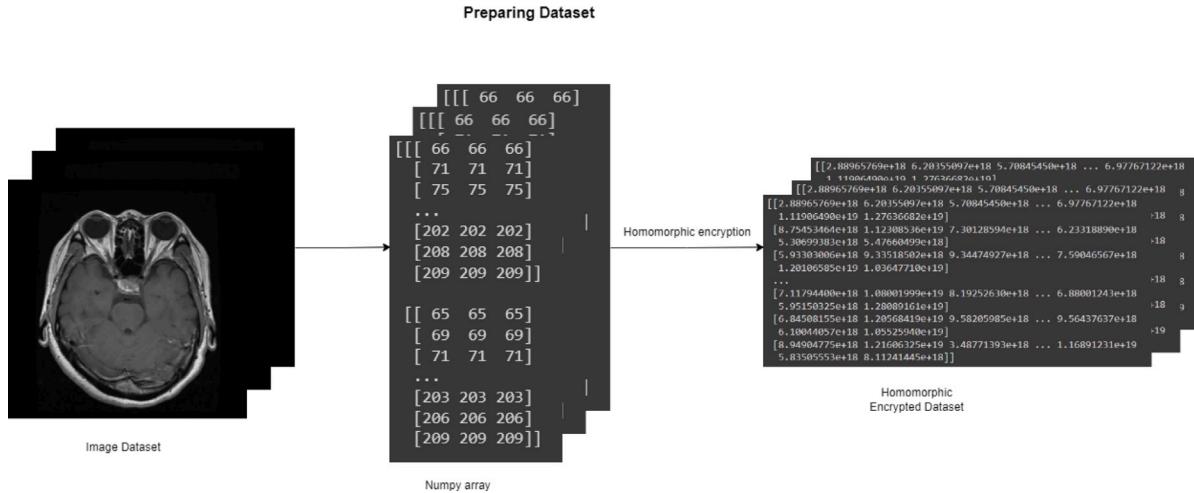
Fig. 11 presents an overview of our proposed data fabric architecture and illustrates the essential operational steps within the architecture.

- (i) **Accessing Data:** Here, initially, to ensure privacy and reduce data volatility, medical data from various users are first selected and then encrypted with Partially Homomorphic Encryption (PHE). Partially Homomorphic Encryption (PHE) enabled us to perform DLOPs on the data securely. Subsequently, the encrypted data is selected to train the model locally for collecting updated model weights, as shown in Fig. 6. For training, various federated learning models like VGG16, VGG19, ResNet50, ResNet152, and our Custom CNN were used which generated model updates for each model. The data of each user is generated and stored locally, without being transferred to the central

<sup>1</sup> <https://www.kaggle.com/datasets/masoudnickparvar/brain-tumor-mri-dataset>



**Fig. 4.** A few samples of the used Brain Tumor dataset.



**Fig. 5.** Overview of the encrypted dataset.

server. Instead, the local model updates were stored and merged for the global model formation.

After combining the local model weights, we selected FedMax as our feature selection algorithm to select important and relevant model weights as it showed more accuracy and efficiency compared to FedAvg and FedMin. The selected data is collected and kept together as “Master Data”, with which we will continue our next works.

- (ii) **Managing Life Cycle:** Selected model weights were differentiated based on the models they were trained on. In this case, for example, updates of the VGG16 model were kept structurally under the “VGG16” name. This enables effective data governance, as the risks of data inconsistency and complicated integration across the whole architecture get lowered. Additionally, organizing data in a structured way ensures proper usage of data and helps strike a balance between data collaboration and privacy mandates. Fig. 7 provides a concise overview of data lifecycle management in our proposed architecture. Furthermore, as the privacy of the data was already ensured in the first step, the collected model weights already comply with existing privacy regulations such as HIPAA, GDPR, etc.

(iii) **Exposing Data:** Since exposing data is a comprehensive approach where appropriate data access controls, data fusion, and cataloging must be implemented, we stored all the collected local model updates in a “Data Lake”. This enabled us to store our raw model weights structurally in their native format. Additionally, since there were huge amounts of model weights, using a Data Lake helped us to store and process it easily.

Most importantly, we can also perform MLOps on the Data Lake directly. On the client side, clients can train their data on the global federated model and can compare the global weights with the local weights of that model stored in the data lake and generate the desired output. Fig. 8 provides a concise overview of exposing data part of our proposed architecture.

#### 4.3.2. Implemented federated learning framework

As we delve into the world of medical imaging and machine learning, we have utilized a cutting-edge approach to store the local weights of our machine learning models using a Data Fabric architecture. This architecture has allowed us to securely store and manage distributed data across multiple environments, providing a consistent and scalable approach to managing data assets.

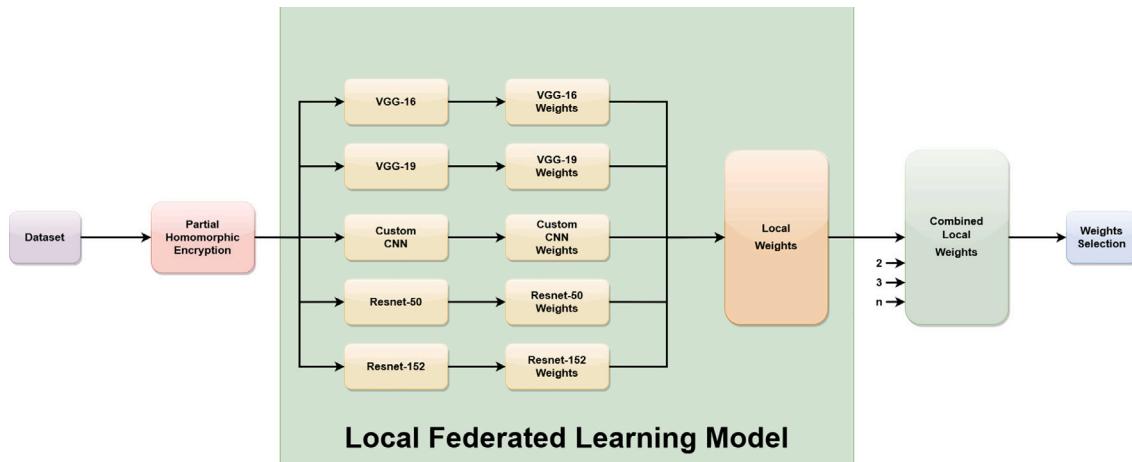


Fig. 6. Flowchart of accessing data section of the proposed data fabric architecture.

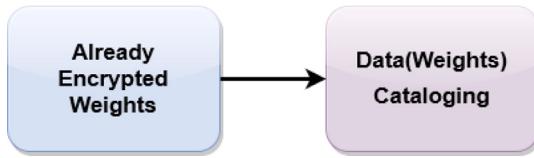


Fig. 7. Flowchart of data lifecycle management in the proposed architecture.

In the federated learning framework, model updates, and the fusion of local and global models are achieved through a collaborative process. Each client trains its own local model using their data, and the model updates are performed locally without sharing raw patient information. The updated local models are then aggregated and combined using the FedMax algorithm to create a global model. This aggregation process ensures that the global model benefits from the knowledge of all clients while preserving data privacy. The resulting global model is stored in a data lake, enabling easy access for evaluation and testing without compromising patient confidentiality.

#### Algorithm 1 Federated Learning Process

```

1: Initialize global model
2: Split data into shards and distribute among clients
3: for i in range(num_communication_rounds) do
4:   Initialize local model
5:   Send local model to each client
6:   for client in clients_data do
7:     client.model.train(client.data)
8:   end for
9:   Clients send updates to the global model
10:  Scale updates by the weight scaling factor
11:  Aggregate updates using FedMax
12:    global_model.update_weights(max(updates)
13:      * weight_scaling_factors)
13: end for

```

Our goal was to develop a pituitary tumor classification model that leverages the decentralized data in a privacy-preserving manner, improving the accuracy of the model while ensuring the privacy of patient information. Using the FedMax algorithm, an improvement on the popular FedAvg algorithm, we employed it to assemble a global model as depicted in Fig. 10, along with local models as illustrated in Fig. 9. This approach was employed for the classification of pituitary tumors from MRI images within our dataset. The FedMax algorithm

considers the model performance of each client and assigns a weighting factor that allows clients with the best performance to contribute more to the global model. We have stored the resulting model weights in a data lake, allowing us to generate a global model and test it using test cases when the user prompts to show results.

#### 4.3.3. Model specification

(i) **VGG16:** A 16-layered convolutional neural network model trained on the ImageNet dataset, it is considered to be one of the best models to date. It is widely regarded for its simple architecture and excellent image classification performance, retaining 92.7% test accuracy in the ImageNet dataset, which consists of almost 14 million training images across a thousand object classes.

As its name suggests, it is composed of 16 layers which include 13 convolutional layers and 3 fully connected layers. It comprises 138 million parameters and uses small convolutional filters and deep architectures to gain a large receptive field and strong discrimination ability, which helps in image classification, object detection, and semantic segmentation. To control overfitting and reduce spatial dimensions, its architecture has five max pooling layers.

The most unique thing about VGG16 is that it is focused on convolution layers of  $3 \times 3$  filter with stride 1 and always used the same padding and max pool layer of  $2 \times 2$  filter with stride 2, and the layers are constantly arranged over the whole architecture. Due to its high-level feature representation, it provides good performance on object detection, fine-grained image classification, etc. [37]. VGG16 is regarded as one of the most accurate pre-trained object detection models while being the complexity on the lower side. Therefore, we found it suitable for our case as we need a balanced model in terms of performance and complexity. Fig. 12 illustrates an overview of the VGG16 model.

(ii) **VGG19:** A variant of the VGG neural network, it is a 19 layer version of the VGG network and similar to the VGG16 architecture. Compared to VGG16, it has 5 convolutional layers and 1 fully connected layer. It has around 143 million parameters and is trained on a dataset with 1.2 million images and 1000 classes. The function accepts an image of shape (128, 128, 1) as input, and the image is passed through concatenate layer which concatenates the image 3 times resulting in a tensor with shape (128, 128, 3), this is passed to a VGG19 model which is a pre-trained convolutional neural network model that is trained on the ImageNet dataset. The VGG19 model serves as a feature extractor and it extracts features from the image by

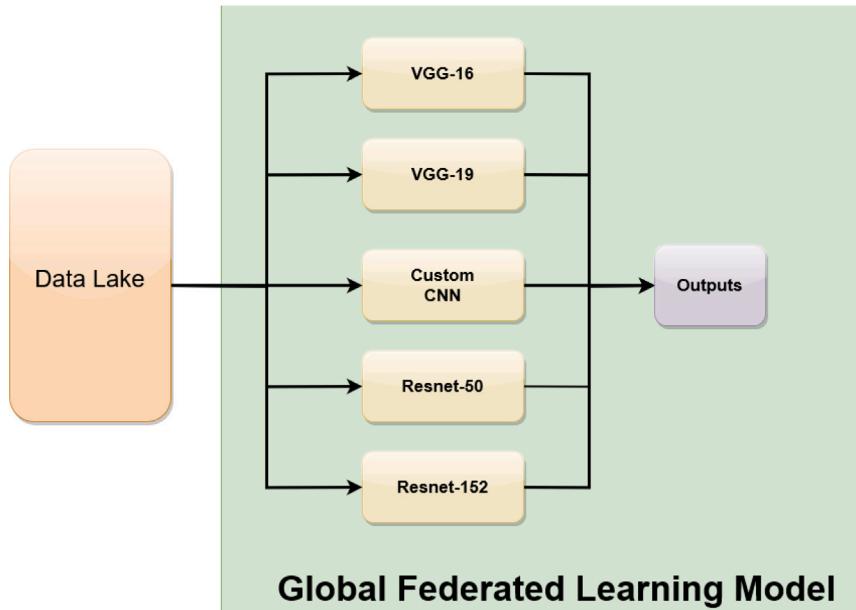


Fig. 8. Flowchart of exposing data part of the proposed architecture.

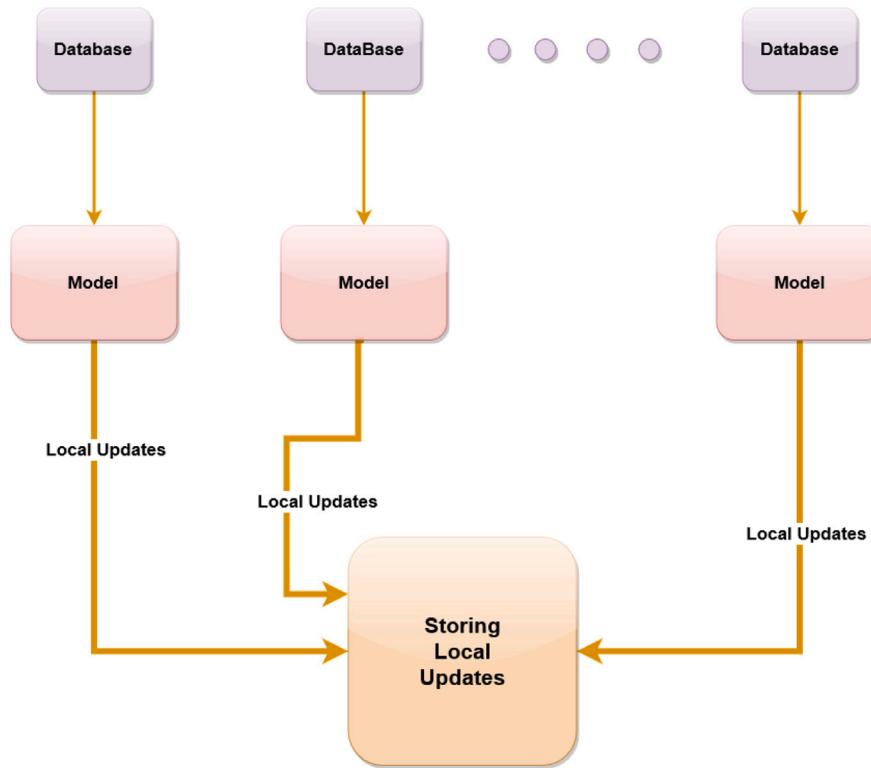
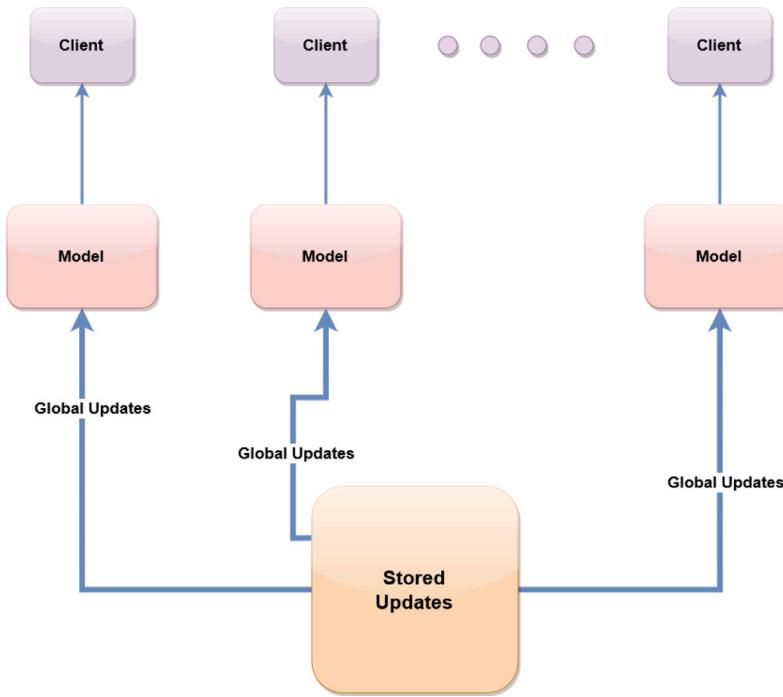


Fig. 9. Workflow and functionality of the local federated learning model of the proposed architecture.

stacking convolutional layers and pooling layers. The output of the final max pooling layer is then passed through a flattened layer which reshapes the output tensor into a 2D array. This flattened layer's output is passed through a dense layer with 1 unit and a sigmoid activation function, it produces the final output of the network which represents the predicted probability of the input image belonging to the target class [37]. VGG19 offers high accuracy in image classification tasks, making it a

compelling choice for classifying homomorphic encrypted images. However, its large number of layers and parameters can lead to increased computational complexity, which should be considered in resource-constrained environments. Despite this, VGG19's robust feature extraction capabilities, especially when utilizing pre-trained weights, make it a strong candidate for accurate encrypted image classification. Fig. 13 provides an overview of the VGG19 model.



**Fig. 10.** Workflow and functionality of the global federated learning model of the proposed architecture.

(iii) **ResNet50:** ResNet50 is a deep convolutional neural network architecture that is trained on more than a million images from the ImageNet dataset. It is known for its use of residual blocks, which address the problem of vanishing gradients in deep networks. These residual blocks have shortcut connections that bypass one or more layers, allowing gradients to flow more easily and making it possible to train very deep networks without the problem of vanishing gradients. The architecture also uses batch normalization layers, which normalize the activation of the layers, making it possible to train the network more quickly and effectively. Additionally, the number of filters increases as the network gets deeper, allowing it to automatically learn more complex features. ResNet50 is widely used in many computer vision tasks and is commonly used as a feature extractor for other tasks.

The ResNet50 model is used as a feature extractor in our experiment, and the output of the final convolutional layer is passed through a flattened layer which reshapes the output tensor into a 2D array. The flattened feature map is then passed through a Dropout layer with a drop rate of 0.5, which is used for regularization to prevent overfitting. The dropout layer is optional and could be removed by commenting out the line. Finally, the output of the dropout layer is passed through a dense layer with 1 unit and a sigmoid activation function, which is used to produce the final output of the network, which represents the predicted probability of the input image belonging to the target class [38]. Its moderate size and efficient performance makes it an attractive option for classifying homomorphic encrypted images. The pre-trained weights of ResNet50 can leverage learned feature representations, enhancing its ability to classify encrypted data accurately. Fig. 14 shows an overview of the ResNet50 model.

(iv) **ResNet152:** We also used a ResNet152 architecture model. The architecture is composed of a stack of convolutional and pooling layers. The model accepts an image of shape (128, 128, 3) as input and the output of the final average pooling layer is passed through a Flatten layer which reshapes the output tensor into a 2D array. Then, there is a Dense layer with classes, number of units, and softmax activation function. This dense layer produces

the final output of the network which represents the predicted probability of the input image belonging to different classes. This architecture does not have a dropout layer and the model is trainable, this means that the model can be fine-tuned with new data for a different task [38]. With 152 layers, ResNet152 offers higher model complexity and potentially increased accuracy for image classification tasks. While it may require more computational resources than ResNet50, ResNet152 can be a suitable choice to test how far PHE can offer accurate results. Its deeper architecture allows it to learn more intricate features, which can be advantageous when working with encrypted data affected by noise. Fig. 15 presents an overview of the ResNet152 model.

(v) **Custom CNN:** We have customized a convolutional neural network (CNN) in a similar structure to existing VGG16 and VGG19 models. It has a total of 16 convolutional layers and 3 fully connected layers. The architecture starts with an image of size (128, 128, 1) as input and predicts a binary output. The first Conv2D layer consists of 64 filters of size  $3 \times 3$ , followed by another Conv2D layer of 64 filters of size  $3 \times 3$ , and a max pooling layer with a pool size  $2 \times 2$ . The next two layers are similar, with 128 filters of size  $3 \times 3$  and another max pooling layer. This is followed by several more pairs of DepthwiseConv2D with a different number of filters and MaxPooling2D layers that are stacked on top of each other and learn increasingly complex features from the images. After two fully connected layers with 4096 units and a dropout of 0.5 applied on each layer, the output is generated with activation "sigmoid" and Dense Layer 1 unit. Compared to VGG16, our customized CNN has a deeper architecture with more convolutional layers and is similar to that of VGG19. It also has a higher number of filters in layers than VGG16 and VGG19. Additionally, Custom CNN uses a dropout layer after each fully connected layer, something that is not present in VGG16 and VGG19. Besides the pre-trained models, this CNN model should provide as good results as VGG models if not better for the encrypted data. The reason for this is the model was designed and tested keeping the limitations of PHE in mind. Fig. 16 depicts an overview of the Custom CNN model.

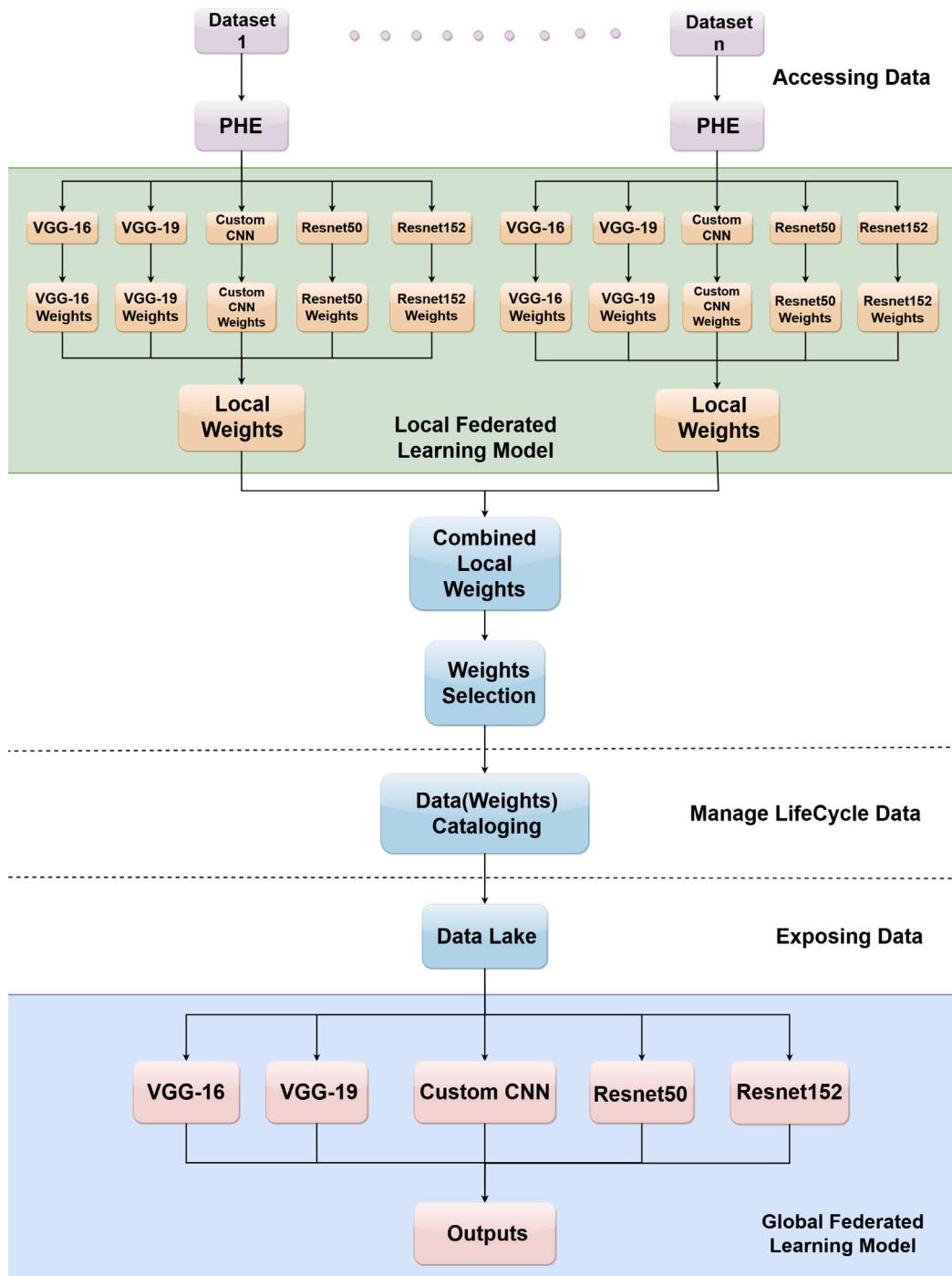


Fig. 11. Comprehensive workflow of the proposed data fabric architecture incorporating federated learning and homomorphic encryption.

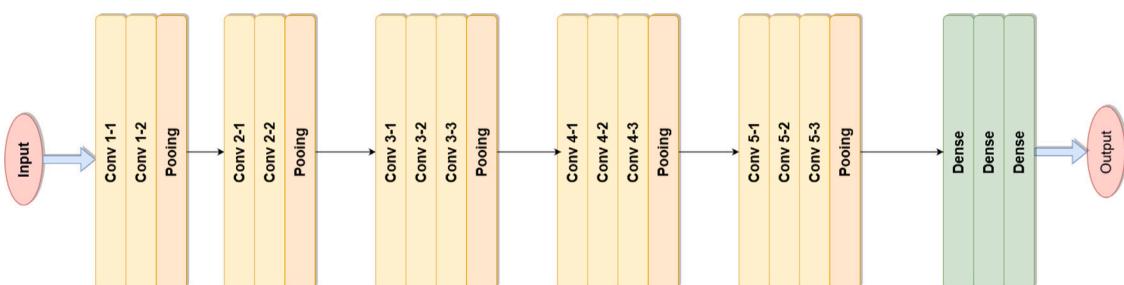


Fig. 12. VGG16 model architecture.

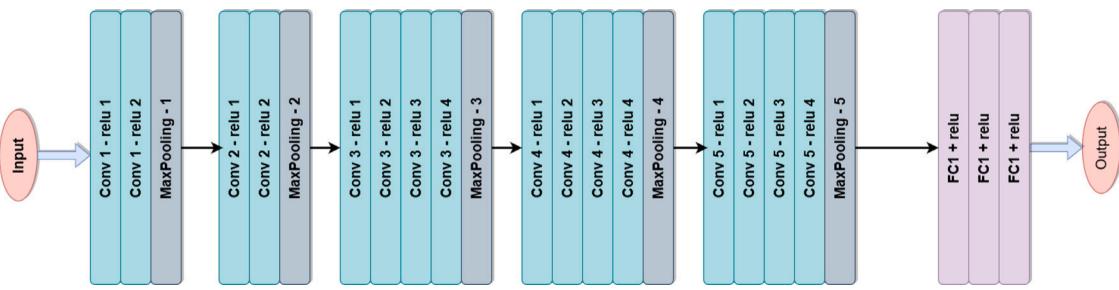


Fig. 13. VGG19 model architecture.

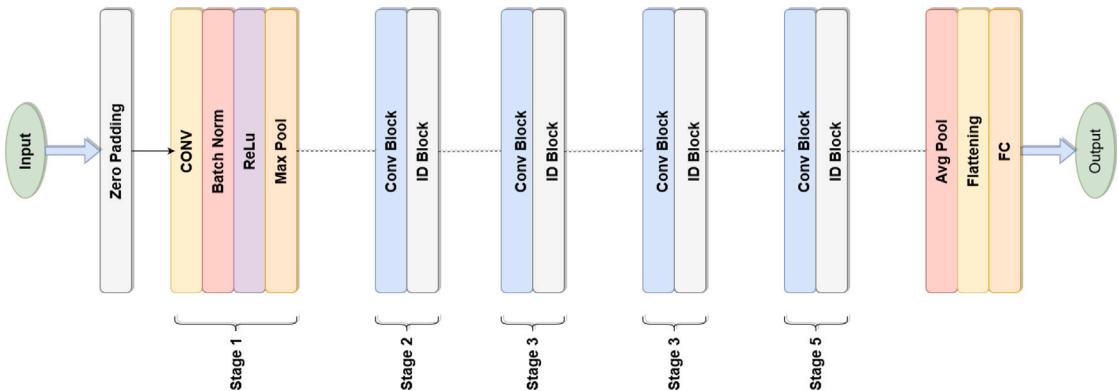


Fig. 14. ResNet50 model architecture.

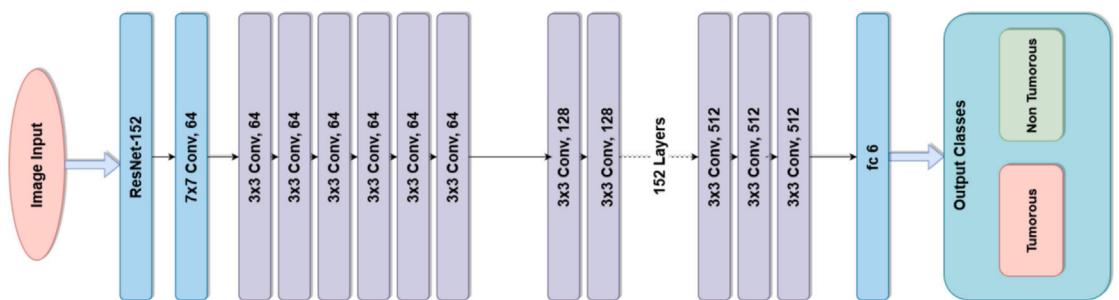


Fig. 15. ResNet152 model architecture.

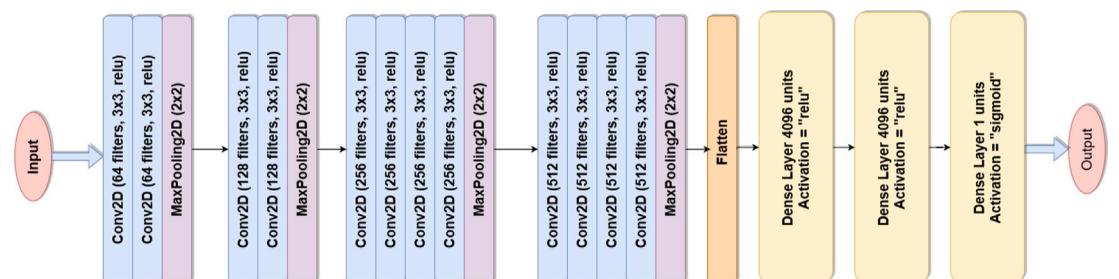


Fig. 16. Custom CNN model architecture.

## 5. Result analysis

We used homomorphic encrypted MRI data for our experiment. In our architecture, this part works as data governance. We tested simple machine-learning algorithms on both encrypted and unencrypted data, as shown in Fig. 17. We saw on unencrypted data we were able to get the highest 97.1% accuracy whereas, on encrypted data, we were able to get the highest 70.12% accuracy. That shows The approach that we followed to encrypt the data is somewhat usable.

For our experiment, we used four different pre-trained deep learning models – VGG16, VGG19, ResNet50, and ResNet152 – as well as a custom CNN model that we developed in-house, to classify pituitary tumor and no tumor from homomorphic encrypted MRI data during federated learning. We evaluated the performance of these models using accuracy and F1-score, as well as precision for each class. These results are from global models of federated learning which we used in our architecture. And these accuracy scores for various models are presented in Table 3.

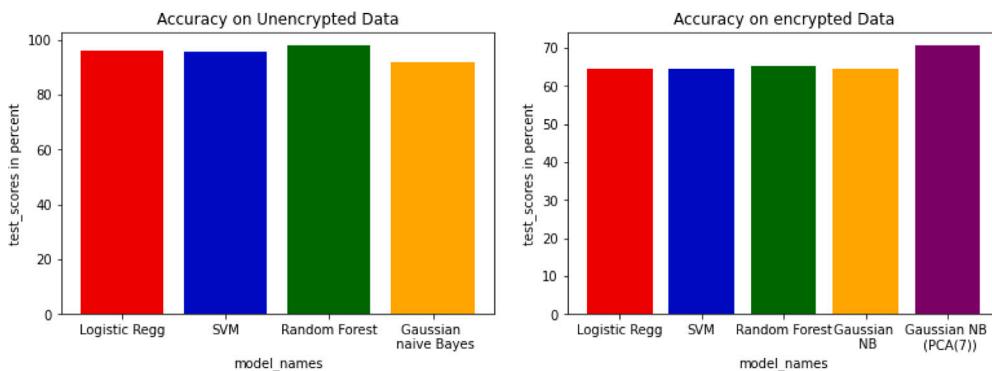


Fig. 17. Performance of tested models on unencrypted and encrypted dataset.

**Table 3**  
Accuracy scores across different models.

Model	Accuracy scores		Precision		Recall	
	Accuracy	F1-Score	Pituitary Tumor	No Tumor	Pituitary Tumor	No Tumor
VGG16	77.25%	77.11%	75.28%	79.22%	80.72%	73.49%
VGG19	78.58%	78.31%	75.82%	81.33%	83.13%	73.49%
ResNet50	61.51%	61.45%	62.38%	60.67%	57.83%	65.06%
ResNet152	65.09%	64.46%	68.18%	62.00%	54.22%	74.70%
Custom CNN	83.31%	83.13%	85.71%	80.90%	79.52%	86.75%

Accuracy is the proportion of correctly classified cases out of all cases. In this study, all five models achieved accuracy greater than 50%, indicating that they performed better than random guessing. The Custom CNN model achieved the highest accuracy of 83.31%, followed by VGG19 with an accuracy of 78.58%, VGG16 with an accuracy of 77.25%, ResNet152 with an accuracy of 65.09%, and ResNet50 with an accuracy of 61.51%.

F1-score is a harmonic mean of precision and recall and is often used as a measure of overall model performance. Looking at the F1-score results, we can see that the Custom CNN model achieved the highest score of 83.13%, followed by VGG19 with an F1-score of 78.31%, VGG16 with an F1-score of 77.11%, ResNet152 with an F1-score of 64.46%, and ResNet50 with an F1-score of 61.45%.

The VGG16 model had a precision of 75.28% for pituitary tumors and 79.22% for no tumors in the binary classification task and 80.72% and 73.49% recall percentages accordingly. This indicates that the model had a moderate ability to correctly identify true positive cases, with slightly higher precision for no tumor cases than for pituitary tumor cases. The associated confusion matrix and training performance are depicted in Fig. 18 and Fig. 19, respectively.

The VGG19 model had the second-highest precision for both classes in the binary classification task. The precision for the pituitary tumor was 75.82% and for no tumor was 81.33%. The recall percentage is 83.13% and 73.49% for the classes. This suggests that the model had a moderate ability to correctly identify true positive cases, with slightly higher precision for no tumor cases than for pituitary tumor cases. Further insights can be gleaned from the confusion matrix presented in Fig. 20 and the training performance depicted in Fig. 21.

The ResNet50 model had the lowest precision for both classes in the binary classification task. The precision for the pituitary tumor was 62.38% and for no tumor was 60.67% and the recall percentage was 57.83% and 65.06% for the classes. This indicates that the model had a lower ability to correctly identify true positive cases, with a higher number of false positives compared to the other models. For a visual representation of these metrics, refer to the confusion matrix in Fig. 22 and the training performance in Fig. 23.

The ResNet152 model had a precision of 68.18% for pituitary tumors and 62.00% for no tumors in the binary classification task while

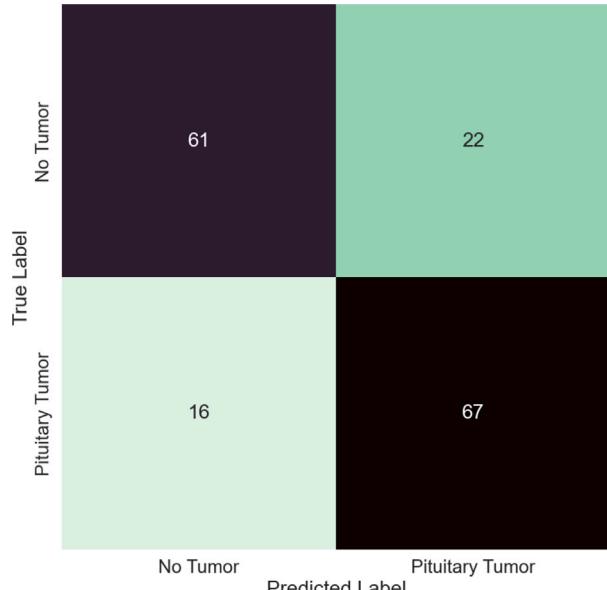


Fig. 18. Confusion matrix of VGG16.

having 54.22% and 74.70% recall. This suggests that the model had a lower ability to correctly identify true positive cases, with a higher number of false positives, particularly for no tumor cases. For a comprehensive view of its performance, refer to the associated confusion matrix in Fig. 24 and the training record in Fig. 25.

The Custom CNN model achieved the highest precision for both classes in the binary classification task. The precision for the pituitary tumor was 85.71% and for no tumor was 80.90%. The recall percentage is 79.52% and 86.75% accordingly. This indicates that the model had a high ability to correctly identify true positive cases, with a relatively low number of false positives. Further insights into its performance can

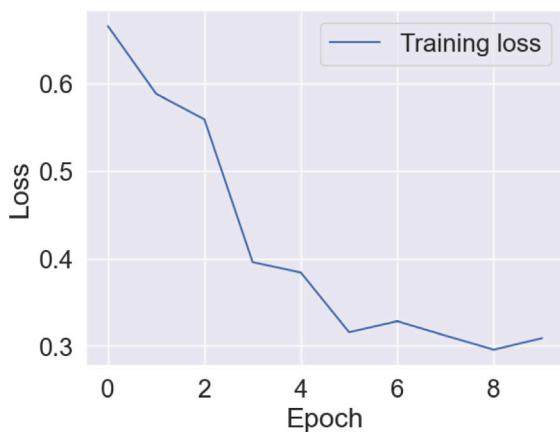


Fig. 19. Training record of VGG16.

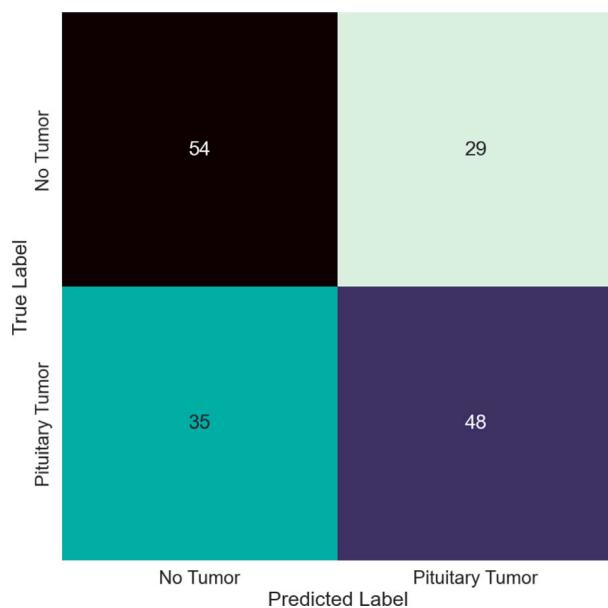


Fig. 22. Confusion matrix of ResNet50.

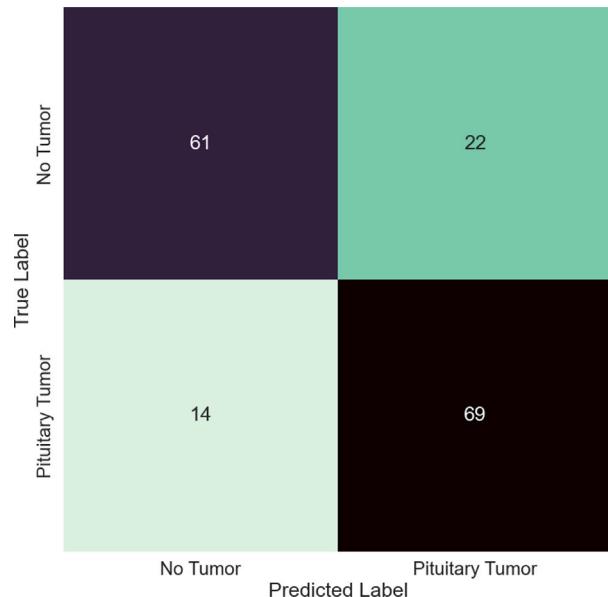


Fig. 20. Confusion matrix of VGG19.

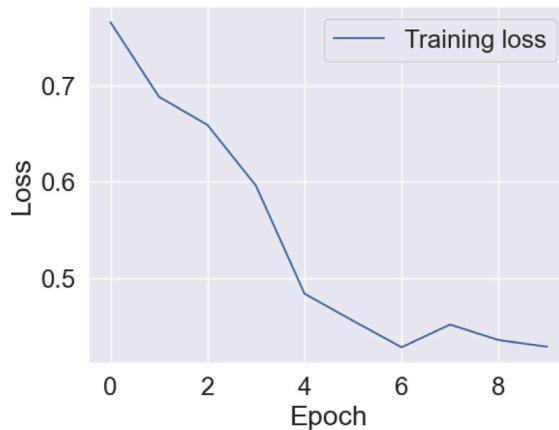


Fig. 23. Training record of ResNet50.

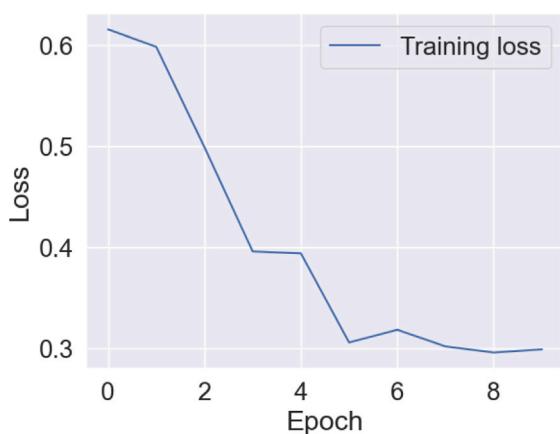


Fig. 21. Training record of VGG19.

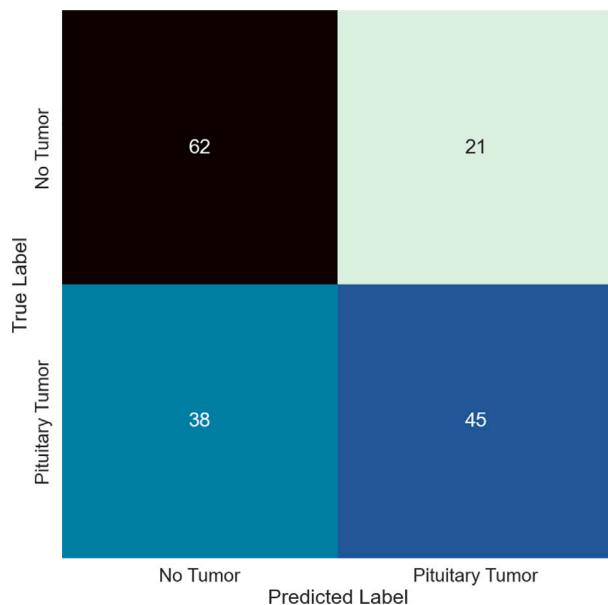


Fig. 24. Confusion matrix of ResNet152.

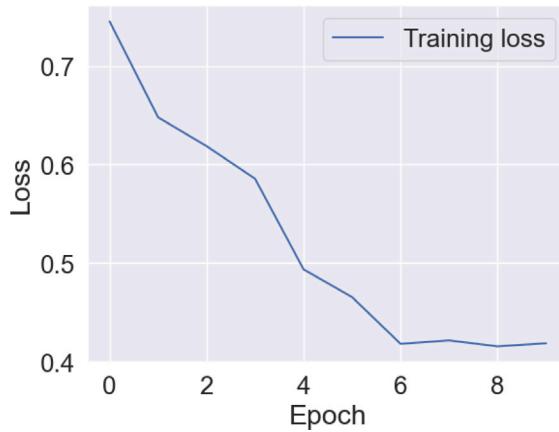


Fig. 25. Training record of ResNet152.

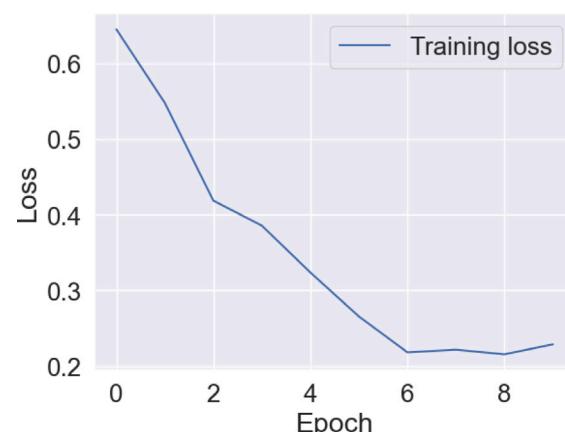


Fig. 27. Training record of custom CNN.

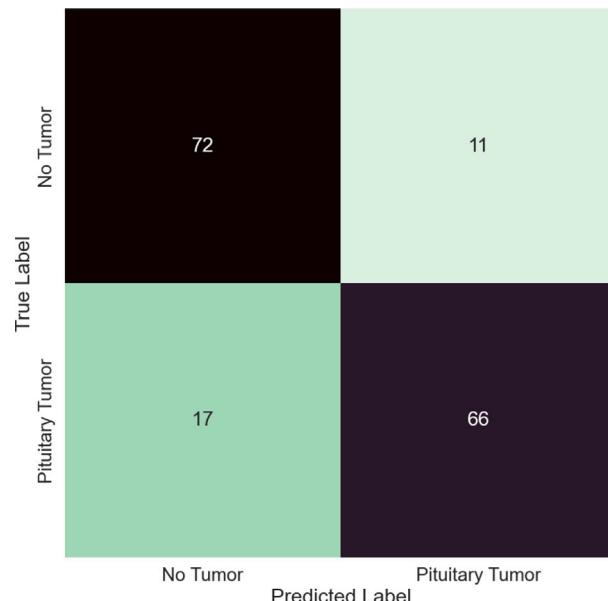


Fig. 26. Confusion matrix of custom CNN.

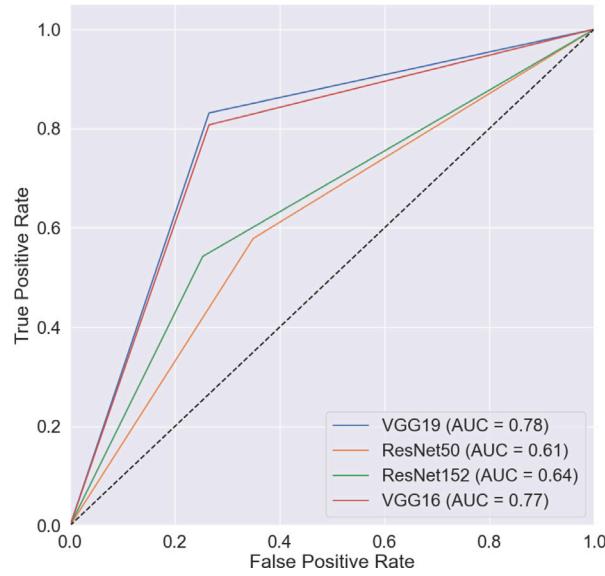


Fig. 28. Receiver operating characteristic of different models.

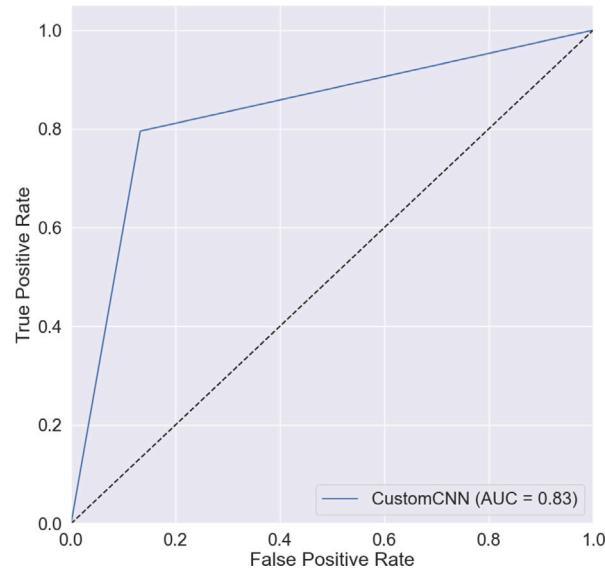


Fig. 29. Receiver operating characteristic of custom CNN.

be gleaned from the associated confusion matrix (Fig. 26) and training records (Fig. 27).

The ROC (Receiver Operating Characteristic) curve is a powerful tool for evaluating the performance of binary classification models. It provides a graphical representation of the trade-off between true positive rate (TPR) and false positive rate (FPR) at various threshold settings, allowing you to choose an appropriate threshold based on your specific needs. A good classifier should have a ROC curve that is as close as possible to the upper left corner of the graph, indicating high TPR and low FPR. By analyzing the ROC curve, you can gain insights into the strengths and weaknesses of your model and make informed decisions about how to improve its performance. The ROC curve is widely used in various fields, including medicine, finance, and machine learning, and is an essential tool for anyone working with binary classification models. The ROC curves of our used DL models are shown in Figs. 28 and 29.

The algorithmic behavior exhibited distinct characteristics across the models. VGG16 and VGG19 demonstrated moderate precision and recall, suggesting a balanced ability to identify true positives, with slightly higher precision for no tumor cases. However, ResNet50 and ResNet152 showcased lower precision and recall, indicating a tendency

for more false positives, particularly in the case of ResNet152 for no tumor instances. The reason for this could be using PHE instead of FHE as ResNet models are quite computationally heavy where PHE falls back. These observations imply that VGG models might have a better trade-off between sensitivity and specificity compared to ResNet models.

On the other hand, our Custom CNN outperformed all of the pre-trained models, aligning with our primary focus on reducing computational complexity, particularly since we utilized PHE as the encryption method. One key factor contributing to the better performance of our Custom CNN model compared to existing models like VGG16, VGG19, and ResNet50 is the number of parameters. In contrast to the 138 million and 144 million parameters of VGG models, Custom CNN utilizes a much smaller set of parameters, approximately 15 million. Additionally, Custom CNN employs a higher number of filters, which significantly aids in achieving greater accuracy. While a higher parameter count typically results in enhanced accuracy for unencrypted data, our architecture, employing Partially Homomorphic Encryption (PHE), contrasts with this outcome. Using PHE can lead to a gradual loss of data integrity when performing scalar multiplication and addition operations on partially encrypted data. This means that in our framework, increasing the complexity of the model does not result in better outcomes due to its effect on data integrity. This highlights the complex interplay between encryption methods, model complexity, and accuracy in our approach to encrypted medical image classification.

Our experiment, which combines homomorphic encryption with federated learning, bears significant promise for advancing the field of medical image analysis and its consequential impact on clinical decision-making. Our data fabric architecture demonstrates potential by ensuring the security of sensitive data through multiple layers, including encryption and federated learning. In terms of clinical decision-making, the impact is substantial. Our test shows a maximum of 83.31% accuracy with satisfactory sensitivity and specificity accordingly 85.71% and 80.90%. In our opinion, these results are quite a milestone for encrypted medical image classification.

## 6. Conclusion

This research demonstrates an advanced data fabric architecture that enables data fusion, integration, and model parameters sharing framework to apply machine learning models without moving the data to a centralized repository. The Partial Homomorphic EncryptionWe and Federated Learning ensured data integrity, privacy, and decentralized learning. We explored the use of pre-trained deep learning models to classify pituitary tumors from homomorphic encrypted MRI data. Our analysis showed that the VGG16 and VGG19 models outperformed the ResNet50 and ResNet152 models in terms of accuracy, precision, recall, and F1-score for both classes. We achieved overall satisfactory accuracy from VGG16 and VGG19. Out of all of these pre-trained models, our custom CNN model performed better with 83.31% accuracy. The reason for the model to perform better is we have made the model similar to VGG16 and VGG19 while making it much less complex. Our model has a total of around 15 million parameters compared to 138 million of VGG16. Despite several significant achievements, the proposed method can be improved by using Fully Homomorphic Encryption. However, it will greatly increase the size of the images, which resulted in significant storage requirements. We have used homogeneous deep learning models for local and global model training in the federated learning approach. The heterogeneous learning models can be used for the robust data fabric architecture,

### 6.1. Assumptions & limitations

We encountered several limitations that hindered our ability to achieve more satisfactory results. One major limitation was the use of Partially Homomorphic Encryption, which only allowed for addition or multiplication operations. However, our data preprocessing, models, and federated process required the ability to perform both operations.

## 6.2. Future directions

A future avenue of exploration encompasses the integration of heterogeneous models to overcome the current user-driven model selection process. This involves developing model fusion strategies that combine diverse models, thus leveraging the collective strengths of each model and eliminating the need for manual selection.

Additionally, the utilization of fully homomorphic encryption is a prospective direction for our advanced data fabric architecture. By adopting this encryption approach, potential constraints on model efficiency attributed to partial homomorphic encryption can be alleviated. In future endeavors, giving priority to essential preprocessing measures has the potential to considerably improve outcomes and overall performance.

## CRediT authorship contribution statement

**Sakib Anwar Rieyan:** Conceptualization, Methodology, Software.  
**Md. Raisul Kabir News:** Conceptualization, Methodology, Software.  
**A.B.M. Muntasir Rahman:** Conceptualization, Methodology, Software.  
**Sadia Afrin Khan:** Conceptualization, Methodology, Software. **Sultan Tasneem Jawad Zaarif:** Conceptualization, Methodology, Software. **Md. Golam Rabiul Alam:** Supervisor, Writing – original draft. **Mohammad Mehedi Hassan:** Software, Writing – review & editing. **Michele Ianni:** Writing – review & editing. **Giancarlo Fortino:** Writing – review & editing.

## Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Data availability

Data will be made available on request.

## Acknowledgments

This work was supported by the King Saud University, Riyadh, Saudi Arabia, through the Researchers Supporting Project under Grant RSP2023R18. We also acknowledge support from the project: PNRR MUR project PE0000013-FAIR;

## References

- [1] A.R. Murray-Watson, Healthcare data breach statistics, HIPAA J. (2023) URL <https://www.hipaajournal.com/healthcare-data-breach-statistics/>.
- [2] P. Qi, D. Chiaro, F. Piccialli, FL-FD: Federated learning-based fall detection with multimodal data fusion, Inf. Fusion (2023) 101890.
- [3] C. Ounoughi, S. Ben Yahia, Data fusion for ITS: A systematic literature review, Inf. Fusion 89 (2023) 267–291.
- [4] A. Imakura, T. Sakurai, Y. Okada, T. Fujii, T. Sakamoto, H. Abe, Non-readily identifiable data collaboration analysis for multiple datasets including personal information, Inf. Fusion 98 (2023) 101826.
- [5] U.D. of Health, H. Services, et al., Public law 104-191: Health insurance portability and accountability act of 1996, 2003, Retrieved November 24 2003.
- [6] R. Viorescu, et al., 2018 reform of eu data protection rules, Eur. J. Law Public Administr. 4 (2) (2017) 27–39.
- [7] L. Zhang, Y. Xie, L. Xidao, X. Zhang, Multi-source heterogeneous data fusion, in: 2018 International Conference on Artificial Intelligence and Big Data, ICAIBD, IEEE, 2018, pp. 47–51.
- [8] M. Al-Hawawreh, M.S. Hossain, A privacy-aware framework for detecting cyber attacks on internet of medical things systems using data fusion and quantum deep learning, Inf. Fusion 99 (2023) 101889.
- [9] K. Wang, C.-M. Chen, Z. Liang, M.M. Hassan, G.M. Sarné, L. Fotia, G. Fortino, A trusted consensus fusion scheme for decentralized collaborated learning in massive IoT domain, Inf. Fusion 72 (2021) 100–109.
- [10] M.M. Hassan, M.G.R. Alam, M.Z. Uddin, S. Huda, A. Almogren, G. Fortino, Human emotion recognition using deep belief network architecture, Inf. Fusion 51 (2019) 10–18.

- [11] C. Savaglio, M. Ganzha, M. Paprzycki, C. Bădică, M. Ivanović, G. Fortino, Agent-based Internet of Things: State-of-the-art and research challenges, *Future Gener. Comput. Syst.* 102 (2020) 1038–1053.
- [12] G. Fortino, W. Russo, C. Savaglio, W. Shen, M. Zhou, Agent-oriented cooperative smart objects: From IoT system design to implementation, *IEEE Trans. Syst. Man. Cybern. Syst.* 48 (11) (2017) 1939–1956.
- [13] L. Fotia, F. Delicato, G. Fortino, Trust in edge-based internet of things architectures: State of the art and research challenges, *ACM Comput. Surv.* 55 (9) (2023) 1–34.
- [14] M.G.R. Alam, M. Haque, M.R. Hassan, S. Huda, M.M. Hassan, F.L. Strickland, S.A. AlQahtani, Feature cloning and feature fusion based transportation mode detection using convolutional neural network, *IEEE Trans. Intell. Transp. Syst.* 24 (4) (2023) 4671–4681.
- [15] A.H. Seh, M. Zarour, M. Alenezi, A.K. Sarkar, A. Agrawal, R. Kumar, R. Ahmad Khan, Healthcare data breaches: Insights and implications, in: *Healthcare. Vol. 8. No. 2, MDPI*, 2020, p. 133.
- [16] H.K. Patil, R. Seshadri, Big data security and privacy issues in healthcare, in: *2014 IEEE International Congress on Big Data*, IEEE, 2014, pp. 762–765.
- [17] Gartner, Definition of data fabric-gartner information technology glossary, 2023, URL <https://www.gartner.com/en/information-technology/glossary/data-fabric>.
- [18] K. Bonawitz, H. Eichner, W. Grieskamp, D. Huba, A. Ingeman, 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* 1 (2019) 374–388.
- [19] L. Morris, Analysis of Partially and Fully Homomorphic Encryption, Vol. 10, Rochester Institute of Technology, 2013, pp. 1–5.
- [20] P. Paillier, Public-key cryptosystems based on composite degree residuosity classes, in: *International Conference on the Theory and Applications of Cryptographic Techniques*, Springer, 1999, pp. 223–238.
- [21] S.B. Belhaouari, Y. Hamdi, A. Hamdi, Revisited Carmichael's reduced totient function, *Mathematics* 9 (15) (2021) 1800.
- [22] C. Stamatellis, P. Papadopoulos, N. Pitropakis, S. Katsikas, W.J. Buchanan, A privacy-preserving healthcare framework using hyperledger fabric, *Sensors* 20 (22) (2020) 6587.
- [23] A. Roehrs, C.A. Da Costa, R. da Rosa Righi, OmniPHR: A distributed architecture model to integrate personal health records, *J. Biomed. Inf.* 71 (2017) 70–81.
- [24] Q. Xia, E.B. Sifah, K.O. Asamoah, J. Gao, X. Du, M. Guizani, MeDShare: Trust-less medical data sharing among cloud service providers via blockchain, *IEEE Access* 5 (2017) 14757–14767.
- [25] Y. Ming, T. Zhang, Efficient privacy-preserving access control scheme in electronic health records system, *Sensors* 18 (10) (2018) 3520.
- [26] J. Fu, J. Xu, S. Zhang, C. Zhang, Research and design of square kilometer array astronomical data management model based on fabric, in: *2020 International Conferences on Internet of Things (IThings) and IEEE Green Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCom) and IEEE Smart Data (SmartData) and IEEE Congress on Cybernetics, Cybermatics*, IEEE, 2020, pp. 513–518.
- [27] M.S. Ram, A. Seeram, M. Poongundran, P. Singh, Y.N. Prajapati, S. Myrzahmetova, Data fusion opportunities in IoT and its impact on decision-making process of organisations, in: *2022 6th International Conference on Intelligent Computing and Control Systems*, ICICCS, IEEE, 2022, pp. 459–464.
- [28] Z. Yin, W. Junli, On traffic management integrated information fusion model, in: *2008 27th Chinese Control Conference*, IEEE, 2008, pp. 546–549.
- [29] F. Wibawa, F.O. Catak, M. Kuzlu, S. Sarp, U. Cali, Homomorphic encryption and federated learning based privacy-preserving cnn training: Covid-19 detection use-case, in: *Proceedings of the 2022 European Interdisciplinary Cybersecurity Conference*, 2022, pp. 85–90.
- [30] W. Ou, J. Zeng, Z. Guo, W. Yan, D. Liu, S. Fuentes, A homomorphic-encryption-based vertical federated learning scheme for rick management, *Comput. Sci. Inform. Syst.* 17 (3) (2020) 819–834.
- [31] I. Kotisuba, A. Velvzhanin, Y. Yanovich, I.S. Bandurova, Y. Dyachenko, V. Zhygulin, Decentralized e-health architecture for boosting healthcare analytics, in: *2018 Second World Conference on Smart Trends in Systems, Security and Sustainability (WorldsS4)*, IEEE, 2018, pp. 113–118.
- [32] C. Zhang, R. Ma, S. Sun, Y. Li, Y. Wang, Z. Yan, Optimizing the electronic health records through big data analytics: a knowledge-based view, *IEEE Access* 7 (2019) 136223–136231.
- [33] F. Rahman, M. Slepian, A. Mitra, A novel big-data processing framwork for healthcare applications: Big-data-healthcare-in-a-box, in: *2016 IEEE International Conference on Big Data, Big Data*, IEEE, 2016, pp. 3548–3555.
- [34] M.J. Kaur, V.P. Mishra, Analysis of big data cloud computing environment on healthcare organizations by implementing Hadoop clusters, in: *2018 Fifth HCT Information Technology Trends*, ITT, IEEE, 2018, pp. 87–90.
- [35] P. Voigt, A. Von dem Bussche, The eu general data protection regulation (GDPR), in: *A Practical Guide*, Vol. 10, No. 3152676, first ed., Springer International Publishing, Cham, 2017, pp. 10–5555.
- [36] M. Nickparvar, Brain tumor MRI dataset, 2021, URL <https://www.kaggle.com/datasets/masoudnickparvar/brain-tumor-mri-dataset>.
- [37] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, 2014, arXiv preprint [arXiv:1409.1556](https://arxiv.org/abs/1409.1556).
- [38] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.