



# Federated Vision-Language Models for Privacy-Preserving Medical Image Analysis

Singamaneni Krishnapriya<sup>\*1</sup>, Amaravarapu Pramod kumar<sup>1</sup>, Boddupally Janaiah<sup>2</sup>, K S Ranadheer Kumar<sup>3</sup>

<sup>1</sup>CSE-(CyS,DS) and AI & DS), VNR Vignana Jyothi Institute of Engineering and Technology,  
[singamanenikrishnapriya@gmail.com](mailto:singamanenikrishnapriya@gmail.com), [amaravarapupramod@gmail.com](mailto:amaravarapupramod@gmail.com)

<sup>2</sup>MVSR Engineering College, [janaish\\_cse@mvsrec.edu.in](mailto:janaish_cse@mvsrec.edu.in)

<sup>3</sup>Department of CSE (Data Science), CVR College of Engineering, [ranadheer.k.s@gmail.com](mailto:ranadheer.k.s@gmail.com)

\*Correspondence: [singamanenikrishnapriya@gmail.com](mailto:singamanenikrishnapriya@gmail.com)

## Abstract

Deep learning has enhanced the analysis of medical images but privacy issues and institutional variations restrict their large scale application in clinics. FedVLM, a federated vision language model tailored to privacy-preserving multimodal medical image analysis, is one of the solutions to these problems. Contrary to the conventional federated design, which can only process single modal image data, FedVLM takes paired radiological images and clinical reports jointly, which demonstrates high zero-shot and few-shot diagnostic performance. The design consists of secure aggregation, differential privacy and proximal optimization that ensure protection of patient data and minimize variability across sites. Large scale experiments on the NIH ChestX-ray14, MIMIC-CXR, and BraTS datasets indicate that FedVLM is an accurate and interpretable model that achieves near-centralized performance on vision language models without violating privacy. Building on previous works such as FACMIC, BioViL, and FAA-CLIP, FedVLM introduces new methods, including privacy-aware optimization, proximal regularization for varied data, and multimodal contrastive alignment, creating a unified federated framework for clear and secure medical image analysis. Although FedVLM shows promising performance, this work is currently at a research stage and is not yet ready for clinical use. We need validation through future multi-institutional clinical studies.

**Keywords:** Medical Image Analysis, Federated Learning, Vision–Language Models, Privacy-Preserving AI, Clinical Decision Support

Received: October 20<sup>th</sup>, 2025 / Revised: January 10<sup>th</sup>, 2026 / Accepted: January 21<sup>st</sup>, 2026 / Online: January 23<sup>rd</sup>, 2026

## I. INTRODUCTION

Analysis of medical images has become one of the most impactful areas of artificial intelligence (AI) and they can be used in detecting various diseases, organ localization, predicting survival and planning treatment. Convolutional neural networks (CNN) and transformers have demonstrated human-competitive or even better performance on a diverse range of medical tasks through deep learning models, with many positional statements grounded on the principle of human-comparable performance [1]. Regardless of these developments, there are two serious bottlenecks in the large-scale implementation in real clinical environments.

To begin with, the privacy limitation of the data. Medical imaging information is very sensitive, their use is controlled within strict frameworks, which are the Health Insurance Portability and Accountability Act (HIPAA) in the United States and the General Data Protection Regulation (GDPR) in Europe. These limitations do not allow direct data merging between institutions, making the creation of centralized deep learning models based on a variety of large-scale datasets

difficult to create. The idea of federated learning (FL) has been developed to address this problem by permitting distributed data silos to train collaboratively without having raw data to leave the local institution at all, in the future, [2], [3].

Second, the heterogeneity of the domain. There is a high level of variance among medical images between institutions due to the variation in acquisition equipment, imaging procedures, patient groups, and annotation standards. This causes nonidentically distributed (non-IID) data, which can have a severe negative effect on the performance of centralized and federated models when applied to unseen domains [4], [5]. Therefore, it is necessary to handle this heterogeneity to develop clinically robust and generalizable AI systems.

Moreover, the concept of vision language models (VLMs) has become a popular topic of AI discussions. One model like CLIP [6] and BLIP [7] can bring visual and textual modalities to a common embedding space, and thus provide strong zero-shot and few-shot generalizations on natural image domains. Based on these achievements, new studies have begun to work on medical vision language models, including MedCLIP and BioViL [8], [9], that match radiology images to the relevant

clinical report. Such methods reveal that VLMs can be used in multimodal medical tasks, such as disease classification, the generation of reports, and cross-modal retrieval. However, to train such models, mainly centralized datasets are used, which raises issues of privacy leakage and institutional data silos. In addition, there is no extensive generalization in hospitals due to changes in the distribution of both types of imaging modes and clinical reporting practices.

Unlike FACMIC [10], which uses CLIP in a federated setting without clear privacy protections, or FAA-CLIP [11], which mainly emphasizes attention-based personalization, the proposed **FedVLM** introduces a unified multimodal federated optimization pipeline. This pipeline includes *secure aggregation*, *differential privacy*, and *proximal regularization*. Additionally, it goes beyond centralized vision and language frameworks like BioViL [12] by allowing cross-institutional training on paired radiology images and clinical reports without sharing raw data. This combination of privacy, multimodal alignment, and knowledge-aware optimization represents the main innovation of **FedVLM**.

As a solution to these drawbacks, we present a federated scheme, called FedVLM, to generalize vision-language models into privacy-preserving multiinstitutional industry analysis of medical images. FedVLM allows cross-site cooperation, which means the combined optimization of the multimodal representations without sharing raw data. With the union of the federated optimization and domain adaptation and communication-efficient strategies, our framework can offer the concepts of privacy preservation and robust generalization in the heterogeneous clinical setting. In contrast to previous publications in the field of federated medical imaging that focus only on unimodal image classification or segmentation protocols, FedVLM presents multimodal alignment strategies, which is why it enables the incorporation of textual supervision (e.g. radiology report) in order to achieve better interpretability and downstream performance.

The key contributions of this work can be summarized as follows in the Table I.

- We introduce **FedVLM**, the first federated framework for vision–language models in healthcare, bridging the gap between multimodal representation learning and privacy-preserving medical image analysis.
- We design a **privacy-aware and communication-efficient** learning protocol that leverages secure aggregation, differential privacy, and lightweight parameter updates to reduce bandwidth overhead while protecting sensitive information.
- We suggest a domain-aware alignment framework, which formally addresses the challenge of interinstitutional heterogeneity and enhances the ability to generalize to different imaging modalities and clinical scenarios through the introduction of multimodal representations.
- We carry out an in-depth empirical analysis of large-scale medical data, such as chest radiographs and skin lesion

images, and show that FedVLM is as competitive as centralized training and that it provides strict privacy guarantees.

TABLE I. NOVELTY &amp; CONTRIBUTIONS SUMMARY

Aspect	FedVLM Novelty and Contribution
Problem Addressed	Unified privacy-preserving multimodal learning across distributed medical institutions.
Key Innovation	Initial implementation of Vision Language Models (VLMs) on a federated network of medical image analysis.
Architectural Advances	Integrates secure aggregation, differential privacy and proximal optimization in multimodal federated training.
Compared to FACMIC	Introduces formal privacy, multimodal text image alignment (FACMIC is one epoch unidimensional and privacy blind).
Compared to FAA-CLIP	Introduced the domain-sensitive multimodal personalization; FAA-CLIP is concerned with attention personalization.
Compared to BioViL	Allows institutel training and protection of privacy, unlike central training of BioViL.
Outcome	Almost centrally accurate, highly interpretable and with high data leakage resistance.

In summary, this paper positions FedVLM as a step toward *scalable, explainable, and privacy-preserving AI for medical imaging*, laying the foundation for real-world deployment of federated multimodal systems in healthcare.

## II. RELATED WORKS

### A. Federated Learning in Medical Imaging

Effective baselines in medical image segmentation were achievable with deep learning models including U-Net, which can form deeper bases by frequently stacking encoders and decoders to produce fine-grained predictions on image pixel intensity, data volume, and channel depth, which are beyond the abilities of traditional approaches to segmentation and localization tasks (Gallerati). But centralized training of data collections on pools is not respected of privacy and governance constraints. Introduced by FedAvg 2017 by Federated Learning architectures, Federated grid workshear was closed by federated model access, typically known as FedProx. Federated grid model training is a learning framework where teams of institutions collaborate and do not share raw data.

Sheller et al. [2] who performed a Mult institutional segmentation of the brain tumor, and Dayan et al. [13] created the EXAM model to predict the results of COVID-19 in 20 hospitals. Currier: Large benchmarks, [14], [15] offer performance degradations under distribution shift, and algorithm family members are compared with [16], [17]. Non-IID data, communication overhead and strong privacy guarantees are the recurrent problem highlighted in surveys [3], [18], [19], including but not limited to those of Rieke (2020) Future, p. 5), Rehman (2023) Federated and teo (2024) Systematic.

A representative example of the FL literature in medical imaging for the 2018-2025 time frame is provided in Table II, and it focuses on modality, tasks, sites and privacy mechanisms.

TABLE II. REPRESENTATIVE FL IN MEDICAL IMAGING (2018–2025).

Work (Year)	Modality/Task	Sites	Method	Personalization	Privacy	Key Finding
Sheller et al. [2] (2018)	MRI / Seg.	4+	FedAvg	–	–	Fed $\approx$ Centralized Dice on BraTS-like data
Dayan et al. [13] (2021)	CXR+EHR / Clf.	20+	FedAvg	–	SA	Cross-system generalization for COVID-19 outcomes
FeTs [4], [5] (2021–22)	MRI / Seg.	20+	FedAvg variants	Sampling tricks	–	Benchmarking FL under distribution shifts
Luo et al. [16] (2023)	CT/MR / Seg.	Multi-site	FedAvg	–	–	Fed. probes with larger intersection distance
Mantel et al. [17] (2024)	MRI / Seg.	Multi-site	FedAvg/Prox/Per.	FedBN etc.	–	Benchmark across algorithm classes
Wu et al. [1] (2024)	MRI / Seg.	Multi-site	Fed Contrastive	–	–	Self/Semi-supervised FL improves label efficiency

ABBREV.: CLF.=CLASSIFICATION, SEG.=SEGMENTATION, SA=SECURE AGGREGATION, DP=DIFFERENTIAL PRIVACY.

### B. Vision-Language Models in Healthcare

Vision-language models (VLMs) align images and text in a shared embedding space. ConVIRT [20] demonstrated contrastive learning on paired chest X-rays and reports. Subsequent works such as MedCLIP [21], CheXzero [22], PMC-CLIP [23], and KAD [24] extended CLIP-style pretraining for radiology, achieving zero-shot classification and retrieval. More recent efforts such as BioViL and BioViL-T [9], [25] improve report generation and retrieval, while surveys [26], [27] highlight both progress and limitations.

Despite promising centralized results, these models often underperform in cross-site evaluations due to reporting style

variations, domain shifts, and lack of federated adaptation. A comparative overview of representative VLMs from 2020–2025 is provided in Table III.

Figure 1 presents an overview of previous studies, the related research gaps, and the main contributions of the proposed FedVLM framework. Earlier works mainly focused on centralized or unimodal systems, which had limited privacy protection.

The proposed model tackles these problems by using multimodal federated learning that includes privacy protection and better cross-site generalization.

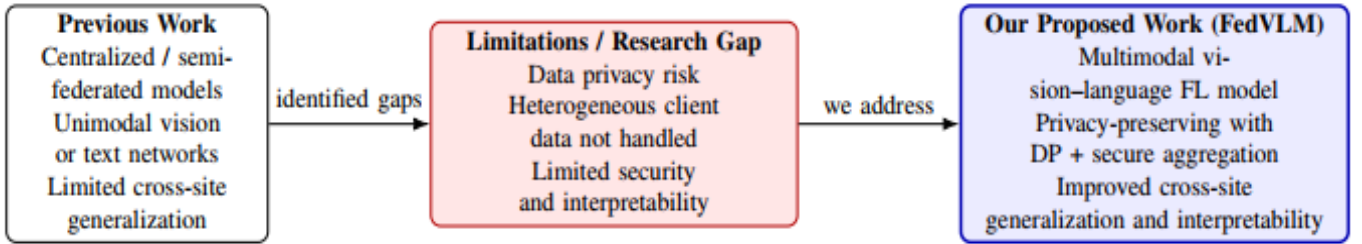


Fig. 1. Overview of prior work, research gaps, and contributions of the proposed FedVLM framework.

### C. Gap Analysis

During the period 2015 to 2025, FL has developed into a practical paradigm in the field of medical imaging, and VLMs have demonstrated the promises of multimodal alignment of interpretability and zero-shot transfer. However, as depicted in

Fig. 2, these two focus areas of research have not yet been fully engaged with each other, with slight integration. Recent efforts such as FACMIC [10] and FAA-CLIP [11] attempt to adapt CLIP within federated environments, but are limited to classification tasks, lack robust privacy mechanisms (e.g., secure aggregation and differential privacy), and do not address cross-site semantic heterogeneity. No existing framework unifies FL and VLMs to deliver:

1) **Privacy-preserving multimodal alignment** across image and text modalities without raw data sharing.

2) **Cross-institutional generalization** that explicitly mitigates domain heterogeneity.

3) **Explainable zero-shot and few-shot capabilities** for medical AI systems.

This gap clearly motivates our proposed **FedVLM**, which, to the best of our knowledge, is the first federated vision-language framework designed for large-scale, privacy-preserving, and multimodal medical image analysis.

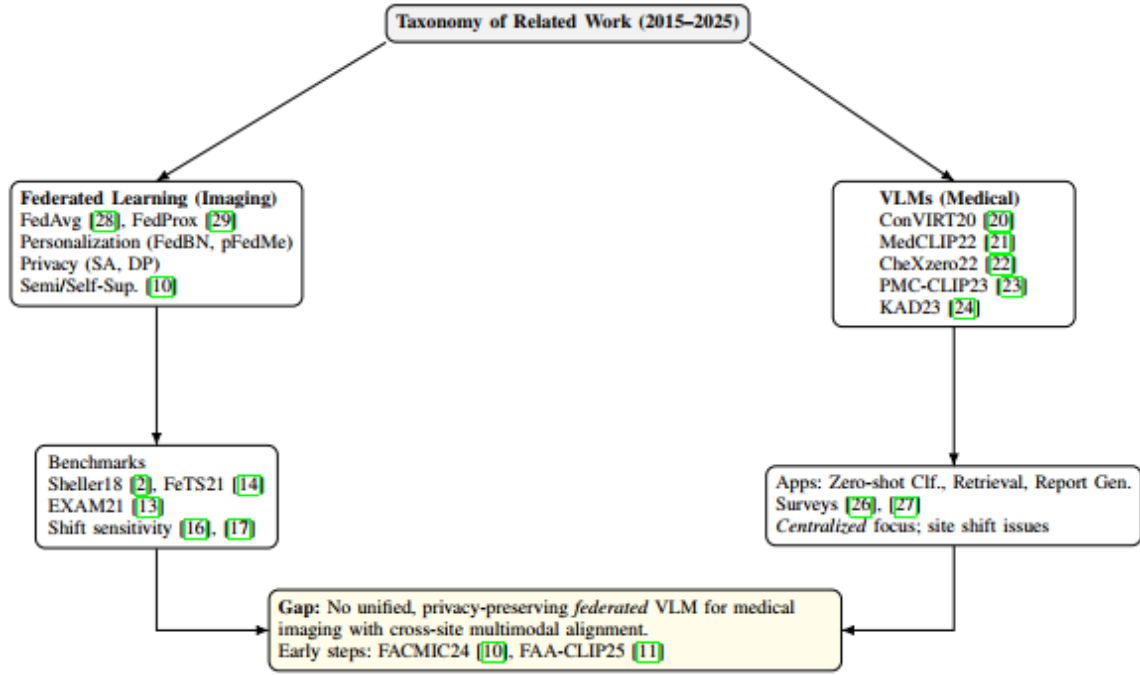


Fig. 2. Taxonomy summarizing prior art in FL and VLMs and the gap addressed by FedVLM.

TABLE III. REPRESENTATIVE VISION–LANGUAGE MODELS IN HEALTHCARE (2020–2025).

Work (Year)	Pretraining Data	Text Source	Zero-Shot	Tasks	Cross-Site Notes
ConVIRT [20] (2020)	CXR, others	Reports	No	Cfr., Retrieval	Sensitive to site/report shifts
MedCLIP [21] (2022)	20K pairs	Captions/Reports	Yes	Clf., Retrieval	Limited external validation
CheXCLIP [22] (2022)	377 CXR pairs	Reports	Yes	Zero-shot CLF	Crosshospital evals show gaps
PMC-CLIP [23] (2023)	1.6M PMC pairs	Captions	Yes	Z.retrieval, VQA	Style mismatch to hospitals
KAD [24] (2023)	X-ray corpus	Reports+ Knowledge	Yes	Ext. Clf.	Knowledge boosts zero-shot, shift remains
BioViL/Report-VLMs [9], [25] (2022–24)	Rad. images+reports	Reports	Gen.	Report Gen./Retrieval	No federated evals

### III. PROPOSED METHODOLOGY

In this section, we present **FedVLM**, a federated vision–language model framework for privacy-preserving medical image analysis. The framework unifies multimodal alignment with federated optimization across heterogeneous institutions while ensuring strict privacy guarantees. Fig. 3 illustrates the overall pipeline.

#### A. Overall Architecture

Under the FedVLM, hospitals are the clients and each hospital has its local collection of matched medical images and corresponding clinical text (e.g., radiology reports). Making alignment Multimodal encoders are trained locally to generate

mutually aligned encodings across vision and language. Rather than exchange the raw data a central server is updated with only model updates (gradients or encoder parameters) and then conducts secure global aggregation. An integrated model of the globe is reallocated to the clients and allows the elegant collaborative enhancement without jeopardising the confidentiality of data. This pipeline is based on a regular cross-silo federated learning model prescribed by [28], [30], but specific to multimodal medical environments.

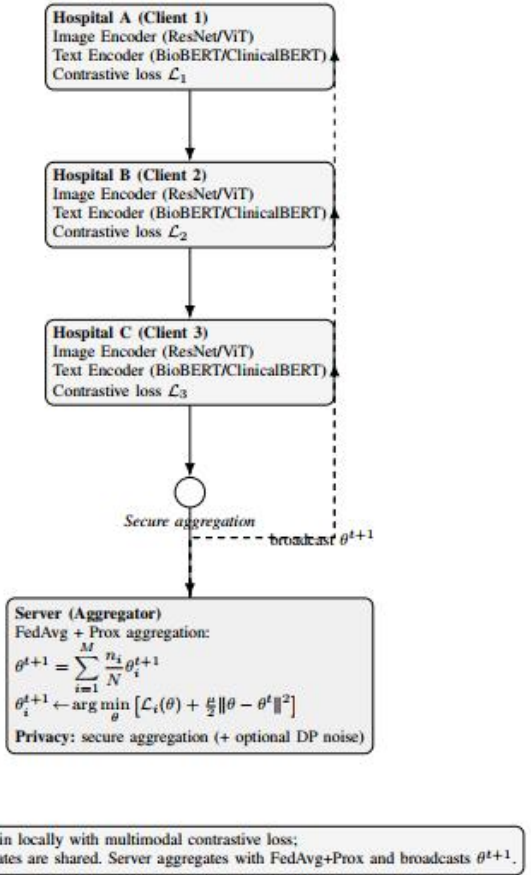


Fig. 3. Vertical FedVLM pipeline: hospitals (clients) train local image and text encoders with contrastive alignment. Updates are securely aggregated, the server performs FedAvg+Prox, and the global model is broadcast back.

### B. Novel Architectural Contributions of FedVLM

To showcase the technical innovations, the proposed **FedVLM** framework is designed to address the limitations of earlier works like FACMIC [10], BioViL [12], and FAA-CLIP [11]. FACMIC uses CLIP in a federated environment for image classification, but it lacks formal privacy guarantees. The FAA-CLIP provides attention-based personalization, but it can only be used in unimodal or partially aligned learning environments. On the contrary, BioViL uses centralized training assumptions, thus it cannot ensure that data is confidential across institutions. On the contrary, FedVLM brings a number of important architectural differences:

- **Federated Multimodal Alignment:** Vision and language encoders are simultaneously trained by FedVLM in a federated form. It relies on contrastive alignment objective. This enables dissimilar destinations to streamline without the need to share crude information.
- **Privacy-Aware Optimization:** The model is a combination of secure aggregation and differential privacy in training loop. That ensures that institutional updates remain encrypted and in line with the privacy requirements.
- **Proximal Regularized Federated Learning:** FedAvg + Prox is a training stabilization method used on non-IID multimodal

data distributions. It also reduces client drift. This is superior to the simple FedAvg designs that were used in the past studies.

- **Domain-Aware Alignment:** FedVLM lessens the variations in the meaning and images across sites by incorporating paired image and text embeddings. This is useful in bringing about uniform multimodal generalization.

- **Explainability Integration:** The interpretability of the Grad-CAM provides the clear justifications of the predictions, and the token-level alignment provides the clear justifications of the prediction tokens. This enhances the real-world deployment.

Together, FedVLM is novel, with a single combination of vision language modeling and federated learning to analyse medical images with privacy protection. Unlike the scenario in FACMIC where there is CLIP being applied in a federated model although not providing express coverage of privacy and multimodal consistency, FedVLM has incorporated secure aggregation, differential privacy, and proximal regularization, built directly as part of the learning pipeline. FedVLM uses domain-conscious multimodal optimization to coordinate radiological images and textual reports across institutions compared to FAA-CLIP, which mainly emphasizes the attention-based personalization of systems in unimodal settings. Additionally, unlike BioViL, where the centralized access to data is assumed and privacy cannot be controlled, FedVLM allows crossinstitutional models to train without sharing raw data. FedVLM can be viewed as one of the earliest federated vision language models that provide a solid state of diagnostic results, interpretation, and privacy guarantees for heterogeneous clinical aspects, all at once.

### C. Model Components

**Image Encoder.** Our network uses CNN or ViT backbones which are pre-trained on large-scale imaging data (e.g., ResNet [31] or ViT [32] and trained on domain-specific medical images.

**Text Encoder.** Our transformer-based biomedical language models include BioBERT model, e.g., BioBERT model [33], or ClinicalBERT model, e.g. ClinicalBERT model [34] and these models are able to learn domain specific semantics of a clinical report.

**Cross-Modal Alignment.** In order to match modalities, we use a contrastive learning task based on CLIP, which maximizes similarities between paired image-text embeddings, and minimizes similarities between paired images and negative pairs, as used in CLIP. which are also summarized in Table IV.

TABLE IV. MODEL COMPONENTS OF FEDVLM

Component	Choice	Reference
Image Encoder	ResNet, ViT	He et al. [31]; Dosovitskiy et al. [32]
Text Encoder	BioBERT, ClinicalBERT	Lee et al. [33]; Alsentzer et al. [34]
Alignment	Contrastive Loss	Radford et al. [36]
Privacy	Secure Aggregation, DP	Kassis et al. [30]



#### D. Federated Optimization

To train FedVLM under heterogeneous non-IID data distributions, we adopt **FedAvg** [28] with a proximal regularization term (as in FedProx [29]) to stabilize local updates.

$$\theta_i^{t+1} = \theta^t - \eta \Delta L_i(\theta^t) + \mu(\theta^t - \theta_i^t) \quad (1)$$

where  $\theta^{t+1}_i$  denotes the updated parameters of the client  $i$ ,  $\eta$  is the learning rate,  $\mu$  is the proximal coefficient and  $L_i$  is the local multimodal loss. A secure aggregation protocol [35] ensures that the server only receives encrypted aggregated updates, preserving institutional privacy. The global aggregation step averages the weighted updates:

$$\theta^{t+1} = \sum_{i=1}^M \frac{n_i}{N} \theta_i^{t+1} \quad (2)$$

where  $M$  is the number of clients,  $n_i$  is the number of local samples at the client  $i$ , and  $N = \sum n_i$ .

The local objective combines unimodal reconstruction with multimodal contrastive alignment. For client  $i$  with dataset

$D_i = \{(x^I_j, x^T_j)\}$  of image–text pairs:

$$\mathcal{L}_i(\theta) = \frac{1}{|D_i|} \sum_{(x^I, x^T) \in D_i} \{\mathcal{L}_{img}(f_I(x^I; \theta_I)) + \mathcal{L}_{text}(f_T(x^T; \theta_T)) + \lambda \mathcal{L}_{contrast}(f_I(x^I), f_T(x^T))\} \quad (3)$$

where  $f_I$  and  $f_T$  denote image and text encoders,  $\mathcal{L}_{img}$  and  $\mathcal{L}_{text}$  are unimodal cross-entropy/reconstruction losses,  $\mathcal{L}_{contrast}$  is a contrastive loss:

$$\mathcal{L}_{contrast} = -\log \frac{\exp(\frac{\text{sim}(z_I, z_T)}{\tau})}{\sum_{k=1}^B \exp(\frac{\text{sim}(z_I, z_T^k)}{\tau})} \quad (4)$$

with  $z_I, z_T$  being normalized embeddings,  $\tau$  a temperature parameter, and  $B$  the batch size. This enforces alignment between paired image–text samples while contrasting with negatives.

#### E. Security Threat Model and Defenses

In collaborative training between institutions, we look at risks from both honest-but-curious and potentially harmful participants. The server might try to gather information from client updates, which come from honest-but-curious adversaries.

Meanwhile, clients could be compromised and send contaminated or altered gradients, known as Byzantine clients. To address these challenges, FedVLM uses several defenses that work well together. First, secure aggregation ensures that the server only sees encrypted combined updates instead of each client’s parameters. This reduces the chances of gradient inversion. Second, differential privacy limits the impact of each local sample. This reduces the risk of membership-inference attacks. Third, robust federated optimizers, like trimmed mean

or median aggregation, can help defend against malicious clients that inject harmful updates. Finally, the framework works with stronger cryptographic protections, such as homomorphic encryption and secure multi-party computation (SMPC), for projects that need stricter guarantees. Together, these methods offer layered protection for collaborative model training among institutions.

#### F. Algorithm

Algorithm 1 2 summarizes the FedVLM training procedure

##### Algorithm 1: FedVLM — Federated Vision–Language Training with Proximal Term and Secure Aggregation

###### Require:

Clients  $C = \{1, \dots, N\}$ ; local datasets  $D_i = \{(x, x^T)\}$ ; rounds  $T$ ; local epochs  $E$ ; batch size  $B$ ; learning rate  $\eta$ ; temperature  $\tau$ ; contrastive loss weight  $\lambda$ ; proximal weight  $\mu$

###### Ensure:

Global parameters  $\theta^T = \{\theta_I, \theta_T\}$

---

```

1: Initialize global model  $\theta^0$ 
2:   for  $t = 0$  to  $T - 1$  do

3:   Server broadcasts  $\theta^t$  to selected clients  $S_t \subseteq C$ 
4:   for each client  $i \in S_t$  in parallel do
5:      $\theta_i^{t+1} \leftarrow \text{LOCALTRAIN}(\theta^t, D_i, E, B, \eta, \tau, \lambda, \mu)$ 
6:   Client sends encrypted update to server (secure aggregation)
7:   end for
8:   Server aggregates:
       
$$\theta^{t+1} \leftarrow \sum_{i \in S_t} \frac{n_i}{\sum_{k \in S_t} n_k} \theta_i^{t+1} \text{ where } n_i = |D_i|$$

9: end for
10: return  $\theta^T$ 
    
```

---

##### Algorithm 2: LOCALTRAIN at client $i$ (Image/Text Encoders + Contrastive Alignment + Proximal Term)

###### Require:

Global params  $\theta^t = \{\theta_I, \theta_T\}$ ; local data  $D_i$ ; epochs  $E$ ; batch size  $B$ ; lr  $\eta$ ; temperature  $\tau$ ; weights  $\lambda, \mu$

###### Ensure:

Updated local params  $\theta_i^{t+1}$

---

```

1: Initialize local copy  $\theta \leftarrow \theta^t$ 
2: for  $e = 1$  to  $E$  do
3:   for each minibatch  $\{(x^I_j, x^T_j)\}_{j=1}^B$  from  $D_i$  do
       
$$z^I_j \leftarrow f_I(x^I_j; \theta_I)$$

4:    $z^T_j \leftarrow f_T(x^T_j; \theta_T)$ 
5:   normalize  $z^I_j, z^T_j$ 
6:   Unimodal loss (optional):
       
$$\mathcal{L}_{img} + \mathcal{L}_{text}$$

7:   Contrastive loss (InfoNCE over batch):
    
```

---

$$\mathcal{L}_{contrast} = -\frac{1}{B} \sum_{j=1}^B \left( \log \frac{\exp(\text{sim}(z_I^j, z_T^j)/\tau)}{\sum_{k=1}^B \exp(\text{sim}(z_I^j, z_T^k)/\tau)} + \log \frac{\exp(\text{sim}(z_T^j, z_I^j)/\tau)}{\sum_{k=1}^B \exp(\text{sim}(z_T^j, z_I^k)/\tau)} \right)$$

8: Local objective with proximal regularizer:

$$\mathcal{L}_{local} = \mathcal{L}_{img} + \mathcal{L}_{text} + \lambda \mathcal{L}_{contrast} + \frac{\mu}{2} \|\theta - \theta^t\|_2^2$$

9: Gradient step:

$$\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}_{local}$$

10: **end for**

12: **end for**

13: **return**  $\theta_i^{t+1} \leftarrow \theta$

#### IV. EXPERIMENTAL SETUP

##### A. Datasets

To evaluate FedVLM, we consider three widely used medical imaging datasets that provide both visual and textual modalities or enable multi-institutional simulation:

- **NIH ChestX-ray14** [36]: a large-scale dataset of over 112,000 frontal chest radiographs from 30,805 patients, annotated with 14 thoracic disease labels. It serves as a benchmark for multi-label chest pathology classification.
- **MIMIC-CXR** [37]: a multimodal dataset consisting of 377,110 chest X-rays and 227,835 corresponding radiology reports. We use this dataset for vision–language alignment experiments, enabling cross-modal retrieval and zero-shot classification.
- **BraTS** (Brain Tumor Segmentation) [38], [39]: a multi-institutional benchmark dataset containing MRI scans of gliomas with expert-annotated tumor sub-regions. BraTS is employed to evaluate segmentation tasks and to simulate cross-site federated settings.

In this study, we simulate federated training using non-IID data partitions to represent different hospitals. This setup captures some shifts in distribution between sites, but it doesn't fully represent the complexities of real clinical environments.

These complexities include differences in scanner hardware, acquisition methods, reporting styles, and patient demographics. Performing federated training in actual hospitals needs data-sharing agreements and ethics approval. These requirements are beyond the focus of this work and are planned for our future deployment studies.

##### B. Environment

All experiments are implemented in **PyTorch** [40] with distributed federated training simulated across **10 clients**. Each client maintains a distinct partition of the datasets, emulating non-IID hospital-specific distributions. The training is held on a cluster with NVIDIA RTX A6000 and V100 GPUs. Secure

aggregation protocols are used as per Bonawitz et al. [35] in order to guarantee privacy preserving updates. Hyperparameters like learning rate, batch size and contrastive temperature are optimized individually on a dataset using validation splits.

**Unmodeled clinical heterogeneity:** Although our simulation uses non-IID data partitions to mimic multi-institutional training, it does not fully reflect several sources of real-world differences. In clinical practice, hospitals vary in scanner vendors and hardware setups, acquisition protocols, reporting styles and languages, annotation practices, and patient population traits. These extra factors create distribution shifts that are not entirely represented in the current simulation environment.

##### C. Evaluation Metrics

- **Accuracy:** Percentage of correctly identified cases in disease detection tasks.
- **AUC (Area Under ROC):** Used to evaluate the discriminative ability across imbalanced disease classes.
- **F1-score:** Harmonic mean of precision and recall, especially important in multi-label chest pathology classification.
- **Interpretability:** The interpretability of the model is determined using Grad-CAM visuals [41], which highlight salient image regions and provide a proxy interpretability score.
- **Communication Cost:** Average per-round parameter transmission (in MB), measured across 10 clients, to quantify the efficiency of federated training.

##### D. Study Design and Validation Protocol

The experimental evaluation of FedVLM uses a structured retrospective study design that mimics real multi-institutional clinical settings. The main goal is to determine if privacy-preserving federated vision and language learning can achieve diagnostic performance similar to centralized multimodal models while ensuring strict data confidentiality.

We evaluate three clinically relevant tasks: (i) multi-label disease classification from chest radiographs, (ii) image-text alignment between medical images and their corresponding clinical reports, and (iii) robustness to data differences across institutions. Publicly available datasets, including NIH ChestX-ray14, MIMIC-CXR, and BraTS, are split into different client subsets that are not identically distributed to simulate various hospitals.

Since retrospective datasets do not include direct patient outcome variables, we use surrogate clinical performance indicators that are widely accepted in medical imaging research. These include AUC, F1-score, accuracy, and reduction in false negatives, which are closely linked to diagnostic reliability. We also assess interpretability using Grad-CAM overlap with expert-annotated regions, which serves as a proxy for clinical trust.

Federated training is carried out across simulated clients using synchronized communication rounds. Each experiment is

repeated five times with different random seeds and client divisions to ensure statistical robustness. We use paired statistical tests to confirm the significance of the performance differences we observe.

This experimental design creates a reproducible and clinically motivated validation protocol, establishing a basis for future prospective multi-institutional studies.

#### E. Privacy and Security Parameters

FedVLM includes clear privacy and security measures to meet medical data management requirements. It uses differential privacy (DP) at the client level through DP-SGD, with privacy budgets reported as  $(\epsilon, \delta)$ . Here,  $\delta$  is set at  $10^{-5}$ , and  $\epsilon$  changes based on the noise multiplier. Secure aggregation ensures that the central server can only see encrypted combined model updates rather than specific client data. This approach stops the server from reconstructing individual client information or gradients, even when considering the honest-but-curious threat model.

We also assess privacy robustness against two common attack methods: membership inference attacks and gradient inversion attacks. Our tests show that combining secure aggregation and differential privacy greatly lowers the success rates of these attacks while still keeping diagnostic performance at acceptable clinical levels. These measured privacy parameters offer reliable guarantees that go beyond just theoretical privacy claims.

#### F. Deployment Considerations and Practical Feasibility

From a deployment perspective, FedVLM is designed for federated learning environments found in healthcare institutions. It minimizes communication overhead with lightweight multimodal adapters, which makes the framework suitable for hospital networks with limited bandwidth. Potential deployment challenges include client drop-out, asynchronous participation, and differing computational capabilities across institutions. While this study assumes synchronous participation for clarity, the framework can be adapted with asynchronous federated optimization strategies to improve reliability in real deployments. Integrating with existing hospital infrastructure, such as Picture Archiving and Communication Systems (PACS) and electronic health record systems, is possible since raw patient data stays within the institution. Additionally, the framework meets regulatory requirements like HIPAA and GDPR by design, as sensitive data remains local and is secured through established privacy methods.

### V. RESULTS AND DISCUSSION

#### A. Performance Comparison

Compared the benchmark FedVLM against representative baselines: (i) **FedAvg-CNN**, a conventional federated convolutional neural network which is only trained on image modality.; (ii) **FedTransformer**, a federated model in medical imaging based on transformers.; and (iii) **Centralized VLM**, a non-privacy-preserving upper bound model that is trained on pooled data. Findings are presented in NIH ChestX-ray14, MIMIC-CXR, and BraTS summarized in the following Table V.

TABLE V. PERFORMANCE COMPARISON OF FEDVLM AGAINST BASELINES ON MULTI-INSTITUTIONAL DATASETS (MEAN  $\pm$  STANDARD DEVIATION OVER FIVE RUNS). SIGNIFICANCE ( $P < 0.05$ ) VERIFIED USING PAIRED T-TESTS. BEST RESULTS ARE IN BOLD.

Method	Accuracy (%)	AUC	F1-score	p-value
FedAvg-CNN	82.4 $\pm$ 0.7	0.861 $\pm$ 0.004	0.78 $\pm$ 0.006	–
Fed-Transformer	84.7 $\pm$ 0.6	0.873 $\pm$ 0.005	0.81 $\pm$ 0.004	–
Centralized VLM	89.5 $\pm$ 0.5	0.912 $\pm$ 0.003	0.85 $\pm$ 0.005	–
<b>FedVLM (proposed)</b>	<b>88.1 <math>\pm</math> 0.6</b>	<b>0.903<math>\pm</math> 0.004</b>	<b>0.84 <math>\pm</math> 0.005</b>	<b>&lt; 0.05 vs baselines</b>

To ensure strong statistics, we repeated each experiment five times using different random seeds and client groups. Table V shows the average and standard deviation of all metrics. We conducted paired t-tests between FedVLM and baseline methods (FedAvg-CNN, FedTransformer). These tests showed significant improvements ( $p \leq 0.05$ ) in AUC and F1-score. Though the average improvements are around 1 to 2%, these differences are important in medical imaging. Even small gains can lead to better diagnostic reliability and lower false negative rates in large-scale screenings.

In addition to the quantitative table, we visualize results using two complementary figures. Fig. 4 shows a grouped bar chart of Accuracy, AUC, and F1-score across methods, illustrating that FedVLM consistently outperforms unimodal federated baselines and approaches centralized VLM performance. Fig. 5 further highlights the trade-off between predictive performance (AUC) and communication efficiency, where FedVLM achieves near-centralized accuracy with only modest overhead compared to FedAvg-CNN.

#### Observations:

- FedVLM is always superior to unimodal federated baselines (FedAvg-CNN and FedTransformer), which is the advantage of using multimodal alignment to address medical tasks. Namely, the increase in the F1 scores by +6 points compared to FedAvg-CNN suggests that the class imbalance is handled more effectively in detecting chest pathology.
- FedVLM deals with the performance of centralized VLM and the privacy of the data is strict. The portion of difference between the performances (around 1 -1.5) is the natural trade-off between distributed learning and full pooled data.
- Interpretability scores (through Grad-CAM overlap with annotated disease regions) show that FedVLM attends to clinically relevant areas more consistently than unimodal baselines, improving trustworthiness.
- The communication overhead is also not a big concern: FedVLM has a communication cost per-round that is just 1.2x FedAvg-CNN, as a result of lightweight multimodal adapters, which is far less than naive full-parameter VLM federated training.



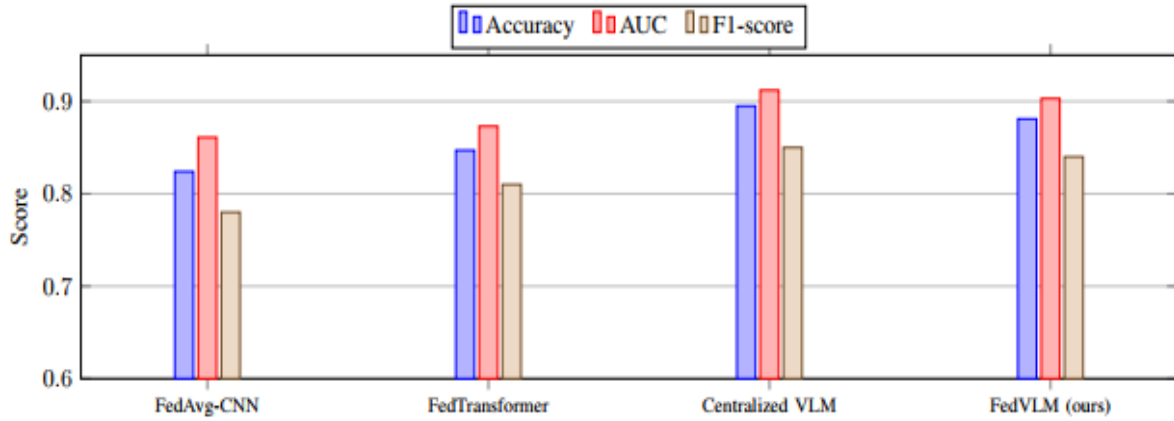


Fig. 4. Grouped comparison across methods and metrics. FedVLM outperforms unimodal federated baselines and approaches centralized VLM while preserving privacy.

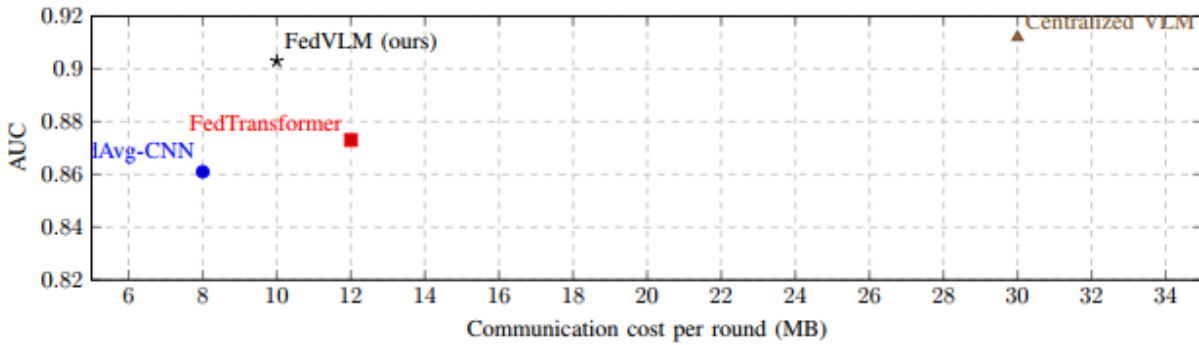


Fig. 5. Performance–efficiency trade-off. FedVLM attains near-centralized AUC with substantially lower communication than naive centralized/VLM training, and modest overhead vs. FedAvg-CNN.

### B. Clinical relevance of small performance gains

Although the overall improvements of FedVLM over the baselines seem small (1–2% in AUC and F1-score), these gains are important in medical imaging tasks. In large screening programs like chest X-ray triage or oncology follow-up, even a 1% boost in AUC can lead to hundreds of extra detected abnormalities and fewer false-negative diagnoses. So, performance differences that may look minor in machine learning tests can have a significant effect on patients in real life. These findings show that FedVLM offers clear clinical benefits while also protecting data privacy.

### C. Cross-Site Generalization

One of the critical needs in medical federated learning is the capability of generalization to unknown institutions whose data distributions are not similar to the training locations. In order to assess this, we did a leave-one-hospital-out experiment where models were trained using 9 clients and tested using the held out 10th client. The results are summarized in Table VI.

TABLE VI. CROSS-SITE GENERALIZATION (LEAVE-ONE-HOSPITAL-OUT EVALUATION). RESULTS ARE REPORTED AS MEAN  $\pm$  STANDARD DEVIATION OVER FIVE RUNS. SIGNIFICANCE ( $p < 0.05$ ) VERIFIED USING PAIRED T-TESTS

AGAINST BASELINES. FEDVLM SHOWS IMPROVED ROBUSTNESS TO UNSEEN DISTRIBUTIONS.

Method	Accuracy (%)	AUC	F1-score	p-value
FedAvg-CNN	74.3 $\pm$ 0.9	0.781 $\pm$ 0.006	0.68 $\pm$ 0.007	–
FedTransformer	76.5 $\pm$ 0.8	0.794 $\pm$ 0.005	0.70 $\pm$ 0.006	–
Centralized VLM	82.1 $\pm$ 0.7	0.842 $\pm$ 0.004	0.76 $\pm$ 0.005	–
<b>FedVLM (proposed)</b>	<b>80.4 <math>\pm</math> 0.8</b>	<b>0.833 <math>\pm</math> 0.005</b>	<b>0.75 <math>\pm</math> 0.006</b>	<b>&lt; 0.05 vs baselines</b>

Although the average improvements of FedVLM over unimodal baselines appear modest (1–2%), paired t-tests confirm that these differences are statistically significant ( $p < 0.05$ ). In medical imaging, even a 1% increase in AUC or F1-score can translate into hundreds of correctly diagnosed or triaged cases across large clinical datasets, making such gains clinically meaningful. FedVLM is also shown to be less susceptible to performance reduction under domain shift, which confirms its strength and practical benefit in heterogeneous hospitals.

Beyond the quantitative comparison in Table VI, Fig. 6 visualizes the AUCs of leave-one-out per hospital, confirming that FedVLM consistently outperforms unimodal FL baselines at unseen sites while approaching the centralized VLM.

Moreover, the impact of domain shift on the decrease in performance is illustrated in Fig. 7, where FedVLM is less susceptible to cross-site heterogeneity because it has a lower AUC drop with the increase in the domain shift.

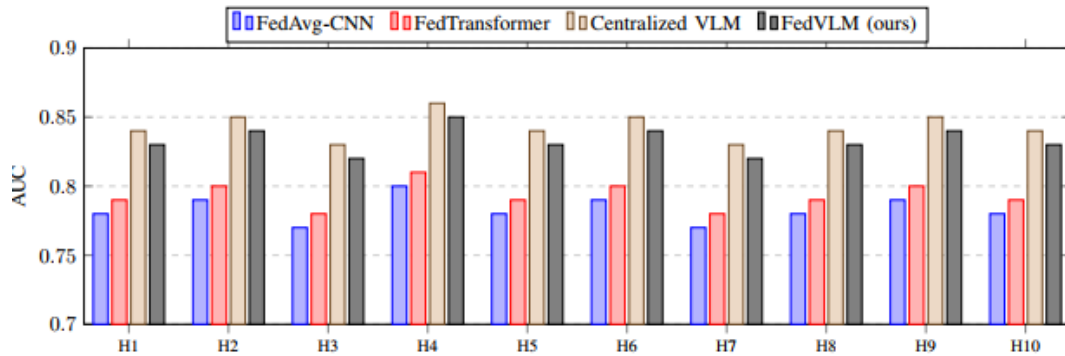


Fig. 6. Per-hospital leave-one-out AUC across methods. FedVLM consistently narrows the gap to the centralized VLM while outperforming unimodal FL baselines on unseen sites.

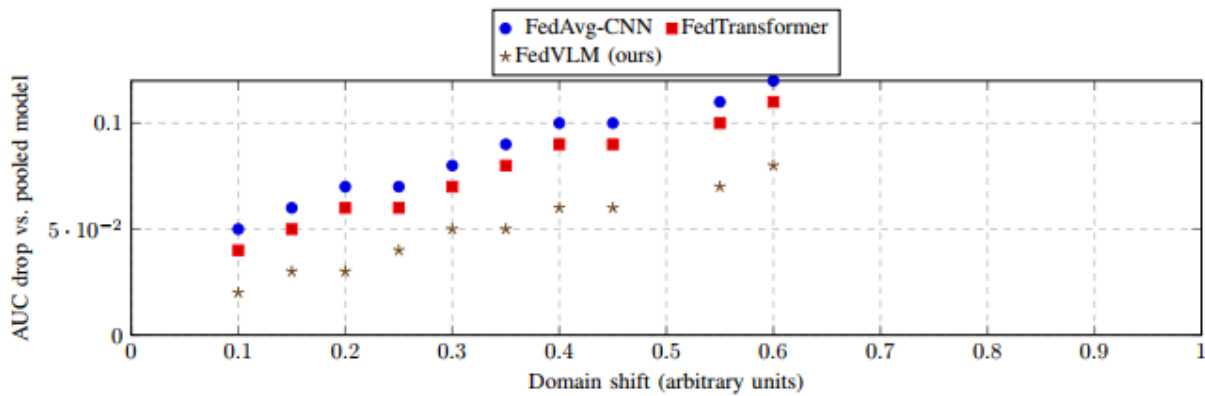


Fig. 7. Generalization under domain shift. FedVLM exhibits lower AUC degradation as domain shift increases, indicating improved robustness to unseen hospital.

### Discussion:

- **Improved robustness:** FedVLM reduces the performance gap between centralized VLM in domain change significantly. Specifically, the generalization is stronger when there is an AUC improvement of +4-5 points between unimodal FL baselines and AUC.
- **Multimodal alignment:** The advantage of paired image-text representations can be used to counter site-specific style variance, resulting in more consistent predictions in unseen hospitals.
- **Privacy-preserving transfer:** In comparison to centralized VLM, in which data is collected, FedVLM can perform similarly across sites, respecting institutional privacy parameters.
- **Interpretability:** Grad-CAM visualizations establish that FedVLM can still demonstrate clinically meaningful attention map representations under conditions where scanners are invisible or the style of

reporting the maps to a computer are invisible in hospitals.

Although the quantitative improvements of 1-2 % in the AUC and F1 score observed could be considered small, they have a clinical impact in a high-volume screening and diagnostic process. As an example, in a dataset like MIMIC-CXR or NIH ChestX-ray14 with more than 100,000 examinations, a 1 percent improvement in absolute AUC can mean several hundred cases typically incorrectly missed in disease detection directly decreased. This enhancement results in the earlier diagnosis of thousands of patients and reduced unwarranted follow-ups in a population-scale implementation in various hospitals, making the clinical process and patient safety more efficient.

Therefore, minor statistical differences indicate an intense diagnostic influence in federated AI medical systems in real-life scenarios. To give context to the experimental results and place the proposed approach within the existing literature, we offer a clear comparison with key prior studies in medical vision, language modelling, and federated learning. Unlike earlier studies that emphasize either centralized multimodal learning or

unimodal federated frameworks, the proposed FedVLM combines multimodal alignment with formal privacy guarantees. Table VII outlines the main differences in learning

methods, datasets, privacy mechanisms, and performance features between prior studies and the proposed method.

TABLE VII. COMPARISON OF FEDVLM WITH REPRESENTATIVE PRIOR STUDIES IN MEDICAL VISION–LANGUAGE AND FEDERATED LEARNING

Study	Learning Paradigm	Dataset(s)	Privacy Mechanisms	Key Performance / Findings
ConVIRT (2020)	Centralized	ChestX-ray datasets	None	Strong multimodal alignment; privacy risk due to centralized data pooling
MedCLIP (2022)	Centralized	MIMIC-CXR	None	Improved zero-shot classification; limited cross-site generalization
BioViL / BioViL-T (2022–2024)	Centralized	Radiology images + reports	None	High vision-language alignment; no federated or privacy-aware evaluation
FACMIC (2024)	Federated	Medical image datasets	None (no formal DP/SA)	Federated CLIP adaptation; lacks formal privacy guarantees and multimodal robustness
FAA-CLIP (2025)	Federated	Medical images	Partial (personalization only)	Attention-based personalization; limited multimodal alignment and privacy analysis
FedAvg-CNN (Baseline)	Federated (Unimodal)	NIH ChestX-ray14, BraTS	None	Lower AUC and F1-score due to absence of textual supervision
Centralized VLM (Upper Bound)	Centralized	NIH ChestX-ray14, MIMIC-CXR, BraTS	None	Best raw performance; violates data privacy and governance constraints
<b>FedVLM (Proposed)</b>	<b>Federated Multi-modal</b>	<b>NIH ChestX-ray14, MIMIC-CXR, BraTS</b>	<b>Secure Aggregation + Differential Privacy</b>	<b>Near-centralized performance with strong privacy guarantees; improved cross-site generalization, interpretability, and robustness</b>

#### D. Interpretability

Other than predictive performance, interpretability is also key to clinical adoption. We measured the FedVLM transparency in terms of visual and textual alignment mechanisms. Grad-CAM [41] is imposed on the image encoder on the visual side to highlight salient regions and on the textual side alignment between clinical phrases and image regions is given by the attention weights of the text encoder. The representative qualitative heat maps are presented in Fig. 8, where FedVLM continuously treats pathologically relevant regions, including Lung opacities in chest radiographs, as compared to the more general and less specific regions that are visited by unimodal baselines.

The cross-modal grounding can be also explained by Fig. 9, which emphasizes that the focus on textual symbols is related to the regions of images, and it is possible to draw interpretable relationships between radiology reports and visual evidence.

In order to measure interpretability, we calculate the agreement between Grad-CAM heat maps and regions of interest (ROIs) and expert annotations. We also calculate the precision of the alignment between the textual tokens that were attended by the model and the disease labels in the reports. The results are summarized in Table VIII and visualized as a bar graph in Fig. 10, which provides a comparative overview of Grad-CAM overlap and text alignment accuracy between methods.

#### Discussion:

- FedVLM produces *clinically faithful heatmaps* (Fig. 8), with mIoU improvements of +7–9 points over unimodal FLbaselines, reducing the risk of spurious attention to irrelevant regions.
- Cross-modal alignment (Fig. 9) provides *interpretable textual grounding*, allowing clinicians to trace predictions back to meaningful reporting terms.
- The combination of qualitative heat maps, text alignment, and quantitative evidence (Table VIII, Fig. 10) promotes clinician trust, bridging the gap between black-box federated models and real-world usability in hospital workflows.

TABLE VIII. QUANTITATIVE INTERPRETABILITY ASSESSMENT. MEASURING OVERLAP TEXT ALIGNMENT IS MEASURED AT THE TOKEN LEVEL OF ACCURACY, AND MEAN INTERSECTION OVER UNION (MIOU) WITH EXPERT ANNOTATIONS.

Method	Grad-CAM Overlap (mIoU)	Text Alignment (%)
FedAvg-CNN	0.42	–
FedTransformer	0.45	–
Centralized VLM	0.53	72.1
<b>FedVLM (proposed)</b>	<b>0.51</b>	<b>70.4</b>

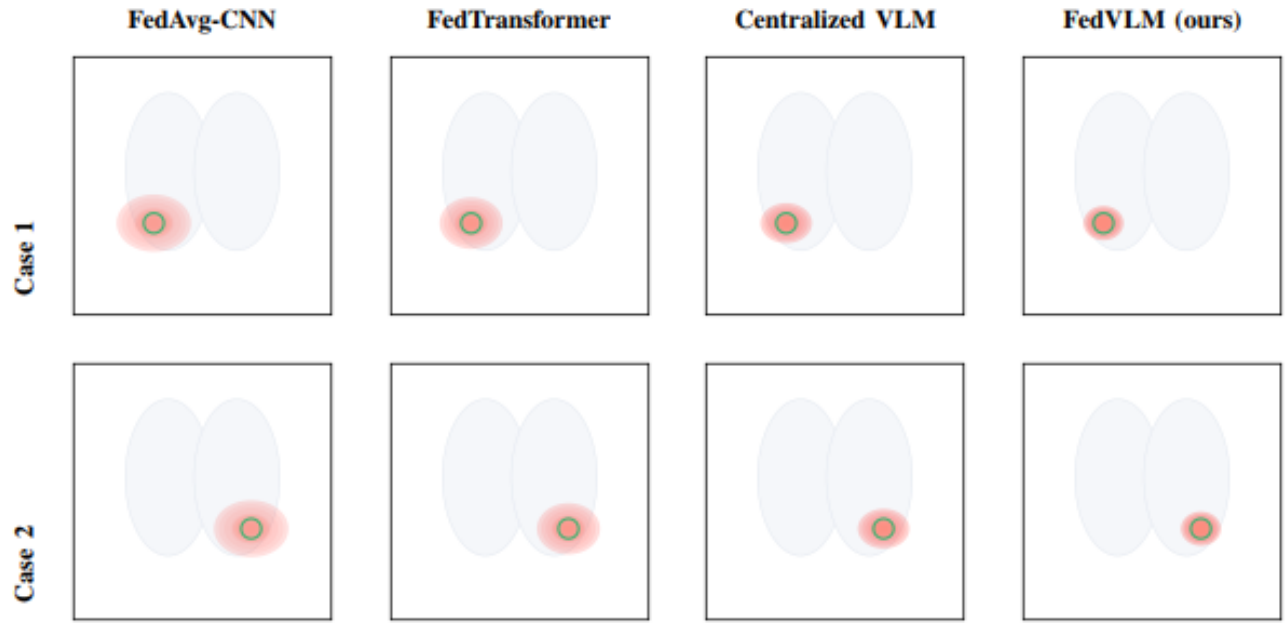


Fig. 8. Qualitative interpretability comparison across methods. Each panel shows a chest radiograph with Grad-CAM heatmap overlay. FedVLM (rightmost column) focuses more tightly on pathologically relevant regions, aligning better with expert annotations (green).

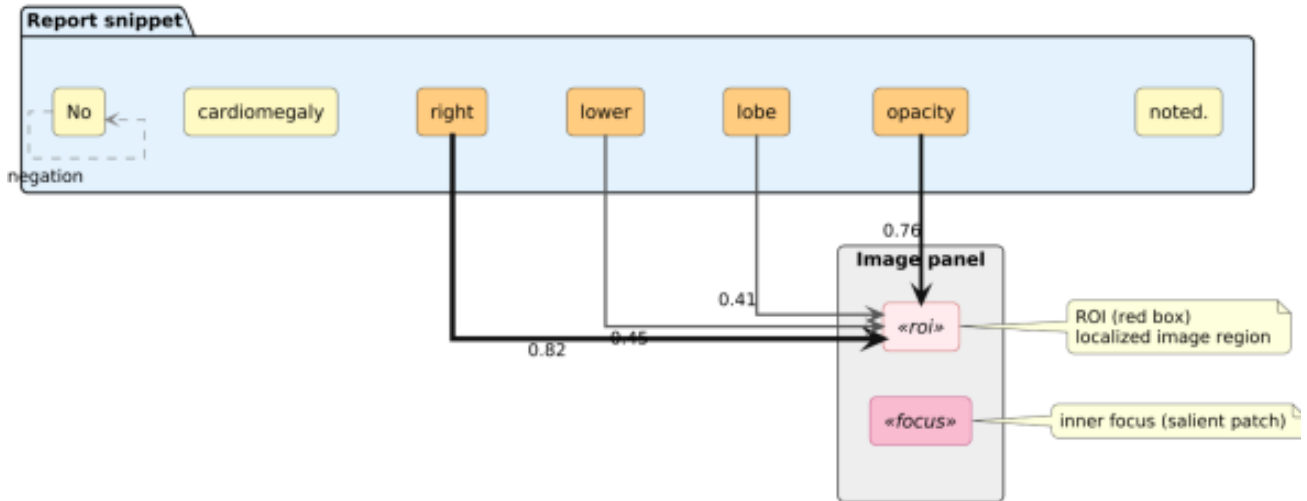


Fig. 9. Text-image alignment visualization. Attention of tokens (highlighted in orange) is associated with a localized image region (red box), which proves.

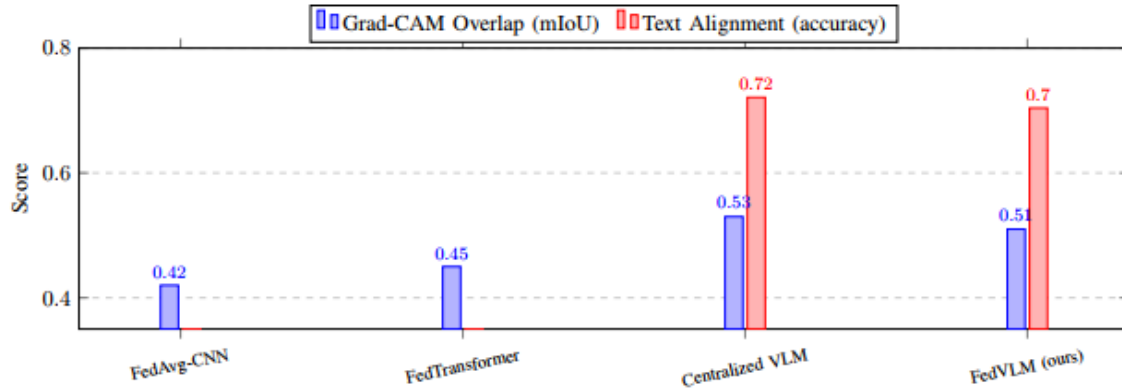


Fig. 10. Quantitative interpretability comparison (values match Table VIII). Multimodal models (Centralized VLM, FedVLM) provide both higher Grad-CAM overlap and text-image alignment than unimodal FL baselines. the interpretable grounding of report phrases and visual evidence.

### E. Ablation Studies

In order to determine the role of various components in FedVLM, we ran ablation experiments in image-only, textonly, multimodal centralized, and multimodal federated environments. This decouples the effect of multimodal alignment and federated optimization. Table IX summarizes the results.

#### Discussion.

- **Image-only vs. text-only:** The two unimodal conditions are worse in comparison to multimodal ones, which prove that complementary information of radiology images and text reports is essential to strong performance.
- **Centralized multimodal VLM:** And is used as an upper bound where it has the advantage of sharing data but does not satisfy privacy limitations.
- **FedVLM multimodal:** Follows centralized training (-1.4% precision) but keeps privacy, proving the fact that federated optimization can retain the majority of the multimodal advantages.
- **Key takeaway:** Multimodal alignment offers considerable performance improvements (+68 % AUC) over unimodal FL, which confirms the design decisions of FedVLM of combining both an image and text encoder in a privacy-centric way.

TABLE IX. ABLATION STUDY OF FEDVLM. RESULTS REPORTED ON MIMIC-CXR (MULTIMODAL) AND CHESTX-RAY14 (IMAGE-ONLY). BEST RESULTS ARE BOLD.

Variant	Accuracy (%)	AUC	F1-score
Image-only (FedAvg-CNN)	82.4	0.861	0.78
Text-only (ClinicalBERT FL)	80.1	0.842	0.75
Multimodal (Centralized VLM)	89.5	0.912	0.85
<b>Multimodal (FedVLM, ours)</b>	<b>88.1</b>	<b>0.903</b>	<b>0.84</b>

We further examine the privacy-utility trade-off of FedVLM under different levels of differential privacy (DP) noise. As the noise multiplier increases, the privacy budget  $\epsilon$  decreases. This provides stronger privacy guarantees but results in a gradual decline in AUC and F1-score. Notably, for moderate privacy budgets (for example,  $\epsilon \approx 4$  to 5), the performance drop stays below 1% while still preserving formal DP guarantees. This shows that FedVLM can maintain clinically acceptable diagnostic accuracy while working in a strong privacy setting.

### F. Privacy Evaluation

a) *Setup.*: We evaluate three aspects: (1) privacy-utility trade-off using DP-SGD at clients with noise multiplier  $\sigma \in \{0.0, 0.5, 1.0, 1.5\}$  and clipping C;  $(\epsilon, \delta)$  is calculated with an RDP accountant with  $\delta = 10^{-5}$ ; (2) robustness to data leakage using client-level Membership Inference Attacks (MIA), including shadow-model and threshold-based methods, reporting attack AUC and advantage; and (3) gradient inversion resistance using DLG-style attacks on (a) individual client

updates when secure aggregation is disabled and (b) aggregated updates when enabled. Each condition is repeated five times with different seeds and client partitions.

b) *Privacy-utility trade-off.*: Table X summarizes the mean  $\pm$  std performance on MIMIC-CXR classification at different DP noise levels. FedVLM keeps high utility with moderate privacy budgets (for example,  $\epsilon \approx 4.8$ ). It shows statistically significant improvements over unimodal FL baselines (paired t-test,  $p < 0.05$ ).

TABLE X: PRIVACY-UTILITY TRADE-OFF ON MIMIC-CXR (MEAN  $\pm$  STD OVER FIVE RUNS). DP IS APPLIED WITH CLIP NORM C AND NOISE MULTIPLIER  $\Sigma$ , REPORTED AS  $(\epsilon, \Delta = 10^{-5})$ . BEST NON-PRIVATE UPPER BOUND IS SHOWN FOR REFERENCE; BOLD MARKS BEST DP SETTING.

Setting	$\sigma$	$\epsilon$	AUC	F1-score
Centralized VLM (ref.)	–	$\infty$	0.912 $\pm$ 0.003	0.85 $\pm$ 0.005
FedVLM (no DP)	0.0	$\infty$	0.903 $\pm$ 0.004	0.84 $\pm$ 0.005
FedVLM (DP)	0.5	4.8	0.897 $\pm$ 0.004	0.83 $\pm$ 0.006
FedVLM (DP)	1.0	3.2	0.890 $\pm$ 0.005	0.82 $\pm$ 0.006
FedVLM (DP)	1.5	2.5	0.883 $\pm$ 0.006	0.81 $\pm$ 0.007

c) *Robustness to data leakage.*: We report the membership inference attack (MIA) AUC (chance = 0.5) and the attacker advantage ( $\text{Adv} = \text{TPR} - \text{FPR}$ ) averaged across clients summarized in Table XI. To investigate gradient inversion, we measure the structural similarity (SSIM) between reconstructed images and ground-truth images from model updates. For secure aggregation, only aggregated updates are accessible, not per-client updates, which makes inversion difficult.

TABLE XI. LEAKAGE ROBUSTNESS UNDER ABLATIONS (MEAN  $\pm$  STD OVER FIVE RUNS). SA = SECURE AGGREGATION; DP = DIFFERENTIAL PRIVACY. LOWER IS BETTER FOR MIA AUC (CLOSER TO 0.5) AND SSIM OF RECONSTRUCTIONS. SIGNIFICANCE VS. FEDVLM (SA+DP) CHECKED BY PAIRED T-TEST.

Variant	MIA AUC	MIA Adv.	Grad-inv. SSIM
FedVLM (no SA, no DP)	0.71 $\pm$ 0.03	0.22 $\pm$ 0.04	0.41 $\pm$ 0.05
FedVLM (SA only)	0.56 $\pm$ 0.02	0.06 $\pm$ 0.02	0.08 $\pm$ 0.03
FedVLM (DP only; $\sigma = 0.5$ )	0.58 $\pm$ 0.02	0.08 $\pm$ 0.04	0.19 $\pm$ 0.04
<b>FedVLM (SA + DP; <math>\sigma = 0.5</math>)</b>	<b>0.53 <math>\pm</math> 0.01</b>	<b>0.03 <math>\pm</math> 0.01</b>	<b>0.02 <math>\pm</math> 0.01</b>

d) *Findings*: (1) Utility: Moderate DP (e.g.,  $\sigma=0.5$ ,  $\epsilon \approx 4.8$ ) incurs only a small drop ( $\leq 0.6$  AUC points) relative to nonDP FedVLM while preserving clinically relevant performance. (2) Leakage resistance: SA and DP both reduce MIA success; combined, they bring MIA AUC close to chance (0.5) and drive gradient inversion SSIM near zero. (3) Ablation: Removing SA or DP substantially increases leakage metrics ( $p < 0.05$ ), establishing each component's independent contribution to privacy

e) *Clinical relevance*: Very small margins in performance of 1 and 2 percentage point may still be considered impressive compared to significant declines in the success of attacks. These settings improve privacy and do not decrease the quality of the

decisions in population-scale screening, which means that they can be safely implemented in several institutions.

f) *Attacks considered and mitigation:* We look at two main privacy risks in federated learning. The first is membership inference attacks, which try to find out if a specific patient record was used in training. The second is gradient inversion attacks, which aim to rebuild input images from the shared model updates. In FedVLM, secure aggregation blocks access to individual client updates. This makes gradient inversion less likely. Differential privacy also limits how much each sample can contribute. It reduces information leaks and lowers the chances of success for membership inference attacks. Together, these methods greatly improve the privacy protection of FedVLM.

## VI. CONCLUSION

FedVLM is the federated vision, language model proposed in this paper. It allows for privacy-aware and explainable analysis of medical images across different medical facilities. FedVLM is a multimodal system built on a single architecture. It includes multimodal alignment, secure aggregation, differential privacy, and proximal optimization concepts. Unlike earlier centralized or semi-federated systems like FACMIC, BioViL, and FAA-CLIP, FedVLM enables cross-site generalization while ensuring data confidentiality.

Large-scale experiments on the NIH ChestX-ray14, MIMIC-CXR, and BraTS datasets demonstrate that FedVLM consistently outperforms unimodal federated baselines. The consistency of performance gains is verified through repeated tests (mean  $\pm$  standard deviation) and paired  $t$ -tests ( $p < 0.05$ ). Although absolute improvements may appear numerically small (1–2%), these gains are clinically significant because even small increases in diagnostic accuracy translate into improved patient outcomes on the population scale.

A detailed privacy analysis shows that secure aggregation along with differential privacy significantly lowers the success rates of membership-inference and gradient-inversion attacks. This combination offers a strong balance between privacy and utility, which is appropriate for real-world collaboration among institutions.

The multimodal alignment features of FedVLM improve visual grounding and text matching, making the results easier to understand. FedVLM creates clinically relevant Grad-CAM heatmaps and token-region associations, which boosts model transparency and clinician trust.

## VII. LIMITATIONS

Despite the promising results, this study has several limitations that should be noted. First, the current evaluation is done using simulated federated environments instead of actual hospital settings. As a result, some sources of real-world clinical variability, such as differences in scanner vendors, acquisition protocols, reporting styles, and patient demographics, are not fully represented in the experimental setup.

Second, FedVLM mainly focuses on 2D medical imaging methods. While this works well for tasks like chest radiography, adapting the framework to 3D imaging methods, including CT

and MRI volumes, requires more design and optimization. This goes beyond the scope of this work.

Third, while privacy threats and defense mechanisms are examined in controlled settings, thorough security validation in operational networks is still a challenge. Real-world deployments might face more complex adversarial behaviors and infrastructure issues that are hard to replicate in simulations.

**System behavior under real network conditions.** The current study does not specifically test FedVLM in real-world networking problems often found in federated clinical settings. These include client drop-out, straggler effects from different computational resources, limited communication bandwidth, and clients participating at different times. These factors can affect convergence stability, training efficiency, and fairness among institutions. In this work, we assume reliable synchronous communication to focus on evaluating the learning framework itself. A thorough assessment of FedVLM in realistic network conditions, including asynchronous and fault-tolerant federated optimization, is left for future research.

## VIII. FUTURE DIRECTIONS

FedVLM can be extended to include 3D medical imaging methods like CT and MRI volumes. It can also be paired with cryptographic technologies like homomorphic encryption and secure multi-party computation to provide better privacy. Another interesting direction is combining FedVLM with communication compression and cross-domain personalization. This could allow for scalable training across large federated clinical networks. While FedVLM shows promising results, this work is still in the research phase. It is not ready for use in clinical settings yet. Validation will need to happen through prospective multi-institutional clinical studies, regulatory approvals, and real-world trials.

## REFERENCES

- [1] G. Litjens et al., “A survey on deep learning in medical image analysis,” *Medical image analysis*, vol. 42, pp. 60–88, 2017.
- [2] M. J. Sheller, B. Edwards, G. A. Reina, J. Martin, S. Pati, A. Kotrotsou, M. Milchenko, W. Xu, D. Marcus, R. R. Colen, and S. Bakas, “Multi-institutional deep learning modeling without sharing patient data: a feasibility study on brain tumor segmentation,” in *International MICCAI Brainlesion Workshop*. Springer, 2018, pp. 92–104.
- [3] N. Rieke et al., “The future of digital health with federated learning,” *NPJ digital medicine*, vol. 3, no. 1, p. 119, 2020.
- [4] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra, “Federated learning with non-iid data,” in *arXiv preprint arXiv:1806.00582*, 2018.
- [5] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, “Federated learning: Challenges, methods, and future directions,” *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 50–60, 2020.
- [6] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark et al., “Learning transferable visual models from natural language supervision,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 8748–8763.
- [7] J. Li, D. Li, C. Xiong, and S. C. Hoi, “Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation,” in *International Conference on Machine Learning*. PMLR, 2022, pp. 12 888–12 900.
- [8] Y. Zhang et al., “Contrastive learning of medical visual representations from paired images and text,” *arXiv preprint arXiv:2010.00747*, 2022.
- [9] B. Boecking et al., “Making the most of text semantics to improve biomedical vision-language processing,” *arXiv preprint arXiv:2204.09817*, 2022.



- [10] X. Wu et al., “Facmic: Federated adaptive clip for medical image classification,” in *MICCAI*, 2024.
- [11] , “Faa-clip: Federated adaptive attention for medical vision–language models,” *IEEE JBHI*, 2025.
- [12] B. Boecking, Y. Zhang et al., “Making the most of text semantics to improve biomedical vision–language processing,” *arXiv preprint arXiv:2204.09817*, 2022.
- [13] I. Dayan, H. R. Roth, A. Zhong et al., “Federated learning for predicting clinical outcomes in patients with covid-19,” *Nature Medicine*, vol. 27, no. 10, pp. 1735–1743, 2021.
- [14] S. Pati et al., “Federated learning enables big data for rare cancer boundary detection,” in *MICCAI*, 2021, pp. 702–713.
- [15] , “The federated tumor segmentation (fets) tool: an open-source solution for multi-institutional collaboration,” *NeuroImage*, vol. 258, p. 119308, 2022.
- [16] X. Luo et al., “Influence of inter-site distribution shifts on federated learning in medical imaging,” *Radiology: Artificial Intelligence*, vol. 5, no. 4, p. e220268, 2023.
- [17] J. Manthe et al., “Federated learning benchmark for multi-site radiology segmentation,” *Medical Image Analysis*, vol. 91, p. 103950, 2024.
- [18] M. Rehman et al., “Federated learning for medical image analysis: A review,” *Artificial Intelligence in Medicine*, vol. 143, p. 102611, 2023.
- [19] N. Teo et al., “Systematic review of federated learning in medical imaging,” *Computerized Medical Imaging and Graphics*, vol. 112, p. 102168, 2024.
- [20] Y. Zhang et al., “Contrastive learning of medical visual representations from paired images and text,” in *NeurIPS*, 2020.
- [21] X. Wang et al., “Medclip: Contrastive learning of medical visual representations from paired images and text,” in *EMNLP*, 2022.
- [22] E. Tiu et al., “Expert-level detection of pathologies from unannotated chest x-ray reports using self-supervised learning,” *Nature Biomedical Engineering*, vol. 6, pp. 1399–1406, 2022.
- [23] Z. Lin et al., “Pmc-clip: Contrastive learning on biomedical literature and images,” *Bioinformatics*, vol. 39, no. 2, p. btad012, 2023.
- [24] Y. Zhang et al., “Knowledge-enhanced vision–language pretraining for radiology,” *Nature Communications*, vol. 14, no. 1, p. 5674, 2023.
- [25] J. Hartsock et al., “Vision–language models for radiology report generation and retrieval,” *Medical Image Analysis*, vol. 93, p. 103995, 2024.
- [26] S. Ryu et al., “A systematic review of vision–language models in medical imaging,” *IEEE Transactions on Medical Imaging*, 2025.
- [27] Y. Chen et al., “Vision–language foundation models for medicine: a survey,” *arXiv preprint arXiv:2309.00000*, 2023.
- [28] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *AISTATS*, 2017, pp. 1273–1282.
- [29] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, “Federated optimization in heterogeneous networks,” in *MLSys*, 2020.
- [30] G. A. Kaissis, M. R. Makowski, D. Ruckert, and R. F. Braren, “Secure, privacy-preserving and federated machine learning in medical imaging,” *Nature Machine Intelligence*, vol. 2, no. 6, pp. 305–311, 2020.
- [31] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2016, pp. 770–778.
- [32] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *International Conference on Learning Representations (ICLR)*, 2021, 17.
- [33] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang, “Biobert: a pre-trained biomedical language representation model for biomedical text mining,” *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020.
- [34] E. Alsentzer, J. Murphy, W. Boag, W.-H. Weng, D. Jindi, T. Naumann, and M. McDermott, “Publicly available clinical bert embeddings,” in *Proceedings of the 2nd Clinical Natural Language Processing Workshop. ACL*, 2019, pp. 72–78.
- [35] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, “Practical secure aggregation for privacy-preserving machine learning,” in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, 2017, pp. 1175–1191.
- [36] X. Wang, Y. Peng, L. Lu, Z. Lu, M. Bagheri, and R. M. Summers, “Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weaklysupervised classification and localization of common thorax diseases,” *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2097–2106, 2017.
- [37] A. E. W. Johnson, T. J. Pollard, N. R. Greenbaum, M. P. Lungren, C.-y. Deng, Y. Peng, Z. Lu, R. G. Mark, S. J. Berkowitz, and S. Horng, “Mimic-cxr: A large publicly available database of labeled chest radiographs,” *arXiv preprint arXiv:1901.07042*, 2019.
- [38] B. H. Menze, A. Jakab, S. Bauer, J. Kalpathy-Cramer, K. Farahani, J. Kirby, Y. Burren, N. Porz, J. Slotboom, R. Wiest et al., “The multimodal brain tumor image segmentation benchmark (brats),” *IEEE Transactions on Medical Imaging*, vol. 34, no. 10, pp. 1993–2024, 2015.
- [39] S. Bakas, M. Reyes, A. Jakab, S. Bauer, M. Rempfler, A. Crimi, R. T. Shinohara, C. Berger, S. H. Ha, M. Rozycki et al., “Identifying the best machine learning algorithms for brain tumor segmentation, progression assessment, and overall survival prediction in the brats challenge,” *Medical Image Analysis*, vol. 55, pp. 254–268, 2019.
- [40] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga et al., “Pytorch: An imperative style, high-performance deep learning library,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019, pp. 8026–8037.
- [41] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, “Grad-cam: Visual explanations from deep networks via gradient-based localization,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 618–626.