# Med-MMFL: A Multimodal Federated Learning Benchmark in Healthcare

Aavash Chhetri[2*]   Bibek Niroula[2*]   Pratik Shrestha[2]   Yash Raj Shrestha[3]   Lesley A Anderson[1]
Prashnna K Gyawali[4]   Loris Bazzani[5]   Binod Bhattarai[1,2,6†]

[1]University of Aberdeen, Aberdeen, UK
[2]NepAl Applied Mathematics and Informatics Institute for research, Nepal
[3]University of Lausanne, Switzerland
[4]West Virginia University, USA
[5]University of Verona, Italy
[6] University College London, UK

## Abstract

*Federated learning (FL) enables collaborative model training across decentralized medical institutions while preserving data privacy. However, medical FL benchmarks remain scarce, with existing efforts focusing mainly on unimodal or bimodal modalities and a limited range of medical tasks. This gap underscores the need for standardized evaluation to advance systematic understanding in medical MultiModal FL (MMFL). To this end, we introduce Med-MMFL, the first comprehensive MMFL benchmark for the medical domain, encompassing diverse modalities, tasks, and federation scenarios. Our benchmark evaluates six representative state-of-the-art FL algorithms, covering different aggregation strategies, loss formulations, and regularization techniques. It spans datasets with 2 to 4 modalities, comprising a total of 10 unique medical modalities, including text, pathology images, ECG, X-ray, radiology reports, and multiple MRI sequences. Experiments are conducted across naturally federated, synthetic IID, and synthetic non-IID settings to simulate real-world heterogeneity. We assess segmentation, classification, modality alignment (retrieval), and VQA tasks. To support reproducibility and fair comparison of future multimodal federated learning (MMFL) methods under realistic medical settings, we release the complete benchmark implementation, including data processing and partitioning pipelines, at* [https://github.com/bhattarailab/Med-MMFL-Benchmark](https://github.com/bhattarailab/Med-MMFL-Benchmark).

## 1. Introduction

Clinicians inherently rely on information from multiple sources and modalities to perform reliable diagnosis, prognosis, and formulating treatment plans. This reliance on heterogeneous information has motivated the recent development of multimodal models for healthcare applications [2, 4, 26]. Like clinicians, the goal of multimodal models is to provide an holistic view of the disease and offer improved and consistent diagnostic performance [42, 45] by integrating information across diverse data sources like medical scans, omics data, and pathology reports. Training such multimodal models in healthcare is challenging because of the fragmentation of medical data owing to their sensitive nature and strict privacy concerns that limit their broad distribution. Federated learning (FL) addresses these challenges by enabling privacy-preserving training of local models for each institution (clients), while aggregating these models at a global level (server) without any exchange or distribution of patient data [36]. More recently, numerous unimodal (e.g., imaging-only) and bimodal (e.g., image–text) FL methods [4, 14, 26, 49] have been proposed to deal with these challenges. However, it still remains open how to scale beyond two modalities and there is a lack of standardized evaluation for different methods in a reproducible manner. To address these limitations, we present a comprehensive benchmark which is composed by multiple modalities, datasets, tasks, evaluation protocols, and baselines with the aim to promote reproducibility and fair comparison and to facilitate fast progress in the domain of multimodal healthcare.

To facilitate FL research, several FL benchmarks have been proposed [6, 15, 27, 30]. However, they remain fragmented across domains and objectives. NIID-Bench [30] focuses on algorithmic diversity under controlled non-IID

---
*Equal Contribution.
†Binod Bhattarai is the corresponding author.

| Features | | NIID-Bench [30] | FLamby [37] | FedLLM-Bench [50] | FedMultimodal [12] | FedVLMBench [55] | Med-MMFL (ours) |
|---|---|---|---|---|---|---|---|
| Modalities | # Modalities (min, max) | (1,1) | (1,1) | (1,1) | (2,2) | (2,2) | **(2,4)** |
| | # Unique Medical Modalities | 0 | 5 | 0 | 2 | 2 | **10** |
| # Multimodal Medical Datasets | | 0 | 0 | 0 | 1 | 2 | **5** |
| # Distinct FL Algorithms | | 4 | 4 | 4 | 4 | 3 | **6** |
| Partitioning Strategies | Real-world | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ |
| | Synthetic IID | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ |
| | Synthetic non-IID | ✓ | ✗ | ✗ | ✓ | ✓ | ✓ |
| Evaluation Tasks | | 1 | 3 | 2 | 1 | 4 | **4** |

Table 1. Comparison of experimental settings across unimodal benchmarks (shaded in grey), multimodal benchmarks (shaded in blue), and our Med-MMFL benchmark (shaded in green). For each feature, **bold** indicates the most comprehensive support and underlined indicates the second best.

settings but is limited to unimodal data. FLamby [37] provides the first medical FL suite, yet without multimodality. In contrast, existing multimodal FL benchmarks such as FedMultimodal [12], FedVLMBench [55], and FedLLM-Bench [50] expand to multiple modalities, yet there remains no comprehensive benchmark that jointly captures the multimodal, medical, and federated aspects of real-world healthcare scenarios.

A closer examination of existing benchmarks (see Tab. 1) reveals several key limitations. **(1) Datasets:** Existing medical FL benchmark, such as [37] focus on unimodal datasets, where as multimodal benchmarks [12, 55] largely ignore medical data. **(2) Modalities:** Prior multimodal works [12, 55] are typically limited to two modalities per dataset, and hence, offer only a limited insight into the characteristics of real-world multimodal clinical data in federated settings. **(3) Tasks:** Most benchmarks cover only a handful of task types (Tab. 1), highlighting the need for broader and more diverse evaluations of medical tasks. **(4) Partitioning:** Partitioning strategies in existing benchmarks remain restricted in scope: [37, 50] focus only on real-world splits, whereas [55] employs synthetic ones, leaving few frameworks that combine both realistic hospital-level silos and synthetic non-IID settings required for comprehensive FL evaluation. **(5) Algorithms coverage:** The set of FL Algorithms studied in existing benchmarks remains narrow, with many variants (e.g., FedAvgM [18], FedAdam, FedAdagrad, FedYogi) reducible to the unified FedOpt formulation [39]. Consequently, existing evaluations offer limited exploration of fundamentally distinct state-of-the-art FL algorithms, urging the need for a broader and more inclusive benchmark.

To address theses limitations, we introduce **Med-MMFL** ( Fig. 1), a comprehensive standardized benchmark designed specifically for multimodal federated learning in healthcare. Our main contributions can be summarized as follows:

1. Med-MMFL integrates **6** distinct FL Algorithms, **4** types of tasks, **3** data partitioning strategies, and **5** medical datasets with varying degree of multimodality (**2 to 4 modalities**).

2. We generalize existing algorithms (e.g., MOON, CreamFL) to handle more than two modalities for fair multimodal evaluation.

3. To foster transparency and facilitate reproducible research, we publicly release our benchmark implementation, including all the dataset processing and partitioning pipelines.

The remainder of this paper is organized as follows. Sec. 2 reviews related work on federated learning and existing benchmarks. Sec. 3 describes the Med-MMFL datasets, including data partitioning and baseline setups. Sec. 4 presents the Med-MMFL framework and the adapted FL algorithms. Finally, Sec. 5 describes the experimental setup and presents the benchmark results.

## 2. Related Works

**FL algorithms.** Most of modern FL approaches are based on FedAvg [36], a foundational work which proposed to train a global model by averaging the models trained locally by each party (clients) while keeping their data private. Subsequent studies [18, 30, 33] built on FedAvg establish that non-Independently and Identically Distributed (non-IID) data degrade the convergence rate of FedAvg. To address this issue of statistical heterogeneity, FedProx [32] adds a $L_2$ regularizer to the loss function that restricts local updates to be close to the global model. SCAFFOLD [23] introduces control variates to correct the drift of local updates, ensuring their directions remain consistent with the global optimization path. FedNova [47] corrects FedAvg's bias toward clients with more local updates by normalizing and scaling their contributions during the aggregation stage. MOON [29] addresses the issue of non-IID data via contrastive learning in model-level by comparing the representations to bridge the gap between the representations learned by the local model and the global model. FedDyn [1] dynamically adjusts each client's objective with a regularization term, ensuring convergence toward stationary points of the global empirical loss. To address the constraint of identical client architectures in all these FedAvg-derived algorithms, knowledge-distillation-based fusion methods [7, 35, 48] have been proposed to
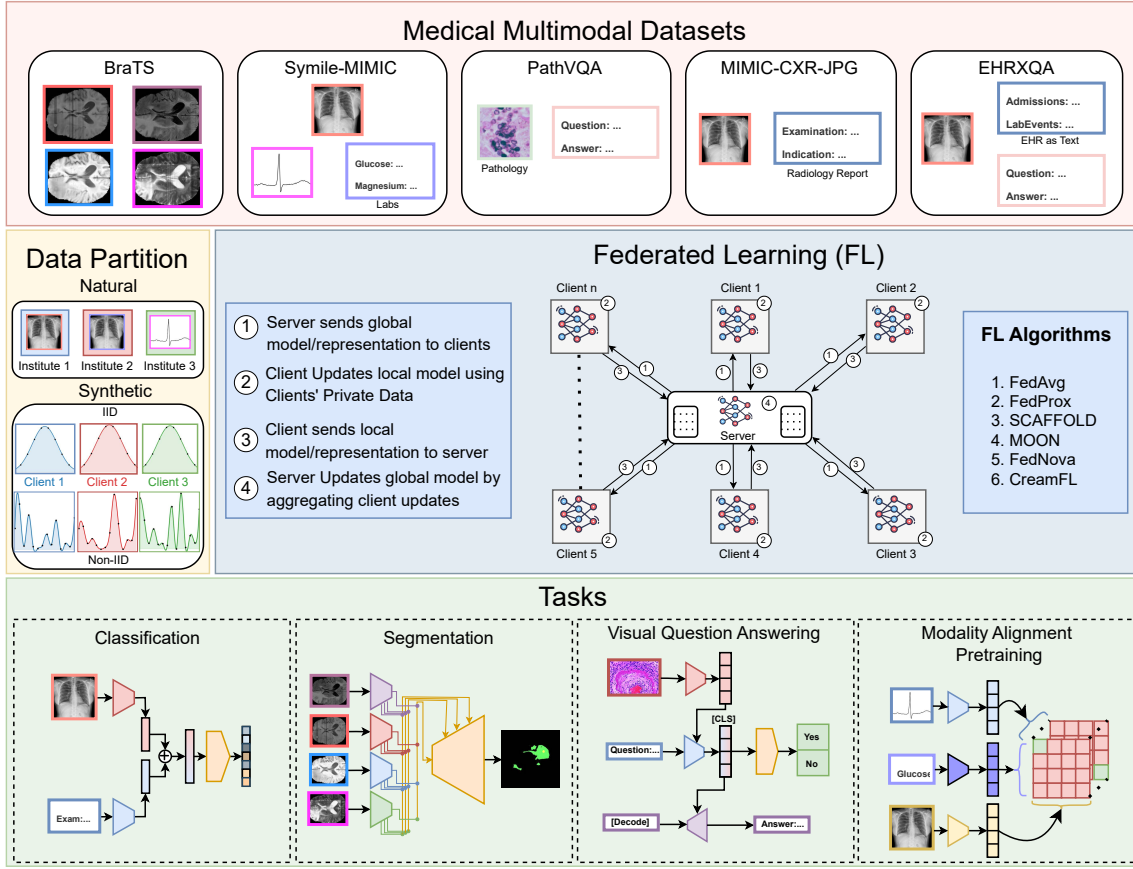
Figure 1. Overview of our proposed Med-MMFL benchmark framework. It spans diverse multimodal medical datasets, task types, and client partitioning strategies, integrating multiple FL algorithms to provide a unified evaluation platform.

enable heterogeneous participation by replacing parameter aggregation with model-agnostic knowledge transfer. CreamFL [52] advances knowledge-distillation fusion toward multimodal learning by adopting representation-level transfer and contrastive objectives to reduce model drift. Although our benchmark focuses on a representative set of fundamentally distinct FL algorithms, some state-of-the-art methods [1, 34, 39] explore complementary approaches and are not directly evaluated in this work. Furthermore, an active and promising research direction of Personalized FL [8, 44] lies outside the scope of our current work and is open for future exploration.

**FL benchmarks.** LEAF [6] was an early benchmark providing federated datasets with realistic natural splits. FedML [15] offers an open research library and benchmark to facilitate FL algorithm development and fair performance comparison. NIID-Bench [30] proposes diverse non-IID partitioning strategies and evaluates multiple FL algorithms, forming a comprehensive unimodal benchmark. However, these benchmarks focus on unimodal and non-medical data, leaving a gap for multimodal medical FL eval-

uation. FLamby [37] marked a major milestone in FL for healthcare research by targeted benchmark with seven medical datasets covering five distinct input modalities. Nevertheless, its unimodal and real-world partitioning design, along with limited representative FL algorithms leaves open opportunities for the community to explore broader multimodal and data heterogeneity scenarios. FedMultimodal [12] introduced a FL benchmark targeting diverse multimodal applications. However, it is limited by its focus on cross-device settings, single evaluation task coverage, as well as datasets with only two modalities. Consequently, the evaluation of federated learning algorithms across scenarios involving more than two modalities remains largely unexplored. Furthermore, its sole healthcare data PTBXL [46] is not inherently a multimodal dataset. While the ECG signals are split and passed as "separate modalities", they all stem from the same underlying modality, merely grouped across electrode groups. FedLLM-Bench [50] establishes a federated benchmark for LLM training (instruction tuning and preference alignment). Building upon it, FedVLM-Bench [55] introduces a systematic benchmark for federated

fine-tuning of Vision-Language Models (VLMs) with three FL algorithms, while incorporating two medical datasets. Despite its contributions, the benchmark still falls short in several aspects: its medical datasets are limited in diversity, the evaluation of FL algorithms remains narrow, and it addresses only two modalities (image and text). Thus, despite prior work, a benchmark that fully integrates multimodality, medical-data diversity, and federated realism remains absent. Thus, we propose Med-MMFL, providing comprehensive coverage of datasets, modalities, partitioning strategies, and FL algorithms for systematic evaluation.

## 3. Med-MMFL Datasets

In this benchmark, we carefully select five publicly available medical datasets covering diverse modalities and clinical tasks, chosen for their widespread adoption and/or representativeness of real-world medical challenges. Since these datasets are not inherently federated, we systematically transform them into FL-compatible versions through realistic client partitioning and adapting baselines to capture broader and clinically meaningful task diversity.
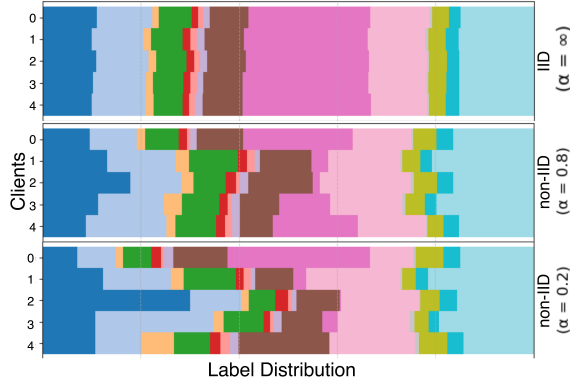


Figure 2. Representative client-level label distribution for the MIMIC-CXR-JPG dataset obtained using our federated partitioning strategy. IID splits produce similar label frequencies across clients, whereas non-IID splits yield heterogeneous distributions where certain labels dominate or are absent on specific clients. Other datasets exhibit analogous distribution patterns under the same protocol (see Sec. 7)

### 3.1. Fed-BraTS-GLI2024

BraTS-GLI2024 [9, 22] is a multimodal 3D Magnetic resonance imaging (MRI) dataset for Brain Tumor Segmentation with four MRI modalities: 1. Pre-contrast T1-weighted (T1), 2. Contrast-enhanced T1-weighted (T1-Gd), 3. T2-weighted (T2), and 4. T2-weighted fluid-attenuated inversion recovery (FLAIR), along with the segmentation mask for each patient. We build a federated version of BraTS-GLI2024 with **5** clients, each corresponding to a unique contributing center in the dataset, and we refer to this setup as a *natural* partition. Additionally, we use pseudo-classes

for *synthetic* data partitioning. Specifically, we apply clustering on the class proportion vectors, computed as the normalized voxel counts per class for each segmentation mask, and treat each cluster as a pseudo-class. To simulate an IID scenario, we uniformly distribute the samples across clients such that each client receives similar proportions of all pseudo-classes. For non-IID synthetic partitions, we follow the de-facto approach from [53] used in many works [12, 18, 30, 55], which is based on sampling from a Dirichlet distribution. Specifically, we sample $p_k \sim Dir_N(\alpha)$ and allocate a $p_{k,j}$ proportion of the instances of pseudo-class $k$ to party $j$. Here, $Dir(\cdot)$ denotes the Dirichlet distribution and $\alpha$ is the concentration parameter which we use to vary the degree of pseudo-class skew across clients, thereby modulating the level of data heterogeneity in the federated setup. As a baseline for the segmentation task, we train the RFNet [11] architecture following its reference implementation. The experiments are evaluated using Dice Score Coefficient (DSC).

### 3.2. Fed-MIMIC-CXR-JPG

The MIMIC-CXR-JPG dataset [21], derived from MIMIC-CXR [13, 19, 20], pairs Chest X-ray (CXR) images with their corresponding textual radiology reports. The dataset provides 14 labels that correspond to radiology findings to form a multi-label classification task measured by the macro Area Under the ROC Curve (AUC). Like Sec. 3.1, we use Dirichlet distribution based sampling to generate *synthetic* federated partitions, Fig. 2. To extend label-distribution skew [30] to multi-label dataset, each label dimension is treated independently: for every class, the corresponding samples are allocated to clients according to proportions drawn from a Dirichlet distribution. In the multi-label scenario, a sample may belong to multiple classes. Hence, when overlap occurs, the final client assignment of such samples effectively follows the allocation determined by the last class processed. To ensure homogeneous distribution in IID partitioning, we use high value ($\alpha \rightarrow \infty$) of the concentration parameter $\alpha$ while sampling from the Dirichlet distribution. As a baseline, we train a multimodal classifier following the implementation of [38].

### 3.3. Fed-Symile-MIMIC

Symile-MIMIC [13, 40, 41] is a multimodal clinical dataset comprising chest X-rays (CXR), electrocardiograms (ECG), and blood laboratory measurements (blood labs). Synthetic federated partitions are derived using the patient and admission metadata. To induce controllable non-IID characteristics, we sample client-wise proportions from a Dirichlet distribution, parameterized by the empirical distribution of the metadata values. Since geographical information of the centers was not available, this partitioning is *synthetic*, nonetheless the most realistic option with

this dataset. Our baseline model follows [40], which performs modality alignment pretraining and we evaluate the results through a zero-shot retrieval setup, wherein CXRs are retrieved based on corresponding ECG and lab representations. We use the ResNet-50 and ResNet-18 architectures [16] for the CXR and ECG encoders, respectively, and a three-layer Multi-layer Perceptron (MLP) to encode the blood labs data following the implementation of [40]. The downstream retrieval task is evaluated using the Accuracy metric.

### 3.4. Fed-PathVQA

PathVQA [17] is a dataset for pathological Visual Question Answering (PVQA) designed to emulate diagnostic reasoning similar to that assessed in the American Board of Pathology examinations. The dataset comprises 4,998 pathology images and 23,700 question-answer pairs. For our experiments, we formulate a close-ended Visual Question Answering (VQA) task by using the yes/no subset of the dataset, which includes 16,334 VQA instances (49.8% of the total). We construct pseudo-classes for partitioning, as described in Sec. 3.1, distribute them uniformly across clients and use Dirichlet sampling to generate IID and non-IID *synthetic* partitions. To infer pseudo-classes, we perform clustering of questions embedded using BioMed-Clip [54], a powerful text model for biomedical applications. As a baseline, we train BLIP [28] and pass the [CLS] token from the text encoder into a 2 layer MLP for binary (yes/no) classification. The quality of the resulting close-ended VQA task is measured by the F1 Score.

### 3.5. Fed-EHRXQA

EHRXQA [5] is a multi-modal question-answering dataset which includes structured Electronic Health Records (EHRs) and chest X-ray images. The original dataset uses SQL queries over the full EHR database to retrieve patient-specific or multi-patient information. In practice, however, clinical reasoning occurs at the individual patient level, where clinicians interpret multimodal data to answer patient-specific questions. Thus, database-wide SQL retrieval is redundant. By converting each patient's structured EHR into text and removing the SQL module, we simplify the VQA task to focus on individual patient data. This preprocessing yields 9,956 QA pairs, each aligned with a single patient record. Consistent with Sec. 3.4, we split the dataset into synthetic IID and non-IID partitions by generating embeddings via BioMedClip [54] and distribute samples based on the clusters of the embeddings. As a baseline, we train the same architecture of BLIP [28], specifically with its generative decoder to generate open-ended answers. The resulting VQA task is evaluated by Token Overlap F1 score [10].

## 4. The Med-MMFL Benchmark Framework

We consider a multimodal federated learning setting with $C$ clients, each having its private dataset $D_C$ with $n_C$ samples. The $i^{th}$ data sample in $D_C$ is represented by tuple $(\{X_m^{(i)}\}_{m=1}^{M_C}, Y^{(i)})$, where $Y^{(i)}$ and $M_C$ represent the label set and the number of modalities in the $C^{th}$ client, respectively. The clients collaboratively train a global model $f_s(\cdot; \mathbf{w}) : \mathbb{R}^n \to \mathbb{R}^d$ parameterized by $\mathbf{w}$, which maps inputs to outputs of dimension $d$. Note that while we explore data heterogeneity in FL settings, the models, however, are homogenous across clients and the server for a dataset.

**Dataset Distribution.** Let $\mathcal{D}$ denote the centralized dataset, which serves as the complete collection of data before federated partitioning. The subsets of this dataset $\{\mathcal{D}_{val}, \mathcal{D}_{test}\} \subset \mathcal{D}$ are used by the server for validation and testing, respectively of the global model. To simulate the federated setting, $\mathcal{D}_{rem} := \mathcal{D} \setminus \{\mathcal{D}_{val}, \mathcal{D}_{test}\}$ is divided into multiple client-specific subsets. If institutional information regarding data collection exists, it is utilized to partition the dataset into $n$ splits, such that $\{\mathcal{D}_1, \mathcal{D}_2, \ldots, \mathcal{D}_n\} \subset \mathcal{D}_{rem}$, constituting a natural partition. In the absence of such information, synthetic partitioning based on metadata and label information is performed to create client datasets under both IID and non-IID settings, with the latter generated using a Dirichlet distribution. Please refer to Sec. 3 to dive deep on how each dataset is partitioned.

### 4.1. FL Algorithms

Med-MMFL implements and evaluates six representative FL algorithms selected to cover diverse state-of-the-art paradigms with distinct optimization and update mechanisms. Notably, most of these methods were originally designed for unimodal or, at best, bimodal applications. In our work, we extend selected algorithms to support multimodal data, enabling a fair and unified evaluation.

**FedAvg [36].** FedAvg introduces a FL framework that allows clients to collectively train a global model keeping its local data private. It performs iterative round-based training where at the start of each round $t + 1$, each client $k$ receives a copy of the global model $w_t$ and updates it to optimize the local objective $\mathcal{L}_{local}$ for a number of local epochs. At the end of each round, the server receives copies of client models $w_k, k = 1, 2..C$ and aggregates them to get new global model $w_{t+1}$. Subsequent studies [18, 30, 33] established that this algorithm's convergence rates may degrade under heterogeneity.

**FedProx [32].** FedProx follows the same model-aggregation framework as FedAvg. However, to tackle heterogeneity in federated settings, it improves the local objective by adding a $L_2$ regularization term to each local training loss, which penalizes the deviation of the

local models from the last global model. Additionally, it introduces a hyper-parameter $\mu$ to control the effect of the regularization.

**SCAFFOLD [23].** When the distribution of each local dataset differs from the global distribution, there exists a *drift* [23] in local updates away from the global optimum. This drift induces bias in the aggregated updates, misguiding the global optimization process and consequently slower or unstable convergence of the global model. To address this, SCAFFOLD introduces control variates that estimates the client drift and corrects the gradients in the direction that compensates for the drift.

**FedNova [47].** FedNova improves FedAvg in the aggregation stage by considering that different parties may conduct different numbers of local steps in each communication round and the parties with a larger number of local steps may significantly bias the global updates towards them. Thus, to ensure a more balanced aggregation, FedNova normalizes and scales the local model updates of each client according to the size of their local steps before updating the global model.

**MOON [29] and m-MOON (ours).** To mitigate drift under non-IID conditions, MOON uses contrastive learning to align local and global representations. It minimizes the distance between features learned by the local and global models while maximizing the distance from the previous local model. Note that MOON only formulates the algorithm for unimodal data. Here, we propose to generalize the model-contrastive learning principle of MOON to more modalities, *i.e.*, multimodal MOON (m-MOON) for evaluation on multimodal datasets. To this end, we perform *modality-wise contrastive alignment* between the local and global representations for each modality $m \in M_C$ (see Eq. (1)).

$$
\mathcal{L}_{m-MOON} = \sum_{m \in M_C} - \log \frac{f(\mathbf{z}_{loc}^m, \mathbf{z}_{glob}^m)}{f(\mathbf{z}_{loc}^m, \mathbf{z}_{glob}^m) + f(\mathbf{z}_{loc}^m, \mathbf{z}_{prev}^m)}
$$
(1)

where $f(\mathbf{z_1}, \mathbf{z_2}) = \exp(\text{sim}(\mathbf{z_1}, \mathbf{z_2})/\tau)$ and $\mathbf{z}_k^m$ is a representation of the modality $m$ from the model $k$. The intuition for this design choice over contrastive learning on a singular fused multi-modal representation stems from the fact that each modality may drift independently during local training due to heterogeneity. To balance the contrastive regularization with the task-specific local objective, m-MOON keeps the hyperparameter $\mu$ from MOON [29], that controls the weight of the model-contrastive loss.

**CreamFL [52] and CreamMFL (ours).** CreamFL enables a server model to be learned from clients with heterogeneous modalities and model architectures. To achieve this, it performs a contrastive representation-level ensemble and a global–local cross-modal aggregation on a small public dataset, allowing the server to fuse client representations while preserving client privacy. The framework additionally introduces inter-modal and intra-modal contrastive regularizers during local training to mitigate local drift that may arise from modality gaps and/or task gaps. However, the formulation of CreamFL is limited to bimodal (image and text) scenarios only. We generalize CreamFL beyond 2 modalities with **CreamMFL**, by simply extending contrastive regularizer for each modality using pairwise inter- and intra-modal losses (see Eq. (2)) and global-local contrastive aggregation by using the available modalities.

$$
\ell_{intra}^{(r)} = \sum_{m \in M_C} - \log \frac{f(\mathbf{z}^{(r,m)}, \mathbf{z}_{glob}^{(r,m)})}{f(\mathbf{z}^{(r,m)}, \mathbf{z}_{glob}^{(r,m)}) + f(\mathbf{z}^{(r,m)}, \mathbf{z}_{prev}^{(r,m)})}
$$
(2a)

$$
\ell_{inter}^{(r)} = \sum_{\substack{m_1, m_2 \in M_C \\ m_1 \neq m_2}} - \log \frac{f(\mathbf{z}^{(r,m_1)}, \mathbf{z}_{glob}^{(r,m_2)})}{\sum_{\substack{n=1 \\ n \neq r}}^{|\mathcal{P}|} f(\mathbf{z}^{(r,m_1)}, \mathbf{z}_{glob}^{(n,m_2)})}
$$
(2b)

where $f(\mathbf{z_1}, \mathbf{z_2}) = \exp(\mathbf{z_1}^\top \cdot \mathbf{z_2})$, $\mathbf{z}_{glob}^{r,m}$ is the representation of $r^{th}$ data point of modality $m$ from $glob$ model and $|\mathcal{P}|$ is the public data available to all clients. For fair comparison across consistent client configurations, CreamMFL trains a separate model on the public data as a virtual client which is also used in aggregation.

## 5. Experiments

### 5.1. Experimental Setup

To evaluate the performance of FL algorithms across medical datasets, we conduct extensive experiments across different partitioning scenarios on five federated multimodal medical datasets, as described in Sec. 3.

**Implementation details.** All datasets are first divided into training, validation, and test subsets prior to federated partitioning, ensuring consistent evaluation data across centralized and FL experiments. Focusing on cross-silo settings, all experiments are performed under full client participation. The Adam optimizer [25] is employed for local optimization, except in the case of FedNova [47]. Since FedNova does not define the normalization vector $\|a\|_1$ [47] when using Adam, we instead use SGD with momentum and tune the learning rate independently. The complete set of hyperparameter values is included in the supplementary material. All algorithms are executed with identical training schedules, including the same number of communication rounds and local epochs. For better reliability, results are averaged over three independent runs with different random seeds. All experiments are performed on NVIDIA A100-PCIE-40GB GPUs.

| Split | Datasets | Clients | Metric ↑ | FedAvg | FedProx | SCAFFOLD | m-MOON | FedNova | CreamMFL |
|---|---|---|---|---|---|---|---|---|---|
| **Natural** | Fed-BraTS-GLI2024 | 5 | DSC | 81.797 | 81.260 | 79.447 | **82.406** | 82.125 | 78.934 |
| **Synthetic IID** | Fed-BraTS-GLI2024 | 3 | DSC | 82.608 | 83.259 | 81.944 | 83.555 | **84.052** | 76.539 |
| | | 5 | DSC | 80.781 | 79.717 | 78.108 | 79.655 | **82.811** | 77.064 |
| | Fed-MIMIC-CXR-JPG | 3 | AUC | 82.996 | 83.245 | **83.539** | 82.016 | 82.582 | 80.607 |
| | | 5 | AUC | **80.202** | 78.983 | 79.978 | 77.186 | 79.905 | 77.871 |
| | Fed-Symile-MIMIC | 3 | Acc | **38.147** | 35.698 | 26.94* | 38.074 | 11.853 | 18.966 |
| | | 5 | Acc | **34.483** | 30.316 | 21.336* | 30.244 | 9.6983 | 16.164 |
| | Fed-PathVQA | 3 | F1 | 86.895 | 87.186 | **87.852** | 87.558 | 84.937 | 84.411 |
| | Fed-EHRXQA | 3 | F1 | 51.008 | 51.044 | **51.552** | 51.4 | 50.706 | 47.524 |
| | # times that performs the best | | | 3 | 0 | 3 | 0 | 2 | 0 |
| **Synthetic non-IID** $\alpha = 0.8$ | Fed-BraTS-GLI2024 | 3 | DSC | 81.536 | 80.518 | **83.941** | 83.260 | 83.164 | 80.592 |
| | | 5 | DSC | 80.292 | 80.466 | 79.381 | 80.432 | **82.325** | 80.069 |
| | Fed-MIMIC-CXR-JPG | 3 | AUC | 83.763 | 83.063 | 83.553 | 82.745 | 82.534 | **84.777** |
| | | 5 | AUC | 81.649 | **81.741** | 79.717 | 79.714 | 79.608 | 71.227 |
| | Fed-Symile-MIMIC | 3 | Acc | 32.902 | **35.345** | 15.086* | 24.856 | 11.207 | 17.026 |
| | | 5 | Acc | **31.824** | 30.244 | 19.181* | 29.597 | 10.991 | 14.224 |
| | Fed-PathVQA | 3 | F1 | **87.022** | 86.99 | 84.978 | 85.863 | 85.194 | 85.764 |
| | Fed-EHRXQA | 3 | F1 | 50.542 | **51.141** | 50.793 | 50.993 | 49.802 | 46.038 |
| | # times that performs the best | | | 2 | 3 | 1 | 0 | 1 | 1 |
| **Synthetic non-IID** $\alpha = 0.2$ | Fed-BraTS-GLI2024 | 3 | DSC | 81.779 | **83.146** | 81.658 | 83.011 | 82.896 | 81.766 |
| | | 5 | DSC | 81.171 | 81.121 | 80.628 | 81.250 | **83.579** | 74.817 |
| | Fed-MIMIC-CXR-JPG | 3 | AUC | **84.108** | 82.731 | 82.884 | 82.157 | 82.447 | 78.827 |
| | | 5 | AUC | 82.745 | **83.793** | 77.252 | 81.755 | 81.098 | 73.969 |
| | Fed-Symile-MIMIC | 3 | Acc | 33.333 | 28.736 | 11.638* | **34.195** | 11.207 | 13.362 |
| | | 5 | Acc | 31.394 | **34.698** | 17.886* | 22.773 | 9.2672 | 18.75 |
| | Fed-PathVQA | 3 | F1 | 87.155 | 87.176 | **87.428** | 85.813 | 85.123 | 86.224 |
| | Fed-EHRXQA | 3 | F1 | 51.445 | **51.549** | 51.170 | 50.993 | 49.751 | 47.990 |
| | # times that performs the best | | | 1 | 4 | 1 | 1 | 1 | 0 |

Table 2. MMFL benchmark results reported for 6 FL algorithms across different multimodal medical datasets and partitions. For a setting, **bold** highlights the best performance, whereas underlined highlights the second-best.

## 5.2. Centralized Baseline

| Dataset | Metric ↑ | # Epochs | Centralized | FL Evaluation (min, max) |
|---|---|---|---|---|
| BraTS-GLI2024 [9] | DSC | 150 | 84.900 | (74.817, 84.052) |
| MIMIC-CXR-JPG [21] | AUC | 40 | 92.380 | (71.227, 84.777) |
| Symile-MIMIC [41] | Acc | 50 | 41.370 | (9.2672, 38.147) |
| PathVQA [17] | F1 | 60 | 87.124 | (84.411, 87.852) |
| EHRXQA [5] | F1 | 30 | 52.385 | (46.038, 51.552) |

Table 3. Centralized baseline performance obtained by aggregating all client data for reference comparison with federated settings.

Tab. 3 reports the centralized baseline obtained by training on the data aggregated from the clients. While this setup is typically unrealistic in federated contexts due to privacy and governance constraints [49], it offers a useful comparison to evaluate how data decentralization and non-IID distributions impact model performance. Across most datasets, the difference between centralized and federated training is minimal, typically less than 1–1.5%. An exception is MIMIC-CXR-JPG, where centralized training achieves nearly 7% higher accuracy than the best-performing federated algorithm. Interestingly, federated algorithms outperform the centralized baseline in certain scenarios, consistent with trends observed in prior studies [3, 51, 55]. This behavior may be attributed to the heterogeneous nature of federated updates. Unlike centralized training, which optimizes the model over the entire global dataset, federated training aggregates updates computed on diverse local data distributions. In non-IID settings, each client's local optimization effectively constitutes a distinct stochastic process, rather than identical copies of the same optimization as assumed under the IID premise. This diversity introduces implicit regularization through stochasticity in the global update, guiding the model toward flatter regions of the loss landscape and improving generalization [24, 31].

## 5.3. Main Results on Med-MMFL

The performance of six FL algorithms, reported across varying partitions and datasets are shown in Tab. 2. Next, we present an in-depth analysis and insights of these results

---

*Gradient scaling was applied to control the gradient correction in SCAFFOLD for Fed-Symile-MIMIC.
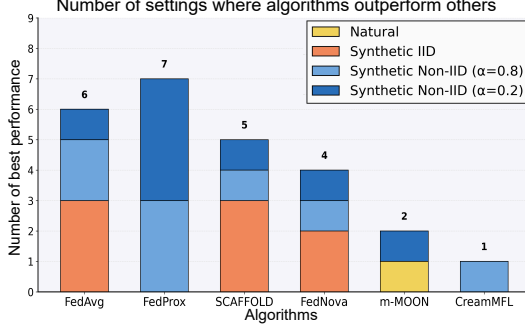
Figure 3. Number of settings in which each algorithm outperforms the others across our Med-MMFL benchmark. The stacked bars are color-coded by data partition type.

from different perspectives.

**Comparison among algorithms.** Overall, no single algorithm consistently outperforms others across all experimental configurations. In Tab. 2, we noticed that Fed-Prox, FedAvg, and SCAFFOLD consistently rank among the top-performing algorithms across diverse combinations of datasets, data splits, and class counts, albeit in different contexts. FedProx achieves the highest performance in 7 out of 16 settings, particularly excelling in non-IID scenarios as highlighted by Fig. 3. In contrast, FedAvg and SCAF-FOLD demonstrate their strongest performance primarily in IID settings. These results suggest that in real-world settings, where datasets are inherently heterogenous, FedProx is likely to offer more robust performance. This also aligns with prior works [18, 33] that establish the degradation in performance of FedAvg with increase in heterogeneity. Interestingly, although FedNova does not rank first or second in most benchmarks, it performs exceptionally well on the Fed-BraTS-GLI2024 dataset, achieving the best results in 4 settings and the second-best in 2 settings, totaling top performances in 6 out of 7 configurations.

**Effect of data partitioning strategy on performance.** The effect of data partitioning strategy on model performance remains generally consistent across most datasets, with notable exceptions. For Fed-BraTS-GLI2024 and Fed-MIMIC-CXR-JPG, performance remains largely consistent across splits, except for CreamMFL, which exhibits a variation of up to 5–6% in the 3-client setting. In contrast, Fed-Symile-MIMIC exhibits a significant drop in performance when moving from IID to non-IID configurations across all algorithms. This decline is likely due to the nature of contrastive learning, which relies on diverse batch samples to effectively push apart embeddings of different instances while pulling together the embeddings of the same instance. Under IID splits, each client receives a representative mix of data, allowing the model to learn discrimina-

tive features across all types of examples. However, in non-IID splits, some data groups may contain only a few samples. In such cases, the model has limited examples to learn from, making it difficult to distinguish instances within that group and reducing overall embedding quality. Conversely, Fed-PathVQA and Fed-EHRXQA demonstrate robustness to data heterogeneity, maintaining comparable performance across different splits.

**Best-performing algorithms across datasets.** The top-performing algorithms vary across datasets, reflecting that each method tends to excel under specific conditions. On Fed-BraTS-GLI2024, FedNova performs best in 4 out of 7 settings. For Fed-MIMIC-CXR-JPG, FedAvg and Fed-Prox lead in 2 settings each, together covering 4 of 6 configurations. In Fed-Symile-MIMIC, FedAvg dominates 3 settings and FedProx 2, while SCAFFOLD and FedNova perform notably worse. As recent work suggests that Self-supervised Learning (SSL) gradients possess predictive capabilities and encode complementary information beyond model embeddings [43], hence, this degradation may stem from instability introduced by frequent gradient corrections or cumulative gradient-based normalization. On Fed-PathVQA, SCAFFOLD achieves the best results in 2 of 3 settings, while on Fed-EHRXQA, FedProx leads in 2 out of 3 cases, highlighting the effectiveness of regularization via weight and gradient control in VQA tasks. Overall, while each dataset favors different algorithms, FedAvg and FedProx emerge as the most consistently strong performers across this experimental setup, ranking highest in 3 out of 5 datasets.

## 6. Conclusion

We present Med-MMFL, the first comprehensive benchmark for multimodal federated learning (MMFL) in healthcare. Our benchmark integrates 6 state-of-the art FL algorithms, 4 types of clinical tasks, 3 data partitioning strategies, and 10 unique medical modalities across 5 datasets. Although Med-MMFL covers a broad spectrum of FL settings, some relevant algorithms remain beyond its current scope and are left for future exploration. Med-MMFL focuses on the federated aspects of multimodal learning, emphasizing the aggregation of client updates rather than the design or evaluation of fusion strategies. A systematic study of fusion methods is therefore beyond the scope of this work. Nonetheless, by publicly releasing the benchmark, Med-MMFL aims to foster transparency, reproducibility, and consistent comparisons for future research. We anticipate that this benchmark will accelerate the development of robust, clinically relevant multimodal FL solutions and serve as a valuable resource for the medical AI community.

## Acknowledgements

## References

[1] Durmus Alp Emre Acar, Yue Zhao, Ramon Matas Navarro, Matthew Mattina, Paul N. Whatmough, and Venkatesh Saligrama. Federated learning based on dynamic regularization. In *International Conference on Learning Representations*, 2021. 2, 3

[2] J.N. Acosta, G.J. Falcone, P. Rajpurkar, and E.J. Topol. Multimodal biomedical ai. *Nature Medicine*, 28(9):1773–1784, 2022. 1

[3] Sultan Alasmari, Rayed AlGhamdi, Ghanshyam G Tejani, Sunil Kumar Sharma, and Seyed Jalaleddin Mousavirad. Federated learning-based multimodal approach for early detection and personalized care in cardiac disease. *Frontiers in Physiology*, 16:1563185, 2025. 7

[4] Rawan AlSaad, Alaa Abd-alrazaq, Sabri Boughorbel, Arfan Ahmed, Max-Antoine Renault, Rafat Damseh, and Javaid Sheikh. Multimodal large language models in health care: Applications, challenges, and future outlook. *J Med Internet Res*, 26:e59505, 2024. 1

[5] Seongsu Bae, Daeun Kyung, Jaehee Ryu, Eunbyeol Cho, Gyubok Lee, Sunjun Kweon, Jungwoo Oh, Lei Ji, Eric Chang, Tackeun Kim, and Edward Choi. Ehrxqa: A multi-modal question answering dataset for electronic health records with chest x-ray images. In *NeurIPS*, pages 3867–3880. Curran Associates, Inc., 2023. 5, 7, 1

[6] Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H. Brendan McMahan, Virginia Smith, and Ameet Talwalkar. Leaf: A benchmark for federated settings, 2019. 1, 3

[7] Yae Jee Cho, Andre Manoel, Gauri Joshi, Robert Sim, and Dimitrios Dimitriadis. Heterogeneous ensemble knowledge transfer for training large models in federated learning. In *International Joint Conference on Artificial Intelligence*, 2022. 2

[8] Qian Dai, Dong Wei, Hong Liu, Jinghan Sun, Liansheng Wang, and Yefeng Zheng. Federated modality-specific encoders and multimodal anchors for personalized brain tumor segmentation. In *AAAI*, 2024. 3

[9] Maria Correia de Verdier, Rachit Saluja, Louis Gagnon, Dominic LaBella, Ujjwall Baid, Nourel Hoda Tahon, Martha Foltyn-Dumitru, Jikai Zhang, Maram Alafif, Saif Baig, Ken Chang, Gennaro D'Anna, Lisa Deptula, Diviya Gupta, Muhammad Ammar Haider, Ali Hussain, Michael Iv, Marinos Kontzialis, Paul Manning, Farzan Moodi, Teresa Nunes, Aaron Simon, Nico Sollmann, David Vu, Maruf Adewole, Jake Albrecht, Udunna Anazodo, Rongrong Chai, Verena Chung, Shahriar Faghani, Keyvan Farahani, Anahita Fathi Kazerooni, Eugenio Iglesias, Florian Kofler, Hongwei Li, Marius George Linguraru, Bjoern Menze, Ahmed W. Moawad, Yury Velichko, Benedikt Wiestler, Talissa Altes, Patil Basavasagar, Martin Bendszus, Gianluca Brugnara, Jaeyoung Cho, Yaseen Dhemesh, Brandon K. K. Fields, Filip Garrett, Jaime Gass, Lubomir Hadjiiski, Jona Hattangadi-Gluth, Christopher Hess, Jessica L. Houk, Edvin Isufi, Lester J. Layfield, George Mastorakos, John Mongan, Pierre Nedelec, Uyen Nguyen, Sebastian Oliva, Matthew W. Pease, Aditya Rastogi, Jason Sinclair, Robert X. Smith, Leo P. Sugrue, Jonathan Thacker, Igor Vidic, Javier Villanueva-Meyer, Nathan S. White, Mariam Aboian, Gian Marco Conte, Anders Dale, Mert R. Sabuncu, Tyler M. Seibert, Brent Weinberg, Aly Abayazeed, Raymond Huang, Sevcan Turk, Andreas M. Rauschecker, Nikdokht Farid, Philipp Vollmuth, Ayman Nada, Spyridon Bakas, Evan Calabrese, and Jeffrey D. Rudie. The 2024 brain tumor segmentation (brats) challenge: Glioma segmentation on post-treatment mri, 2024. 4, 7

[10] Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, Eric Lehman, Caiming Xiong, Richard Socher, and Byron C. Wallace. ERASER: A benchmark to evaluate rationalized NLP models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4443–4458, Online, 2020. Association for Computational Linguistics. 5

[11] Yuhang Ding, Xin Yu, and Yi Yang. Rfnet: Region-aware fusion network for incomplete multi-modal brain tumor segmentation. In *ICCV*, pages 3975–3984, 2021. 4

[12] T. Feng, D. Bose, T. Zhang, R. Hebbar, A. Ramakrishna, R. Gupta, M. Zhang, S. Avestimehr, and S. Narayanan. Fedmultimodal: A benchmark for multimodal federated learning. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, pages 4035–4045, 2023. 2, 3, 4

[13] A. Goldberger, L. Amaral, L. Glass, J. Hausdorff, P. C. Ivanov, R. Mark, and H. E. Stanley. Physiobank, physiotoolkit, and physionet: Components of a new research resource for complex physiologic signals. *Circulation [Online]*, 101(23):e215–e220, 2000. RRID:SCR_007345. 4

[14] Hao Guan, Pew-Thian Yap, Andrea Bozoki, and Mingxia Liu. Federated learning for medical image analysis: A survey. *Pattern Recognition*, 151:110424, 2024. 1

[15] Chaoyang He, Songze Li, Jinhyun So, Xiao Zeng, Mi Zhang, Hongyi Wang, Xiaoyang Wang, Praneeth Vepakomma, Abhishek Singh, Hang Qiu, Xinghua Zhu, Jianzong Wang, Li Shen, Peilin Zhao, Yan Kang, Yang Liu, Ramesh Raskar, Qiang Yang, Murali Annavaram, and Salman Avestimehr. Fedml: A research library and benchmark for federated machine learning, 2020. 1, 3

[16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5

[17] Xuehai He, Yichen Zhang, Luntian Mou, Eric Xing, and Pengtao Xie. Pathvqa: 30000+ questions for medical visual question answering, 2020. 5, 7

[18] Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data distribution for federated visual classification. In *Neurips Workshop on Federated Learning*, 2019. 2, 4, 5, 8

[19] A. Johnson, T. Pollard, R. Mark, S. Berkowitz, and

S. Horng. Mimic-cxr database (version 2.1.0), 2024. RRID:SCR_007345. 4

[20] Alistair E. W. Johnson, Tom J. Pollard, Seth J. Berkowitz, Nathaniel R. Greenbaum, Matthew P. Lungren, Chih-ying Deng, Roger G. Mark, and Steven Horng. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific Data*, 6(1):317, 2019. 4

[21] Alistair E. W. Johnson, Tom J. Pollard, Nathaniel R. Greenbaum, Matthew P. Lungren, Chih ying Deng, Yifan Peng, Zhiyong Lu, Roger G. Mark, Seth J. Berkowitz, and Steven Horng. Mimic-cxr-jpg, a large publicly available database of labeled chest radiographs, 2019. 4, 7, 1

[22] A. Karargyris, R. Umeton, M. J. Sheller, A. Aristizabal, J. George, A. Wuest, S. Pati, et al. Federated benchmarking of medical artificial intelligence with medperf. *Nature Machine Intelligence*, 5:799–810, 2023. 4

[23] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J. Reddi, Sebastian U. Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning, 2020. 2, 6

[24] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima, 2017. 7

[25] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2015. 6

[26] Felix Krones, Umar Marikkar, Guy Parsons, Adam Szmul, and Adam Mahdi. Review of multimodal machine learning approaches in healthcare. *Information Fusion*, 114:102690, 2025. 1

[27] Fan Lai, Yinwei Dai, Sanjay S. Singapuram, Jiachen Liu, Xiangfeng Zhu, Harsha V. Madhyastha, and Mosharaf Chowdhury. Fedscale: Benchmarking model and system performance of federated learning at scale, 2022. 1

[28] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. BLIP: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *Proceedings of the 39th International Conference on Machine Learning*. PMLR, 2022. 5

[29] Qinbin Li, Bingsheng He, and Dawn Song. Model-contrastive federated learning, 2021. 2, 6, 3

[30] Qinbin Li, Yiqun Diao, Quan Chen, and Bingsheng He. Federated learning on non-iid data silos: An experimental study. In *IEEE International Conference on Data Engineering*, 2022. 1, 2, 3, 4, 5

[31] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine*, 37(3):50–60, 2020. 7

[32] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks, 2020. 2, 5

[33] Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data, 2020. 2, 5, 8

[34] Xiaoxiao Li, Meirui Jiang, Xiaofei Zhang, Michael Kamp, and Qi Dou. Fedbn: Federated learning on non-iid features

via local batch normalization. In *International Conference on Learning Representations*, 2021. 3

[35] Tao Lin, Lingjing Kong, Sebastian U. Stich, and Martin Jaggi. Ensemble distillation for robust model fusion in federated learning. In *NeurIPS*, 2020. 2

[36] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B.A. y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 1273–1282. PMLR, 2017. 1, 2, 5

[37] Jean Ogier du Terrail, Samy-Safwan Ayed, Edwige Cyffers, Felix Grimberg, Chaoyang He, Regis Loeb, Paul Mangold, Tanguy Marchand, Othmane Marfoq, Erum Mushtaq, Boris Muzellec, Constantin Philippenko, Santiago Silva, Maria Teleńczuk, Shadi Albarqouni, Salman Avestimehr, Aurélien Bellet, Aymeric Dieuleveut, Martin Jaggi, Sai Praneeth Karimireddy, Marco Lorenzi, Giovanni Neglia, Marc Tommasi, and Mathieu Andreux. Flamby: Datasets and benchmarks for cross-silo federated learning in realistic healthcare settings. In *Advances in Neural Information Processing Systems*, pages 5315–5334. Curran Associates, Inc., 2022. 2, 3

[38] Pranav Poudel, Prashant Shrestha, Sanskar Amgain, Yash Raj Shrestha, Prashnna Gyawali, and Binod Bhattarai. CAR-MFL: Cross-Modal Augmentation by Retrieval for Multimodal Federated Learning with Missing Modalities . In *Medical Image Computing and Computer Assisted Intervention*. Springer Nature Switzerland, 2024. 4

[39] Sashank Reddi, Zachary Burr Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and Brendan McMahan. Adaptive federated optimization. In *International Conference on Learning Representations*, 2021. 2, 3

[40] Adriel Saporta, Aahlad Puli, Mark Goldstein, and Rajesh Ranganath. Contrasting with symile: Simple model-agnostic representation learning for unlimited modalities. In *NeurIPS*. Curran Associates, Inc., 2024. 4, 5

[41] Adriel Saporta et al. Symile-mimic: a multimodal clinical dataset of chest x-rays, electrocardiograms, and blood labs from mimic-iv (version 1.0.0). https://doi.org/10.13026/3vvj-s428, 2025. RRID:SCR_007345. 4, 7

[42] P. Shrestha, S. Amgain, B. Khanal, C.A. Linte, and B. Bhattarai. Medical vision language pretraining: A survey. *arXiv preprint arXiv:2312.06224*, 2023. 1

[43] Walter Simoncini, Spyros Gidaris, Andrei Bursuc, and Yuki M. Asano. No train, all gain: Self-supervised gradients improve deep frozen representations, 2024. 8

[44] Alysa Ziying Tan, Han Yu, Lizhen Cui, and Qiang Yang. Towards personalized federated learning. *IEEE Transactions on Neural Networks and Learning Systems*, 34(12): 9587–9603, 2023. 3

[45] J. Venugopalan, L. Tong, H.R. Hassanzadeh, and M.D. Wang. Multimodal deep learning models for early detection of alzheimer's disease stage. *Scientific Reports*, 11(1):3254, 2021. 1

[46] Patrick Wagner, Nils Strodthoff, Ralf-Dieter Bousseljot, Dieter Kreiseler, Fatima I. Lunze, Wojciech Samek, and Tobias

Schaeffter. Ptb-xl, a large publicly available electrocardiography dataset. *Scientific Data*, 7(1):154, 2020. 3

[47] Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H. Vincent Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization, 2020. 2, 6

[48] Chuhan Wu, Fangzhao Wu, Lingjuan Lyu, Yongfeng Huang, and Xing Xie. Communication-efficient federated learning via knowledge distillation. *Nature Communications*, 13(1), 2022. 2

[49] J. Xu, B.S. Glicksberg, C. Su, P. Walker, J. Bian, and F. Wang. Federated learning for healthcare informatics. *Journal of Healthcare Informatics Research*, 5(1):1–19, 2021. Epub 2020 Nov 12. 1, 7

[50] Rui Ye, Rui Ge, Xinyu Zhu, Jingyi Chai, Yaxin Du, Yang Liu, Yanfeng Wang, and Siheng Chen. Fedllm-bench: Realistic benchmarks for federated learning of large language models. In *NeurIPS*, 2024. 2, 3

[51] Jong Chan Yeom, Jae Hoon Kim, Young Jae Kim, Jisup Kim, and Kwang Gi Kim. A comparative study of performance between federated learning and centralized learning using pathological image of endometrial cancer. *Journal of Imaging Informatics in Medicine*, 37(4):1683–1690, 2024. Comparative Study; Research Support, Non-U.S. Gov't; Epub 2024 Feb 21. 7

[52] Qiying Yu, Yang Liu, Yimu Wang, Ke Xu, and Jingjing Liu. Multimodal federated learning via contrastive representation ensemble. In *International Conference on Learning Representations*, 2023. 3, 6

[53] Mikhail Yurochkin, Mayank Agarwal, Soumya Ghosh, Kristjan Greenewald, Nghia Hoang, and Yasaman Khazaeni. Bayesian nonparametric federated learning of neural networks. In *International Conference on Machine Learning*. PMLR, 2019. 4

[54] Sheng Zhang, Yanbo Xu, Naoto Usuyama, Hanwen Xu, Jaspreet Bagga, Robert Tinn, Sam Preston, Rajesh Rao, Mu Wei, Naveen Valluri, Cliff Wong, Andrea Tupini, Yu Wang, Matt Mazzola, Swadheen Shukla, Lars Liden, Jianfeng Gao, Angela Crabtree, Brian Piening, Carlo Bifulco, Matthew P. Lungren, Tristan Naumann, Sheng Wang, and Hoifung Poon. A multimodal biomedical foundation model trained from fifteen million image–text pairs. *NEJM AI*, 2(1), 2024. 5

[55] Weiying Zheng, Ziyue Lin, Pengxin Guo, Yuyin Zhou, Feifei Wang, and Liangqiong Qu. Fedvlmbench: Benchmarking federated fine-tuning of vision-language models, 2025. 2, 3, 4, 7

# Med-MMFL: A Multimodal Federated Learning Benchmark in Healthcare

## Supplementary Material

In this supplementary material, we first provide additional information about each dataset individually, including details on the corresponding data partitioning strategy, in Sec. 7. In Sec. 8, we describe further details regarding the experimental setup, including hyperparameter configurations. Sec. 9 presents additional results, covering both quantitative and qualitative components.

## 7. Datasets and Partitioning Strategy

**Fed-BraTS-GLI2024.** We evaluate the methods on the multimodal brain tumor segmentation task in Fed-BraTS-GLI2024. The objective is to segment three tumor sub-regions: **whole tumor**, **tumor core**, and **enhancing tumor** from multimodal MRI scans. The whole tumor comprises all three constituent sub-regions: the necrotic and non-enhancing tumor core (NCR/NET), the peritumoral edema (ED), and the GD-enhancing tumor (ET). The tumor core includes NCR/NET and ET, whereas the enhancing tumor corresponds exclusively to ET.

**Fed-MIMIC-CXR-JPG.** Fed-MIMIC-CXR-JPG is evaluated on a multi-label classification task. Fig. 4 illustrates the distribution of the 14 pathology labels in the dataset, showing the frequency of instances annotated with: Atelectasis, Cardiomegaly, Consolidation, Edema, Enlarged Cardiomediastinum, Fracture, Lung Lesion, Lung Opacity, No Finding, Pleural Effusion, Pleural Other, Pneumonia, Pneumothorax, and Support Devices.
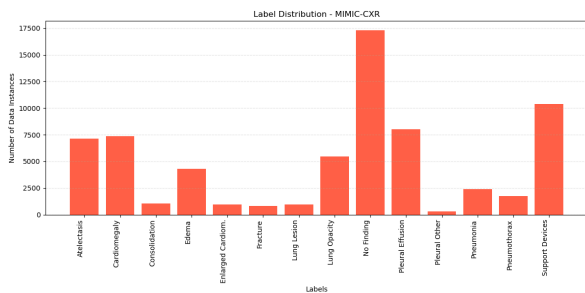


Figure 4. The distribution of labels in MIMIC-CXR-JPG [21]

**Fed-EHRXQA.** As described in Sec. 3, the original EHRXQA dataset [5] contains multi-patient SQL-style queries that retrieve information across the entire EHR database. Our goal, however, is to focus the evaluation on tasks that naturally reflect clinical reasoning, where a model interprets multimodal information for a single patient. Population-level SQL retrieval such as counting, listing, or filtering patients, does not align with this objective. To align the dataset with this objective, we removed templates requiring database-wide aggregation (e.g., listing or counting all patients matching some condition). Representative examples of discarded templates are shown in Tab. 4, where both the template and an actual sample question are included.

To ensure that the task aligns with realistic, single-patient clinical VQA, we retain only those templates explicitly anchored to a `patient_id` or `study_id`. These can be answered using a single patient's EHR and corresponding chest X-ray. Examples of retained templates and their sample questions are shown in Tab. 5. This allows us to preserve rich, per-patient EHR data and support meaningful multimodal reasoning between EHRs and CXRs. The retained instances naturally constrain the task to single-patient, multimodal clinical VQA (see Fig. 7), eliminating the need for database-wide SQL retrieval. For a full description of the original EHRXQA schema and query design, we refer readers to the dataset's accompanying publication [5].
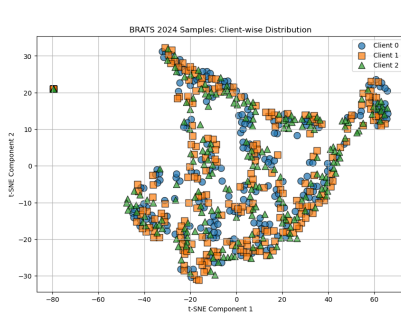
**Synthetic Partitioning.** For our experiments with synthetic partitions, we follow the Dirichlet-based partitioning strategy described in Sec. 3. When $\alpha$ is very high (approaching $\infty$), the sampled probabilities $p_k$ are nearly uniform across clients, resulting in roughly equal proportions of each class being assigned to every client. This produces nearly homogeneous splits, which we refer to as *synthetic IID* partitions. Conversely, as $\alpha$ decreases, the sampling becomes increasingly skewed, with certain clients receiving disproportionately more instances of some pseudo-classes than others. This leads to heterogeneous, non-IID distributions across clients, which we denote as *synthetic non-IID*. In our experiments, $\alpha = 0.2$ represents a more heterogeneous partition than $\alpha = 0.8$. For a qualitative understanding of the distributional differences introduced by our synthetic Dirichlet partitions in Fed-BraTS-GLI2024, we reduce the class-proportion vectors to two dimensions using t-SNE and visualize them in Fig. 5. The plots highlight the transition from balanced to highly skewed client distributions as the Dirichlet parameter decreases. Similarly, client-level distributions of Fed-EHRXQA and Fed-PathVQA resulting from our federated partitioning strategy are illustrated in Fig. 6.

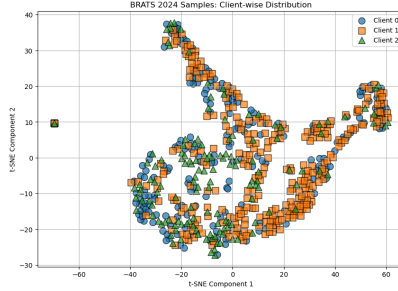| Template | Sample Question |
|---|---|
| `list the ids of patients who had any chest x-ray study indicating $attribute [time_filter_global1].` | `identify the patients by their ids who had chest x-ray studies indicating low lung volumes since 03/2103.` |
| `count the number of patients aged [age_group] who had a chest x-ray study during hospital visit indicating any $category in the $object [time_filter_global1].` | `how many patients in the 20s age group had chest x-ray findings of any anatomical findings in the right lung until 05/2101?` |

Table 4. Examples of removed EHRXQA question templates that require multi-patient SQL-style retrieval.

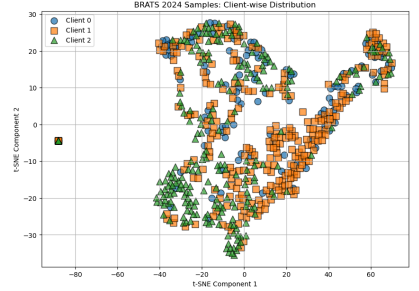| Template | Sample Question |
|---|---|
| `has_verb patient {patient_id} been prescribed with {drug_name} [time_filter_global1] and also had a chest x-ray study indicating $attribute within the same period?` | `was patient 11887414 prescribed milk of magnesia since 91 months ago and had a chest x-ray showing a chest port within the same timeframe?` |
| `given the {study_id1} study, list all anatomical locations related to any $category that are $comparison compared to the previous study?` | `enumerate all anatomical locations related to any diseases still present in the 50978999 study relative to the previous one.` |

Table 5. Examples of retained EHRXQA templates suitable for patient-level clinical VQA.



(a) Synthetic IID partition: Class proportion vectors are uniformly distributed across clients.

(b) Synthetic non-IID partition ($\alpha = 0.8$): moderate heterogeneity yields partially separated clusters, indicating uneven pseudo-class exposure across clients.

(c) Synthetic non-IID partition ($\alpha = 0.2$): strong heterogeneity produces distinct, well-separated client distributions, consistent with highly skewed pseudo-class allocation.

Figure 5. t-SNE visualization of class-proportion vectors across clients for the Fed-BraTS-GLI2024 dataset under three partitioning strategies: synthetic IID, synthetic non-IID ($\alpha = 0.8$), and synthetic non-IID ($\alpha = 0.2$), each with 3 clients. Class-proportion vectors are derived from normalized voxel counts per class within each segmentation mask. As the Dirichlet concentration parameter $\alpha$ decreases, client distributions become increasingly separated, illustrating the controlled progression from balanced to highly skewed label distributions.

# 8. Experimental Settings

For each dataset, we first implement the centralized baseline. To achieve this, we either tune the hyperparameters in the centralized setting or adopt the hyperparameters used in the corresponding original reference implementations (see Sec. 3). Except for FedNova [47], where the local optimizer differs, we use a single shared set of hyperparameters for each dataset across all FL algorithms. In addition, we tune the algorithm-specific hyperparameters for each FL method using one data partition per dataset and then apply the same tuned values to all remaining partitions. This procedure is intended to ensure consistency and allow us to reliably assess how variations in data partitions influence the performance of FL algorithms under different multimodal medical settings. The complete hyperparameter configura-

(a) Fed-EHRXQA pseudo-label distribution. IID splits yield similar client-level frequencies, while non-IID splits result in heterogeneous distributions.



(b) Fed-PathVQA pseudo-label distribution under the same synthetic partitioning strategy.

Figure 6. Client-level pseudo-label distributions for Fed-EHRXQA and Fed-PathVQA obtained using our synthetic federated partitioning strategy.

| Dataset | Local epochs | Communication rounds | Learning rate | $\mu_{\text{fedprox}}$ | $\mu_{\text{moon}}$ | $\tau_{\text{moon}}$ | $\gamma_{\text{creamfl}}$ | $\alpha_{\text{creamfl}}$ |
|---|---|---|---|---|---|---|---|---|
| Fed-BraTS-GLI2024 | 3 | 50 | 0.0002 | 0.1 | 0.1 | 0.5 | 0.002 | 0.03 |
| Fed-MIMIC-CXR-JPG | 3 | 30 | 0.0001 | 0.1 | 0.1 | 0.5 | 0.002 | 0.03 |
| Fed-Symile-MIMIC | 3 | 30 | 0.001 | 0.1 | 0.1 | 0.5 | 0.002 | 0.03 |
| Fed-PathVQA | 3 | 20 | 0.00001 | 0.1 | 10 | 0.5 | 0.002 | 0.03 |
| Fed-EHRXQA | 5 | 30 | 0.00001 | 0.1 | 0.1 | 0.5 | 0.002 | 0.03 |

Table 6. Hyperparameters used for FL evaluations.

tion is provided in Tab. 6. For FedNova, the learning rate is tuned separately over the set $\{0.1, 0.01, 0.001\}$. In contrast, the search spaces for both $\mu_{\text{moon}}$ and $\mu_{\text{fedprox}}$ are $\{0.1, 1, 10\}$. We directly set $\tau_{\text{moon}} = 0.5$ without additional tuning, following the original work [29], which reported relatively low sensitivity to this hyperparameter.

# 9. Additional Results

This section provides supplementary results corresponding to the experiments reported in Sec. 5.

## 9.1. Quantitative Results

**Fed-BraTS-GLI2024.** Following the definitions of the whole tumor, tumor core, and enhancing tumor regions provided in Sec. 7, we present in Tab. 7, the per-subregion Dice scores for all six FL algorithms. These results correspond to the same experimental settings reported in the main paper. Notable observations from the per–tumor-subregion Dice scores include the following: in the synthetic IID (3-client) setting, FedProx reports higher Dice scores for the tumor core and enhancing regions, while FedNova achieves the highest overall average. Similarly, in the synthetic non-IID $\alpha = 0.2$ case, FedProx achieves the highest overall average but has a lower whole-tumor score than four other algorithms, with its average largely influenced by stronger tumor-core and enhancing-region performance again. Under the synthetic non-IID configuration with $\alpha = 0.8$, m-MOON attains a higher whole-tumor score than SCAFFOLD, despite SCAFFOLD yielding the better overall aver-

age. Although CreamMFL does not outperform any method in terms of overall Dice averages, it records a higher whole-tumor score than FedNova in the synthetic non-IID $\alpha = 0.8$ setting. These patterns indicate that overall Dice averages can obscure sub-region–specific behavior, and that algorithmic performance may vary substantially across tumor sub-regions.

## 9.2. Qualitative Results

To provide a clearer sense of the underlying tasks, we present qualitative predictions from centralized baseline models on PathVQA (Fig. 8) and EHRXQA (Fig. 7). These examples illustrate the nature of the multimodal reasoning required in both datasets. Note that these examples are intended solely to demonstrate the task characteristics rather than to compare different federated training methods. Aside from these examples, we also include qualitative results for Fed-BraTS-GLI2024, showcasing the input MRI modalities (T2-FLAIR, T1Gd, T1, and T2), the corresponding ground-truth tumor segmentation mask, and the predicted masks produced by six federated learning algorithms. These visualizations are drawn from the synthetic non-IID setting with Dirichlet $\alpha = 0.2$ across three clients, and are provided to offer a representative view of the segmentation outputs across methods.

| Split | Clients | Sub-region | FedAvg | FedProx | SCAFFOLD | m-MOON | FedNova | CreamMFL |
|---|---|---|---|---|---|---|---|---|
| Natural | 5 | Whole | 81.902 | 82.000 | 78.904 | **83.224** | 82.664 | 80.958 |
| | | Core | 81.606 | 80.730 | 79.556 | **81.974** | 81.704 | 77.928 |
| | | Enhancing | 81.883 | 81.049 | 79.880 | **82.020** | 82.007 | 77.919 |
| Synthetic IID | 3 | Whole | 79.778 | 80.972 | 81.062 | 83.595 | **84.210** | 76.293 |
| | | Core | 83.628 | **84.081** | 82.146 | 83.215 | 81.865 | 76.435 |
| | | Enhancing | 84.419 | **84.722** | 82.623 | 83.853 | 82.359 | 76.889 |
| | 5 | Whole | 80.392 | 79.878 | 78.502 | 79.635 | **84.210** | 77.393 |
| | | Core | 80.415 | 79.399 | 77.731 | 79.471 | **81.865** | 76.723 |
| | | Enhancing | 81.537 | 79.873 | 78.092 | 79.858 | **82.359** | 77.075 |
| Synthetic non-IID $\alpha = 0.8$ | 3 | Whole | 81.135 | 79.935 | 83.686 | **84.022** | 83.982 | 78.346 |
| | | Core | 81.487 | 80.591 | **83.820** | 82.707 | 82.524 | 81.346 |
| | | Enhancing | 81.985 | 81.027 | **84.318** | 83.052 | 82.988 | 82.085 |
| | 5 | Whole | 79.459 | 80.449 | 77.275 | 81.407 | 84.320 | **84.743** |
| | | Core | 80.531 | 80.250 | 80.194 | 79.725 | **81.169** | 77.911 |
| | | Enhancing | 80.886 | 80.699 | 80.674 | 80.164 | **81.487** | 78.153 |
| Synthetic non-IID $\alpha = 0.2$ | 3 | Whole | 81.258 | 81.678 | 82.668 | 83.011 | **83.225** | 79.716 |
| | | Core | 81.965 | **83.738** | 81.066 | 82.923 | 82.671 | 82.689 |
| | | Enhancing | 82.116 | **84.023** | 81.241 | 83.362 | 82.850 | 82.894 |
| | 5 | Whole | 80.006 | 80.834 | 80.419 | 80.898 | **84.499** | 74.015 |
| | | Core | 81.333 | 81.034 | 80.501 | 81.207 | **82.997** | 74.875 |
| | | Enhancing | 82.172 | 81.494 | 80.964 | 81.645 | **83.240** | 75.562 |

Table 7. Supplementary per–tumor-subregion Dice scores for the Fed-BraTS-GLI2024 dataset. These metrics complement the average Dice scores reported in Sec. 5 by detailing the whole tumor, tumor core, and enhancing tumor regions.
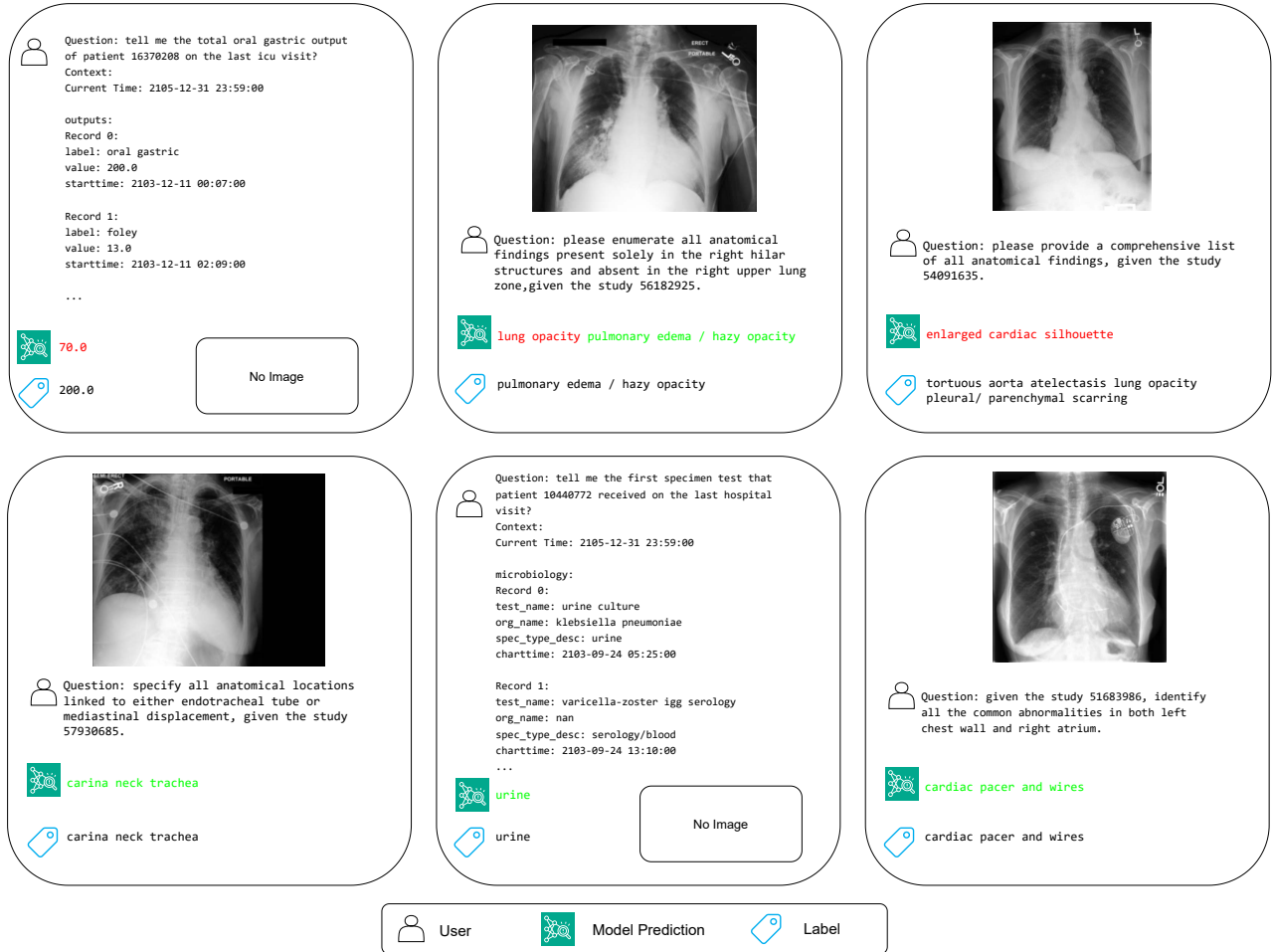
**Panel 1 (top-left):**
Question: tell me the total oral gastric output of patient 16370208 on the last icu visit?
Context:
Current Time: 2105-12-31 23:59:00

outputs:
Record 0:
label: oral gastric
value: 200.0
starttime: 2103-12-11 00:07:00

Record 1:
label: foley
value: 13.0
starttime: 2103-12-11 02:09:00

...

Model Prediction: 70.0

Label: 200.0

No Image

**Panel 2 (top-middle):**
Question: please enumerate all anatomical findings present solely in the right hilar structures and absent in the right upper lung zone,given the study 56182925.

Model Prediction: lung opacity pulmonary edema / hazy opacity

Label: pulmonary edema / hazy opacity

**Panel 3 (top-right):**
Question: please provide a comprehensive list of all anatomical findings, given the study 54091635.

Model Prediction: enlarged cardiac silhouette

Label: tortuous aorta atelectasis lung opacity pleural/ parenchymal scarring

**Panel 4 (bottom-left):**
Question: specify all anatomical locations linked to either endotracheal tube or mediastinal displacement, given the study 57930685.

Model Prediction: carina neck trachea

Label: carina neck trachea

**Panel 5 (bottom-middle):**
Question: tell me the first specimen test that patient 10440772 received on the last hospital visit?
Context:
Current Time: 2105-12-31 23:59:00

microbiology:
Record 0:
test_name: urine culture
org_name: klebsiella pneumoniae
spec_type_desc: urine
charttime: 2103-09-24 05:25:00

Record 1:
test_name: varicella-zoster igg serology
org_name: nan
spec_type_desc: serology/blood
charttime: 2103-09-24 13:10:00
...

Model Prediction: urine

Label: urine

No Image

**Panel 6 (bottom-right):**
Question: given the study 51683986, identify all the common abnormalities in both left chest wall and right atrium.

Model Prediction: cardiac pacer and wires

Label: cardiac pacer and wires

**Legend:**
User    Model Prediction    Label

Figure 7. Qualitative examples from the centralized EHRXQA baseline on patient-level open-ended VQA task.
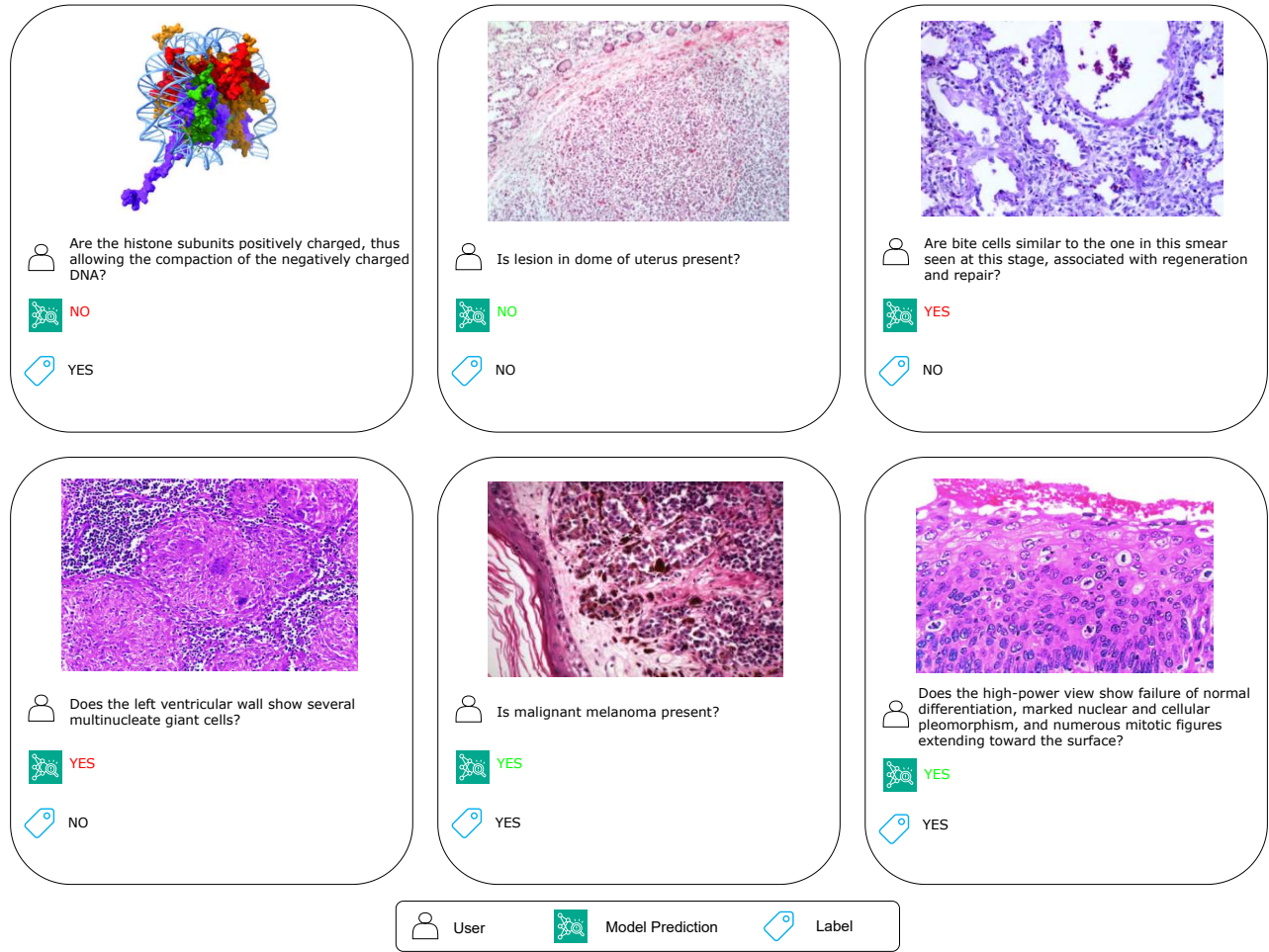
5

Figure 8. Qualitative examples from the centralized PathVQA baseline on close-ended VQA task.
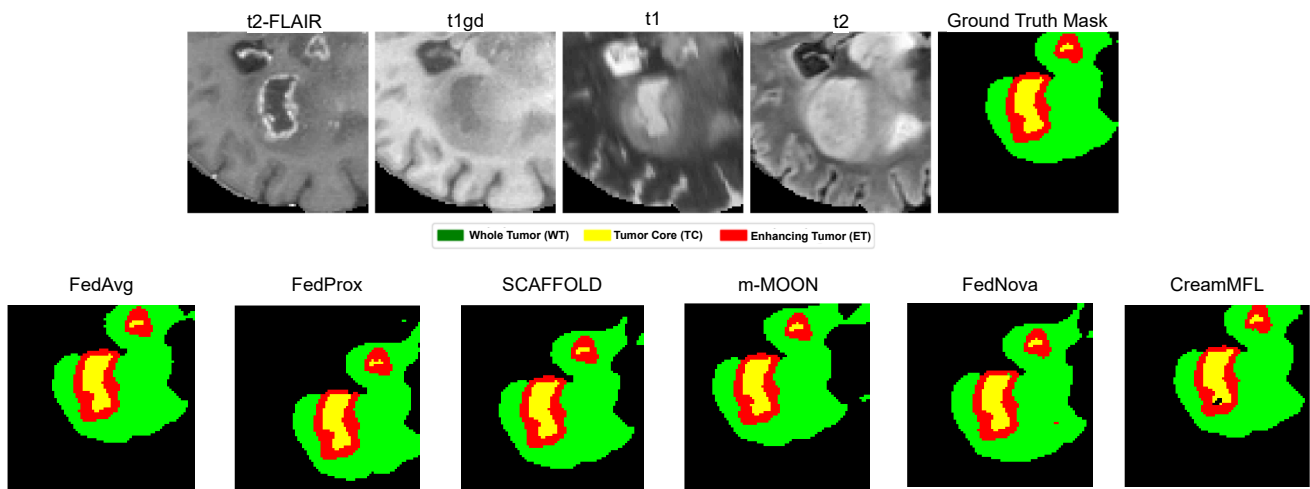


Figure 9. Multimodal Brain Tumor Segmentation results on Fed-BraTS-GLI2024. Top row shows the four MRI input modalities (T2-FLAIR, T1Gd, T1, T2), followed by the ground-truth segmentation mask , and the bottom row shows the predictions from six federated learning algorithms under the synthetic non-IID setting ($\alpha = 0.2$, three clients).

6